

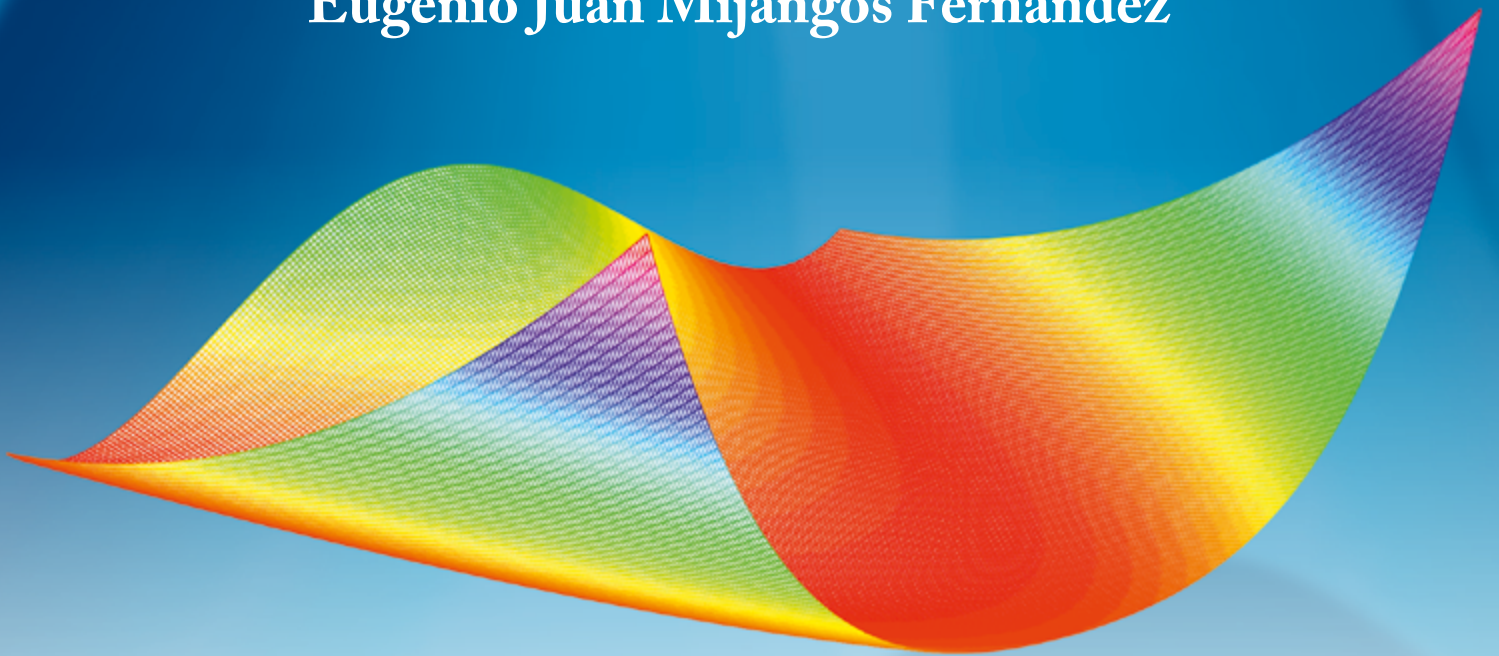
Zenbakizko metodoak

MATLAB<sup>®</sup>

erabiliz

2.  
edizioa

Eugenio Juan Mijangos Fernández



eman ta zabal zazu



Universidad  
del País Vasco

Euskal Herriko  
Unibertsitatea

UPV/EHUko Euskara Zerbitzuak sustatua eta zuzendua, Euskarazko ikasmaterialgintza sustatzeko deialdiaren bitartez.

© Servicio Editorial de la Universidad del País Vasco  
Euskal Herriko Unibertsitateko Argitalpen Zerbitzua  
eISBN: 978-84-9860-545-7

**Zenbakizko metodoak**

---

**MATLAB<sup>®</sup>**

---

**erabiliz**

**2.**  
edizioa

**Eugenio Juan Mijangos Fernández**

eman ta zabal zazu



Universidad  
del País Vasco

Euskal Herriko  
Unibertsitatea

Eugenio Juan Mijangos Fernández

Matematika Saila

Zientzia eta Teknologia Fakultatea

`eugenio.mijangos@ehu.eus`

# Gaien Aurkibidea

<b>1. Sarrera</b>	<b>1</b>
1.1. Problema bat duzu . . . . .	1
<b>2. MATLABi buruzko oinarriko nozioak</b>	<b>5</b>
2.1. Eragiketa aritmetikoak . . . . .	5
2.2. Programari gehitutako funtzioak . . . . .	6
2.3. Esleipen-instrukzioak . . . . .	7
2.4. Matrizeak . . . . .	7
2.5. Eragiketa matrizialak . . . . .	10
2.6. Gaiez gai egiten diren eragiketak . . . . .	11
2.6.1. MATLABek emandako funtzioak . . . . .	12
2.7. Grafikoak . . . . .	13
2.8. MATLABez programatzen: M fitxategiak . . . . .	17
2.8.1. Instrukzioen fitxategiak . . . . .	17
2.8.2. Funtzio-fitxategiak . . . . .	18
2.8.3. Azpifuntzioak . . . . .	20
2.8.4. Input-output . . . . .	21
2.8.5. Fitxategiak sortu eta fitxategietan sartu . . . . .	23

2.8.6.	Programazio egituratua . . . . .	24
2.8.7.	M fitxategiei funtzioak igorri . . . . .	32
2.9.	Problemak . . . . .	37
<b>3.</b>	<b>Ordenagailuaren aritmetika eta errorearen analisisa</b>	<b>43</b>
3.1.	Algoritmoak eta erroreak . . . . .	43
3.1.1.	Erroreak . . . . .	44
3.1.2.	Diskretizazio-erroreak . . . . .	45
3.1.3.	Biribiltze-errorearen eragin kaltegarria . . . . .	47
3.1.4.	Algoritmo bat ebaluatzeko irizpideak . . . . .	48
3.2.	Puntu higikorreko sistemak . . . . .	49
3.2.1.	Inaustea eta biribiltzea . . . . .	49
3.2.2.	Biribiltzearen unitatea . . . . .	51
3.3.	Zenbaki bitarrak . . . . .	51
3.3.1.	Zenbaki oso bitarrak . . . . .	52
3.3.2.	Zatiki bitarrak . . . . .	53
3.3.3.	Ordenagailuko zenbakiak . . . . .	55
3.3.4.	Ordenagailu baten zehaztasuna . . . . .	56
3.3.5.	Ordenagailuko puntu higikorreko zenbakiak . . . . .	56
3.4.	Errorearen analisisa . . . . .	60
3.4.1.	Zifra esanguratsuen ezeztapena . . . . .	61
3.4.2.	Trunkatze-errorea . . . . .	64
3.4.3.	$O(h^n)$ hurbiltze-ordena . . . . .	64
3.4.4.	Segida baten hurbiltze-ordena . . . . .	67
3.4.5.	Errorearen hedapena . . . . .	67

<i>GAIEN AURKIBIDEA</i>	iii
3.4.6. Datuen ziurgabetasuna . . . . .	70
3.5. Problemak . . . . .	71
<b>4. Ekuazio ez-linealen ebazpena</b>	<b>81</b>
4.1. Metodo grafikoak . . . . .	81
4.2. Bakartze-metodoak . . . . .	83
4.2.1. Erroak bakartzea . . . . .	83
4.2.2. Bisekzio-metodoa . . . . .	85
4.2.3. <i>Regula falsi</i> metodoa . . . . .	88
4.3. Metodo irekiak . . . . .	90
4.3.1. Puntu finkoaren metodoa . . . . .	90
4.3.2. Konbergentziaren ordena . . . . .	94
4.3.3. Newton-Raphson-en metodoa . . . . .	95
4.3.4. Ebakitzailaren metodoa . . . . .	100
4.3.5. Muller-en metodoa . . . . .	102
4.3.6. Alderantzizko interpolazio koadratikoa . . . . .	104
4.3.7. Zeroin algoritmoa . . . . .	105
4.4. Algoritmoak gelditzeko irizpideak . . . . .	106
4.5. Problemak . . . . .	107
<b>5. Sistema linealak: metodo zuzenak eta iteratiboak</b>	<b>115</b>
5.1. Sistema linealen ebazpena . . . . .	115
5.2. 3x3 adibide bat . . . . .	116
5.3. Permutazio- eta triangelu-matrizeak . . . . .	118
5.4. Pibotatzearen beharra . . . . .	120

5.4.1.	Pibotatzte partziala . . . . .	121
5.5.	<i>LU</i> faktORIZAZIOA . . . . .	122
5.5.1.	Pibotatzte baztergarria . . . . .	125
5.6.	Matematikako problema baten baldintza . . . . .	125
5.7.	Matrizeen normak . . . . .	126
5.7.1.	Bi-norma eta espektro-erradioa . . . . .	129
5.8.	Sistema lineal baten baldintzazko zenbakia . . . . .	129
5.9.	Cholesky-ren faktORIZAZIOA . . . . .	133
5.10.	Metodo iteratiboak . . . . .	138
5.10.1.	Jacobi-ren iterazioa . . . . .	139
5.10.2.	Gauss-Seidel-en iterazioa . . . . .	141
5.10.3.	Metodo egonkorren konbergentzia . . . . .	142
5.11.	Problemak . . . . .	147
<b>6.</b>	<b><i>QR</i> faktORIZAZIOA eta minimo karratu linealak</b>	<b>155</b>
6.1.	Householder-en islapenak . . . . .	155
6.2.	<i>QR</i> faktORIZAZIOA . . . . .	158
6.2.1.	<b>A</b> matrize karratu ez-singularra . . . . .	158
6.2.2.	Sistema lineal karratu determinatu baten ebazpena . . . . .	161
6.3.	Givens-en biraketak . . . . .	161
6.4.	Minimo karratu linealak: sistema gaindeterminatuak . . . . .	162
6.5.	<i>QR</i> faktORIZAZIOAREN propietateak . . . . .	169
6.6.	Hein urriko <i>QR</i> faktORIZAZIOA . . . . .	171
6.7.	Deskonposizio ortogonal osoa . . . . .	175
6.8.	Balio singularretako deskonposizioa . . . . .	176



6.8.1.	Moore-Penrose-ren sasiaderantzizko matrizea . . . . .	178
6.8.2.	Baldintzazko zenbaki orokortua . . . . .	179
6.8.3.	Minimo karratuen problemaren ebazpen egonkorrena . . . . .	180
6.9.	Problemak . . . . .	181
<b>7.</b>	<b>Ekuazio ez-linealen sistemen ebazpena</b>	<b>187</b>
7.1.	Newtonen metodoa . . . . .	187
7.1.1.	Newtonen metodoaren konbergentzia lokala . . . . .	190
7.2.	Newtonen metodoaren aldaketak, sistema ez-linealak ebazteko . . . . .	192
7.2.1.	Diferentzia finituzko Newtonen metodoa . . . . .	192
7.2.2.	Newtonen metodo aldatua . . . . .	195
7.2.3.	Jacobiren aldaera . . . . .	195
7.2.4.	Gauss-Seidelen aldaera . . . . .	195
7.3.	Quasi-Newton metodoak . . . . .	195
7.3.1.	Broyden-en metodoa . . . . .	196
7.4.	Murrizketarik gabeko optimizazioa . . . . .	201
7.4.1.	Oinarrizko metodoak . . . . .	202
7.4.2.	Metodo globalak . . . . .	204
7.4.3.	Minimo karratu ez-linealak . . . . .	212
7.5.	Problemak . . . . .	215



# Irudien Zerrenda

1.1. Abiadura/denbora grafikoa. . . . .	2
1.2. Zenbakizko soluzioa eta soluzio analitikoa. . . . .	4
2.1. Zirkulu bidimentsional bat eta helize tridimentsional bat, bi zatitako irudi batean. . . . .	16
3.1. Kurba urdin jarraituak diskretizazio-errorea ematen du, eta puntu-marra zuzen gorriak biribiltze-errererik gabeko errore hori (MATLABeko <code>loglog</code> funtzioak egindako grafikoa). . . . .	47
3.2. Zehaztasun bikoitzean adieraz daitezkeen zenbakien heina . . . . .	58
3.3. Zenbaki baten biltzea, notazio bitarrean, IEEE-754 estandarrean. . . . .	58
3.4. 22.5 zenbakiaren biltegia, notazio bitarrean, IEEE-754 estandarrean. . . . .	60
4.1. $(mp, fp)$ grafikoa. . . . .	82
4.2. Bisekzio-metodoa. . . . .	85
4.3. $f(x) = 8 - 4.5(x - \sin(x))$ -ren grafikoa. . . . .	87
4.4. <i>Regula falsi</i> metodoa. . . . .	89
4.5. (a) P. finko erakargarria ( $ F'(x)  < 1$ ). (b) P. finko alderagarria ( $ F'(x)  > 1$ ). . . . .	92
4.6. Newtonen metodoa. . . . .	97
4.7. Ebakitzailearen metodoa. . . . .	101
4.8. Mullerren metodoa. . . . .	103

6.1. Householderren islapena. . . . .	156
6.2. Minimo karratu linealen adibidea. . . . .	162
6.3. Minimo karratu linealen soluzioa. . . . .	163
7.1. Gradiente metodoaren urratsak. . . . .	205
7.2. Armijoren baldintza. . . . .	209
7.3. Armijo-Goldsteinen baldintzak. . . . .	210

# Taulen Zerrenda

- 2.1. Funtzio marratzaileen laburpen bat. . . . . 16
  
- 3.1.  $\{x_n\} = \{1/3^n\}$  segida eta bere hurbilpenak. . . . . 69
  
- 3.2. Erroreen segidak. . . . . 70
  
  
- 4.1. Bisekzio-metodoa. . . . . 86
  
- 4.2. Puntu finkoaren iterazioak. . . . . 93
  
- 4.3. Newtonen iterazioak. . . . . 99
  
- 4.4. Newtonen iterazioak. . . . . 99
  
  
- 5.1. Jacobiren iterazioaren konbergentzia 5.3. adibidean. . . . . 139
  
- 5.2. Gauss-Seidelen iterazioaren konbergentzia 5.3. adibidean. . . . . 142
  
  
- 6.1. Matrize-faktorizazioen kostu aritmetikoa. . . . . 160



## Hitzaurrea

Zenbakizko metodoak I irakasgaietan Matematika Gradu 2. mailako euskarazko ikasleek duten hutsune bibliografikoa liburu honen bidez betetzea da testugile honen helburu nagusia. Egia esan, dakidanez, gai honetaz ez dago liburu bat ere euskaraz. Liburu hau oso erabilgarria izan daiteke bai Ingeniaritzako ikasketetan, eta baita Zientziako beste ikasketa batzuetan ere; hau da, oso erabilgarria izan daiteke.

Liburu honen bidez zenbakizko metodoen oinarriak eskainiko dira, bai ikuspuntu teoriko batetik, bai ikuspuntu aplikatu batetik. Iraganean, zenbakizko analisisian eta teorian oinarritzen zen zenbakizko metodoei buruzko ikasturte bat, baina, gaur egun, irakasgai horren edukia garatzean, kontuan hartu behar dugu ordenagailu ahaltsuak eta kalkulu-software eraginkorrak eskura izatea. Oraingo joera, gero eta gehiago, aplikazioetara eta inplementazioetara joatea da, jadanik prest dagoen kalkulu-tresneria erabiliz. Ikasturte batean, ikasleek zenbakizko metodoen oinarriak ikasiko dituzte. Gainera, ordenagailu-lengoaia batean programatzen ikasiko dute, eta software aurreratua erabiliko problemak ebazteko tresna gisa. MATLAB da horrelako softwarearen adibide egoki bat. Ikasleek beren programak idazteko erabil dezakete pakete informatiko hori, eta haren funtzioak problemak ebazteko tresna gisa.

Dezimalekin erabili beharreko puntuazioari dagokionez, testuan koma eta programazio-adibideetan puntua ez erabiltzearen, puntuz bereiziko dira dezimalak, nahiz eta euskaraz dezimalak komaz banandu.

Bigarren argitalpeneko oharra:

Argitalpen honetan akats batzuk zuzendu ditut, eta hobekuntza arin batzuk sartu.

Nire gurasoei eskaintzen diet.





# 1. kapitulua

## Sarrera

### 1.1. Problema bat duzu

Demagun zubi-jauzi enpresa batek zure zerbitzuak kontratatu nahi dituela; hau da zure lana: jauzilariaren erortzeko abiadura, denboraren funtzioan (erorketa librean), aurretik jakitea. Informazio hori erabiliko da zehazteko (beste analisi handiago batean) masa desberdineko jauzilarienezako kordaren luzera eta erresistentzia.

Zuk badakizu, fisika-ikasketei esker, Newtonen bigarren legearen arabera  $a = F/m$  erlazioa bete behar dela. Baina, fluidoaren mekanikari buruzko ezagupenen arabera, eredu matematiko hau garatzen duzu:

$$\frac{dv}{dt} = g - \frac{c_d}{m}v^2, \quad (1.1)$$

non  $v$  = abiadura bertikala baita,  $t$  = denbora (s),  $g$  = grabitate-azelerazioa ( $\approx 9.81$  m/s<sup>2</sup>),  $c_d$  = airearen erresistentzia-koefizientea (kg/m) eta  $m$  = jauzilariaren masa (kg).

Ekuzio diferentzial horren  $v = v(t)$  soluzio analitiko edo zehatz bat kalkula dezakegu eta,  $t = 0$  denean  $v = 0$  dela hartzen badugu, hau da soluzio analitikoa:

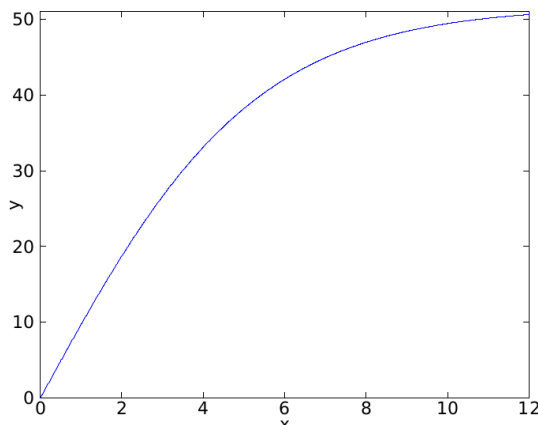
$$v(t) = \sqrt{\frac{gm}{c_d}} \tanh\left(\sqrt{\frac{gc_d}{m}}t\right), \quad (1.2)$$

non tangente hiperbolikoa ( $\tanh$ ) zuzen kalkula baitezakegu, edo funtzio esponentzialak erabiliz, adierazpen honetan:

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}.$$

**1.1. adibidea.** *68.1 kg-ko zubi-jauzilari batek zubi batetik salto egiten badu, erabili (1.2)*

$t$ (s)	$v$ (m/s)
0	0
2	18.7292
4	33.1118
6	42.0762
8	46.9575
10	49.4214
12	50.6175
$\infty$	51.6938



**1.1. irudia.** Abiadura/denbora grafikoa.

*adierazpena, bere abiadura erorketa librean 12 s-ko unean kalkulatzeko. Zein izango da azken abiadura, kordaren luzera infinitua bada? Erabili  $c_d = 0.25$  kg/m.*

*Ebazpena.* Parametroen balioak (1.2) adierazpenean ordezkatzan baditugu hau ematen du:

$$v(t) = \sqrt{\frac{9.81 \cdot 68.1}{0.25}} \tanh\left(\sqrt{\frac{9.81 \cdot 0.25}{68.1}} t\right) = 51.6938 \tanh(0.18977t),$$

eta hori erabil dezakegu taula hau kalkulatzeko:

Eredu honen arabera, jauzilariak azkar azeleratzen du; 49.4214 m/s-ko abiadura hartzen du 10. segunduan. Kontuan izan, baita ere, denbora nahiko handia denean, 51.6983 m/s-ko azken abiadura hartzera jotzen duela. Abiadurak konstante izatera jotzen du, zeren grabitatearen indarra orekatu egingo baita airearen erresistentziarekin. Beraz, indar garbia zero izango da eta azelerazioa gelditu egingo da.  $\square$

(1.2) ekuazioari *soluzio analitikoa* (edo *zehatza*) deritzogu, zehazki betetzen duelako jatorrizko ekuazio diferentziala. Zorritzarez, badaude zehazki ebatzi ezin diren eredu matematikoak. Kasu askotan, dagoen aukera bakarra *zenbakizko soluzio* bat garatzea da, eta horrek soluzio zehatza hurbiltzen du.

*Zenbakizko metodoetan*, problema matematiko bat birformulatzen da eragiketa aritmetikoen bidez ebazteko moduan. Hori argitu daiteke, (1.1) ekuaziorako honako hau kontuan hartuz:

$$\frac{dv}{dt} \approx \frac{\Delta v}{\Delta t} = \frac{v(t_{i+1}) - v(t_i)}{t_{i+1} - t_i}, \quad (1.3)$$

Kontuan izan  $dv/dt \approx \Delta v/\Delta t$  hurbilpena dugula  $\Delta t$  finitua delako. Gogoratu kalkulu infinitesimalen hau betetzen dela:

$$\frac{dv}{dt} = \lim_{\Delta t \rightarrow 0} \frac{\Delta v}{\Delta t}.$$

(1.3) ekuazioari  $t$  uneko deribatuaren *diferentzia finituen hurbilpena* deritzo. Hura ordezkatu dezakegu, (1.1) ekuazioan hau emateko:

$$\frac{v(t_{i+1}) - v(t_i)}{t_{i+1} - t_i} = g - \frac{c_d}{m}v(t_i)^2.$$

Ekuazio hori honela ordena dezakegu:

$$v(t_{i+1}) = v(t_i) + \left[ g - \frac{c_d}{m}v(t_i)^2 \right] (t_{i+1} - t_i). \quad (1.4)$$

Kontuan izan makoen arteko gaia (1.1) ekuazio diferentzialaren eskuineko gaia dela. Hau da, ekuazio horrek ematen digu bide bat  $v$ -ren malda kalkulatzeko. Ondorioz, ekuazioa honela berriidatz dezakegu:

$$v(t_{i+1}) = v(t_i) + \frac{dv}{dt}(t_i) \cdot \Delta t. \quad (1.5)$$

Orain, ekuazio diferentzialaren bidez beste ekuazio bat lortu dugu, zeina  $t_{i+1}$  uneko abiadura kalkulatzeko erabil baitezakegu aurreko  $t$  eta  $v$ -ren balioak erabiliz. Eta lortutako balioak,  $t_{i+1}$  eta  $v_{i+1}$ , erabil ditzakegu  $v_{i+2}$  kalkulatzeko, elkarren segidan. Hau da:

$$\text{Balio berria} = \text{balio zaharra} + \text{malda} \times \text{urratsaren luzera}.$$

Metodo horri *Eulerren metodoa* deritzogu.

**1.2. adibidea.** *Aurreko adibidean, orain (1.5) ekuazioa erabiliko dugu abiadura kalkulatzeko. Urratsaren luzera 2 s izango da.*

*Ebazpena.* Hasieran  $t_0 = 0$  hartuko dugu; orduan, jauzilariaren abiadura 0 da. Aurreko adibideko parametroen balioak erabiliz,  $t_1 = 2$  s unean abiadura kalkulatu dugu:

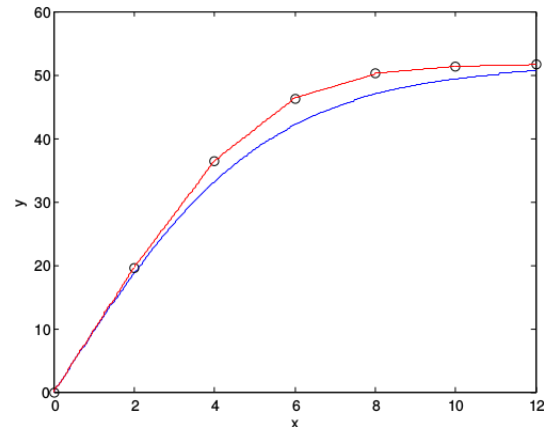
$$v = 0 + \left[ 9.81 - \frac{0.25}{68.1}0^2 \right] \times 2 = 19.62 \text{ m/s}.$$

Hurrengo tarterako ( $t = 2$ -tik  $t = 4$ -ra), kalkuluak errepikatu eta hau lortuko dugu:

$$v = 19.62 + \left[ 9.81 - \frac{0.25}{68.1}19.62^2 \right] \times 2 = 36.4137 \text{ m/s}.$$

Kalkuluaz horrela jarraituz, taula hau lortuko dugu:

$t$ (s)	$v$ (m/s)
0	0
2	19.6200
4	36.4137
6	46.2983
8	50.1802
10	51.3123
12	50.6008
$\infty$	51.6938



**1.2. irudia.** Zenbakizko soluzioa eta soluzio analitikoa.

## 2. kapitulua

# MATLABi buruzko oinarrizko nozioak

Kapitulu honetan ikusiko dugu MATLAB paketea programa matematikoen multzo bat dela eta matrizeen erabileran oinarritzen dela. MATLAB izenak MATrix LABoratory esan nahi du, zeren bere oinarrizko datu-gaia matrize bat baita. MATLAB kalkulu teknikorako lengoiaia ahaltzua da, eta oso erabilia da unibertsitateetan, Matematika, Zientzia eta Ingeniaritza ikasketetan. Pakete horrek zenbakizko programen eta marrazketa bidimentsionala eta tridimentsionala egiteko programa grafikoen bilduma zabal bat dauka. Gainera, goi-mailako lengoiaia erabiliz, programa gehigarriak idaztea onartzen du.

MATLAB paketearen instrukzioak ~~idazmakina-letraz~~ idatziko dira. Adibideetan ikusiko dugu MATLABeko lan-leiho batek (ingelesez: *command window*) erakusten diguna. Honelako `>>` ikur baten jarraian agertzen da sartutako datua edo instrukzioa. Nahi duguna idatzi eta gero, «sartu» tekla sakatu behar dugu; orduan, ordenagailuak eragiketa egingo du eta emaitza erakutsiko du, honela: `ans =`. Programarekin erabiltzeko gidak eta erreferentziako gidak datoz; eta laguntza-leihoan ere aurki daitezke instrukzioei buruzko eta haien aukerei buruzko informazio gehiago eta adibideak.

### 2.1. Eragiketa aritmetikoak

---

+	Batuketa
-	Kenketa
*	Biderketa
/	Eskuineko zatiketa
\	Ezkerreko zatiketa
^	Berreketa
pi, e, i	Konstanteak

---

### 2.1. adibidea.

```
>>(2+3*pi)/2
ans =
    5.7124
```

## 2.2. Programari gehitutako funtzioak

Jarraian, MATLAB paketeko funtzio erabilgarri batzuen zerrenda labur bat eskaintzen da:

---

abs(#)	cos(#)	exp(#)	log(#)	log10(#)	cosh(#)	sum(#)	length(#)
sin(#)	tan(#)	sqrt(#)	acos(#)	tanh(#)	floor(#)	ceil(#)	round(#)

---

Programak berak funtzio erabilgarriei buruzko informazioa eskaintzen du. Funtzio horien zerrenda izateko, idatzi «`help elfun`». Honako adibide honek argitzen du nola erabiltzen eta nola konbinatzen diren eragiketa aritmetikoak eta funtzioak.

### 2.2. adibidea.

```
>>3*cos(sqrt(4.7))
ans=
   -1.6868
```

Normalean, emaitzan bost zifra hamartar esanguratsu erakusten dira; `format long` instrukzioari esker, hamabost zifra hamartar esanguratsu arte lor ditzakegu.

### 2.3. adibidea.

```
>> format long
>> 3*cos(sqrt(4.7))
ans=
  -1.68686892236893
```

Jarraian, formatu-instrukzioekin agertzen da taula bat. Haiek `format` mota sintaxiarekin idazten dira.

Mota	Emaitza	Adibidea
<code>short</code>	Puntu finkoa 5 digiturekin	3.1416
<code>long</code>	Puntu finkoa 15 digiturekin	3.14159265358979
<code>short e</code>	Puntu higikorren formatuan 5 digiturekin	3.1416e+000
<code>long e</code>	Puntu higikorren formatuan 15 digiturekin	3.14159265358979e+000

## 2.3. Esleipen-instrukzioak

Berdintza ikurraren bidez, adierazpen baten kalkuluaren emaitzari izen bat eslei diezaiokegu.

### 2.4. adibidea.

```
>>a=3-floor(exp(2.9))
a=
    -15
```

Adierazpen baten bukaeran puntu eta koma idazten dugunean, konputagailuak dagozkion eragiketak egiten ditu, eta bere emaitza gordetzen du guk esleitutako izenpean. Hurrengo adibidean, ez da erakusten `b`-ren balioa.

### 2.5. adibidea.

```
>>b=sin(a);
>>2*b^2
ans=
    0.8457
```

## 2.4. Matrizeak

MATLAB paketea, aldagai guztiak matrizeak dira. Matrizeak zuzen sartzen dira.

**2.6. adibidea.**

```
>>A=[1 2 3;4 5 6;7 8 9]
A=
    1 2 3
    4 5 6
    7 8 9
```

«;» ikurrak matrizeko lerroak bereizten ditu, eta lerro bateko gaiak hutsune-espazio baten bidez (edo koma batez) bereizi behar ditugu. Matrizeak lerroz lerro ere sar ditzakegu.

**2.7. adibidea.**

```
>>A=[1 2 3
     4 5 6
     7 8 9]
A =
    1 2 3
    4 5 6
    7 8 9
```

Gehitutako funtzio batzuk erabiliz, matrize berezi batzuk sor ditzakegu.

**2.8. adibidea.**

```
>>Z=zeros(3,5);           (zerozko 3 x 5 dimentsioko matrize bat sortzen du)
>>X=ones(3,5);           (batezko 3 x 5 dimentsioko matrize bat sortzen du)
>>Y=0:0.5:2              (1 x 5 dimentsioko matrize hau sortzen du:)
Y=
0 0.5000 1.0000 1.5000 2.0000

>>sin(Y)                 (Y matrizeko gai bakoitzaren sinua hartuz, 1 x 5
                           dimentsioko matrize hau sortzen du:)
ans=
    1.0000 0.8776 0.5403 0.0707 -0.4161
```



`linspace(a,b,n)` instrukzioak  $a$  eta  $b$  puntuen artean tarte berak bereizitako  $n$  puntu sortzen ditu, muturrak sartuta. Aurreko adibideko  $Y$  bektore berdina ematen du, hau idazten badugu:

```
» linspace(0,2,5)
```

Ez badugu jartzen  $n$ , instrukzioak 100 puntu sortuko ditu automatikoki.

`logspace(a,b,n)` instrukzioak `linspace(10a,10b,n)` instrukzioaren berdina ematen du.

Ez badugu jartzen  $n$ , instrukzioak 50 puntu sortuko ditu automatikoki.

Matrize bateko gaiekin hainbat erataraz jola dezakegu.

## 2.9. adibidea.

```
>>A(2,3)                (A-ren (2,3) lekuko gaia aukeratzen du)
ans=
     6
>>A(1:2,2:3)           (A-ren azpimatriz bat aukeratzen du)
ans=
     2 3
     5 6
>>A([1 3],[1 3])      (beste era bat A-ren azpimatriz bat aukeratzeko)
ans=
     1 3
     7 9
>>A(2,2)=tan(7.8);    (beste balio bat esleitzen dio A-ren (2,2) lekuan
                        dagoen gaiari)

>>A(2,:)
ans=
     4 5 6
>>A(:,3)
ans=
     3
     6
     9
```

Laguntza-leihoak beste funtzio matritzialei buruzko informazioa eskaintzen digu.

## 2.5. Eragiketa matrizialak

---

+	Batuketa
-	Kenketa
/	Eskuineko zatiketa
\	Ezkerreko zatiketa
*	Biderketa
^	Berreketa
'	Irauli konjugatua

---

### 2.10. adibidea.

```
>>B=[1 2;3 4];
>>C=B'                (C B-ren iraulia da)
C=
     1     3
     2     4
>>3*(B*C)^3
ans=
    13080    29568
    29568    66840
```

Azken hori  $3(BC)^3$  da.

### 2.11. adibidea.

```
>>a=[1 2 3];
>>b=[2;4;6];
>>a*b
ans=
    28
>>b*a
ans=
     2     4     6
     4     8    12
     6    12    18
>>a*A
ans=
```

```

    30 36 42
>>c=[4 5 6]';
>>A*c
ans=
    32
    77
   122
>>A*a
??? Error using ==> mtimes
Inner matrix dimensions must agree.
>>A/pi
ans=
    0.3183    0.6366    0.9549
    1.2732    1.5915    1.9099
    2.2282    2.5465    2.8648
>>2\[1 2;3 4]
ans=
    0.5    1
    1.5    2

```

$A \setminus B$  idazten badugu,  $A^{-1}B$  emango digu, eta  $B/A$  idatziz gero,  $BA^{-1}$ . Ondorioz,  $4 \setminus 1$ -ek eta  $1/4$ -ek emaitza bera dute: 0.25.

Zer emango du MATLABek  $2/[1 \ 2;3 \ 4]$  idazten badugu?

## 2.6. Gaiez gai egiten diren eragiketak

MATLAB paketearen ezaugarri erabilgarrienetariko bat da matrize batez gaiez gai eragiten duen funtzio kopuru handi bat daukala. Aurreko adibide batean ikusi dugu  $1 \times 5$  matrize baten gai bakoitzeko sinua kalkulatzen duela. Gaiez gai egiten dira batuketa, kenketa eta eskalar bateko biderketa eragiketa matrizialak; aldiz, hori ez da gertatzen biderketa, zatiketa eta berreketa eragiketa matrizialekin. Azken hiru eragiketak gaiez gai egin ditzakegu eragiketaren aurrean puntu bat idazten badugu, hau da:  $*$ ,  $/$ ,  $\setminus$  eta  $\wedge$ . Oso garrantzitsua da jakitea nola eta noiz erabili behar ditugun eragiketa horiek, zeren gaiez gaiko eragiketak oso garrantzitsuak baitira zenbakizko programak eta grafiko programak MATLAB paketearekin eraginkorki diseinatzeko eta implementatzeko orduan.

### 2.12. adibidea.

```
>>A=[1 2;3 4];
```

```

>>A^2          (AA biderketa kalkulatzen du.)
ans=
     7 10
    15 22
>>A.^2        (A-ren gai bakoitza karratura jasotzen du.)
ans=
     1 4
     9 16
>>2.\A
ans=
    0.5000 1.0000
    1.5000 2.0000
>>cos(A./2)   (A-ren gai bakoitza 2rekin zatitu eta gero,
               kosinua kalkulatzen du.)
ans=
    0.8776 0.5403
    0.0707 -0.4161

```

### 2.6.1. MATLABek emandako funtzioak

Ikus 2.2. atalean MATLABeko oinarrizko funtzioak. Funtzio horietako propietate garrantzitsu bat da haietako gehienek bektoreen eta matrizeen gainean eragiten dutela, aurreko puntua jarri gabe.

#### 2.13. adibidea.

```

>>log(A)
ans=
    0.0000    0.6931
    1.0986    1.3863
>>B=sqrt(A)
B=
    1.0000    1.4142
    1.7321    2.0000
>>round(B)   (B-ren gaiak zenbaki oso hurbilenera biribiltzen ditu.)
ans=
     1     1
     2     2
>>ceil(B)    (B-ren gaiak goiko zenbaki oso hurbilenera biribiltzen ditu.)

```

```

ans =
     1     2
     2     2
>>floor(B) (B-ren gaiak beheko zenbaki oso hurbilenera biribiltzen ditu.)
ans=
     1     1
     1     2
>>F=[3 5 4 6 1];
>>sum(F)
ans =
    19
>>min(F),max(F),mean(F),prod(F),sort(F)
ans =
     1
ans =
     6
ans =
    3.8000
ans =
    360
ans =
     1     3     4     5     6

```

## 2.7. Grafikoak

MATLAB paketeak kurben eta gainazalen marrazketa bidimentsionalak eta tridimentsionalak egin ditzake. Laguntza-paketeen kontsulta daitezke instrukzio-grafikoen aukera eta alderdi gehigarriak.

`plot` instrukzioaz, kurba lauen grafikoak sor daitezke. Adibide honetan,  $[0, \pi]$  tarteko  $y = \cos(x)$  eta  $y = \cos^2(x)$  funtzioen grafikoak lor ditzakegu.

### 2.14. adibidea.

```

>>x=0:0.1:pi;
>>y=cos(x);
>>z=cos(x).^2;
>>plot(x,y,x,z,'o')

```

Lehenengo lerroan zehazten dira eremua eta 0.1 urratseko tamaina. Hurrengo bi lerroetan funtzioak definitzen dira. Kontuan izan lehenengo hiru lerroak puntu eta koma batez bukatzen direla; puntu eta koma hori erabiltzen da ez agertzeko pantailan  $\mathbf{x}$ ,  $\mathbf{y}$  eta  $\mathbf{z}$  matrize bakoitzeko 32 gaiak. Laugarren lerroak grafikoak ematen duen marrazketa-instrukzioa dauka. Lehenengo bi gaiak,  $\mathbf{x}$  eta  $\mathbf{y}$ ,  $y = \cos(x)$  funtzioa marrazten dute. Hirugarren eta laugarren gaiak,  $\mathbf{x}$  eta  $\mathbf{z}$ ,  $z = \cos^2(x)$  funtzioa marrazten dute. Azken gaiak, 'o', behartzen du marraztera 'o' puntu hauetan:  $(x_k, z_k)$ , non  $z_k = \cos^2(x_k)$ .

Hirugarren lerroan, oinarrizkoa da «. $\wedge$ » gaiez gaiko eragiketa-adierazlea erabiltzea; izan ere, lehenik gai bakoitzaren kosinua kalkulatzen da, eta, gero,  $\cos(\mathbf{x})$  matrizearen gai bakoitza karratura jasotzen da  $\wedge$  instrukzioa erabiliz.

`fplot` marrazketa-instrukzioa `plot` instrukziorako aukera erabilgarria da. Instrukzio horren sintaxia `fplot('izena', [a,b], n)` da. Honek `izena.m` funtzioaren grafikoa ematen digu, haren balioa zehaztuz  $[a, b]$  tarteko  $n$  puntutan. Ez bada ematen  $n$ -ren balioa,  $n = 25$  da.

**2.15. adibidea.** *Honek  $[-2, 2]$  tartean  $y = \tanh$  marrazten du:*

```
>>fplot('tanh', [-2,2])
```

`plot` eta `plot3` instrukzioak erabiltzen dira kurba parametrizatu bidimentsionalak eta tridimentsionalak marrazteko.

**2.16. adibidea.**  $c(t) = (2 \cos(t), 3 \sin(t))$  *elipsearen marrazketa,  $0 \leq t \leq 2\pi$ , instrukzio hauekin lortzen da:*

```
>>t=0:0.2:2*pi;
>>plot(2*cos(t),3*sin(t))
```

**2.17. adibidea.**  $c(t) = (2 \cos(t), t^2, 1/t)$  *kurbaren marrazketa,  $0.1 \leq t \leq 4\pi$ , instrukzio hauekin lortzen da:*

```
>>t=0.1:0.1:4*pi;
>>plot3(2*cos(t),t.^2,1./t)
```

`meshgrid` instrukzioaz marrazketa tridimentsionalak lortzeko funtzioaren eremu angeluzuzen bat zehaztu behar dugu, eta, gero, grafikoa lortu `mesh` edo `surf` instrukzioekin.

### 2.18. adibidea.

```
>>x=-pi:0.1:pi;
>>y=x;
>>[x,y]=meshgrid(x,y);
>>z=sin(cos(x+y));
>>mesh(z)
```

`hold on` instrukzioak datu eta propietate guztiekin gordetzen ditu marrazkiak, eta, orain, beste instrukzio batzuk gehi daitezke.

`hold off` instrukzioak leiho grafiko hori bukatzen du.

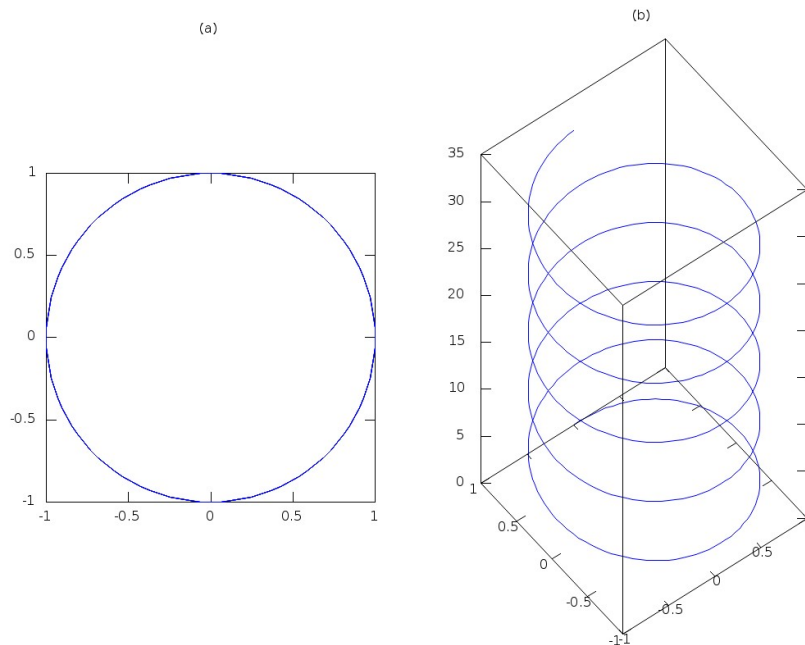
`subplot(m,n,p)` instrukzioak leiho grafikoa zatitzen du  $m \times n$  zatitan, eta  $p$ . azpicleihoa (ezkerretik eskuinera eta goitik behera zenbatuta) aukeratzen du, oraingo grafikoa marrazteko.

### 2.19. adibidea.

```
>>t=0:pi/50:10*pi;
>>subplot(1,2,1);
>>plot(sin(t),cos(t));
>>axis square
>>title('(a)')
>>subplot(1,2,2);
>>plot3(sin(t),cos(t),t);
>>title('(b)')
```

Irudia gorde dezakegu pdf eta jpg formatuetan zera idatziz:

```
>>print -dpdf irudiizena.pdf;
>>print -djpg irudiizena.jpg;
```



**2.1. irudia.** Zirkulu bidimentsional bat eta helize tridimentsional bat, bi zatitako irudi batean.

Agindua	Ezaugarria
plot	$x - y$ grafiko soila
loglog	Grafikoa logaritmikoki eskalatutako ardatzekin
semilogx	Grafikoa logaritmikoki eskalatutako $x$ ardatzarekin
semilogy	Grafikoa logaritmikoki eskalatutako $y$ ardatzarekin
polar	Grafikoa koordenatu polarretan
fplot	Funtzio marratzailea
plot3	$x - y - z$ grafiko soila
contour	Sestra-kurben grafikoa
mesh	Sare-begizko gainazala
meshc	Sare-begizko gainazala sestra-kurbekin
surf	Gainazal solidoa
surf	Gainazal solidoa sestra-kurbekin

**2.1. taula.** Funtzio marratzaileen laburpen bat.

Beste aukera batzuk ere badaude; begiratu laguntza-leihoan.



## 2.8. MATLABez programatzen: M fitxategiak

### 2.8.1. Instrukzioen fitxategiak

MATLABen instrukzioen segida bat fitxategi batean gordetzen dugunean, horrelako instrukzioen M fitxategi bat dugu. Fitxategi horiek erabilgarriak dira, behin baino gehiagotan erabiltzen badugu instrukzioen segida hori. Era hauetan egikari daiteke:

- Aginduen leihoan fitxategiaren izena idatziz.
- Fitxategia editatuz eta editorearen leihotik exekutatzuz.

**2.20. adibidea.** *Garatu instrukzioen fitxategi bat, puenting-jauzilariaren abiadura kalkulatzeko.*

*Ebazpena:*

```
g=9.81; m=68.1; t=12; cd=0.25;  
v=sqrt(g*m/cd)*tanh(sqrt(g*cd/m)*t)
```

Hau `instrukadib.m` izeneko fitxategi batean gordetzen badugu eta, gero, aginduen leihoan hau idazten badugu:

```
>>instrukadib
```

emaitza hau izango da:

```
v=  
50.6175
```

Parametro baten balioa jakin nahi badugu, adibidez `g`-rena, hau idatziko dugu:

```
>>g  
g=  
9.8100
```

## 2.8.2. Funtzio-fitxategiak

Guk funtzio berriak defini ditzakegu MATLAB paketearekin erabiltzeko. Horretarako, `m` luzapeneko testu-fitxategi bat idatzi behar dugu MATLABek berak duen editorea erabiliz: `M` fitxategi bat. Behin definituz gero, beste edozein funtzio bezala erabil daiteke.

**2.21. adibidea.**  $f(x) = -x^2/4 + x + 1$  funtzioa definituko dugu, eta `fun.m` izeneko `M` fitxategi batean gordeko dugu. Horretarako, testu-editoreaz hau idatziko dugu:

*Ebazpena:*

```
function y=fun(x)
% Funtzio honen bidez f(x)=a*x^2+b*x+c polinomioa kalkula dezakegu.
a=-1/4; b=1; c=1;
y=a*x^2+b*x+c
```

Geroago, funtzio hori baliozta dezakegu aginduen leihoan `x` argumentuaren balio baterako. Hots,  $f(2)$  kalkulatzeko, nahikoa da honako hau idaztea:

```
>>fun(2)
ans=
    2.0000
```

Gainera, funtzio honi buruz laguntza eska dezakegu:

```
>>help fun
Funtzio honen bidez f(x)=a*x^2+b*x+c polinomioa kalkula dezakegu.
```

Aldagaietarako letra desberdinak erabil ditzakegu eta funtzioari izen desberdina eman diezaiokegu, baina formatu berdina izan behar dute. Funtzio hori `fun.m` fitxategi batean gorde eta gero, MATLABeko beste edozein funtzio bezala erabil dezakegu.

```
>>cos(fun(3))
ans=
   -0.1782
```

Bestalde, `feval` instrukzioaz baliozta ditzakegu funtzioak era erabilgarri eta eraginkor batean.

```
>>feval('fun',4)
ans=
     1
```

Geroago, funtzioaren izena ahazten badugu, guk honela erabil dezakegu `lookfor` instrukzioa:

```
>>lookfor polinomioa
```

eta honako informazio hau jasoko dugu:

```
fun.m: % Funtzio honen bidez f(x)=a*x^2+b*x+c polinomioa kalkula dezakegu.
```

Aldiz, hau gertatuko da:

```
>>a
??? Undefined function or variable 'a'.
```

Horren arrazoia zera da: `a` aldagaiak  $-1/4$  balioa hartzen du M fitxategiaren barnean, hots, aldagai *lokala* da, eta aldagai lokalak ezabatu egiten dira fitxategia exekutatu eta gero. Aldiz, instrukzio M fitxategi baterako aldagaiak lan-espazioan gordetzen dira, banan-banan instrukzio bakoitza aginduen leihoan eskuz sartuko bazenu bezala; hori gertatzen da, adibidez, `instrukadib.m` M fitxategiarekin. Ariketa gisa, egiazta ezazu hori.

M fitxategiek emaitza bat baino gehiago itzul ditzakete. Adibidez, `estatistikoak.m` fitxategi honek bektore bati dagozkion gaien batezbestekoa eta desbideratze estandarra kalkulatzeko dituzte:

```
function [bbest,dest]=estatistikoak(x)
n=length(x);
bbest=sum(x)/n;
dest=sqrt(sum((x-bbest).^2/n));
```

Orduan, honela joka dezakegu:

```
>>y=[8 5 10 12 6 7.5 4];
>>[b,d]=estatistikoak(y)
b =
    7.5000
d =
    2.6049
```

Instrukzio M fitxategiak gutxi erabiltzen dira; funtzio M fitxategiak erabiliko ditugu batez ere eta, ondorioz, M fitxategiak izen laburtuarekin adieraziko ditugu hemendik aurrera horiek.

### 2.8.3. Azpifuntzioak

Funtzio batek beste funtzio batzuk barnera ditzake. Horrelako funtzioak M fitxategi berezitu moduan existitu arren, M fitxategi bakar batean sartuta egon daitezke.

Har dezagun, berriro, jauzilariaren 2.20. adibidea; kalkula dezagun abiadura ere M fitxategi honen bitartez:

```
function v=jauzilari(t,m,cd)
v=abi(t,m,cd);
end
```

```
function v=abi(t,m,cd)
g=9.81;
v=sqrt(g*m/cd)*tanh(sqrt(g*cd/m)*t)
end
```

end funtzio bakoitzaren muga adierazteko erabiliko dugu, baina ez da beharrezkoa. Fitxategi hori `jauzilari.m` izenarekin gordeko dugu. Lehenengo funtzioari *funtzio nagusia* (*main function*) deritzogu. Funtzio hori aginduen leihoan bakarrik eskura dezakegu. Beste funtzioak, kasu honetan `abi`, *azpifuntzioak* (*subfunctions*) dira. Azpifuntzio bat funtzio nagusirako M fitxategiaren barnean bakarrik da eskuragarria. Aginduen leihotik exekutatzuz, hau dugu:

```
>>jauzilari(12,68.1,0.25)
ans=
    50.6175
```

Hala ere, `abi` azpifuntzioa exekutatzan saiatzan bagara, honelako errore-mezu bat aterako da:

```
>>abi(12,68.1,0.25)
??? Undefined command/function 'abi'.
```

### 2.8.4. Input-output

Funtzioan, aginduen leihoaren bidez soilik ez dira sartzen edo ateratzen parametroen balioak edo beste informazio bat; badaude beste bi bide ere lan hori egiteko.

- `input` funtzioa. Funtzio honen bidez, erabiltzaileak balio batzuk zuzenean eman ditzake; adibidez, honela idatziz M fitxategian:

```
m=input('Masa (kg): ')
```

Lerro hau exekutatzen denean, monitorean hau agertuko da:

```
Masa (kg):
```

Erabiltzaileak balio bat sartuko du, eta orduan hura esleituko zaio `m` aldagaiari.

Funtzio honen bidez hitzak (*strings*) ere sar ditzakegu; horretarako, 's' gehituko diogu funtzioaren argumentuen zerrendari. Adibidez,

```
izena=input('Sartu zure izena: ','s')
```

- `disp` funtzioa. Funtzio honek era praktikoa ematen du parametro baten balioa edo hitza (*string*) ikusi ahal izateko. Bere sintaxia `disp(parametroa)` da.

#### 2.22. adibidea.

```
function jauzilaria
% jauzilaria: jauzilariaren abiadura kalkulatzeko era interaktibo bat
g=9.81;
m=input('Masa (kg): ');
cd=input('Erresistentzia-koefizientea (kg/m): ');
t=input('Denbora (s): ');
disp(' ')
disp('Abiadura (m/s):')
disp(sqrt(g*m/cd)*tanh(sqrt(g*cd/m)*t))
```

`jauzilaria.m` izeneko fitxategian gordez gero, eta hau idatziz:

```
>> jauzilaria
```

```
Masa (kg): 68.1
```

```
Erresistentzia-koefizientea (kg/m): 0.25
```

```
Denbora (s): 12
```

```
Abiadura (m/s):
```

```
50.6175
```

- `fprintf` funtzioa. Funtzio honen bidez kontrolatzen dugu monitorean ateratzen den informazioa. Haren sintaxia, labur emanda, hau da:

`fprintf('formatua',x,...)` non `formatua`-ren bidez adierazten baitugu nola ikusi nahi dugun `x`-ren balioa. Adibidez,

```
>>fprintf('Abiadura %8.4f m/s da\n',abiadura)
Abiadura    50.6175 m/s da
```

Normalean, hauek dira `fprintf` funtzioarekin erabiltzen diren formatuen eta kontrolen kodeak:

Formatuen kodeak	Deskripzioa
%d	Zenbaki osoaren formatuan
%e	Formatu zientifikoa e minuskularekin
%E	Formatu zientifikoa e maiuskularekin
%f	Formatu hamarrenarekin
%g	%e edo %f formatuen trinkoena
Kontrolen kodea	Deskripzioa
\n	Hasi lerro berri bat
\t	Tabuladorea

Ikus ditzagun, orain, bi adibide funtzio honen erabilera hobe ulertzeko.

### 2.23. adibidea.

```
>> fprintf('%5d %10.3f %8.5e\n',100,2*pi,pi);
    100      6.283   3.14159e+000
```

### 2.24. adibidea.

```
function fprintfdemo
x=[1 2 3 4 5];
y=[20.4 12.6 17.8 88.7 120.4];
z=[x;y];
fprintf('      x      y\n');
fprintf('%5d %10.3f\n',z);
```

Hau da horren emaitza:

```
>> fprintfdemo
```

```

x      y
1      20.400
2      12.600
3      17.800
4      88.700
5      120.400

```

### 2.8.5. Fitxategiak sortu eta fitxategietan sartu

MATLABek datu-fitxategiak irakur eta idatz ditzake. Modu errazena fitxategi bitar berezi bat erabiltzen du, MAT fitxategi deritzona; hori espresuki gauzatzen da MATLABen barnean inplementatzeko. Horrelako fitxategiak `save` instrukzioaz sortzen dira; `load` instrukzioaz fitxategian sar daitezke, eta datuak irakurri. `fitxategiizena.mat` datu-fitxategia sortzeko edo irakurtzeko, sintaxi hau erabiliko dugu:

```
save fitxategiizena var1 var2 ... varn
```

```
load fitxategiizena var1 var2 ... varn
```

Ez baditugu `var1 var2 ... varn` aldagaiak idazten, lan-espazioko aldagai guztiak gordeko dira (`save` kasuan) edo kargatuko dira (`load` kasuan).

Adibidez,

```

>> g=9.81;m=80;t=5;
>> cd=[.25 .267 .1245 .28 .273]';
>> v=sqrt(g*m/cd)*tanh(sqrt(g*cd/m)*t);

```

Orduan, abiadura- eta erresistentzia-koefizienteak gorde ditzakegu `abikoef.mat` fitxategian, hau idatziz:

```
>> save abikoef v cd
```

`clear` instrukzioa idazten badugu, lan-espazioko aldagai guztiak ezabatuko dira. Aldagai horiek berreskuratzeke, nahikoa da hau idaztea:

```
>> load abikoef v cd
```

Lan-espazioan zein aldagai ditugun jakiteko, hau idatziko dugu:

```
>> who
```

Your variables are:

```
cd    v
```

Horrelako fitxategiak oso erabilgarriak dira MATLABekin lan egiteko; baina, agian, beste programa batzuekin hobe izango da testu-fitxategi bat izatea, hau da, ASCII fitxategi bat. Horretarako, nahikoa da `save abikoef v cd -ascii` idaztea, eta `abikoef.txt` fitxategi batean gordeko du. Datuak zehaztasun bikoitzarekin gorde nahi baditugu, `-ascii -double` idatziko dugu. Adibidez,

```
>> A=[5 7 9 2;3 6 3 9];
>> save simpmatrix.txt -ascii -double
```

Horrelako fitxategi bat beste programa batzuek irakur dezakete; esate baterako, Excel-ek eta Word-ek. Bestalde, `simpmatrix.txt` ez denez MAT fitxategi bat, behin MATLABek irakurriz gero, MATLABek zehaztasun bikoitzeko matrize bat sortuko du, honela:

```
>> load simpmatrix.txt
>> simpmatrix
simpmatrix =
     5     7     9     2
     3     6     3     9
```

Gainera, `load` instrukzioa funtzio gisa erabil dezakegu. Adibidez,

```
>> A = load(simpmatrix.txt)
```

## 2.8.6. Programazio egituratua

### Begiztak eta adarkadurak

Erlazio eragileak:

Eragilea	Erlazioa	Adibidea
<code>==</code>	Zerbaiten berdin	<code>x == 0</code>
<code>~=</code>	Zerbaiten desberdin	<code>unitate ~= 'm'</code>
<code>&lt;</code>	Zerbait baino txikiago	<code>a &lt; 0</code>
<code>&gt;</code>	Zerbait baino handiago	<code>s &gt; t</code>
<code>&lt;=</code>	Zerbait baino txikiago edo berdin	<code>2.8 &lt;= a/3</code>
<code>&gt;=</code>	Zerbait baino handiago edo berdin	<code>r &gt;= 0</code>



Eragile logikoak:

~	Ez	(Egiazkoa, baldin eta soilik baldin, proposizioa faltsua bada)
&	Eta	(Egiazkoa, baldin bi proposizioak egiazkoak badira)
	Edo	(Egiazkoa, baldin bi proposizioetako bat egiazkoa bada)

Booleko balioak:

1	Egiazkoa
0	Faltsua

`for`, `if`, `switch` eta `while` egiturek MATLAB paketeen eragiten dute beste programazio-lengoaia batzuetan egiten duten antzeko era batean. Instrukzio horiek oinarriko sintaxi hau hartzen dute:

```
for (begiztaren aldagaia=begiztaren heina)
```

```
    (adierazpenak)
```

```
end
```

Baina MATLABen, batzuetan, ez dugu `for` begizta erabiltzeko beharrik. Adibidez,

```
i=0;
for t=0:0.02:50
    i=i+1;
    y(i)=cos(t);
end
```

adieraz daiteke *era bektorizatu* batean, honela:

```
t=0:0.02:50
y=cos(t);
```

MATLABek automatikoki handitzen du bektoreen eta matrizeen luzera. Adibidez,

```
t=0:0.01:5;
for i=1:length(t)
    if t(i)>1
        y(i)=1/t(i);
    end
end
```

```
    else
      y(i)=1;
    end
end
```

```
if (baldintza)
    (adierazpenak)
```

```
end
```

```
if (baldintza)
    (adierazpenak)
```

```
    else
```

```
        (adierazpenak)
```

```
end
```

Baldintza bat baino gehiago daudenean, `if...elseif...else...end` egitura erabiltzen dugu:

```
if (baldintza)
    (adierazpenak)
```

```
elseif (baldintza)
```

```
    (adierazpenak)
```

```
elseif (baldintza)
```

```
    (adierazpenak)
```

```
    .
```

```
    .
```

```
    .
```

```
else
```

```
    (adierazpenak)
```

```
end
```

“error” funtzioa erabiltzeko, if oso egokia da. Adibidez, hau dugu:

```
function f=errortest(x)
if x==0
    error('zero balioa aurkitu du');
end
f=1/x
```

Hau da, zatiketa kalkulatu da, argumentua zero ez denean bakarrik.

```
>> errortest(10)
ans =
    0.1000
```

Bestela, mezu hau emango digu:

```
>> errortest(0)
??? Error using ==> errortest
zero balioa aurkitu du
```

switch egitura if...elseif...else...end egituraren antzekoa da. Hala ere, banakako baldintzak jarri beharrean, adarkatzea test-adierazpen txikiagoetan oinarritzen da:

```
switch (test-adierazpena)
    case (balioa);
        (adierazpenak)
    case (balioa);
        (adierazpenak)
        .
        .
        .
    otherwise
        (adierazpenak)
```

```
end
```

Adibidez,

```
kalifikazioa = 'B';
switch kalifikazioa
    case 'A'
        disp('Bikain')
    case 'B'
        disp('Oso ondo')
    case 'C'
        disp('Ondo')
    case 'D'
        disp('Nahiko')
    case 'E'
        disp('Txarto')
    otherwise
        disp('Oso txarto')
end
```

Kode hori exekutatzen badugu, «Oso ondo» emango digu.

```
while egituran begizta bat errepikatzen da, baldintza logikoa egia bada, hots:
```

```
while (baldintza)
    (adierazpenak)
end
```

Adibidez,

```
x=8
while x>0
    x=x-3;
    disp(x)
end
```

Kode hori exekutatzen denean, emaitza hau da:

```
x=
```

```
8
```

```

5
2
-1

```

Hurrengo adibidean erakusten da nola sar ditzakegun begizta batzuk elkarren gainka, matrize bat sortzeko. Testu-lerroak `habia.m` fitxategi batean gordez, orduan, MATLAB paketeko lan-leihoan `habia` idazten dugun bakoitzean, `A` matrizea lortuko dugu. Ohartaraziko dugu `A` matrizeko gaiek, goiko ezkerreko izkinatik hasiz, Pascal-en triangelua sortzen dutela.

### 2.25. adibidea.

```

for i=1:5
    A(i,1)=1;
    A(1,i)=1;
end
for i=2:5
    for j=2:5
        A(i,j)=A(i,j-1)+A(i-1,j);
    end
end
A

```

`break` instrukzioa azken (`for` edo `while`) begiztatik ateratzeko, hau bete baino lehen erabiliko dugu.

### 2.26. adibidea.

```

for k=1:100
    x=sqrt(k);
    if ((k>10)&(x-floor(x)==0))
        break
    end
end
k

```

`disp` instrukzioa erabiliko dugu testu-lerro bat edo matrize bat erakusteko.

**2.27. adibidea.**

```
n=10;
k=0;
while k<=n
    x=k/3;
    disp([x x^2 x^3]);
    k=k+1;
end
```

`pause` instrukzioa erabiltzen da, programa baten exekuzioa gelditu nahi badugu leku egoki batera heltzen denean. `pause(n)` idazten badugu, `n` segundotan zehar geldi egongo da. `tic` instrukzioak oraingo denbora gordetzen du, eta `toc` instrukzioak monitorean erabilitako denbora emango digu. `beep` instrukzioak, berriz, ordenagailuaren «beep» soinu bat emango digu. Adibidez,

```
tic
beep
pause(5)
beep
toc
```

Kode hau exekutatzen denean, «beep» bat emango digu. Bost segundo geroago beste «beep» bat entzungo dugu, eta monitorean mezu hau izango dugu:

```
Elapsed time is 5.006306 seconds.
```

Programa baten exekuzioa moztu nahi badugu, nahikoa da **Ctrl+C** jotzea.

**2.28. adibidea. Egitura habiaratuak.** *Adierazpen honen bidez kalkula ditzakegu  $f(x) = ax^2 + bx + c$  ekuazio koadratikoaren erroak:*

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

*Garatu funtzio bat, formula hori implementatzeko koefizienteen balioak finkatu eta gero.*

*Ebazpena:*

```

function errokoad(a,b,c)
% errokoad: ekuazio koadratikoaren erroak
% errokoad(a,b,c): ekuazio koadratikoaren erro erreal eta konplexuak
% input:
% a=2. mailako koefizientea
% b=1. mailako koefizientea
% c=0. mailako koefizientea
% output:
% r1=1. erroaren zati errealak
% i1=1. erroaren zati irudikaria
% r2=2. erroaren zati errealak
% i2=2. erroaren zati irudikaria
if a==0
    % kasu bereziak
    if b~=0
        % erro erreal bakuna
        r1=-c/b
    else
        % soluzio nabaria
        error('Soluzio nabaria. Saiatu berriro')
    end
else
    % formula koadratikoa
    d=b^2-4*a*c; % diskriminantea
    if d>=0
        % erro errealak
        r1=(-b+sqrt(d))/(2*a)
        r2=(-b-sqrt(d))/(2*a)
    else
        % erro konplexuak
        r1=-b/(2*a)
        i1=sqrt(abs(d))/(2*a)
        r2=r1
        i1=-i1
    end
end
end

```

Orain, kalkulu hauek egin ditzakegu:

```

>> errokoad(1,1,1)
r1=
    -0.5000
i1=

```

```

    0.8660
r2=
   -0.5000
i2=
   -0.8660
>> errokoad(1,5,1)
r1=
   -0.2087
r2=
   -4.7913
>> errokoad(0,5,1)
r1=
   -0.2000
>> errokoad(0,0,0)
??? Error using ==> errokoad
Soluzio nabaria. Saiatu berriro

```

### 2.8.7. M fitxategiei funtzioak igorri

M funtzio bat garatu daiteke ekuazio berri bakoitzerako, baina aukera hobea da funtzio generiko bat diseinatzea eta analizatu nahi dugun ekuazioa argumentu gisa pasatzea. MATLABen hizkeran, horrelako funtzioek `funtzio-funtzioak` izena dute.

#### Funtzio anonimoak

*Funtzio anonimoak* erabiliz funtzio soil bat sor dezakegu, M fitxategi bat sortu gabe. Aginduen leiho baten barnean defini daitezke haiek, honelako sintaxi baten bidez:

```
feskuleku = @(argumzerrenda) adierazpena
```

non `feskuleku` = zuk erabil dezakezun funtzio-eskulekua baita, `argumzerrenda` funtzioaren argumentuen zerrenda, komaz bereizia, eta `adierazpena` MATLABeko adierazpen bat. Adibidez,

```

>> f1=@(x,y) x^2 + y^2;
>> f1(3,4)
ans =
    25

```



Beste adibide bat:

```
>> a=4;
>> b=2;
>> f2=@(x) a*x^b;
>> f2(3)
ans =
    36
```

a eta b-rako balio berriak sartzen baditugu, ez da aldatzen funtzio anonimoaren balioa:

```
>> a=3;
>> f2(3)
ans =
    36
```

Baina, guk birsortzen badugu funtzio bera, balioa aldatuko da:

```
>> f2=@(x) a*x^b;
>> f2(3)
ans =
    27
```

Antzeko emaitzak lortzeko beste bide bat da `inline` funtzioa erabiltzea. Hau da:

```
>> f1=inline('x^2 + y^2','x','y');
```

## Funtzio-funtzioak

*Funtzio-funtzioak* beste funtzio batzuetan eragiten duten funtzioak dira, eta funtzioak argumentu gisa pasatzen dira. Adibide argi bat `fplot` funtzioarena da; honek funtzioen grafikoak marrazten ditu. Hau da bere idazkera labur bat:

```
fplot(fun,lims)
```

non *fun* funtzio matematikoa marraztuko baita *lims*=[*xmin,xmax*] mugen artean. Adibidez:

```
>> vel=@(t) ...
sqrt(9.81*68.1/0.25)*tanh(sqrt(9.81*0.25/68.1)*t)
>> fplot(vel,[0 12])
```

Honek  $t = 0$ tik  $t = 12$ rako marrazki bat sortuko du.

**2.29. adibidea.** *Sortu funtzio-funtzioaren  $M$  fitzategi bat, tarte bateko funtzio baten balioen batezbestekoa kalkulatzeko. Hau ikusiko dugu  $t \geq 0$ tik  $t = 12$ rako tarteko puenting-jauzilariaren abiadura erabiliz.*

*Ebazpena.* Hau da hori lortzeko bide bat:

```
>> t=linspace(0,12);
>> v=sqrt(9.81*68.1/0.25)*tanh(sqrt(9.81*0.25/68.1)*t);
>> mean(v)
ans =
    36.0870
```

Beste bide bat:

```
function fbb = funcbb(a,b,n)
%
% funcbb: batezbestekoa kalkulatzeko du
% input:
%   a = tartearen ezker muturra
%   b = tartearen eskuin muturra
%   n = tartearen puntuen kopurua
% output:
%   fbb = funtzioaren balioen batezbestekoa
x = linspace(a,b,n);
y = func(x);
fbb = mean(y);
end

function f = func(t)
f=sqrt(9.81*68.1/0.25)*tanh(sqrt(9.81*0.25/68.1)*t);
end
```

eta gero,

```
>> funcbb (0,12,60)
ans =
    36.0127
```

## Funtzio argumentua

Emaitza berdina lortzeko beste era bat da abiadura kalkulatzeko *funtzioa argumentu moduan* sartzea; hau da:

```
function fbb = funcbb(f,a,b,n)
%
% funcbb: batezbestekoa kalkulatzeko du
% input:
%   f = balioztatu nahi dugun funtzioa
%   a = tartearen ezker muturra
%   b = tartearen eskuin muturra
%   n = tartearen zatien kopurua
% output:
%   fbb = funtzioaren balioen batezbestekoa
x = linspace(a,b,n);
y = f(x);
fbb = mean(y);
end
```

Gero, hau idatziko dugu:

```
>> abi=@(t) ...
sqrt(9.81*68.1/0.25)*tanh(sqrt(9.81*0.25/68.1)*t);
>> funcbb(abi,0,12,60)
ans =
    36.0127
```

## Funtzioen parametroak aldatzea

Ikertzeko orduan oso arrunta da ikustea nola aldatzen den menpeko aldagaiaren balioa funtzioen parametroak aldatzen ditugunean. Parametroei buruzko ikerketa horri ereduaren *sentsibilitatearen ikerketa* deritzogu.

Aurreko 2.29. adibideko `funcbb` funtzioaren sentsibilitate-analisi bat egin nahi badugu, M-file hau idatz dezakegu:

```
function fbb=funcbb(f,a,b,n,varargin)
x = linspace(a,b,n);
y = f(x,varargin(:));
fbb = mean(y);
```

Geroago, funtzio anonimo hau idatziko dugu:

```
>> abi=@(t,m,cd) sqrt(9.81*m/cd)*tanh(sqrt(9.81*cd/m)*t);
```

eta gero,  $m$  eta  $c_d$  parametroen balio desberdinetarako abiadura kalkula dezakegu. Adibidez,  $m = 68.1$  eta  $c_d = 0.25$  badira, hau ematen du:

```
>> funcbb(abi,0,12,60,68.1,0.25)
ans =
    36.0127
```

Gero,  $m = 100$  eta  $c_d = 0.28$  badira, hau ematen du:

```
>> funcbb(abi,0,12,60,100,0.28)
ans =
    38.9345
```

## 2.9. Problemak

1. Defini ditzagun  $a$ ,  $b$ ,  $c$  eta  $d$  aldagaiak, honela:  $a = 14.75$ ,  $b = -5.92$ ,  $c = 61.4$  eta  $d = 0.6(ab - c)$ . MATLABen bidez kalkula itzazu:

$$(a) \quad a + \frac{ab(a+d)^2}{c\sqrt{|ab|}} \qquad (b) \quad de^{d/2} + \frac{(ad+cd)/\left(\frac{25}{a} + \frac{35}{b}\right)}{a+b+c+d}.$$

2. Sortu tarte berak bereizitako lerro-bektore bat 16 gai dauzkana, non lehenengo gaia 4 baita eta azkena 61.
3. Sortu zutabe-bektore bat non lehenengo gaia 31 baita, hurrengo gaiak txikiagoak egiten baitira -4 gehituz, eta azkena -9 baita. (Zutabe-bektore bat sor daiteke lerro-bektore bat irauliz).
4. Sortu jarraian ematen den matrizea. Lerroak sartzean, erabili bektoreen notazioa tarte berak bereizitako bektoreak sortzeko. (Hots, ez sartu gaiak banan-banan.)

$$A = \begin{bmatrix} 0 & 1.0000 & 2.0000 & 3.0000 & 4.0000 & 5.0000 & 6.0000 \\ 3.0000 & 9.1667 & 15.3333 & 21.5000 & 27.6667 & 33.8333 & 40.0000 \\ 28.0000 & 27.7500 & 27.5000 & 27.2500 & 27.0000 & 26.7500 & 26.5000 \\ 6.0000 & 5.0000 & 4.0000 & 3.0000 & 2.0000 & 1.0000 & 0 \end{bmatrix}.$$

Lehenengo bi ariketetan bi puntuen ikurra (:) erabiliz, egin hau:

- (a) Sortu 4 gaiko lerro-bektore bat,  $\mathbf{va}$  izenekoa,  $A$ -ren bigarren lerroko azken lau gaiak dauzkana.
- (b) Sortu 4 gaiko zutabe-bektore bat,  $\mathbf{vb}$  izenekoa,  $A$ -ren seigarren zutabeko gaiak dauzkana.
- (c) Sortu  $3 \times 4$  matrizea,  $\mathbf{B}$  izenekoa,  $A$  matrizearen 1., 2. eta 4. lerroetako eta 1., 2. eta 7. zutabeetako gaiak erabiliz.
- (d) Sortu  $2 \times 3$  matrizea,  $\mathbf{C}$  izenekoa,  $A$  matrizearen 2. eta 4. lerroetako eta 2., 5. eta 6. zutabeetako gaiak erabiliz.
5. Demagun  $y = \frac{(x^3 + 1)^2}{x^2 + 2}$  funtzioa. Kalkula ezazu  $y$ -ren balioa  $x$ -ren balio hauetarako: -1.6, -1.2, -0.8, -0.4, 0, 0.4, 0.8, 1.2. Ebatzi problema hau  $x$  bektore bat sortuz eta gero  $y$  bektore bat sortuz, gaiez gaiko kalkuluak erabiliz. Bikote horietarako egin grafiko bat, non puntuak izartxo batez adierazita agertzen baitira, eta puntuak lotuta lerro beltzez. Etiketatu ardatzak.
6. Defini dezagun  $a = 0.8$  eskalarra eta  $x = -3, -2.8, -2.6, \dots, 2.6, 2.8, 3$ . Orduan, erabili aldagai horiek honela,  $y$  kalkulatzeko:  $y = \frac{8a^2}{x^2 + 4a^2}$ . Marraztu  $y$ ,  $x$ -rekiko.

7. Erabili MATLAB frogatzeko  $\sum_{n=0}^{\infty} (-1)^n \frac{1}{2n+1}$  serie infinituaren batura,  $\pi/4$ ra jotzen duena. Egin batura, hori kalkulatzuz,  $n$ -ren balio hauetarako:

- (a)  $n = 100$
- (b)  $n = 1000$
- (c)  $n = 5000$ .

Sortu  $n$  izeneko bektore bat alde bakoitzean, non lehenengo gaia 0 baita, gehikuntza 1, eta azken gaia 100, 1000 edo 5000 baita. Orduan, erabili gaiez gaiko kalkulua bektore bat sortzeko, non gaiak  $(-1)^n \frac{1}{2n+1}$  baitira. Azkenik, erabili `sum` funtzioa seriearen gaiak batzeko. Konparatu (a), (b) eta (c) ataletan lortutako balioak  $\pi/4$  balioarekin.

8. St. Louis-eko Sarrerako Arkua osatzen da ekuazio honen arabera:

$$y = 693.8 - 68.8 \cosh\left(\frac{x}{99.7}\right)$$

(oinetan neurtuta). Egin Arkuaren marrazki bat, non  $-299.25 \leq x \leq 299.95$  oin.



Daniel Schwen CC BY-SA 3.0

9. Bektore bat honela emango dugu:  $\mathbf{x}=[15 \ 85 \ 72 \ 59 \ 100 \ 80 \ 44 \ 60 \ 91 \ 38]$ . Baldintzazko egiturak eta begiztak erabiliz, idatzi programa bat 59 baino handiagoak diren  $\mathbf{x}$ -ko gaien batezbestekoa kalkulatzeko.
10. Kosinu-funtzioaren balioa kalkula dezakegu serie infinitu honen bidez:

$$\cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots$$

Sortu  $M$  fitxategi bat, formula hori implementatuz, hau egiteko:  $\cos x$ -ren balioak kalkulatu, eta monitorean seriearen gai bakoitza erakutsi gehitzen dituen neurrian. Alegia, kalkulatu balio hauek elkarren segidan eta erakutsi monitorean:

$$\begin{aligned} \cos x &= 1 \\ \cos x &= 1 - \frac{x^2}{2!} \\ \cos x &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} \\ &\vdots \end{aligned}$$

eta segi horrela zuk aukeratutako ordenaraino. Aurreko balio hurbildu bakoitzerako, kalkulatu ehuneko errore erlatiboa, eta erakutsi monitorean hau kontuan hartuz:

$$\text{errorea} = \frac{\text{egiazkoa} - \text{seriearen hurbilpena}}{\text{egiazkoa}} \times 100.$$

Proba bezala, erabili programa  $\cos(1.5)$  kalkulatzeko, batuketan 8 batugai sartuta; hau da, seriea garatuz  $x^{14}/14!$  batugairaino.

11. Planoko puntu bat kokatzeko, bi koordenatu behar ditugu:

- Koordenatu kartesiarretan,  $(x, y)$  ardatzekiko distantzia horizontala eta bertikala.
- Koordenatu polarretan,  $(r, \theta)$  erradioa eta angelua.

Koordenatu polarrak ezagutuz gero, nahiko erraza da koordenatu kartesiarrek kalkulatzea. Aldiz, alderantziz ez da hain erraza. Erradioa adierazpen honen bidez kalkulatu dugu:

$$r = \sqrt{x^2 + y^2}.$$

Puntua lehenengo edo laugarren koadrantean dagoenean (hots,  $x > 0$ ), erraza da  $\theta$  kalkulatzeko formula:

$$\theta = \arctan\left(\frac{y}{x}\right).$$

Beste kasuetarako agertzen da zailtasuna. Honako taula honek laburtzen ditu aukerak:

$x$	$y$	$\theta$
$<0$	$>0$	$\arctan(y/x) + \pi$
$<0$	$<0$	$\arctan(y/x) - \pi$
$<0$	$=0$	$\pi$
$=0$	$>0$	$\pi/2$
$=0$	$<0$	$-\pi/2$
$=0$	$=0$	$0$

Idatzi M fitxategi ondo egituratu bat  $r$  eta  $\theta$  kalkulatzeko,  $x$ -ren eta  $y$ -ren funtzioan. Adierazi  $\theta$ -rako azken emaitzak, gradutan. Proba ezazu zure programa, kasu hauek balioztatuz:

$x$	$y$	$r$	$\theta$
1	0		
1	1		
0	1		
-1	1		
-1	0		
-1	-1		
0	0		
0	-1		
1	-1		

12. Garatu M fitxategi bat, non 0tik 100erako zenbakizko balio bat igortzen badiogu, hark taula honi jarraituz letrazko kategoria bat bueltatuko baitigu:

Letra	Irizpidea
A	$90 \leq \text{zenbakizko balioa} \leq 100$
B	$80 \leq \text{zenbakizko balioa} < 90$
C	$70 \leq \text{zenbakizko balioa} < 80$
D	$60 \leq \text{zenbakizko balioa} < 70$
E	zenbakizko balioa $< 60$

13. Manning-en ekuazio hau erabil dezakegu uraren abiadura kalkulatzeko ubide ireki laukizuzen batean:

$$U = \frac{\sqrt{S}}{n} \left( \frac{BH}{B + 2H} \right)^{2/3},$$

non  $U$  = abiadura (m/s),  $S$  = ubidearen malda,  $n$  = marruskadura koefizientea,  $B$  = zabalera (m) eta  $H$  = sakonera (m). Datu hauek eskuragarriak dira bost ubidetarako:

$n$	$S$	$B$	$H$
0.035	0.0001	10	2
0.020	0.0002	8	1
0.015	0.0010	20	1.5
0.030	0.0007	24	3
0.022	0.0003	15	2.5

Idatzi M fitxategi bat ubide horietarako abiadura kalkulatzeko duena. Sartu balio horiek matrize batean, non zutabe bakoitzak parametro bat adierazten baitu eta lerro bakoitzak ubide bat. Monitorean taula itxura izan behar du M fitxategiaren irteerak; baina handituta, bosgarren zutabe batekin, non ubide bakoitzari dagokion abiadura agertzen baita, zutabeak etiketatzeko taularen goiburuak sartuta.

14. Habe baten desplazamendua honelako funtzio batek neurtzen du:

$$y(x) = \frac{-5}{6} [\langle x-0 \rangle^4 - \langle x-5 \rangle^4] + \frac{15}{6} \langle x-8 \rangle^3 + 75 \langle x-7 \rangle^2 + \frac{57}{6} x^3 - 238.25x,$$

non  $x$  habearen zeharkako distantzia baita, eta *singulartasun-funtzioa* honela definitzen baita:

$$\langle x - a \rangle^n = \begin{cases} (x - a)^n, & x > a \text{ denean} \\ 0, & x \leq a \text{ denean} \end{cases}$$

Garatu M fitxategi bat,  $x$  abszisa eta  $y$  ordenatua dituen grafiko bat egiteko.

15. Likido baten  $B$  bolumena,  $r$  erradiodun eta  $L$  luzeradun zilindro horizontal baten barnean, likidoaren  $h$  sakoneraren menpean dago formula honen bidez:

$$B = \left[ r^2 \arccos \left( \frac{r-h}{r} \right) - (r-h) \sqrt{2rh - h^2} \right] L.$$

Garatu M fitxategi bat, bolumena/sakonera grafiko bat marrazteko. Proba ezazu programa,  $r = 2$  m eta  $L = 5$  m denean.



16. «Zatitu eta erdibanatu» metodoa,  $a$  zenbaki positibo baten erro karratua hurbiltzeko metodo zaharra, honela formula daiteke:

$$x = \frac{x + a/x}{2}.$$

Idatzi `while ... break` begiztako egituran oinarritutako M fitxategi ondo egituratu bat, algoritmo hori inplementatzeko. Erabili koska egoki bat, egitura argi bat izateko moduan. Urrats bakoitzean, balioetsi zure hurbilpenaren errorea adierazpen honen bidez:

$$\varepsilon = \left| \frac{x_{berria} - x_{zaharra}}{x_{berria}} \right|.$$

Errepikatu begizta  $\varepsilon$ -ren balioa, zuk emandako  $\varepsilon_e$  tolerantzia bat baino txikiago izan arte. Diseinatu zure programa, emaitza eta errorea emateko. Ziurta ezazu zuk kalkula ditzakezula zero edo zero baino txikiagoak diren zenbaki guztien erro karratuak. Azken kasu horretan, erakutsi emaitza zenbaki irudikari bat bezala. Esate baterako, monitorean -4ren erro karratuaren itzultzeak  $2i$  izan behar du. Proba ezazu zure programa  $a = 0, 2, 4$  eta  $-9$  balioztatuz,  $\varepsilon_e = 10^{-4}$  hartuz.

17. *Zatikako funtzioak*, batzuetan, oso erabilgarriak dira menpeko aldagaiaren eta aldagai askearen arteko erlazioa ekuazio bakar batean adierazi ezin dugunean. Adibidez, honela deskriba liteke espazio-suziri baten abiadura:

$$v(t) = \begin{cases} 11t^2 - 5t, & 0 \leq t \leq 10 \\ 1100 - 5t, & 10 < t \leq 20 \\ 50t + 2(t - 20)^2, & 20 < t \leq 30 \\ 1520e^{-0.2(t-30)}, & t > 30 \\ 0, & \text{beste kasu batean.} \end{cases}$$

Garatu M fitxategi bat  $v$  kalkulatzeko  $t$ -ren funtzioan. Orduan, erabili  $v = v(t)$  funtzio hori, grafiko bat sortzeko  $t = -5$ etik  $t = 50$  arte ( $t$  abszisa eta  $v$  ordenatua erakutsiz).



## 3. kapitulua

# Ordenagailuaren aritmetika eta errorearen analisisia

Soluzio bat lortzeko erabiltzen dugun ordenagailua tresna inperfektua da. Izan ere, mugatuta dago hark duen gaitasuna zenbakiak zehaztasunez adierazteko. Ondorioz, makinak berak lortzen dituen emaitzek erroreak dituzte.

Bestalde, metodo hurbilduak erabiltzen baditugu, erroreak sortzen dira. Adibidez, abiaduraren deribatua kalkulatzeko (punting-jauzilariaren adibidean bezala) diferentzia finituak erabiltzean:

$$\frac{dv}{dt} \approx \frac{\Delta v}{\Delta t} = \frac{v(t_{i+1}) - v(t_i)}{t_{i+1} - t_i}.$$

Auzia da, beraz, nola maneiatu horrelako ziurgabetasuna.

### 3.1. Algoritmoak eta erroreak

Ordenagailu bat erabiltzen denean problema baten zenbakizko soluzio bat lortzeko, programak gauzatzen ditu erabilitako zenbakizko metodoari elkartutako eragiketak. Zenbakizko metodo batzuk errazak dira inplementatzeko, baina batzuetan zenbakizko prozedurak zailak dira programatzeko.

Zenbakizko metodo bat programatu baino lehen, oso onuragarria da zenbakizko metodoa inplementatzeko jarraitu behar ditugun urrats guztiak planifikatzea. Horrelako plan bati *algoritmo* deritzogu, eta soluziora heltzeko urratsez urratseko instrukzioen bilduma da. Algoritmoak xehetasun-maila batzuetan idatz daitezke. Adibidez, jo dezagun  $ax^2 + bx + c = 0$

ekuazio koadratikoa ebatzi nahi dugula, soluzio hauek kalkulaturaz, algoritmo batez:

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}.$$

### 3.1. algoritmoa. Ekuazio koadratiko baten erro errealak ebazteko algoritmoa.

Ekuazio koadratikoaren  $a$ ,  $b$  eta  $c$  konstanteak emanda daude.

1. Kalkulatu diskriminantearen balioa,  $D = b^2 - 4ac$ .
2.  $D > 0$  bada, kalkulatu bi erroak goiko adierazpenak erabiliz.
3.  $D = 0$  bada, kalkulatu  $x = -b/2a$  erroa, eta erakutsi mezu hau: «Ekuazioak erro bakar bat du».
4.  $D < 0$  bada, erakutsi mezu hau: «Ekuazioak ez du erro errealik».

Behin algoritmoa asmatuz gero, inplementa daiteke ordenagailuko programa batean.

#### 3.1.1. Erroreak

Zenbakizko kalkuluaren praktikan, kontuan hartu behar dugu ordenagailuak lortutako emaitzak ez direla soluzio matematiko zehatzak. Hainbat faktorerengatik, zenbakizko soluzioaren zehaztasuna txikitu daiteke, eta zailtasun hori ulertzeak maiz gida gaitzake zenbakizko algoritmo egokiak eraikitzen.

**3.1. Definizioa.** Demagun  $\hat{p}$  balioa  $p$ -ren hurbilpen bat dela. Hurbilpenaren **errore absolutua**  $E_p = |p - \hat{p}|$  da eta **errore erlatiboa**  $R_p = |p - \hat{p}|/|p|$ ,  $p \neq 0$ .

**3.1. adibidea.** Izan bitez  $x = 3.141592$  eta  $\hat{x} = 3.14$ ; orduan, errore absolutua hau da:

$$E_x = |x - \hat{x}| = |3.141592 - 3.14| = 0.001592$$

eta errore erlatiboa hau da:

$$R_x = \frac{|x - \hat{x}|}{|x|} = \frac{0.001592}{3.141592} = 0.00507.$$

Puntu higikorreko adierazpenetan, nahiago da errore erlatiboaz lan egitea, hau mantisarekin zuzen erlazionatuta baitago.

## Errore motak

Atal honetan aztertuko dugu zein diren erroreen iturriak.

### 1. Ebatzi behar den problemako erroreak.

*Eredu matematikoaren hurbiltze-erroreak* izan daitezke. Adibidez, sarritan zeruko gorputzak hurbiltzen ditugu esferen bidez haien propietateak kalkulatzeko; esate baterako, jakiteko asteroide batek talka egingo duen planeta batekin emandako data bat baino lehen.

Problema bateko beste ohiko errore-iturri bat *datuetako erroreak* dira. Hau neurri fisikoetatik etor daiteke, horiek ez baitira inoiz zeharo zehatzak.

### 2. Hurbiltze-erroreak.

Horrelako erroreak sortzen dira benetako funtzio bat balioztatzeko adierazpen hurbildu bat erabiltzen denean. Sarritan, bi errore mota hauek aurkitzen ditugu:

- *Diskretizazio-erroreak* diskretizatze prozesu batetik sortzen dira; esate baterako, interpolazioa, diferentziazioa eta integrazioa.
- *Konbergentzia-erroreak* metodo iteratiboetan sortzen dira. Esate baterako, problema ez-linealak orokorki prozesu iteratibo baten bidez ebazten dira. Horrelako prozesuak konbergitu egingo dira infinitu iterazioetan, baina moztu egiten da iterazio kopuru finitu bat egin ondoren. Metodo iteratiboak aljebra linealean sortzen dira.

### 3. Biribiltze-erroreak.

Ordenagailu batez (edo kalkulagailu batez) zenbaki errealekin egindako edozein kalkuluk biribiltze-errore bat izaten du. Hori gertatzen da edozein ordenagailutan zenbaki errealeen zehaztasun finituko adierazpenagatik; horrek eragiten du bai datuen adierazpenetan, bai ordenagailuaren aritmetikan.

Diskretizazio- eta konbergentzia-erroreak balioztatu behar dira erabilitako metodoaren analisi batez. Aldiz, biribiltze-erroreek egitura leunago bat dute, eta hori, batzuetan, ustiatu daiteke. Gure oinarrizko hipotesia izango da hurbiltze-erroreen tamaina biribiltze-erroreena baino handiagoa dela kalkulu erreal arrakastatsuetan.

#### 3.1.2. Diskretizazio-erroreak

Ikus dezagun, adibide bat argitzeko, diskretizazio-erroreen portaera.

**3.2. adibidea.** Izan bedi  $x = x_0$  puntu bateko  $f(x)$  funtzio baten  $f'(x_0)$  deribatua hurbiltzearen problema. Esate baterako, izan bedi  $f(x) = \cos(x)$  eta  $x_0 = 0.5$ . Orduan,  $f(x_0) = \cos(0.5) = 0.87758\dots$

Jo dezagun egoera berezi batean ezin dugula  $f'(x_0)$  kalkulatu, edo bere kalkulua konputazionalki oso garestia dela, baina  $f(x)$  balioa kalkulatu dezakegula  $x_0$ -tik  $x$  hurbil badago.

*Ebazpena.* Algoritmo erraz bat eraiki dezakegu Taylorren seriea erabiliz. Hots, izan bedi  $h > 0$  balio txiki bat; orduan:

$$f(x_0 + h) = f(x_0) + hf'(x_0) + \frac{h^2}{2}f''(x_0) + \frac{h^3}{6}f'''(x_0) + \frac{h^4}{24}f''''(x_0) + \dots$$

Beraz,

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} - \left( \frac{h}{2}f''(x_0) + \frac{h^2}{6}f'''(x_0) + \frac{h^3}{24}f''''(x_0) + \dots \right).$$

Gure algoritmo  $f'(x_0)$  hurbiltzeko, hau kalkulatzeko datza:

$$\frac{f(x_0 + h) - f(x_0)}{h}.$$

Lortutako hurbilpenak errore hau dauka:

$$\left| f'(x_0) - \frac{f(x_0 + h) - f(x_0)}{h} \right| = \left| \frac{h}{2}f''(x_0) + \frac{h^2}{6}f'''(x_0) + \frac{h^3}{24}f''''(x_0) + \dots \right|.$$

Baldin ezagutzen badugu  $f''(x_0)$  eta hura ez bada zero, orduan  $h$  txiki baterako errore hori, *diskretizazio-errorea*, hurbil dezakegu adierazpen honen bidez:

$$\left| f'(x_0) - \frac{f(x_0 + h) - f(x_0)}{h} \right| \approx \left| \frac{h}{2}f''(x_0) \right|.$$

Hurbiltzearen errorea taula honetan ikus daiteke,  $h$ -ren balio hauetarako:

$h$	Errorea
0.1	4.3044e-2
0.01	4.3799e-3
0.001	4.3871e-4
1.e-4	4.3878e-5
1.e-5	4.3879e-6
1.e-6	4.3887e-7
1.e-7	4.3963e-8

Ikus daitekeenez, ematen du diskretizazio errorea txikiago egiten dela  $h$ -rekin batera. Bigarren deribatua kalkulatu,  $\frac{1}{2}f''(x_0) \approx -0.43879$ . Hots,  $0.43879h$ -k erroreak zehaztasunez hurbiltzen dituela ematen du.  $\square$

### 3.1.3. Biribiltze-errorearen eragin kaltegarria

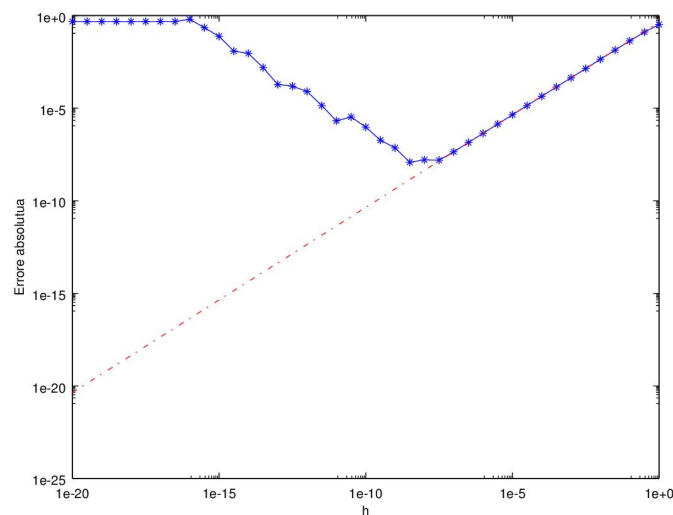
Itzul gaitezen 3.2. adibidera, ikusteko zelako ondorioak ekar ditzakeen biribiltze-erroreak.

**3.3. adibidea.** 3.2. adibideko taulako zenbakiak iradokitzen dute algoritmoak edozein zehaztasun lor dezakeela, betiere  $h$  nahiko txiki bat hartuz. Esate baterako, demagun hau nahi dugula:

$$\left| -\sin(0.5) - \frac{\cos(0.5 + h) - \cos(0.5)}{h} \right| < 10^{-10}.$$

*Ebazpena.* Adibide haren arabera, ematen du nahikoa dela  $h \leq 10^{-10}/0.43879$  hartzea. Baina honako taula honen arabera, zer gertatzen den ikus dezakegu:

$h$	Errorea
1.e-8	1.6207e-08
1.e-9	7.1718e-08
1.e-10	9.5990e-07
1.e-11	2.0701e-06
1.e-12	7.9786e-05
1.e-13	1.9081e-04
1.e-15	0.075686
1.e-16	0.63080



**3.1. irudia.** Kurba urdin jarraituak diskretizazio-errorea ematen du, eta puntu-marra zuzen gorriak biribiltze-errore gabeko errore hori (MATLABeko `loglog` funtzioak egindako grafikoa).

Diskretizazio-errorea txikiago egiten da era ordenatu batean  $h$  txikiago egiten den ahala, eta biribiltze-errorea menperatzen du  $h$  handi samarra den bitartean. Baina,  $h$   $10^{-8}$  baliora gutxi gorabehera jaisten denean, diskretizazioaren erroreak oso balio txikia hartzen du, eta biribiltze-errorea diskretizazioarena menperatzen hasten da (errorerik handiena bihurtzen baita).  $\square$

### 3.1.4. Algoritmo bat ebaluatzeko irizpideak

Algoritmo baten kalitatearen eta erabilgarritasunaren ebaluazio bat irizpide batzuetan oinarri dezakegu.

#### Zehaztasuna

Hau elkarri lotuta dago errore motekin (geroago aztertuko dugu). Algoritmo baten ebaluazioan zenbakizko algoritmo baten zehaztasuna parametro erabakigarria dela kontuan hartzea oso garrantzitsua da, eta zenbakizko algoritmo bat diseinatzean beharrezkoa da kalkulu bat burutzean zelako errore tamaina espero dezakegun jakitea; ikus 3.2. adibidea.

#### Eraginkortasuna

Propietate teoriko handiak dituen zenbakizko algoritmo bat alferra da, bere kalkulua egiteko denbora gehiegi erabiltzen badu. Eraginkortasuna CPU denboraren eta memoria-espazioaren eskakizunen menpe dago. Algoritmo baten implementazio xehetasun batzuek garrantzi handia izan dezakete eraginkortasunean. Beste propietate teoriko batzuk ere eraginkortasunaren adierazle izan daitezke; esate baterako, konbergentzia-ordena (ratioa).

Sarritan egin behar diren eragiketa aritmetikoen estimazio batek algoritmoaren eraginkortasunaren informazioa ematen digu. Adibidez, kalkulu arruntean,  $n$  mailako polinomio baten balioa kalkulatzeko  $n^2$  ordenako eragiketak egiten dira:

$$p_n(x) = c_0 + c_1x + \dots + c_nx^n.$$

Aldiz, Horner-en metodoa erabiltzen badugu,  $n$  ordenako eragiketak:

$$p_n(x) = (\dots((c_nx + c_{n-1})x + c_{n-2})x \dots)x + c_0.$$

#### Sendotasuna

Sarritan, zenbakizko softwarea idazteko ahalegin handiena ez da gastatzen algoritmo baten mamia implementatzeko, baizik eta ziurtatzeko berak lan egingo duela baldintza ahul guztien



menpean. Hots, errutinak, tolerantzia-maila baten barne, emaitza zuzen bat eman beharko luke, edo, huts egitekotan, dotoreziaz (esate baterako, abisu batez) egin beharko luke.

## 3.2. Puntu higikorreko sistemak

Euskaraz dezimalak komaz banantzen dira, baina testu honetan koma eta programazio-adibideetan puntua ez erabiltzearen puntuz bereiziko dira dezimalak. Hori dela eta, ingelesez *floating point* dena *puntu higikor* deituko dugu hemendik aurrera testu honetan.

Puntu higikorreko sistema bat lau baliok definitzen dute:  $(\beta, t, L, U)$ , non

$$\begin{aligned}\beta &= \text{zenbaki-sistemaren oinarria;} \\ t &= \text{zehaztasuna (digitu kopurua);} \\ L &= e \text{ berretzailearen behe-bornea;} \\ U &= e \text{ berretzailearen goi-bornea.}\end{aligned}$$

Orokorrean, zera dugu:

$$\text{fl}(x) = \pm \left( \frac{\tilde{d}_0}{\beta^0} + \frac{\tilde{d}_1}{\beta^1} + \dots + \frac{\tilde{d}_{t-1}}{\beta^{t-1}} \right) \times \beta^e,$$

$\beta$  oinarria 1 baino handiagoa den zenbaki oso bat da, eta  $\tilde{d}_i$  zenbaki osoak  $0 \leq \tilde{d}_i \leq \beta - 1$  tartean daude.  $\text{fl}(x)$  zenbakia  $x$ -ren hurbilpen bat da.

Adierazpen horren bakartasuna ziurtatzeko,  $\tilde{d}_0 \neq 0$  betetzeko normalizatzen dugu,  $e$  berretzailea doituz aurreko zeroak kentzeko. Kontuan hartu,  $\beta = 2$  kasuan izan ezik,  $\tilde{d}_0$  finkatu gabe dagoela, eta, ondorioz, gorde egin behar da, baita ere. Gainera,  $e$  berretzaileak  $L \leq e \leq U$  bete behar du.

### 3.2.1. Inaustea eta biribiltzea

Nola gorde  $x = \pm(d_0.d_1d_2d_3 \dots d_{t-1}d_t d_{t+1} \dots) \times \beta^e$  zenbakia  $t$  digitu bakarrik erabiliz? Bi estrategia hauek erabil daitezke:

- inaustea:  $d_t, d_{t+1}, d_{t+2}, d_{t+3} \dots$ , digituak baztertzea,  $\tilde{d}_i = d_i$  izanik, eta

$$\text{fl}(x) = \pm(d_0.d_1d_2d_3 \dots d_{t-1}) \times \beta^e;$$

- biribiltzea: aztertu  $d_t$  hurbilpena zehazteko:

$$\text{fl}(x) = \begin{cases} \pm(d_0.d_1d_2d_3 \dots d_{t-1}) \times \beta^e, & d_t < \beta/2 \text{ bada,} \\ \pm(d_0.d_1d_2d_3 \dots d_{t-1} + \beta^{1-t}) \times \beta^e, & \text{bestela.} \end{cases}$$

**3.4. adibidea.** *Inaustearen eta biribiltzearen emaitzak,  $\beta = 10$  eta  $t = 3$  direnean.*

$x$	inausia	biribildua
6.552	6.55	6.55
-6.552	-6.55	-6.55
6.557	6.55	6.56
-6.557	-6.55	-6.56
6.592	6.59	6.59
6.595	6.59	6.60

**3.5. adibidea.** *Izan bedi sistema hamartarra, hots,  $\beta = 10$ , eta sistema horretako zenbaki hau:*

$$\frac{8}{3} = 2.6666\dots = \left( \frac{2}{10^0} + \frac{6}{10^1} + \frac{6}{10^2} + \frac{6}{10^3} + \frac{6}{10^4} + \dots \right) \times 10^0.$$

*Zenbaki horrek 10 oinarriko berretura-serie infinitu bat ematen du, nahiz eta  $\beta = 3$  oinarrian seriea finitua izan (hain zuzen,  $2 \cdot 3^0 + 2 \cdot 3^{-1}$ ).*

*Adierazi zenbaki hamartar hori  $t = 4$  erabiliz.*

*Ebazpena.* Inausiz, zera dugu:

$$\frac{8}{3} \approx \left( \frac{2}{10^0} + \frac{6}{10^1} + \frac{6}{10^2} + \frac{6}{10^3} \right) \times 10^0 = 2.666 \times 10^0.$$

Bestalde, biribiltzean  $d_t = 6 \geq \beta/2 = 5$  denez, zenbaki inausiari  $\beta^{1-t} = 10^{-3}$  gehituta  $2.667 \times 10^0$  ematen du.

Puntu higikorreko adierazpen hau ez da bakarra; adibidez, hau dugu:

$$2.667 \times 10^0 = 0.2667 \times 10^1.$$

Beraz, adierazpena normalizatzen dugu  $d_0 \neq 0$  eta  $1 \leq d_0 \leq 9$ ,  $0 \leq d_i \leq 9$ ,  $i = 1, \dots, t-1$  nahitaez hartuz, edozein anbiguetate ezabatuz. Horrela,  $x = 2.6666\dots$  biribilduz,  $\text{fl}(x) = 2.667 \times 10^0$  dugu.  $\square$

0 zenbakia ezin da adierazi era normalizatu batean, eta  $\pm\infty$  adierazten dira biten konbinazio berezi batzuen bidez, zenbaki-sistemaren arabera.

**3.6. adibidea.** *Izan bedi (jostailuzko) puntu higikorreko sistema hamartar bat  $t = 4$ ,  $U = 1$  eta  $L = -2$  hartuz. Zein dira zenbaki handienak eta txikienak? Zenbat zenbaki ditu?*

*Ebazpena.* Horrela, hain zuzen, 2.666 zenbaki hamartarra adierazgarria da bere mantisak lau digitu dituelako eta  $L \leq e \leq U$ .

Kasu honetan, zenbaki handiena  $99.99 < 10^2 = 100$  da, txikiena  $-99.99 > -100$ , eta zenbaki positibo txikiena  $10^{-2} = 0.01$ .

Zenbat zenbaki desberdin ditugu? Lehenengo digituak 9 balio desberdin har ditzake, beste digituetako bakoitzak 10 balio. Hortaz,  $9 \times 10 \times 10 \times 10 = 9000$  zatiki normalizatu desberdin posible daude. Berretzaileak har ditzake  $U - L + 1 = 4$  balio desberdin, beraz, guztira  $4 \times 9000 = 36000$  zenbaki positibo desberdin daude. Beste horrenbeste zenbaki negatibo daude, eta, gainera, 0 zenbakia dugu. Ondorioz, 72001 zenbaki desberdin daude puntu higikorreko sistema honetan.  $\square$

### 3.2.2. Biribiltzearen unitatea

Puntu higikorreko adierazpenari dagokion errore erlatiboa hau da:

$$\frac{|\text{fl}(x) - x|}{|x|}.$$

Zenbaki batetik hurrengorako mantisen arteko distantzia  $\beta^{1-t}$  da. Askotan,  $t - 1$  zenbakiari *digitu esanguratsuen kopurua* izendatzen diogu.

Ondorioz,  $(\beta, t, L, U)$  sistema orokor baterako *biribiltzearen unitatea*  $\eta = \frac{1}{2}\beta^{1-t}$  da eta  $\frac{|\text{fl}(x) - x|}{|x|} \leq \eta$  betetzen da. Adibidez,  $\beta = 10$ ,  $t = 4$  eta  $x = 12743.25$  badira,  $\text{fl}(x) = 1.274 \times 10^4$  dugu eta, ondorioz,

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{|1.274 \times 10^4 - 1.274325 \times 10^4|}{|1.274325 \times 10^4|} = \frac{0.000325 \times 10^4}{1.274325 \times 10^4} = \frac{0.000325}{1.274325} < \frac{10^{-3}}{2} = \eta.$$

Zehaztasun bikoitzeko sistemaren kasuan, non  $t - 1 = 52$  baita,  $\eta = \frac{1}{2}2^{-52} = 1.1 \times 10^{-16}$  dugu.

## 3.3. Zenbaki bitarrak

Nahiz eta gizakiok zenbaki-sistema hamartarra erabili kalkulu aritmetikoetan, ordenagailu gehienek zenbaki-sistema bitarra erabiltzen dute. Ordenagailuak datu guztiak zenbaki bitar bihurtzen ditu. Kalkulu aritmetikoak 2 oinarrian egiten ditu, eta gero 10 oinarria itzultzen ditu. Zehaztasuneko bederatzi zifra dezimal dauzkan ordenagailu batek emaitza hau eman zuen:

$$\sum_{i=1}^{100000} 0.1 = 9999.99447.$$

Batura 10000 izan behar zenez, gure lanetariko bat izango da hura gertatzeko arrazoia jakitea.

### 3.3.1. Zenbaki oso bitarrak

Sistema hamartarrean, 1563 zenbakia era garatuan idatzita hau da:

$$1563 = (1 \times 10^3) + (5 \times 10^2) + (6 \times 10^1) + (3 \times 10^0).$$

Baina, zenbaki bera sistema bitarrean idatzita hau da:

$$\begin{aligned} 1563 = & (1 \times 2^{10}) + (1 \times 2^9) + (0 \times 2^8) + (0 \times 2^7) + (0 \times 2^6) + (0 \times 2^5) \\ & + (1 \times 2^4) + (1 \times 2^3) + (0 \times 2^2) + (1 \times 2^1) + (1 \times 2^0). \end{aligned}$$

Beraz,  $1563 = 11000011011_{bi}$  dugu.

Hau da beste adibide bat:

$$(111 \dots 11)_2 = 2^{n-1} + 2^{n-2} + 2^{n-3} + \dots + 2^1 + 2^0 = (2^n - 1).$$

Azken berdintza lortzeko, kontuan hartu behar dugu  $1 + x + x^2 + \dots + x^n = \frac{x^{n+1} - 1}{x - 1}$  betetzen dela,  $x \neq 1$  denean.

Nola lor dezakegu sistema bitarrera igarotzea?

1563 zenbakia birekin zatitzean, hau ateratzen dugu:

$$1563 = 2 \times 781 + 1.$$

Orain, 781 zenbakia birekin zatitzean, hau ateratzen dugu:

$$781 = 2 \times 390 + 1.$$

Eragiketa hori errepikatuz, zera lortzen da:

$$390 = 2 \times 195 + 0$$

$$195 = 2 \times 97 + 1$$

$$97 = 2 \times 48 + 1$$

$$48 = 2 \times 24 + 0$$

$$24 = 2 \times 12 + 0$$

$$12 = 2 \times 6 + 0$$

$$6 = 2 \times 3 + 0$$

$$3 = 2 \times 1 + 1$$

$$1 = 2 \times 0 + 1$$

eta behetik gorako hondarrak hartuz, zenbaki bitarra eraikitzen dugu:  $1563 = 11000011011_{bi}$ .

### 3.3.2. Zatiki bitarrak

$R$  zenbaki erreal bat bada, non  $0 < R < 1$ , orduan  $\{0, 1\}$  multzoan badago  $d_1, d_2, \dots, d_n, \dots$  zifra segida bat hau betetzen duena:

$$R = (d_1 \times 2^{-1}) + (d_2 \times 2^{-2}) + \dots + (d_n \times 2^{-n}) + \dots \quad (3.1)$$

Eskuineko adierazpena laburki emanda, hots, notazio zatikiar bitarrean, honela da:

$$R = 0.d_1d_2\dots d_n\dots_{bi}.$$

Askotan, zenbaki errealek 1 zifraren kopuru infinitu bat behar dute adierazpen bitarrean. Adibidez,  $7/10$ , oinarri hamartarrean  $0.7$ ren bidez adierazten da, baina bere adierazpen bitarrak 1 zifraren kopuru infinitu bat behar du:

$$\frac{7}{10} = 0.1\overline{0110}_{bi}.$$

Zatiki hori periodikoa da:  $0110$  lau zifrako taldea errepikatu egiten da, bukaerarik gabe.

Orain, algoritmo eraginkor bat garatu dezakegu 2 oinarriko adierazpenak aurkitzeko. Birekin (3.1) adierazpena biderkatzen badugu, hau dugu:

$$2R = d_1 + (d_2 \times 2^{-1}) + \dots + (d_n \times 2^{-n+1}) + \dots,$$

non  $d_1$   $2R$ -ren zati osoa baita,  $d_1 = \lfloor 2R \rfloor$ . Adierazpen horren zatikia,  $zat(2R)$ , hau da:

$$Z_1 = zat(2R) = (d_2 \times 2^{-1}) + \dots + (d_n \times 2^{-n+1}) + \dots$$

eta  $(0,1)$  tartean dago. Azken adierazpen hori birekin biderkatuz, zera lortzen dugu:

$$2Z_1 = d_2 + (d_3 \times 2^{-1}) + \dots + (d_n \times 2^{-n+2}) + \dots$$

Berdintza horren zati osoa hartuz, hots,  $d_2 = \lfloor 2Z_1 \rfloor$ . Prozesua aurrera joango da, seguraski bukaerarik gabe ( $R$ -ren adierazpena 2 oinarrian ez bada finitua, ezta periodikoa ere), eta era errepikarian honela sortzen ditu  $\{d_k\}$  eta  $\{Z_k\}$  bi segidak:

$$\begin{aligned} d_k &= \lfloor 2Z_{k-1} \rfloor \\ Z_k &= zat(2Z_{k-1}), \end{aligned}$$

non  $d_1 = \lfloor 2R \rfloor$  eta  $Z_1 = zat(2R)$  baitira. Ondorioz,  $R$ -ren adierazpen bitarra serie konbergente honek ematen digu:

$$R = \sum_{i=1}^{\infty} d_i(2)^{-i}$$

eta hau  $1/2$  arrazoiko serie geometriko baten azpiseria da.

**3.7. adibidea.** *Kalkulatu  $7/10$ en adierazpen bitarra.*

*Ebazpena.*  $R = 7/10 = 0.7$  denez, honela kalkulatu dugu:

$$\begin{array}{lll}
 2R = 1.4 & d_1 = \lfloor 1.4 \rfloor = 1 & Z_1 = \text{zat}(1.4) = 0.4 \\
 2Z_1 = 0.8 & d_2 = \lfloor 0.8 \rfloor = 0 & Z_2 = \text{zat}(0.8) = 0.8 \\
 2Z_2 = 1.6 & d_3 = \lfloor 1.6 \rfloor = 1 & Z_3 = \text{zat}(1.6) = 0.6 \\
 2Z_3 = 1.2 & d_4 = \lfloor 1.2 \rfloor = 1 & Z_4 = \text{zat}(1.2) = 0.2 \\
 2Z_4 = 0.4 & d_5 = \lfloor 0.4 \rfloor = 0 & Z_5 = \text{zat}(0.4) = 0.4 \\
 2Z_5 = 0.8 & d_6 = \lfloor 0.8 \rfloor = 0 & Z_6 = \text{zat}(0.8) = 0.8 \\
 2Z_6 = 1.6 & d_7 = \lfloor 1.6 \rfloor = 1 & Z_7 = \text{zat}(1.6) = 0.6.
 \end{array}$$

Kontuan izan  $2Z_2 = 1.6 = 2Z_6$  dela. Beraz,  $d_k = d_{k+4}$  eta  $Z_k = Z_{k+4}$  patroiak errepikatuko dira  $k = 2, 3, 4, \dots$ -tarako. Ondorioz,  $7/10 = 0.\overline{10110}_{bi}$ .  $\square$

Serie geometrikoa erabil dezakegu, adierazpen bitar bati dagokion zenbaki arrazional hamartarra kalkulatzeko.

**3.8. adibidea.** *Izan bedi  $0.\overline{01}_{bi}$  zenbaki bitarra. Aurkitu zenbaki horri dagokion zenbaki arrazional hamartarra.*

*Ebazpena.*  $0.\overline{01}_{bi}$  era garatuan idatziz, zera dugu:

$$\begin{aligned}
 0.\overline{01}_{bi} &= (0 \times 2^{-1}) + (1 \times 2^{-2}) + (0 \times 2^{-3}) + (1 \times 2^{-4}) + \dots \\
 &= \sum_{i=1}^{\infty} (2^{-2})^i = -1 + \sum_{i=0}^{\infty} (2^{-2})^i \\
 &= -1 + \frac{1}{1 - \frac{1}{4}} = -1 + \frac{4}{3} = \frac{1}{3}. \quad \square
 \end{aligned}$$

**3.9. adibidea. Desplazamendu bitarra.** *Demagun  $S = 0.00000\overline{11000}_{bi}$  dela. Kalkula ezazu  $S$  zenbakia sistema hamartarrean.*

*Ebazpena.* Hau dugu:

$$2^5 S = 32S = 0.\overline{11000}_{bi}$$

Bestalde,

$$2^{10} S = 1024S = 11000.\overline{11000}_{bi}$$

Bigarren berdintzari lehenengoa kenduz:

$$1024S - 32S = 11000.\overline{11000}_{bi} - 0.\overline{11000}_{bi} = 11000_{bi} = 1 \times 2^4 + 1 \times 2^3 = 24,$$

eta orduan  $992S = 24$ , azkenik  $S = 24/992 = 8/33$ .  $\square$

### 3.3.3. Ordenagailuko zenbakiak

Ordenagailuek zenbaki errealetarako puntu higikorren adierazpen bitarra erabiltzen dute. Zenbaki bitarrak 0 eta 1 digituek osatuta daudenez, puntuaren ezker aldeko zenbakia 1 izateko normalizatzen dira beti. Beraz, *bit* (0 edo 1 digitu bitar bakoitza) hori ez da gorde behar (beti 1 baita). Ondorioz, zero ez diren zenbaki bitarrak honela gordeko dira:

$$\text{fl}(x) = \text{zeinu}(x)(1 + f) \times 2^e. \quad (3.2)$$

$f$  zenbakiari *mantisa* deritzogu eta adierazpen bitar finitua da;  $e$  zenbakiari *berretzaile* deritzogu. Ordenagailuek zenbaki errealeen azpimultzo txiki bat soilik erabiltzen dute, zeren  $f$  eta  $e$  izan ditzaketen zifra bitarren kopurua murriztea beharrezkoa baita. Adibidez, demagun era honetako zenbaki erreal positiboen multzoa:

$$1.d_1d_2d_3d_4 \text{ }_{bi} \times 2^e,$$

non  $d_1, d_2, d_3, d_4 \in \{0, 1\}$  eta  $e \in \{-3, -2, -1, 0, 1, 2, 3, 4\}$ . Mantisarako  $2^4 = 16$  aukera dugu, eta berretzailerako 8 aukera ere bai, horrek 128 zenbakien multzo bat ematen digu:

$$\{1.0000_{bi} \times 2^{-3}, 1.0001_{bi} \times 2^{-3}, \dots, 1.1110_{bi} \times 2^4, 1.1111_{bi} \times 2^4\},$$

non, esate baterako:

$$\begin{aligned} 1.0000_{bi} \times 2^{-3} &= 1 \times 2^{-3} = 0.1250 \\ 1.1111_{bi} \times 2^4 &= (1 + 1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4}) \times 2^4 \\ &= 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = 16 + 8 + 4 + 2 + 1 = 31. \end{aligned}$$

Ordenagailu batek 4 zifrako mantisa batez soilik  $(1/10 + 1/5) + 1/6$  eragiketa egin beharko balu, ordenagailuak zenbaki bitar hurbilenari biribilduko dio zenbaki erreal bakoitza; kasu honetan, hauek batzen ditu:

$$\begin{aligned} \frac{1}{10} = 0.1 &= \frac{1.6}{2^4} = 1.6 \times 2^{-4} = 1.\overline{1001}_{bi} \times 2^{-4} \approx 1.1010_{bi} \times 2^{-4} = 0.11010_{bi} \times 2^{-3} \\ \frac{1}{5} = 0.2 &= \frac{1.6}{2^3} = 1.6 \times 2^{-3} = 1.\overline{1001}_{bi} \times 2^{-3} \approx 1.1010_{bi} \times 2^{-3} = 1.1010_{bi} \times 2^{-3}, \\ \text{batuz hau lortzen da:} \\ \frac{3}{10} &\approx 10.01110_{bi} \times 2^{-3}. \end{aligned}$$

Ordenagailuak erabakitzen du nola gorde behar duen  $10.01110_{bi} \times 2^{-3} = 1.001110_{bi} \times 2^{-2}$  zenbakia; demagun  $1.0100_{bi} \times 2^{-2}$  zenbakian biribiltzen duela. Hurrengo urratsa hau da:

$$\begin{aligned} \frac{3}{10} = 0.3 &\approx 1.0100_{bi} \times 2^{-2} = 1.0100_{bi} \times 2^{-2} \\ \frac{1}{6} = \frac{1.\overline{3}}{2^3} &= 1.\overline{3} \times 2^{-3} = 1.0\overline{1}_{bi} \times 2^{-3} \approx 1.0101_{bi} \times 2^{-3} = 0.10101_{bi} \times 2^{-2}, \\ \text{batuz hau lortzen da:} \\ \frac{7}{15} &\approx 1.11101_{bi} \times 2^{-2}. \end{aligned}$$

Ordenagailuak erabakitzen du nola gorde behar duen  $1.11101_{bi} \times 2^{-2}$ . Biribiltzen du  $1.1111_{bi} \times 2^{-2}$  gordez. Beraz, hau da ordenagailuak batuketara problemari ematen dion soluzioa:

$$\frac{7}{15} \approx 1.1111_{bi} \times 2^{-2}.$$

Ordenagailuak sortutako errorea hau izan da:

$$\frac{7}{15} - 1.1111_{bi} \times 2^{-2} \approx -0.017708,$$

eta hori 7/15en %3.79 da.

### 3.3.4. Ordenagailu baten zehaztasuna

Mantisak 32 zifra baditu, 9 zifrara arteko zenbakiak gorde daitezke. Itzul gaitezen atalaren hasierara, ordenagailu baten bidez 1/10 elkarren ondoan 100000 bider batu nahi genuen lekura.

Demagun  $f$  mantisak 32 zifra bitar dauzkala. Beraz,

$$1 + f = 1.d_1d_2d_3 \dots d_{31}d_{32bi}.$$

Zatiki bat era bitarrean adierazten badugu, adierazpen hori periodikoa izango da; adibidez:

$$\frac{1}{10} = 0.\overline{00011}_{bi}.$$

Baina 32 zifrako mantisa bat erabiltzen dugunean, ordenagailuak trunkatzen du, eta barneko hurbilpen gisa hau erabiltzen du:

$$\frac{1}{10} = 1.10011001100110011001100110011001_{bi} \times 2^{-4},$$

eta horren errorea (hau da, aurreko bi zenbaki bitarren arteko diferentzia) hau da:

$$0.\overline{1001}_{bi} \times 2^{-36} \approx 8.731149137 \times 10^{-12}.$$

Arrazoi horregatik, ordenagailuak errore bat egin behar du, 1/10 zenbakia 100000 bider batzea eskatzen diogunean. Errore hori  $(100000)(8.731149137 \times 10^{-12}) = 8.731149137 \times 10^{-7}$  izan behar zen gutxienez. Batura handitzen doan heinean, batura partzialak ere biribiltzen dira, eta batugaiak txikiak dira une horretan dagoen batura partzialarekin konparatuta. Ondorioz, trunkaketa gogorragoa da batugaiaren ekarpenerako. Errore horien guztien eragin konposatuak azken errore hau ematen du:  $10000 - 9999.99447 = 5.53 \times 10^{-3}$ .

### 3.3.5. Ordenagailuko puntu higikorreko zenbakiak

IEEE 754 puntu higikorreko aritmetikarako IEEEko estandarra da. Ordenagailuek bi modu dituzte zenbakiak adierazteko: *modu osoa* eta *puntu higikorreko modua*. Modu osoa erabiltzen da emaitza osoa izateko ziurtasuna duten eragiketak egiteko. Baina, normalean,



zientzietan eta ingeniartzan, puntu higikorren adierazpenak erabiltzen dira; eta (3.2) adierazpena erabiltzeak mugak jartzen ditu  $f$  mantisako zifren kopuruaren eta  $e$  berretzailearen heinaren gainean.

Ordenagailuak *zehaztasun bikoitzarekin* zenbaki errealak adierazteko 64 zifra erabiltzen ditu; orduan, lehenengo bita zenbakiaren zeinua gordetzeko izango da (0 + da eta 1 - da), hurrengo 11 bitak berretzailerako, eta azken 52 bitak mantisarako. Beraz, esanguratsutasunak bit implizitu 1 dauka (finko alde osoan) eta 52 bit esplizitu, guztira 53 bit; ondorioz, zehaztasun osoak 53 bit dauzka, hau da  $\approx 16$  digitu dezimal (16 digitu esanguratsu),  $\log_{10}(2^{53})$ . Balio handiena, **realmax**, hau da:

$$+1.1111 \dots 1111_{bi} \times 2^{+1023} = (2 - 2^{-52}) \times 2^{+1023} = 1.7977 \times 10^{308}.$$

Eta balio positibo txikiena, **realmin**, hau da:

$$+1.0000 \dots 0000_{bi} \times 2^{-1022} = 2.2251 \times 10^{-308}.$$

Zenbakien arteko tartearen luzera haien tamainaren mendekoa da. Mantisaren balio txikiena  $2^{-52} = 2.2204 \times 10^{-16}$ . Balio horrek makinaren zehaztasuna, **eps**, (edo makinaren errorea,  $\varepsilon_M$ ) ematen du 1 tamainako zenbakiatarako; hori da bi zenbakien artean egon daitekeen diferentzia txikiena. Zenbakiaren tamaina handiagoa denean, diferentzia hori handitu egiten da.

MATLABek, berez, zehaztasun bikoitza erabiltzen du; hots,  $(\beta, t, L, U) = (2, 52, -1022, 1023)$  puntu higikorreko sistema (IEEE sisteman, zenbakiaren 1 lehenengo digitua ez da zenbatzen  $t$  emateko).

Puntu higikorreko adierazpenari dagokion errore erlatiboa hau da:

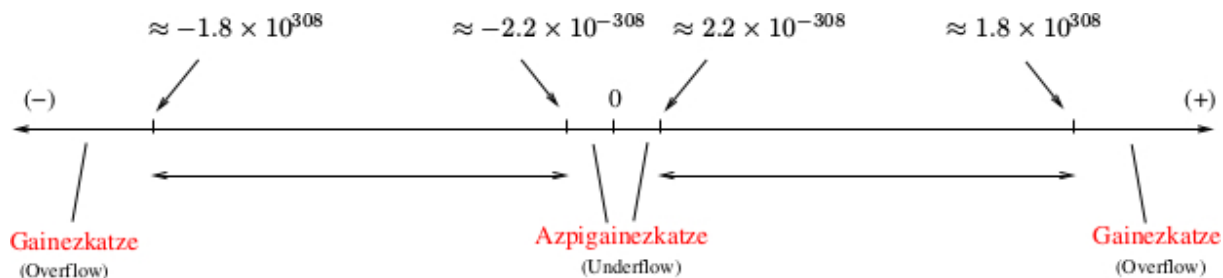
$$\frac{|\text{fl}(x) - x|}{|x|}.$$

Gaurko puntu higikorreko sistemek, IEEEko estandarrek bezala, bermatzen dute errore erlatibo hori bornatuta dagoela zenbaki honetaz:

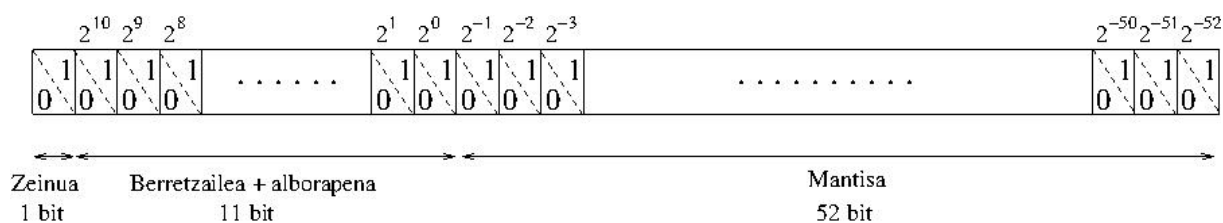
$$\eta = \frac{1}{2} \times 2^{-t}.$$

Zenbaki horri *biribiltzearen unitatea* deritzogu. Zehaztasun bikoitzaren kasuan,  $\eta = \frac{1}{2} \times 2^{-52} \approx 1.1 \times 10^{-16}$  (beraz, 52 digitu bitar esanguratsu daude, edo 16 digitu hamartar).

Mantisaren balioa era bitarrean sartzen da. Berretzailearen balioa *alborapen* batekin sartzen da. Hemen, alborapenak esan nahi du berretzailearen balioari konstante bat batzen zaiola. Alborapena erabiltzen da berretzailearen zeinuan bit bat ez gastatzeko. Zehaztasun bikoitzean, notazio bitarrean 11 bitekin idatz dezakegun zenbaki handiena 2047 da (hots, 11 digituak 1 direnean). Erabiltzen den alborapena 1023 da (berretzaile txikiena 1 da). Hots, berretzailea 4 bada, gordetako balioa  $4+1023=1027$  izango da (1 gordeta badago, benetako



**3.2. irudia.** Zehaztasun bikoitzean adieraz daitezkeen zenbakien heina



**3.3. irudia.** Zenbaki baten biltzea, notazio bitarrean, IEEE-754 estandarrean.

berretzailea  $1-1023=-1022$  dugu). Beraz, ordenagailuan gordetako berretzaile txikiena  $-1022$  izango da (1 bezala gordez), eta handiena  $1023$  ( $2046$  bezala gordez).

*Zehaztasun bakunarekin* zenbaki errealak adierazteko, 32 zifra bitar (bit) erabiltzen dituzte ordenagailuek. Lehenengo bitak zenbakiaren zeinua gordetzen du. Hurrengo 8 bitak berretzailea gordetzeko erabiltzen dira. Eta azken 23 bitak mantisa gordetzeko erabiltzen dira. Horrek, zeroaz gain,  $1.1755 \times 10^{-38}$ tik  $3.4028 \times 10^{38}$ ra arteko zenbaki errealak adierazten uzten du; hots,  $2^{-126}$ tik  $(2 - 2^{-23})2^{127}$ ra arte. Zenbakien arteko tartearen luzera haien tainaren mendekoa da. Gorde dezakegun mantisaren balio txikiena  $2^{-23} = 1.1921E - 7$  da. Ondorioz, zehaztasun bakunak  $(\beta, t, L, U) = (2, 23, -126, 127)$  puntu higikorreko sistema da.

Zehaztasun bakunean, non 8 bit erabiltzen baititugu, haiekin idatz dezakegun zenbaki handiena  $255$  da (hots,  $11111111_{bi}$ ) eta txikiena  $1$ ; ondorioz,  $127$  da berretzailearen balioa gordetzeko alborapena.

Zehaztasun bikoitzaren kasuan, *biribiltzearen unitatea*  $\eta = \frac{1}{2} \times 2^{-23} \approx 6.0 \times 10^{-8}$  da (beraz, 23 digitu bitar esanguratsu daude, edo 7 digitu hamartar).

## IEEEren biribiltze zehatza

Nahiz eta gure puntu higikorreko sisteman zenbakizko adierazpenak zehatzak izan, zenbaki horien eragiketa aritmetikoek biribiltze-erroreak sartzen dituzte. Errore horiek erlatiboki

nahiko handiak izan daitezke, *eskolta-digituak* erabiltzea izan ezik. Aparteko digitu horiek interim kalkuletan erabiltzen dira. IEEE estandarrak *biribiltze zehatza* behar du, eta, horri esker, eragiketa aritmetiko bakoitzean errore erlatiboa bornatuta dago  $\eta$  balioaz.

**3.10. adibidea.** *Izan bedi puntu higikorreko sistema bat, non  $\beta = 10$  oinarria eta  $t = 4$  digituak baitira. Ondorioz, biribiltze unitatea  $\eta = \frac{1}{2} \times 10^{-3}$  da. Aztertu biribiltze zehatzaren eragina.*

*Ebazpena:* Izan bitez

$$x = 0.1103 = 1.103 \times 10^{-1}, \quad y = 9.963 \times 10^{-3}.$$

Bi zenbaki horiek kenduz, balio zehatza  $d = x - y = 0.100337$  da. Beraz, biribiltze zehatzak  $\hat{d} = 0.1003$  ematen du. Jakina, errore erlatiboa hau da:

$$\frac{|d - \hat{d}|}{|d|} = \frac{|(x - y) - \text{fl}(x - y)|}{|x - y|} = \frac{|0.100337 - 0.1003|}{0.100337} = 0.37 \times 10^{-3} < \eta.$$

Aldiz, baldin guk kentzen baditugu bi zenbaki horiek eskolta-digiturik gabe,  $\bar{d} = 0.1103 - 0.0099 = 0.1004$  lortuko genuke. Orain, lortutako errore erlatiboa ez da  $\eta$  baino txikiago, zeren

$$\frac{|d - \bar{d}|}{|d|} = \frac{|0.100337 - 0.1004|}{0.100337} = 0.63 \times 10^{-3} > \eta.$$

Hortaz, eskolta-digituak erabili behar dira biribiltze zehatza lortzeko.  $\square$

## Balio bereziak

Ez dira erabiltzen berretzaile posible guztiak. Zehaztasun bikoitzerako  $2^{11} = 2048$  berretzaile desberdin ditugu, eta zehaztasun bakunerako  $2^8 = 256$ . Baina, 2046 eta 254 bakarrik, hurrenez hurren, erabiltzen ditugu. Hori da IEEE estandarrak berretzaileen muturreko puntuak ( $b = 0000000000_{bi} = 0$  eta  $b = 1111111111_{bi} = 2047$  zehaztasun bikoitzean eta  $b = 00000000_{bi} = 0$  eta  $b = 11111111_{bi} = 255$  zehaztasun bakunean) helburu berezietarako erreserbatzen dituelako.

Nola gordetzen dira 0 eta  $\infty$ ? Hemen erabiltzen dira berretzailearen zenbaki berezi horiek.

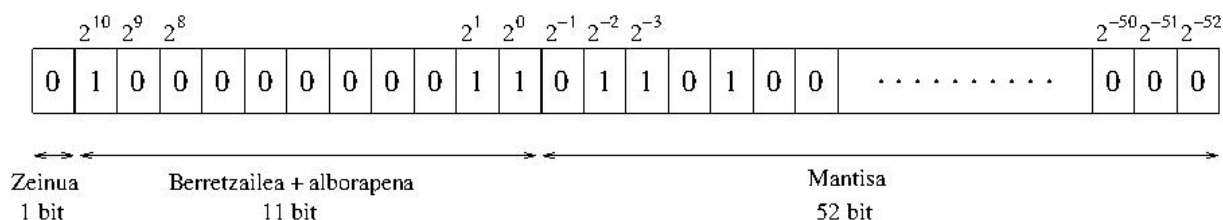
- 0-rako,  $b = 0, f = 0$  jartzen dugu, zeinua edozein izanik.
- $\pm\infty$ -rako,  $b = 1 \dots 1, f = 0$ .

- $b = 1 \dots 1, f \neq 0$  patroia konbentzioz NaN da (Not a Number); emaitza bat kalkula ezina denean agertzen da ( $1/0, 0/0, \dots$ ).

**3.11. adibidea.** *Ikus dezagun 22.5 zenbakia nola gordetzen den zehaztasun bikoitzean, IEEE-754 estandarren arabera. Lehenengo, zenbakia normalizatzen da (hots, alde osoan 1 zenbakia jartzen da) honela:*

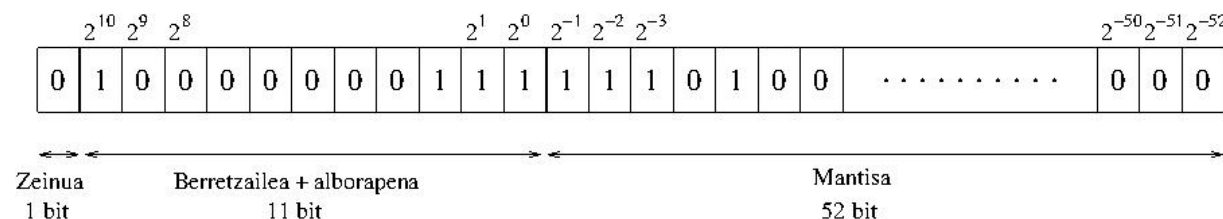
$$\frac{22.5}{2^4} 2^4 = 1.40625 \times 2^4.$$

*Zehaztasun bikoitzean, berretzailea alborapenarekin  $4+1023=1027$  da, zeina era bitarrean  $10000000011_{bi}$  gordetzen baita. Mantisa  $0.40625$  da, zeina era bitarrean  $.01101000_{bi} \dots 000$  gordetzen baita. Zenbaki horren biltzea hau da:*



**3.4. irudia.** 22.5 zenbakiaren biltegia, notazio bitarrean, IEEE-754 estandarrean.

**3.12. adibidea.** *Zein zenbaki hamartar dago gordeta zehaztasun bikoitzeko biltegi honetan:*



*Ebazpena.*  $b = 10000000111_{bi} = 1031$ , beraz,  $e = 1031 - 1023 = 8$ .

Bestalde,  $f = 1110100 \dots 0_{bi} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{32} = 0.90625$ . Eta zeinua + da.

Ondorioz,  $+(1 + f)2^e = +1.90625 \cdot 2^8 = 488$ .  $\square$

## 3.4. Errorearen analisia

**3.2. Definizioa.**  $\hat{p}$  zenbakia  $d$  zifra dezimal esanguratsuko  $p$ -ren hurbilpen bat dela esango dugu, hau betetzen duen zenbaki arrunt handiena  $d$  zenbakia bada:

$$\frac{|p - \hat{p}|}{|p|} < \frac{10^{-d}}{2}.$$

**3.13. adibidea.**  $x = 3.141592$  eta  $\hat{x} = 3.14$  badira, orduan  $|x - \hat{x}|/|x| = 0.000507 < 10^{-2}/2$ . Beraz,  $\hat{x}$  bi zifra esanguratsuko  $x$ -ren hurbilpena da.

*Gainezkatze* bat lortzen da zenbaki bat handiegia denean adierazteko erabilitako puntu higikor sisteman; esate baterako,  $e > U$  denean. *Azpigainezkatze* bat lortzen da  $e < L$  denean. Kalkulu bat egitean gainezkatze bat gertatzen denean, hori oso txarra da. Aldiz, hori ez da gertatzen azpigainezkatzearekin; sistemak orokorki jartzen du zenbakia 0-ra eta jarraitu egiten du (MATLABek hori egiten du, ezer esan gabe).

### 3.4.1. Zifra esanguratsuen ezeztapena

Errore hau gertatzen da bi zenbaki ia berdinak kentzen direnean.

Demagun  $z = x - y$ , non  $x \approx y$ . Orduan

$$|z - \text{fl}(z)| \leq |x - \text{fl}(x)| + |y - \text{fl}(y)|,$$

eta, hortaz, errore erlatiboak hau betetzen du:

$$\frac{|z - \text{fl}(z)|}{|z|} \leq \frac{|x - \text{fl}(x)| + |y - \text{fl}(y)|}{|x - y|}.$$

Zenbakitzailarekin ez da egon behar inolako arazorik puntu higikorreko sistemak  $x$  eta  $y$  ondo adierazten baditu. Baina, izendatzailea zerotik oso hurbil dago  $x \approx y$  bada, eta orduan  $z$ -ren errore erlatiboa oso handia bihur daiteke.

Jo dezagun  $p = 3.1415926536$  eta  $q = 3.1415957341$  zenbakiak berdintsuak direla, eta 11 zifra dezimaleko zehaztasunez adierazita daudela. Haien kendura kalkulatzen badugu,  $p - q = -0.0000030805$ , ikusiko dugu  $p$ -ren eta  $q$ -ren lehenengo sei zifrak berdinak direla; bere diferentziak  $p - q$  bost zifra dezimal bakarrik dauzka; fenomeno horri ***zifra esanguratsuen ezeztapena*** edo ***galera*** deritzogu, eta kontuz ibili behar dugu, zeren konturatu gabe kalkulatutako azken emaitzaren zehaztasuna txikiagotu egin baitezake.

**3.14. adibidea.** Izan bitez  $f(x) = x(\sqrt{x+1} - \sqrt{x})$  eta  $g(x) = x/(\sqrt{x+1} + \sqrt{x})$  algebratikoki funtzio baliokideak (egiaztatu). Konparatuko ditugu  $f(500)$  eta  $g(500)$ , sei zifra esanguratsuko biribiltzea erabiliz.

*Ebazpena.* Lehenengo funtzioarekin hau lortzen dugu:

$$f(500) = 500(\sqrt{501} - \sqrt{500}) = 500(22.3830 - 22.3607) = 500(0.0223) = 11.1500.$$

Orain,  $g(x)$ -rekin zera dugu:

$$g(500) = \frac{500}{\sqrt{501} + \sqrt{500}} = \frac{500}{22.3830 + 22.3607} = \frac{500}{44.7437} = 11.1748.$$

Bigarrenaren errore absolutua txikiagoa da. Hori da lortuko genukeena 11.174755300747198... emaitza zehatza sei zifra esanguratsuetara biribilduz.  $\square$

**3.15. adibidea.** Demagun  $y = \sqrt{x+1} - \sqrt{x}$  kalkulatu nahi dugula  $x = 100000$  bost digituzko aritmetika hamartarrean ( $\beta = 10, t = 5$ ).

*Ebazpena.* Bistan dago 100001 zenbakia ezin dela zehazki adierazi puntu higitokorren sistema horretan, eta bere adierazpena (inausiz edo biribilduz) 100000 da. Beste hitzez,  $x$ -ren balio horretarako sistema horretan  $x+1 = x$  dugu, eta kalkuluak sistema horretan  $\sqrt{x+1} - \sqrt{x} = 0$  ematen du.

Askoz hobe egin dezakegu identitate hau erabiltzen badugu:

$$\frac{(\sqrt{x+1} - \sqrt{x})(\sqrt{x+1} + \sqrt{x})}{(\sqrt{x+1} + \sqrt{x})} = \frac{1}{\sqrt{x+1} + \sqrt{x}}.$$

Formula horren eskuineko adierazpena erabiliz 5 digitu esanguratsurekin,  $1.5811 \times 10^{-3}$ , eta hori da hain zuzen, balio zuzena 5 digitu hamartarrekin.  $\square$

Badago beste era bat ere ezeztatze-erroreak saihesteko, eta teknika ezagun bat Taylorren garapena erabiltzea da.

**3.16. adibidea.** Demagun hau kalkulatu nahi dugula:

$$y = \sinh(x) = \frac{1}{2}(e^x - e^{-x}).$$

*Guk ziurtatu nahi dugu  $x$  guztietarako emaitza zehatz (doi) bat ematen duela (ez derrigorrez zuzena).*

*Ebazpena.* Goiko formula zuzenean aplikatu ahal izan arren, guk lor dezakegu formula sendoago bat  $x$  ia zero den kasurako, eta Taylorren garapen honetatik lortzen da:

$$\sinh(x) = x + \frac{x^3}{6} + \frac{\psi^5}{120},$$

$\psi$  baterako, non  $|\psi| \leq |x|$ . Orduan,  $x + x^3/6$  formulak hurbilpen efikaza emango luke  $|x|$  nahiko txikia denean, zeren trunkatze-errorea borna baitaiteke  $\frac{x^5}{120}$  balioaz, eta biribiltze-errorea jadanik ez da arazo bat.  $\square$

## Horner-en metodoa

Polinomio bat balioztatzeko orduan, emaitza hobeak lortuko ditugu Hornerren metodoa (biderkadura ahokatu) erabiltzen badugu.

Izan bedi  $p(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots + a_nx^n$  polinomioa; demagun haren balioa  $x_0$  baliorako,  $p(x_0)$ , kalkulatu nahi dugula. Hori lortzeko, segida hau definituko dugu:

$$\begin{aligned} b_n &= a_n \\ b_{n-1} &= a_{n-1} + b_nx_0 \\ b_{n-2} &= a_{n-2} + b_{n-1}x_0 \\ &\vdots \\ b_1 &= a_1 + b_2x_0 \\ b_0 &= a_0 + b_1x_0, \end{aligned}$$

orduan,  $p(x_0)$ -ren balioa  $b_0$  da.

Metodoa ulertzeko, kontuan izan polinomioa biderkadura ahokatu hau bezala idatz dezakegula:

$$p(x) = a_0 + x(a_1 + x(a_2 + x(a_3 + \dots + x(a_{n-1} + a_nx) \dots))).$$

Gero,  $b_i$  banan-banan ordezkaturaz, zera dugu:

$$\begin{aligned} p(x_0) &= a_0 + x_0(a_1 + x_0(a_2 + x_0(a_3 + \dots + x_0(a_{n-1} + b_nx_0) \dots))) \\ &= a_0 + x_0(a_1 + x_0(a_2 + x_0(a_3 + \dots + x_0(b_{n-1}) \dots))) \\ &\vdots \\ &= a_0 + x_0(a_1 + x_0b_2) \\ &= a_0 + x_0b_1 \\ &= b_0. \end{aligned}$$

**3.17. adibidea.** Izan bitez  $P(x) = -1 + 3x - 3x^2 + x^3$  eta  $Q(x) = -1 + x(3 + x(-3 + x))$ . Hiru zifra esanguratsutarako biribiltzea erabiliz,  $P(2.19)$  eta  $Q(2.19)$  kalkulatu ditugu, eta bi hurbilpen horiek konparatu ditugu benetako balioekin,  $P(2.19) = Q(2.19) = 1.685159$ .

*Ebazpena:*

$$\begin{aligned} P(2.19) &= (2.19)^3 - 3(2.19)^2 + 3(2.19) - 1 \\ &\approx -1 + 6.57 - 14.4 + 10.5 = 1.67. \end{aligned}$$

Hornerren metodoaz, honela kalkulatu da:

$$\begin{aligned} b_3 &= a_3 = 1 \\ b_2 &= a_2 + b_3x_0 = -3 + 1(2.19) = -0.81 \\ b_1 &= a_1 + b_2x_0 = 3 - 0.81(2.19) \approx 1.23 \\ b_0 &= a_0 + b_1x_0 \approx -1 + 1.23(2.19) \approx 1.69, \end{aligned}$$

hots,  $Q(2.19) \approx 1.69$ .

Errore absolutuak 0.015159 eta 0.004841 dira, hurrenez hurren. Beraz,  $Q(2.19) \approx 1.69$  hurbilpena hobea da.  $\square$

### 3.4.2. Trunkatze-errorea

Trunkatze-errorearen nozioa orokorki erlazionatuta dago adierazpen zail bat beste formula errazago batez ordezkatzearekin. Adibidez, funtzio bat Taylorren polinomio batez ordezkatzeko dugunean; esate baterako, adierazpen honen kasuan:

$$e^{x^2} = 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \frac{x^8}{4!} + \dots,$$

bere integrala zenbakizko kalkuluen bitartez aurkitzeko orduan, lehenengo bost gaien batu-  
rarekin ordezkatu dezakegu:  $1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \frac{x^8}{4!}$ .

### 3.4.3. $O(h^n)$ hurbiltze-ordena

Jakin badakigu  $\left\{\frac{1}{n^2}\right\}_{n=1}^{\infty}$  eta  $\left\{\frac{1}{n}\right\}_{n=1}^{\infty}$  segidak konbergenteak direla; hala ere, lehenengoak bigarrenak baino azkarrago jotzen du zerora. Jarraian, segida baten konbergentziaren azkartasuna deskribatzeko beharrezko terminologia eta notazioa sartuko ditugu.

**3.3. Definizioa.**  $f(h)$  funtzioa  $h \rightarrow 0$  denean,  $g(h)$  ordenakoa dela esango dugu, eta  $f(h) = O(g(h))$  idatziko,  $C$  eta  $c$  konstanteak existitzen direnean, non hau betetzen baita:

$$|f(h)| \leq C|g(h)|, \quad |h| \leq c \text{ bada.}$$

**3.18. adibidea.** Izan bitez  $f(x) = x^3 + 2x^2$  eta  $g(x) = x^2$ . Izan ere,  $|x| \leq 1$ -erako  $x^3 \leq x^2$  da;  $|x| \leq 1$ -erako  $x^3 + 2x^2 \leq 3x^2$  lortzen dugu. Beraz,  $f(x) = O(g(x))$ .  $\square$

Landau-ren  $O(\cdot)$  notazioa deritzo, eta infinituetarako limiteetan ere erabiltzen da. Oso erabilgarria da, funtzio baten handitzearen abiadura deskribatzeko, funtzio ezagun batekin konparatzea ( $x^n$ ,  $x^{1/n}$ ,  $a^x$ ,  $\log_a x$ , eta abar).

Segiden konbergentziaren abiadura ere deskriba dezakegu antzeko era batean.

**3.4. Definizioa.** Izan bitez  $\{x_n\}_{n=1}^{\infty}$  eta  $\{y_n\}_{n=1}^{\infty}$  segidak.  $\{x_n\}$  segida  $\{y_n\}$  ordenakoa dela esaten da, eta  $x_n = O(y_n)$  idazten,  $C$  eta  $N$  konstanteak existitzen badira, non hau betetzen baita:

$$|x_n| \leq C|y_n|, \quad n \geq N \text{ bada.}$$

**3.19. adibidea.**  $\frac{n^2 - 1}{n^3} = O\left(\frac{1}{n}\right)$ ; izan ere,  $\frac{n^2 - 1}{n^3} \leq \frac{n^2}{n^3} = \frac{1}{n}$ ,  $n \geq 1$  bada.  $\square$



**3.5. Definizioa.** Demagun  $p(h)$  funtzioak beste  $f(h)$  funtzio bat hurbiltzen duela, eta existitzen direla  $M > 0$  zenbaki erreal bat eta  $n$  zenbaki arrunt bat non hau betetzen baita:

$$\frac{|f(h) - p(h)|}{|h^n|} \leq M, \quad h \text{ nahiko txiki baterako.}$$

Orduan,  $p(h)$  funtzioak  $f(h)$  hurbiltzen duela esaten da  $O(h^n)$  hurbiltze-ordenarekin, eta honela idazten da:

$$f(h) = p(h) + O(h^n).$$

Goiko desberdintza  $|f(h) - p(h)| \leq M|h^n|$  bezala idatziz,  $O(h^n)$  adierazpenak  $M|h^n|$  errore-bornearen lekua betetzen duela dakusagu.

**3.1. teorema.** Demagun  $f(h) = p(h) + O(h^n)$  eta  $g(h) = q(h) + O(h^m)$ , eta izan bedi  $r = \min\{m, n\}$ . Orduan:

$$\begin{aligned} f(h) + g(h) &= p(h) + q(h) + O(h^r), \\ f(h)g(h) &= p(h)q(h) + O(h^r), \\ f(h)/g(h) &= p(h)/q(h) + O(h^r), \quad g(h) \neq 0 \text{ eta } q(h) \neq 0 \text{ badira.} \end{aligned}$$

Oso interesgarria da aintzat hartzea  $f(x)$  funtzioaren Taylorren polinomioen bidezko  $n$ -garren hurbilpena  $p(x)$  den kasua. Orduan Taylorren formularen hondarra  $O(h^{n+1})$ -ren bidez adierazten da, eta idatzi gabeko gai guztiak ordezkatzen ditu, hau da,  $h^{n+1}$  berreketa daukana eta goi-ordenakoak. Taylorren formularen hondarrak jotzen du zerora  $h \rightarrow 0$  denean,  $h^{n+1}$ -en abiadura berdinarekin, adierazpen honek erakusten duen bezala:

$$O(h^{n+1}) \approx Mh^{n+1} \approx \frac{f^{(n+1)}(c)}{(n+1)!} h^{n+1},$$

zeina baliozkoa baita  $h$  nahiko txikia denean. Beste hitz batzuetan esanda,  $O(h^{n+1})$  gaiak  $Mh^{n+1}$  kopurua ordezkatzen du, non  $M$  konstantea baita.

**3.2. Taylorren teorema.** Demagun  $f \in C^{n+1}[a, b]$ . Baldin  $x_0$  eta  $x = x_0 + h$   $[a, b]$  tartean badaude, orduan, hau betetzen da:

$$f(x_0 + h) = \sum_{i=0}^n \frac{f^{(i)}(x_0)}{i!} h^i + O(h^{n+1}).$$

Kalkuluetan propietate hauek erabiltzen dira:

- (i)  $O(h^p) + O(h^p) = O(h^p)$ .
- (ii)  $O(h^p) + O(h^q) = O(h^r)$ , non  $r = \min\{p, q\}$ .

(iii)  $O(h^p)O(h^q) = O(h^s)$ , non  $s = p + q$ .

Honako adibide honek 3.1. teorema argitzen du.

**3.20. adibidea.** *Izan bitez Taylorren garapen hauek:*

$$e^h = 1 + h + \frac{h^2}{2!} + \frac{h^3}{3!} + O(h^4) \quad \text{eta} \quad \cos(h) = 1 - \frac{h^2}{2!} + \frac{h^4}{4!} + O(h^6).$$

*Aurki itzazu baturaren eta biderkaduraren hurbiltze-ordenak.*

*Ebazpena:*

$$\begin{aligned} e^h + \cos(h) &= 1 + h + \frac{h^2}{2!} + \frac{h^3}{3!} + O(h^4) + 1 - \frac{h^2}{2!} + \frac{h^4}{4!} + O(h^6) \\ &= 2 + h + \frac{h^3}{3!} + O(h^4) + \frac{h^4}{4!} + O(h^6). \end{aligned}$$

Eta  $O(h^4) + \frac{h^4}{4!} = O(h^4)$  eta  $O(h^4) + O(h^6) = O(h^4)$  direnez, zera dugu:

$$e^h + \cos(h) = 2 + h + \frac{h^3}{3!} + O(h^4),$$

eta, ondorioz,  $O(h^4)$  hurbiltze-ordena da.

Biderkaduran, antzeko eran egiten da.

$$\begin{aligned} e^h \cos(h) &= \left(1 + h + \frac{h^2}{2!} + \frac{h^3}{3!} + O(h^4)\right) \left(1 - \frac{h^2}{2!} + \frac{h^4}{4!} + O(h^6)\right) \\ &= \left(1 + h + \frac{h^2}{2!} + \frac{h^3}{3!}\right) \left(1 - \frac{h^2}{2!} + \frac{h^4}{4!}\right) \\ &\quad + \left(1 + h + \frac{h^2}{2!} + \frac{h^3}{3!}\right) O(h^6) + \left(1 - \frac{h^2}{2!} + \frac{h^4}{4!}\right) O(h^4) \\ &\quad + O(h^4)O(h^6) \\ &= 1 + h - \frac{h^3}{3} - \frac{5h^4}{24} - \frac{h^5}{24} + \frac{h^6}{48} + \frac{h^7}{144} \\ &\quad + O(h^6) + O(h^4) + O(h^4)O(h^6). \end{aligned}$$

Izan ere,  $O(h^4)O(h^6) = O(h^{10})$  eta hau betetzen da:

$$-\frac{5h^4}{24} - \frac{h^5}{24} + \frac{h^6}{48} + \frac{h^7}{144} + O(h^6) + O(h^4) + O(h^{10}) = O(h^4),$$

honako erlazio hau lortzen da:

$$e^h \cos(h) = 1 + h - \frac{h^3}{3} + O(h^4).$$

Ondorioz, biderkadura  $O(h^4)$  hurbiltze-ordenakoa da.  $\square$

### 3.4.4. Segida baten hurbiltze-ordena

**3.6. Definizioa.** Demagun  $\lim_{n \rightarrow \infty} x_n = x$  dugula eta  $\{r_n\}_{n=1}^{\infty}$  segida bat dela, non  $\lim_{n \rightarrow \infty} r_n = 0$ . Orduan, esango dugu  $\{x_n\}_{n=1}^{\infty}$  segidak  $x$ -ra jotzen duela  $O(r_n)$  hurbiltze-ordenarekin,  $K > 0$  konstante bat existitzen bada hau betetzen duena:

$$\frac{|x_n - x|}{|r_n|} \leq K \quad n \text{ nahiko handi baterako.}$$

Hori  $x_n = x + O(r_n)$  idatziz adieraziko dugu, edo  $x_n \rightarrow x + O(r_n)$  hurbiltze-ordenarekin.

**3.21. adibidea.** Izan bitez  $x_n = \cos(n)/n^2$  eta  $r_n = 1/n^2$ , orduan  $\lim_{n \rightarrow \infty} x_n = 0$  dugu  $O(1/n^2)$  hurbiltze-ordenarekin. Hori erlazio honetatik ateratzen da:

$$\frac{|\cos(n)/n^2|}{|1/n^2|} = |\cos(n)| \leq 1 \quad n \text{ guztietarako.} \quad \square$$

### 3.4.5. Errorearen hedapena

Atal honetan aztertuko dugu nola heda daitezkeen erroreak segidako eragiketen kate batean. Jo ditzagun  $p$  eta  $q$  zenbakien batura (balio zehatzak izanik), haien  $\hat{p}$  eta  $\hat{q}$  balio hurbilduak, eta azken horien erroreak,  $\varepsilon_p$  eta  $\varepsilon_q$ , hurrenez hurren. Hots,  $p = \hat{p} + \varepsilon_p$  eta  $q = \hat{q} + \varepsilon_q$ . Batura horiek gaiez gai batuz, hau dugu:

$$p + q = (\hat{p} + \varepsilon_p) + (\hat{q} + \varepsilon_q) = (\hat{p} + \hat{q}) + (\varepsilon_p + \varepsilon_q).$$

Beraz, batuketa batean, errorea da batugaien erroreen batura.

Biderketa batean, errorearen hedapena konplexuagoa da. Biderkadura hau da:

$$pq = (\hat{p} + \varepsilon_p)(\hat{q} + \varepsilon_q) = \hat{p}\hat{q} + \hat{p}\varepsilon_q + \hat{q}\varepsilon_p + \varepsilon_p\varepsilon_q.$$

Beraz,  $|\hat{p}|$  eta  $|\hat{q}|$  1 baino handiagoak badira,  $\hat{p}\varepsilon_q$  eta  $\hat{q}\varepsilon_p$  gaiek handitu ditzakete jatorrizko  $\varepsilon_p$  eta  $\varepsilon_q$  erroreak. Berrordenatuz, hau lortuko dugu:

$$pq - \hat{p}\hat{q} = \hat{p}\varepsilon_q + \hat{q}\varepsilon_p + \varepsilon_p\varepsilon_q.$$

Demagun  $p \neq 0$  eta  $q \neq 0$ ; orduan, zera betetzen da:

$$R_{pq} = \frac{pq - \hat{p}\hat{q}}{pq} = \frac{\hat{p}\varepsilon_q + \hat{q}\varepsilon_p + \varepsilon_p\varepsilon_q}{pq} = \frac{\hat{p}\varepsilon_q}{pq} + \frac{\hat{q}\varepsilon_p}{pq} + \frac{\varepsilon_p\varepsilon_q}{pq}. \quad (3.3)$$

Gainera, suposatzen badugu  $\hat{p}$  eta  $\hat{q}$   $p$ -ren eta  $q$ -ren hurbilpen onak direla, orduan,  $\hat{p}/p \approx 1$ ,  $\hat{q}/q \approx 1$  eta  $R_p R_q = (\varepsilon_p/p)(\varepsilon_q/q) \approx 0$  ( $R_p$  eta  $R_q$   $\hat{p}$ -ren eta  $\hat{q}$ -ren errore erlatiboak dira). Hurbilpen horiek (3.3) adierazpenean ordezkatzuz, hau dugu:

$$R_{pq} = \frac{pq - \hat{p}\hat{q}}{pq} \approx \frac{\varepsilon_q}{q} + \frac{\varepsilon_p}{p} = R_q + R_p.$$

Beraz,  $\widehat{p}\widehat{q}$  biderkadurari dagokion errore erlatiboa, gutxi gorabehera,  $\widehat{p}$  eta  $\widehat{q}$  faktoreei dagozkien errore erlatiboen batura da.

Normalean, datuen hasierako erroreak hedatzen dira eragiketen kate batean zehar. Edozein zenbakizko prozesutan, kualitate desiragarria da hasierako baldintzetako errore txiki batek azken emaitzan errore txikiak eragitea. Algoritmo batek propietate hori badauka, *egonkorra* dela esango dugu; bestela, *ezeگونkorra* deritzo. Ahal den neurrian, metodo egonkorak aukeratuko ditugu.

**3.7. Definizioa.** Demagun  $\varepsilon$  hasierako errore bat dela, eta  $\varepsilon(n)$  adierazpenak  $n$  eragiketa ondoren errore horren hazkuntza adierazten duela. Baldin  $|\varepsilon(n)| \approx n\varepsilon$  bada, hazkuntza lineala dela esaten da. Baldin  $|\varepsilon(n)| \approx K^n\varepsilon$  bada, hazkuntza esponentziala dela esaten da. Baldin  $K > 1$  bada, orduan,  $n \rightarrow \infty$  denean, errore esponentziala bornerik gabe hazten da; baina,  $0 < K < 1$  bada, orduan,  $n \rightarrow \infty$  denean, errore esponentzialak zerora jotzen du.

Ondorengo bi adibideen bitartez ikusten da hasierako errore bat modu egonkorrean edo ezeگونkorrean heda daitekeela. Lehenengoan hiru algoritmo aurkezten dira, eta ikusiko dugu hiru horien bidez segida berdina lortuko genukeela, eragiketak era zehatzean egingo bagenitu. Bigarrean, hastapen-baldintzetako errore txikien hedapena aztertuko dugu.

**3.22. adibidea.** Erakutsiko dugu  $\{1/3^n\}_{n=0}^\infty$  segidaren gaiak hiru eskema desberdin erabiliz gara ditzakegula.

*Ebazpena.* Izan bitez hiru adierazpen hauek:

$$r_0 = 1 \qquad r_n = \frac{1}{3}r_{n-1} \qquad n = 1, 2, \dots$$

$$p_0 = 1, \quad p_1 = \frac{1}{3} \quad p_n = \frac{4}{3}p_{n-1} - \frac{1}{3}p_{n-2} \quad n = 2, 3, \dots$$

$$q_0 = 1, \quad q_1 = \frac{1}{3} \quad q_n = \frac{10}{3}q_{n-1} - q_{n-2} \quad n = 2, 3, \dots,$$

$\{r_n\}$ -rena bistan dago. Jarraian egiaztatuko dugu  $\{p_n\}$  segidan diferentzien kendurak  $p_n = A(1/3^n) + B$  soluzio orokorra daukala.

$$\begin{aligned} \frac{4}{3}p_{n-1} - \frac{1}{3}p_{n-2} &= \frac{4}{3} \left( \frac{A}{3^{n-1}} + B \right) - \frac{1}{3} \left( \frac{A}{3^{n-2}} + B \right) \\ &= \left( \frac{4}{3^n} - \frac{3}{3^n} \right) A + \left( \frac{4}{3} - \frac{1}{3} \right) B = A \frac{1}{3^n} + B = p_n. \end{aligned}$$

$A = 1$  eta  $B = 0$  hartzen baditugu,  $p_0 = 1$  eta  $p_1 = 1/3$  ditugu, eta hori da gure segida.

Jarraian egiaztatuko dugu  $\{q_n\}$  segidan diferentzien kendurak  $q_n = A(1/3^n) + B3^n$  soluzio orokorra daukala.

$$\begin{aligned} \frac{10}{3}q_{n-1} - q_{n-2} &= \frac{10}{3} \left( \frac{A}{3^{n-1}} + B3^{n-1} \right) - \left( \frac{A}{3^{n-2}} + B3^{n-2} \right) \\ &= \left( \frac{10}{3^n} - \frac{9}{3^n} \right) A + (10 - 1)3^{n-2}B \\ &= A \frac{1}{3^n} + B3^n = q_n. \end{aligned}$$

$A = 1$  eta  $B = 0$  hartzen baditugu,  $q_0 = 1$  eta  $q_1 = 1/3$  ditugu, eta horrek gure segida garatzen du.  $\square$

**3.23. adibidea.**  $\{x_n\} = \{1/3^n\}$  segidaren hurbilpenak garatuko ditugu adierazpen hauek erabiliz:

$$r_0 = 0.99996 \quad r_n = \frac{1}{3}r_{n-1} \quad n = 1, 2, \dots$$

$$p_0 = 1, \quad p_1 = 0.33332 \quad p_n = \frac{4}{3}p_{n-1} - \frac{1}{3}p_{n-2} \quad n = 2, 3, \dots$$

$$q_0 = 1, \quad q_1 = 0.33332 \quad q_n = \frac{10}{3}q_{n-1} - q_{n-2} \quad n = 2, 3, \dots,$$

$\{r_n\}$  kasuan  $r_0$ -ren errorea  $0.00004$  da.  $\{p_n\}$  eta  $\{q_n\}$  kasuetan,  $p_1$ -en eta  $q_1$ -en errorea  $0.00001\bar{3}$  da. Hastapen-baldintzetako errore txiki horien hedapena aztertuko dugu.

$n$	$x_n$	$r_n$	$p_n$	$q_n$
0	1.0000000000	0.9999600000	1.0000000000	1.0000000000
1	0.3333333333	0.3333200000	0.3333200000	0.3333200000
2	0.1111111111	0.1111066667	0.1110933330	0.1110666667
3	0.0370370370	0.0370355556	0.0370177778	0.0369022222
4	0.0123456790	0.0123451852	0.0123259259	0.0119407407
5	0.0041152263	0.0041150617	0.0040953086	0.0029002469
6	0.0013717421	0.0013716872	0.0013517695	-0.0022732510
7	0.0004572474	0.0004572291	0.0004372565	-0.0104777503
8	0.0001524158	0.0001524097	0.0001324188	-0.0326525834
9	0.0000508053	0.0000508032	0.0000308063	-0.0983641945
10	0.0000169351	0.0000169344	-0.0000030646	-0.2952280648

**3.1. taula.**  $\{x_n\} = \{1/3^n\}$  segida eta bere hurbilpenak.

$n$	$x_n - r_n$	$x_n - p_n$	$x_n - q_n$
0	0.0000400000	0.0000000000	0.0000000000
1	0.0000133333	0.0000133333	0.0000133333
2	0.0000044444	0.0000177778	0.0000444444
3	0.0000014815	0.0000192593	0.0001348148
4	0.0000004938	0.0000197531	0.0004049383
5	0.0000001646	0.0000199177	0.0012149794
6	0.0000000549	0.0000199726	0.0036449931
7	0.0000000183	0.0000199909	0.0109349977
8	0.0000000061	0.0000199970	0.0328049992
9	0.0000000020	0.0000199990	0.0984149998
10	0.0000000007	0.0000199997	0.2952449999

### 3.2. taula. Erroreen segidak.

$\{r_n\}$ -ren errorea egonkorra da eta era esponentzian txikitzen da.  $\{p_n\}$ -ren errorea egonkorra da.  $\{q_n\}$ -ren errorea ezegonkorra da eta abiadura esponentzialarekin handitzen da. Nahiz eta  $\{p_n\}$ -ren errorea egonkorra izan, haren gaiek  $n \rightarrow \infty$  denean  $p_n \rightarrow 0$  betetzen dutenez, epe luzean errorea menperatu egiten da, eta  $p_8$ -tik aurrerako zifra esanguratsuak ez daude ados  $x_n$ -ri dagozkienekin. Ikus 3.4.5. eta 3.4.5. taulak.

### 3.4.6. Datuen ziurgabetasuna

Errealitatean agertzen diren problemen datuek ziurgabetasunak edo erroreak dauzkate. Errore mota hori *zarata* izenarekin ezagutzen da, eta datu horietan oinarritzen den edozein zenbakizko kalkuluren zehaztasunari eragiten dio. Ezin dugu hobetu kalkuluen zehaztasuna, zaratak erasandako datuekin eragiketak gauzaten baditugu. Beraz,  $d$  zifra esanguratsu dauzkaten datuekin hasten bagara, orduan, datu horien bidez lortutako emaitzak  $d$  zifra esanguratsuekin erakutsi beharko lirateke. Adibidez, demagun  $p_1 = 4.152$  eta  $p_2 = 0.07931$  datuek lau zifrako zehaztasuna daukatela; orduan, tentagarria izango litzateke kalkulagailu baten pantailan agertzen diren zifra guztiak ematea, esate baterako, haien batuketa egitean:  $p_1 + p_2 = 4.23131$ . Baina hori ez da zuzena, ez genituzke emaitzak erabili beharko, horiek jatorrizko datuek baino zifra esanguratsu gehiago baldin badituzte. Hau izango litzateke emaitza egokia egoera horretan:  $p_1 + p_2 = 4.231$ .

## 3.5. Problemak

### Eskuz ebazteko problemak:

1. Egin 3.3. adibidearen pareko kalkuluak  $x_0 = 0.5$  puntuko  $f(x) = e^{-2x}$  funtzioaren deribatua hurbiltzeko. Behatu antzekotasunak eta desberdintasunak konparatzean zure taula adibideko taularekin.
2. Egin 3.2. adibidearen antzeko kalkuluak  $f'(x_0)$  hurbiltzeko, adierazpen hau erabiliz:

$$\frac{f(x_0 + h) - f(x_0 - h)}{2h}.$$

Erakutsi errorea  $O(h^2)$  ordenakoa dela. Hain zuzen, erroreaken gai nagusia ( $h$  nahiko txikia denean)  $-\frac{h^2}{3!}f'''(x_0)$  dela  $f'''(x_0) \neq 0$  bada.

3. Egin 3.3. adibidearen pareko kalkuluak 2. problemaren hurbilpena erabiliz. Behatu antzekotasunak eta desberdintasunak konparatzean zure taula adibideko taularekin.
4. Bihurtu zenbaki bitar hauek zenbaki dezimal:
  - (a)  $10101_{bi}$
  - (b)  $11111110_{bi}$
  - (c)  $0.11011_{bi}$
  - (d)  $0.1010101_{bi}$
  - (e)  $1.0110101_{bi}$
  - (f)  $11000101.101_{bi}$ .
5. Bihurtu zenbaki dezimal hauek zenbaki bitar:
  - (a) 23
  - (b) 378
  - (c) 0.6
  - (d)  $7/16$
  - (e)  $23/32$
  - (f)  $1/7$ .
6. Idatzi 81, 66.25 eta -0.625 zenbakiak era hauetan:
  - (a) Era bitarrean.
  - (b) Puntu higikorreko era normalizatuan.
  - (c) 32 biteko zehaztasun bakuneko kate baten bidez (IEEE-754 estandarrean).
7. Idatzi 256.1875, -30952 eta -0.0032 zenbakiak era hauetan:

- (a) Era bitarrean.
- (b) Puntu higikorreko era normalizatuan.
- (c) 64 biteko zehaztasun bikoitzeko kate baten bidez (IEEE-754 estandarrean).
8. Aurkitu  $R = 0.d_1d_2d_3d_4d_5d_6d_7$  hurbiltze-errorea, zazpi zifra esanguratsuko hurbilpen bitar hauetan:
- (a)  $R = 1/10 \approx 0.0001100_{bi}$ .
- (b)  $R = 1/7 \approx 0.0010010_{bi}$ .
9. Serie geometriko konbergenteen baturaren formula erabiliz, frogatu  $1/7 = 0.\overline{001}_{bi}$  garapen bitarra eta  $\frac{1}{7} = \frac{1}{8} + \frac{1}{64} + \frac{1}{512} + \dots$  baliokideak direla.
10. Serie geometriko konbergenteen baturaren formula erabiliz, frogatu  $1/5 = 0.\overline{0011}_{bi}$  garapen bitarra eta  $\frac{1}{5} = \frac{3}{16} + \frac{3}{256} + \frac{3}{4096} + \dots$  baliokideak direla.
11. Zehaztu zer gertatzen den, lau zifrako mantisa duen ordenagailu batek eragiketa hauek egiten dituenean:
- (a)  $\left(\frac{1}{3} + \frac{1}{5}\right) + \frac{1}{6}$
- (b)  $\left(\frac{1}{10} + \frac{1}{3}\right) + \frac{1}{5}$
- (c)  $\left(\frac{3}{17} + \frac{1}{9}\right) + \frac{1}{7}$
- (d)  $\left(\frac{7}{10} + \frac{1}{9}\right) + \frac{1}{7}$ .
12. (a) Zenbat zenbaki positibo desberdin adieraz daitezke  $(\beta, t, L, U) = (10, 2, -9, 10)$  puntu higikorreko sisteman?
- (b) Zenbat zenbaki normalizatu desberdin adieraz daitezke  $(\beta, t, L, U)$  puntu higikorreko sistema orokorrean?
13. (a)  $8/7 = 1.14285714285714\dots$  zenbakiak, noski, ez du adierazpen zehatzik sistema hamartarrean. Aurkitu  $\beta$  oinarriko eta  $t$  mantisako (zehaztasuneko) puntu higikorreko sistema bat, non zenbaki horrek adierazpen zehatza daukan.
- (b)  $\pi$  zenbakirako al dago horrelako sistema?
14. Aurkitu  $E_x$  errore absolutua eta  $R_x$  errore erlatiboa, eta zehaztu hurbilpenaren zifra esanguratsuen kopurua kasu hauetan:
- (a)  $x = 2.71828182$ ,  $\hat{x} = 2.7182$ .
- (b)  $y = 98350$ ,  $\hat{y} = 98000$ .
- (b)  $z = 0.000068$ ,  $\hat{z} = 0.00006$ .



15. Bete ezazu kalkulu hau:

$$p = \int_0^{1/4} e^{x^2} dx \approx \int_0^{1/4} \left( 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} \right) dx = \hat{p}.$$

Esan zein errore mota gertatzen den, eta konparatu  $\hat{p}$  emaitza egiazko balioarekin:  $p = 0.2553074606$ .

16. (a) Jo dezagun  $p_1 = 1.414$  eta  $p_2 = 0.09125$  datuak lau zifra esanguratsuko zehaztasunarekin emanda daudela. Egoera horretan, aurkitu  $p_1 + p_2$  eta  $p_1 p_2$  eragiketei dagozkien emaitzak.

(b) Jo dezagun  $p_1 = 31.415$  eta  $p_2 = 0.027182$  datuak bost zifra esanguratsuko zehaztasunarekin emanda daudela. Egoera horretan, aurkitu  $p_1 + p_2$  eta  $p_1 p_2$  eragiketei dagozkien emaitzak.

17. Bukatu kalkulu hauek eta esan nolako errorea gertatzen den:

$$(a) \frac{\sin(\pi/4 + 0.00001) - \sin(\pi/4)}{0.00001} = \frac{0.70711385222 - 0.70710678119}{0.00001} = \dots$$

$$(b) \frac{\ln(2 + 0.00005) - \ln(2)}{0.00005} = \frac{0.69317218025 - 0.69314718056}{0.00005} = \dots$$

18. Hiru zifrako eta biribiltzeko puntu higikorreko aritmetika erabiliz, kalkula itzazu batura hauek (batu adierazten den ordenan):

$$(a) \sum_{k=1}^6 \frac{1}{3^k}.$$

$$(b) \sum_{k=1}^6 \frac{1}{3^{7-k}}.$$

19. Izan bitez Taylorren garapen hauek:

$$\frac{1}{1-h} = 1 + h + h^2 + h^3 + O(h^4)$$

eta

$$\cos(h) = 1 - \frac{h^2}{2!} + \frac{h^4}{4!} + O(h^6).$$

Aurkitu horien baturaren eta biderkaduraren hurbiltze-ordenak.

20. Izan bitez Taylorren garapen hauek:

$$\sin(h) = h - \frac{h^3}{3!} + \frac{h^5}{5!} + O(h^7)$$

eta

$$\cos(h) = 1 - \frac{h^2}{2!} + \frac{h^4}{4!} + O(h^6).$$

Aurkitu horien baturaren eta biderkaduraren hurbiltze-ordenak.

21.  $f_1(x_0, h) = \cos(x_0 + h) - \cos(x_0)$  funtzioa eralda dezakegu beste  $f_2(x_0, h)$  bat lortzeko, formula trigonometriko hau erabiliz:

$$\cos(a) - \cos(b) = -2 \sin\left(\frac{a+b}{2}\right) \sin\left(\frac{a-b}{2}\right).$$

Beraz,  $f_1$  eta  $f_2$  funtzioek balio berdinak dituzte, aritmetika zehatzean,  $x_0$  eta  $h$  aldagaien balio guztietarako.

- (a) Aurkitu  $f_2(x_0, h)$ -ri dagokion adierazpena.
- (b)  $(f(x_0 + h) - f(x_0))/h$  adierazpenaren bidezko  $f'(x_0)$ -ren kalkulu hurbilduan, iradoki adierazpen bat ezeztapenaren errorea saihesteko. Idatzi MATLABeko programa bat zure adierazpena inplementatzen duena, eta kalkula itzazu  $f'(0.5)$ -ren hurbilpenak,  $h = 1.e - 20, 1.e - 19, \dots, 1$  balioetarako.
- (c) Azaldu zure emaitzen eta 3.3. adibideko emaitzen artean dagoen diferentzia.
22.  $f_1(x, \delta) = \sin(x + \delta) - \sin(x)$  funtzioa eralda dezakegu beste  $f_2(x, \delta)$  bat lortzeko, formula trigonometriko hau erabiliz:

$$\sin(a) - \sin(b) = 2 \cos\left(\frac{a+b}{2}\right) \sin\left(\frac{a-b}{2}\right).$$

Beraz,  $f_1$  eta  $f_2$  funtzioek balio berdinak dituzte, aritmetika zehatzean,  $x$  eta  $\delta$  aldagaien balio guztietarako.

- (a) Aurkitu  $f_2(x, \delta)$ -ri dagokion adierazpena.
- (b) Froga ezazu, analitikoki,  $f_1(x, \delta)/\delta$  eta  $f_2(x, \delta)/\delta$  adierazpenak,  $\delta$  nahiko txikia denean,  $\cos(x)$  funtzioaren hurbilpenak direla.
- (c) Idatzi MATLABeko funtzio bat  $g_1(x, \delta) = f_1(x, \delta)/\delta - \cos(x)$  eta  $g_2(x, \delta) = f_2(x, \delta)/\delta - \cos(x)$  kalkulatzeko,  $x = 3$  eta  $\delta = 1.e - 11$  hartuz.
- (d) Azaldu bi kalkuluen arteko emaitzen diferentziak.
23. Jo dezagun lehenengo deribatuaren hurbilpena, hots:

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}.$$

Formula honetarako trunkatze (edo diskretizazio) errorea  $O(h)$  da. Demagun  $f$  balioztatzean errore absolutua  $\varepsilon$ -k bornatzen duela, eta baztertu egingo ditugu oinarriko eragiketa aritmetikoak egitean sortutako erroreak.

- (a) Froga ezazu konputazio-errore osoa (trunkatzearen eta biribiltzearen konbinazioa) bornatuta dagoela balio honetaz:

$$\frac{Mh}{2} + \frac{2\varepsilon}{h},$$

non  $M = |f''(x)|$ -ren borne bat baita.

- (b) Zein  $h$ -ren baliotarako minimizatzen da aurreko bornea?
- (c) Erabiltzen dugun biribiltze-errorea, gutxi gorabehera,  $10^{-16}$  da. Erabili aurreko galderaren emaitza, 3.3. adibideko grafikoaren portaera azaltzeko. Grafikoaren itxura azaltzean, azaldu, baita ere, non espero dezakegun minimo hori aurkitzea.
24. Demagun  $f(x) = \frac{1 - \cos(x)}{\sin(x)}$  funtzioa.
- (a) Erabili formatu hamartarra, sei digitu esanguratsurekin (biribilduz),  $f(0.007)$  kalkulatzeko (kalkulagailua erabiliz).
- (b) Erabili MATLAB (`format long` erabiliz)  $f(x)$ -ren balioa kalkulatzeko, eta benetako errore erlatiboa, biribiltzearen ondorioz, (a) atalean lortutako  $f(x)$ -ren balioarekiko.
- (c) Biderkatu  $f(x)$  funtzioa  $\frac{1 + \cos(x)}{1 + \cos(x)}$  adierazpenaz, biribiltze-errorea izateko joera gutxiago duen  $f(x)$ -ren era bat lortzeko. Era berrian, erabili sei digitu esanguratsutako formatu hamartarra (biribilduz)  $f(0.007)$  kalkulatzeko (kalkulagailua erabiliz). Konparatu balioa, (a) eta (b) ataletan lorturiko balioekin.
25. Demagun  $f(x) = \frac{\sqrt{9+x} - 3}{x}$  funtzioa.
- (a) Erabili formatu hamartarra sei digitu esanguratsurekin (biribilduz)  $f(0.005)$  kalkulatzeko (kalkulagailua erabiliz).
- (b) Erabili MATLAB (`format long` erabiliz)  $f(x)$ -ren balioa kalkulatzeko, eta benetako errore erlatiboa, biribiltzearen ondorioz, (a) atalean lortutako  $f(x)$ -ren balioarekiko.
- (c) Biderkatu  $f(x)$  funtzioa  $\frac{\sqrt{9+x} + 3}{\sqrt{9+x} + 3}$  adierazpenaz, biribiltze-errorea izateko joera gutxiago duen  $f(x)$ -ren era bat lortzeko. Era berrian, erabili sei digitu esanguratsutako formatu hamartarra (biribilduz)  $f(0.005)$  kalkulatzeko (kalkulagailua erabiliz). Konparatu balioa (a) eta (b) ataletan lorturiko balioekin.
26. Honako kasu hauetan, aurkitu formula baliokide bat zifra esanguratsuen galera saihesteko:
- (a)  $\ln(x+1) - \ln(x)$ ,  $x$  handi baterako.
- (b)  $\sqrt{x^2+1} - x$ ,  $x$  handi baterako.
- (c)  $\cos^2(x) - \sin^2(x)$ ,  $x \approx \pi/4$  baterako.
- (d)  $\sqrt{\frac{1 + \cos(x)}{2}}$ ,  $x \approx \pi$  baterako.

Kasu bakoitzean, azaldu zergatik aukeratu duzun formula hori, eta eman zenbakizko adibide bat zehaztasunaren diferentzia agerian uzteko.

27. (a) Froga ezazu  $\ln(x - \sqrt{x^2 - 1}) = -\ln(x + \sqrt{x^2 - 1})$  betetzen dela  $|x| \geq 1$  bada.  
 (b) Bi formula horietako zein da egokiena zenbakizko kalkulurako? Azaldu zergatik, eta eman zenbakizko adibide bat zehaztasunaren diferentzia agerian jartzeko.
28. Azaldu adierazpen hauetarako ager daitezkeen zailtasunak, eta berridatzi formula horiek zenbakizko kalkulurako era egokiago batean:

(a)  $\sqrt{x + \frac{1}{x}} - \sqrt{x - \frac{1}{x}}$ ,  $x \gg 1$  denean.

(b)  $\sqrt{\frac{1}{a^2} + \frac{1}{b^2}}$ ,  $a \approx 0$  eta  $b \approx 1$  direnean.

Kasu bakoitzean, eman zenbakizko adibide bat zehaztasunaren diferentzia nabaria izateko.

29. Izan bitez adierazpen hauek:

$$p(x) = x^3 - 3x^2 + 3x - 1, \quad q(x) = ((x - 3)x + 3)x - 1, \quad r(x) = (x - 1)^3.$$

- (a) Lau zifrako puntu higikorreko aritmetika erabiliz eta biribilduz, kalkula itzazu  $p(2.72)$ ,  $q(2.72)$  eta  $r(2.72)$ . Kontuan izan  $p(x)$ -ren kalkuluan, aritmetika horrekin,  $(2.72)^3 = 20.12$  eta  $(2.72)^2 = 7.398$  direla.
- (b) Lau zifrako puntu higikorreko aritmetika erabiliz eta biribilduz, kalkula itzazu  $p(0.975)$ ,  $q(0.975)$  eta  $r(0.975)$ . Kontuan izan  $p(x)$ -ren kalkuluan, aritmetika horrekin,  $(0.975)^3 = 0.9268$  eta  $(0.975)^2 = 0.9506$  direla.
30. Demagun  $(\beta, t, L, U) = (10, 8, -50, 50)$  puntu higikorreko sistema duen makina bat erabiltzen dugula ekuazio koadratiko honen erroak kalkulatzeko:

$$ax^2 + bx + c = 0,$$

non  $a$ ,  $b$  eta  $c$  zenbaki erreal ezagunak baitira.

Ondorengo kasu bakoitzerako, esan zelako zenbakizko zailtasunak gerta daitezkeen formula estandarra erabiltzen badugu erroak kalkulatzeko. Azaldu, baita ere, nola gainditu zailtasun horiek (ahal denean):

(a)  $a = 1$ ;  $b = -10^5$ ;  $c = 1$ .

(b)  $a = 6 \cdot 10^{30}$ ;  $b = 5 \cdot 10^{30}$ ;  $c = -4 \cdot 10^{30}$ .

(c)  $a = 10^{-30}$ ;  $b = -10^{30}$ ;  $c = 10^{30}$ .

31. *Bigarren mailako ekuazioaren ebazpenaren formula hobetua.* Demagun  $a \neq 0$  eta  $b^2 - 4ac > 0$  ditugula, eta jo dezagun  $ax^2 + bx + c = 0$  ekuazioa. Haren erroak adierazpen ezagun hauen bidez kalkulatzeko dira:

(a)  $x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$  eta  $x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$ .

Frogatu erro horiek kalkula daitezkeela adierazpen baliokide hauen bidez:

$$(b) \ x_1 = \frac{-2c}{b + \sqrt{b^2 - 4ac}} \quad \text{eta} \quad x_2 = \frac{-2c}{b - \sqrt{b^2 - 4ac}}.$$

*Oharra.*  $|b| \approx \sqrt{b^2 - 4ac}$  denean, kontuz ibili behar dugu ezeztatzearen bidezko zehaztasunaren galera saihesteko. Baldin  $b > 0$  bada,  $x_1$  erroa (b) formularen bidez kalkulatu beharko genuke, eta  $x_2$  erroa (a) formulaz. Aldiz,  $b < 0$  bada,  $x_1$  erroa (a) formularen bidez kalkulatu beharko genuke, eta  $x_2$  erroa (b) formula erabiliz.

32. Demagun gure zenbaki-sistema hamartarrak 8 digitu esanguratsu erabiltzen dituela (hots,  $t = 8$ ). Erabili formula egokia  $x_1$  eta  $x_2$  kalkulatzeko, aurreko ariketan azaltzen den bezala, bigarren mailako ekuazio hauen erroak aurkitzeko:

(a)  $x^2 - 1000.001x + 1 = 0.$

(b)  $x^2 - 100000.0001x + 1 = 0.$

(c)  $x^2 + 100000.00001x + 1 = 0.$

33. Aztertu eragiketa hauen erroen hedapena (begira ezazu 3.4.5. azpiatala):

- (a) Hiru zenbakiren batura:

$$p + q + r = (\hat{p} + \varepsilon_p) + (\hat{q} + \varepsilon_q) + (\hat{r} + \varepsilon_r).$$

- (b) Zero ez den zenbaki baten alderantzizkoa:

$$\frac{1}{q} = \frac{1}{\hat{q} + \varepsilon_q}.$$

- (c) Zero ez den zenbaki batez zatitzea:

$$\frac{p}{q} = \frac{\hat{p} + \varepsilon_p}{\hat{q} + \varepsilon_q}.$$

### MATLABez ebazteko problemak:

34. (a) Marraztu  $(\beta, t, L, U) = (2, 3, -2, 3)$  puntu higikorren sistemako zenbaki guztien adierazpen hamartarrak.  $d_0.d_1d_2$  mantisako zenbaki txikiena 1.00 da eta zenbaki posible handiena 1.11 (sistema hamartarrean 1.75). Hortaz, begizta batean 1etik 1.75era joango gara, eta gehikuntzak  $\beta^{1-t} = 2^{-2} = 0.25$  izango dira. Hori da begizta baterako oinarria. Bestalde, berretzailea -2tik 3ra doa; hori da aurreko begiztaren barneko begizta baten tartea. Begizta bikoitz horrek eraikitzen du bektore bat sistema horretako zenbaki positibo guztiekin. Baina, zenbaki negatiboak eta zeroa ere marraztu behar dituzu. Marrazteko orduan,  $\mathbf{x}$  bektorean zenbaki hamartar horiek guztiak gorde behar dituzu, eta  $y_i = 0$  hartu  $(x_i, y_i)$  bikote guztietarako.
- (b) Zer gertatzen da sistemako zenbakien arteko distantziarekin? Zergatik?

35. Formula koadratiko klasikoak dio  $ax^2 + bx + c = 0$  ekuazioaren bi erroak honela kalkulatzeko direla:

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad \text{eta} \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}.$$

- (a) Sortu MATLABeko M fitxategi bat (`[x1,x2]=fkklas(a,b,c)` motakoa) erroak kalkulatzeko, eta erabili funtzio hori bi erroak kalkulatzeko kasu honetan:

$$a = 1, \quad b = -100000000, \quad c = 1.$$

- (b) Konpara itzazu emaitza horiek MATLABeko funtzio honek ematen duenarekin:  
`roots([a b c])`
- (c) Zer gertatzen da erroak eskuz edo kalkulagailu arrunt batez kalkulatzeko? Ikusi beharko duzu formula klasikoa ona dela erro bat kalkulatzeko, baina ez bestea. Beraz, erabili ezazu zuk sortutako funtzio hura zehaztasunez erro bat aurkitzeko, eta gero erabili propietate hau:

$$x_1 x_2 = \frac{c}{a}$$

beste erroa kalkulatzeko.

36. *Formula hobetua.* Erabili aurreko 31. problemaren informazioa ariketa hauek egiteko.

- (a) Idatzi algoritmo bat erroak era horretan kalkula ditzan eta  $|b| \approx \sqrt{b^2 - 4ac}$  de-nean, kontuz ibil dadila ezeztatzearen bidezko zehaztasunaren galera saihesteko. Alegia, zera kontuan hartuz:
- Baldin  $b > 0$  bada,  $x_1$  erroa azken formularen bidez kalkulatu beharko genuke, eta  $x_2$  erroa formula klasikoaz.
  - Aldiz,  $b < 0$  bada,  $x_1$  erroa formula klasikoaren bidez kalkulatu beharko genuke, eta  $x_2$  erroa azken formula erabiliz.
- (b) Gero, inplementatu algoritmo hori MATLABeko M fitxategi bat eraikiz (adibidez, `[x1,x2]=fkhobe(a,b,c)` motakoa) eta probatu ekuazio koadratiko egoki batzuekin (esate baterako,  $x^2 - 5x + 6 = 0$  edo  $x^2 - 100000000x + 1 = 0$ ).

37. Formula klasikoa eta formula hobetua erabiliz, kalkula itzazu bigarren mailako ekuazio hauen erroak:

- (a)  $x^2 - 1000.001x + 1 = 0$ .
- (b)  $x^2 - 100000.0001x + 1 = 0$ .
- (c)  $x^2 - 100000.00001x + 1 = 0$ .
- (d)  $x^2 - 1000000.000001x + 1 = 0$ .
- (e)  $x^2 - 100000000.00000001x + 1 = 0$ .

Egiaztatu emaitzak ( $x_1 + x_2 = -b/a$  eta  $x_1x_2 = c/a$  bete behar dute), eta konpara itzazu kasu bakoitzeko formula klasikoak eta formula hobetuak lorturiko emaitzak. Zein da metodo egokiena? Zergatik huts egiten du besteak?

38. Hau da  $f(x) = e^x$  funtzioaren Taylorren seriearen garapena:

$$f(x) = e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \frac{x^5}{5!} + \dots$$

Erabili ekuazio hori  $e^{-2}$  kalkulatzeko, kasu hauetan:

- (a) Lehenengo lau batugaiak erabiliz.
- (b) Lehenengo sei batugaiak erabiliz.
- (c) Lehenengo zortzi batugaiak erabiliz.

Kasu bakoitzean, kalkula ezazu trunkeze-errore absolutua eta errore erlatiboa. Erabili MATLAB, `format long` instrukzioaz,  $e^{-2}$ -ren benetako balioa kalkulatzeko. Erabili sei zenbaki esanguratsuko zenbaki hamartarrak (biribilduz), eragiketa guztietan.

39. Hau da  $\sin(x)$ -rako berretura-seriea:

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

MATLABeko funtzio honek seriea erabiltzen du  $\sin(x)$  kalkulatzeko:

```
function s = berresin(x)
% sin(x) berretura seriearen bidez kalkulatzeko du.
s = 0;
t = x;
n = 1;
while s+t ~= s;
s = s+t;
t = -x.^2/((n+1)*(n+2)).*t;
n = n+2;
end
```

- (a) Zergatik bukatuko da `while` begizta?
- (b) Aldatu `berresin.m`-ren lehenengo lerroa honela:

```
function [s,tmax,tkop]=berresin(x)
```

Gero, sartu lerro hauek `while` begizta hasi baino lehen:

```
tmax = abs(t);
tkop = 0;
```

eta lerro hauek begiztaren bukaeran:

```
tmax = max(tmax,abs(t));
tkop = tkop+1;
```

Erantzun galdera hauek  $x = \pi/2, 11\pi/2, 21\pi/2$  eta  $31\pi/2$  balioetarako:

- Zein da emaitzaren zehaztasuna?
- Zenbat gai erabili behar izan ditu?
- Zein izan da seriearen tamaina handieneko gaia?

Galdera horiek erantzuteko, komeni da honelako taula bat egitea:

x	$\pi/2$	$11\pi/2$	$21\pi/2$	$31\pi/2$
sin(x)-berresin(x)				
tmax				
tkop				

- (c) Zer ondorioztatzen da puntu higikorreko aritmetikaren eta berretura-serieen erabilera funtzioak balioztatzeke orduan?

40. 3.22. eta 3.23. adibideak kontuan hartuz, frogatu  $\{1/2^n\}_{n=1}^{\infty}$  segidaren gaiak garatzeko hiru metodo hauek erabil ditzakegula:

- (a)  $r_0 = 1$  eta  $r_n = \frac{1}{2}r_{n-1}$ ,  $n = 2, 3, \dots$  izanik.
- (b)  $p_0 = 1$ ,  $p_1 = 0.5$  eta  $p_n = \frac{3}{2}p_{n-1} - p_{n-2}$ ,  $n = 2, 3, \dots$  izanik.
- (c)  $q_0 = 1$ ,  $q_1 = 0.5$  eta  $q_n = \frac{5}{2}q_{n-1} - q_{n-2}$ ,  $n = 2, 3, \dots$  izanik.

Orain, kasu bakoitzean hasierako errore txiki bat sortzen da:

- (a)  $r_0 = 0.994$ .
- (b)  $p_0 = 1$ ,  $p_1 = 0.497$ .
- (c)  $q_0 = 1$ ,  $q_1 = 0.497$ .

MATLAB erabiliz, sortu metodo bakoitzerako lehenengo 10 iterazioak, eta aurkeztu emaitzak 3.4.5. eta 3.4.5. tauletan bezala (MATLABen bidez, hori ere).

Nolako errorea (egonkorra/ezegonkorra) gertatzen da kasu bakoitzean? Non txikitzen da errore hori?



## 4. kapitulua

# Ekuazio ez-linealen ebazpena

Medikuntza-ikasketek baieztatu dute puenting-jauzilari batek bizkarrezurreko kalte garrantzitsu bat jasateko duen probabilitatea handitu egiten dela baldin erorketa libreko abiadura 36 m/s baino handiagoa bada, erorketa librean 4 s egon ondoren. Orduan, zure nagusiak nahi du puenting-enpresan zuk zehaztea zein masatarako gainditzen den irizpide hori, airearen erresistentzia-koefizientea 0.25 kg/m bada.

Zuk badakizu, aurreko ikasketen bidez, adierazpen hau erabil dezakegula denboraren menpeko erorketaren abiadura jakiteko:

$$v(t) = \sqrt{\frac{g \cdot m}{c_d}} \tanh\left(\sqrt{\frac{g \cdot c_d}{m}} t\right).$$

Ezin da bakandu  $m$ . Aukera bat da  $m$  aurkitzeko funtzio berri honen erroa kalkulatzeko:

$$f(m) = \sqrt{\frac{g \cdot m}{c_d}} \tanh\left(\sqrt{\frac{g \cdot c_d}{m}} t\right) - v(t). \quad (4.1)$$

Kapitulu honetan ikusiko dugu nola erabil dezakegun ordenagailua horrelako soluzioak kalkulatzeko.

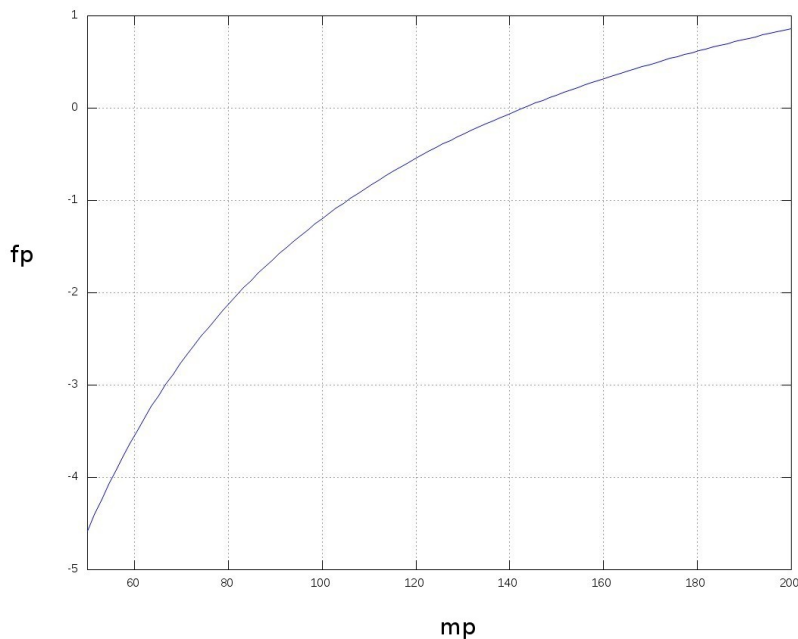
### 4.1. Metodo grafikoak

Metodo erraz bat da  $f(x) = 0$  ekuazioko erroaren hurbilpen bat lortzeko funtzioaren marrazketa egitea, eta aztertzea non zeharkatzen duen  $x$  ardatza. Puntu horrek erroaren saihurbilpen bat ematen digu.

**4.1. adibidea.** *Hurbilpen grafikoa erabiliz, aurkitu puenting-jauzilariaren masa,  $c_d = 0.25$  kg/m bada, erorketa libreko 4 s igaro ondoren abiadura 36 m/s-koa izateko ( $g = 9.81$  m/s<sup>2</sup>).*

*Ebazpena:*  $m$ -rekiko (4.1) funtzioa marrazteko, MATLAB erabiliz, hau egingo dugu:

```
>> cd=0.25; g=9.81; v=36; t=4;
>> mp=linspace(50,200);
>> fp=sqrt(g*mp/cd).*tanh(sqrt(g*cd./mp)*t)-v;
>> plot(mp,fp), grid
```



**4.1. irudia.**  $(mp, fp)$  grafikoa.

Funtzio horrek 140 eta 150 kg artean zeharkatzen du  $m$  ardatza. Marrazketaren azterketak erroaren 145 kg-ko sasihurbilpen bat ematen digu. Hurbilpen grafiko horren baliotasuna probatzeko, (4.1) adierazpenean ordezkatu dugu:

```
>> sqrt(g*145/cd).*tanh(sqrt(g*cd./145)*t)-v
ans =
    0.0456
```

zeina zerotik hurbil baitago.  $\square$

Teknika grafikoek erabilera praktikoa mugatua dute, oso zehatzak ez direlako. Hala ere, metodo grafikoak erabil daitezke erroen sasihurbilpenak aurkitzeko, eta horiek zenbakizko metodoetan erabil daitezke, hasierako hurbiltze-puntuak bezala.

## 4.2. Bakartze-metodoak

Erroari buruzko informazioa lortzeko metodo nagusiak hauek dira:

- *Bakartze-metodoak*. Metodo hauetan, hasierako bi hurbiltze-puntuk erro bakoitza bakartzen du tarte batean.
- *Metodo irekiak*. Metodo hauetan, hasierako hurbiltze-puntu bat (edo gehiago) erabiltzen da, baina ez da beharrezkoa erroa tarte batean bakartzea.

Bakartze-metodoen konbergentzia motela da. Aldiz, metodo irekien konbergentzia, oro har, azkarragoa da; baina haiek dibergenteak izan daitezke.

Bi kasuetan, beharrezkoak dira hasierako hurbiltze-puntuak. Aztertzen ari garen testuinguru fisikotik atera ditzakegu horiek. Hala ere, batzuetan ez da erraza horrelako hastapen-datu onak iragartzea.

Ekuazio batek erro anitz baldin badauzka tarte batean, erroak bakartzeko tarte txikiagoak definituko ditugu. Hori funtzioa marraztuz egin dezakegu, edo funtzioa balioztatuz segidako puntuen multzo bat hartuz, eta aztertuz funtzioaren zeinua zein tartetan aldatzen den.

Ondorengo atalean kasu berezi bat ikusiko dugu, non erroak analitikoki bakartu baititzakegu.

### 4.2.1. Erroak bakartzea

Izan bitez  $[a, b]$  tarte eta  $f : [a, b] \subset \mathbb{R} \rightarrow \mathbb{R}$  funtzio ondo definitua eta jarraitua tarte horretan. Funtzioaren erroak bilatzea eta  $f(x) = 0$  ekuazioaren soluzioak asmatzea problema bera da. Berdintza hori betetzen duen  $p$  soluzioari *funtzioaren erro* esango diogu.

Batzuetan,  $f'(x)$  funtzioaren existentzia funtsezkoa izango da helburua lortzeko eta, zenbaitetan,  $f''(x)$ -rena ere bai. Jo dezagun  $f(x) = 0$  ekuazioaren erro guztiak desberdinak direla. Hortaz, erro bakoitzaren ingurune nahiko txiki baterako ez dago beste errorik. Orduan, bi urrats hauek beteko ditugu:

1. Erroen bereizketa: erro bakoitzaren  $[\alpha, \beta]$  ingurune ahalik eta txikiena aurkitzea; non tartearen barnean ez baita ekuazioaren beste errorik egon behar.
2. Erro hurbilduen balioak doitzea: emaitza gero eta hobeak hurbiltzea.

Jarraian, kontuan hartu beharko ditugu bi propietate hauek:

- (i) Izan bedi  $f$  jarraitua  $[a, b]$  tartean eta  $[a, b]$ -ko muturretan zeinu desberdinekin (hots,  $f(a) \cdot f(b) < 0$ ); orduan,  $f$  funtzioak gutxienez  $c$  erro bat izango du  $(a, b)$  tartean; hots,  $\exists c \in (a, b)$ , non  $f(c) = 0$  baita (Bolzano-ren teorema).
- (ii) Baldin  $f$  deribagarria bada  $(a, b)$  tartean, eta  $f'$ -k tarte horretan zeinua aldatzen ez bada (hau da,  $f$  monotonoa bada), orduan,  $f$  funtzioak ez du  $[a, b]$  tartean erro bat baino gehiago izango.

Ondorioz,  $f(x) = 0$  ekuazioaren erroak kalkulatzeko urrats hauei jarraituko diegu:

1. Funtzioen erroak tarte desberdinetan bakartuko ditugu (ii) aplikatuz, ahal bada. Lehendabizi, emandako tartearen partizio bat egiten da:  $a < \alpha_1 < \alpha_2 < \dots < \alpha_n < b$  (adibidez,  $f$  deribagarria bada, deribatuaren zeinu-aldaketak erabiliz lor dezakegu partizioa; hau da,  $f'(x)$  funtzioaren erroak kalkulatzuz).
2. Ondoz ondoko puntu-bikoteak bilatuko ditugu, (i) aplikatuz: hala,  $f$  funtzioak zeinu desberdinak hartuko ditu. Hau da,  $\alpha_k, \alpha_{k+1}$  izango ditugu, non  $f(\alpha_k) \cdot f(\alpha_{k+1}) < 0$  baita. Horrek adierazten digu ezen tarte horretan  $f(x)$ -ren erro bat existitzen dela.

#### 4.2. adibidea. Ebatzi $x^4 - 4x - 1 = 0$ ekuazioa.

*Ebazpena.* Izan bedi  $f(x) = x^4 - 4x - 1$ . Funtzio horrek, 4. mailako polinomio bat denez, gehienez lau erro izango ditu. Funtzio horren zeinu-aldaketak finkatzeko, aski dira emandako tartearen muturrak eta funtzioaren maximo eta minimo lokalak. Kasu horretan,  $f'(x) = 4(x^3 - 1)$  aztertzen da, eta, ondorioz,  $x = 1$  da  $f'(x)$ -ren zeinu-aldaketa bakarra. Hain zuzen ere:

$$\begin{aligned} x \in (-\infty, 1) &\Rightarrow f'(x) < 0 \\ x \in (1, \infty) &\Rightarrow f'(x) > 0. \end{aligned}$$

Hortaz,  $(-\infty, 1)$  tartean,  $f(x)$  beherakorra da, eta,  $(1, \infty)$  tartean,  $f(x)$  gorakorra. Gainera, tarte horien muturretan zera gertatzen da:

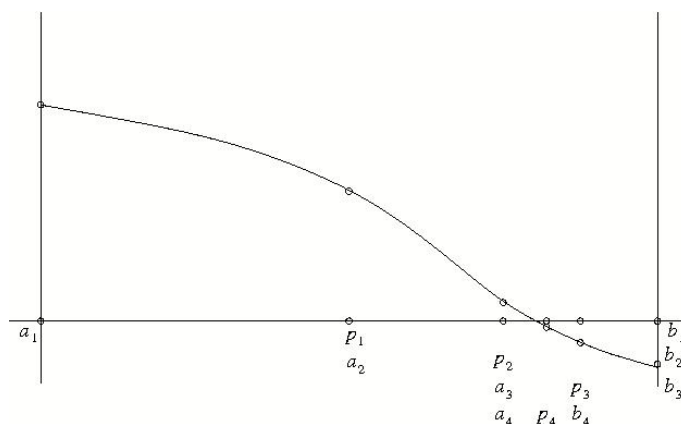
$$f(-\infty) \cdot f(1) < 0 \quad \text{eta} \quad f(1) \cdot f(\infty) < 0,$$

Ondorioz, bi tarte horietan  $f(x)$ -k 0 balioa hartzen du (i) aplikatuz. Hau da, tarte bakoitzean erro bat dago.  $\square$

Erroak tarte desberdinetan bakartzeko beste metodo bat da funtzioaren grafikoa aztertzea (esate baterako, MATLAB erabiliz). Geroago, erro bat duen tarte bakoitza txikitzen da, adibidez, bisekzio-metodoaren bitartez.

### 4.2.2. Bisekzio-metodoa

Demagun  $f$  funtzioa jarraitua dela  $[a, b]$  tartean, eta  $f(a) \cdot f(b) < 0$  betetzen dela; orduan, (i) propietatearen arabera  $p$  balio bat existitzen da  $(a, b)$  tartean, non  $f(p) = 0$  baita. Suposa dezagun, orobat, tarte horretan erro hori baino ez dagoela.



4.2. irudia. Bisekzio-metodoa.

Lehendabizi,  $a_1 = a$  eta  $b_1 = b$  izendatuko dira, eta emandako tarte bi zati berdinetan bereiziko da;  $[a, b] = [a_1, p_1] \cup [p_1, b_1]$ , non  $p_1 = (a_1 + b_1)/2$  erdiguneko balioa baita. Orain, hasierako tartearen luzera-erdiko azpitarte batean egongo da  $p$  erroa, eta, zein den jakiteko,  $f(a_1) \cdot f(p_1)$  balioztatzen da. Baldin azken balio hori negatiboa bada, orduan lehenengo azpitartean egongo da, eta, ondorioz, prozesua  $[a_2, b_2]$  tartean errepikatu behar da, tarte horretarako  $a_2 = a_1$  eta  $b_2 = p_1$  hartuz; bestela,  $a_2 = p_1$  eta  $b_2 = b_1$  izendatu beharko dira. Iterazio bakoitzean dagoen tartearen luzera erdibitu egiten da, eta,  $n$  iterazio gauzatu ondoren,  $[a_n, b_n]$ -ren luzera  $(b - a)/2^{n-1}$  izango da. Noski,  $p \in [a_n, b_n]$  dagoenez, errorea  $|p_n - p| \leq (b - a)/2^n$  izango da. Beraz,  $n$  handituz, errorea gero eta txikiagoa egiten da. Hori formaliza daiteke teorema honetan:

**4.1. teorema.** Demagun  $f \in C[a, b]$  eta  $f(a)f(b) < 0$  dela. Izan bedi  $\{p_n\}_{n=1}^{\infty}$  segida, non  $p_n = (a_n + b_n)/2$  baita. Orduan, badago  $p \in (a, b)$  puntu bat, non  $f(p) = 0$ , hau betetzen duena:

$$\varepsilon_n = |p - p_n| \leq \frac{b - a}{2^n}, \quad n = 1, 2, \dots \quad (4.2)$$

eta, bereziki,  $\lim_{n \rightarrow \infty} p_n = p$ .

*Frogapena.* Metodoaren deskripziotik (4.2) lortzen da. Bestalde, hau betetzen da:

$$\lim_{n \rightarrow \infty} |p - p_n| \leq \lim_{n \rightarrow \infty} \frac{b - a}{2^n} = 0,$$

eta horrek  $\lim_{n \rightarrow \infty} p_n = p$  inplikatzeko du.  $\square$

Aurreko teoremaren arabera, hasierako  $\tau > 0$  tolerantzia bat (errorearen goi-borne bat) jar dezakegu, eta  $\varepsilon_n < \tau$  denean bukatutzat eman prozesua. Bestalde, prozesua mozteko beti jar dezakegu iterazioen  $n_{max}$  kopuru maximo bat; hots,  $n = n_{max}$  denean, prozesua bukatzen da eta  $p_n$ -rekin geratzen gara.

**4.3. adibidea.** *Bisekzio-algoritmoa erabiliz, kalkulatu  $e^x = \sin x$  ekuazioaren soluzioa  $x \in [-4, -3]$  tartean, errorearen tolerantzia 0.01 izanik.*

*Ebazpena.* Izan bedi  $f(x) = e^x - \sin x$ ; orain, ekuazioaren soluzio bat  $f$  funtzioaren erro bat izango da. Hasteko,  $[a_1, b_1] = [-4, -3]$ ,  $f(a_1) = f(-4) < 0$  eta  $f(b_1) = f(-3) > 0$  direnez,  $[-4, -3]$  tartean gutxienez erro bat dago.

Jarraian,  $[-4, -3]$  tartea erdibituko dugu. Erdigunea  $p_1 = -3.5$  da;  $[-4, -3] = [-4, -3.5] \cup [-3.5, -3]$ ,  $f(p_1) = f(-3.5) < 0$  denez,  $[a_2, b_2] = [-3.5, -3]$  tartean egongo da erroa.

Orain, tartearen erdigunea  $p_2 = -3.25$  da;  $[-3.5, -3] = [-3.5, -3.25] \cup [-3.25, -3]$ ,  $f(p_2) = f(-3.25) < 0$  denez,  $[a_3, b_3] = [-3.25, -3]$  tartean egongo da erroa.

Eta horrela jarraituz, nahi dugun bezainbeste hurbilduko gara erroantz.

$n$	$a_n$	$b_n$	$p_n$	$ f(p_n) $
1	-4	-3	-3.5	0.3206
2	-3.5	-3	-3.25	0.0694
3	-3.25	-3	-3.125	0.0605
4	-3.25	-3.125	-3.1875	0.0046

**4.1. taula.** Bisekzio-metodoa.

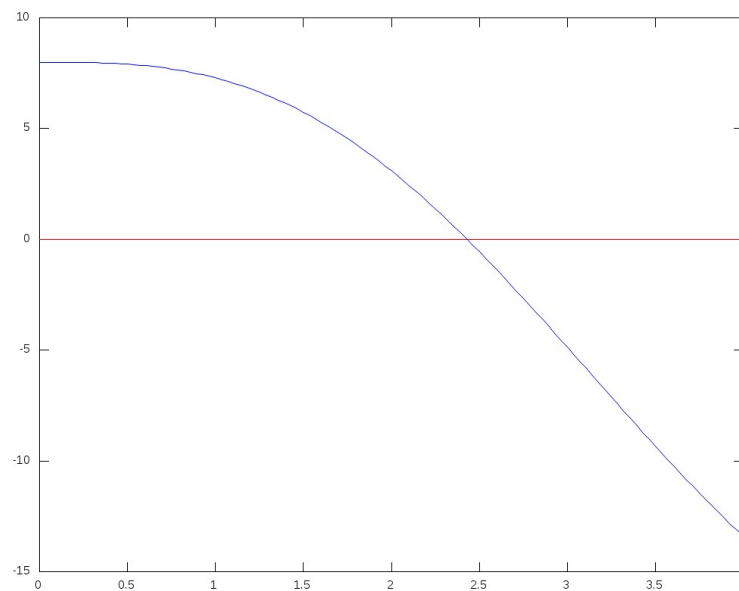
Beraz, lau iterazio egin ondoren lorturiko soluzio hurbildua  $x_4 = -3.1875$  da, eta  $|f(x_4)| = 0.0046 < 0.01$  sortzen den erroa.  $\square$

### Bisekzio-metodoari buruzko oharra

Alde batetik, bisekzio-metodoak bi eragozpen nabari ditu, eta hauexek dira: bata, metodoa astiro hurbiltzen dela soluziorantz, eta bestea, hurbildu bitartean alde batean utz ditzakeela hurbilketa onak. Beste aldetik, metodo honek soluziorantz jotzen du, eta, horregatik, erabilgarria da hasierako estimazio bat asmatzeko, eta, jarraian, beste metodo eraginkorrago bat aplikatu ahal izateko. Metodo horrek huts egin dezake  $p$  puntuan, baldin abszisa  $y = f(x)$  kurbaren zuzen ukitzaila bada eta kurbak abszisa ez badu zeharkatzen.

**4.4. adibidea.** Idatzi MATLAB programa bat, instrukzioen  $M$  fitxategi bat,  $8 - 4.5(x - \sin(x)) = 0$  ekuazioaren soluzioa aurkitzeko, bisekzio-metodoa erabiliz. Soluzioaren errore-tolerantziak  $0.001$  rad izan behar du (hots, errorea  $< 0.001$ ). Sortu taula bat, non bisekzio-prozesuko iterazio bakoitzean  $a$ ,  $b$ ,  $p_n$ ,  $f(p_n)$  eta errorea agertzen baitira.

*Ebazpena.* Soluzioaren hasierako hurbiltze bat egiteko,  $f(x) = 8 - 4.5(x - \sin(x))$  marrazten dugu, MATLABen `fplot` instrukzioa erabiliz. Horrek erakusten digu  $x = 2$ ren eta  $x = 3$ ren artean soluzioa dagoela. Beraz,  $a = 2$  eta  $b = 3$  muturrekin aukeratuko da hasierako tartea.



**4.3. irudia.**  $f(x) = 8 - 4.5(x - \sin(x))$ -ren grafikoa.

MATLAB programa honek ebazten du problema hori:

```
clear all
f=inline('8-4.5*(x-sin(x))');
a=2; b=3; nmax=20; tol=0.001;
fa=f(a); fb=f(b);

if fa*fb>0
    disp('Errorea: Funtzioak a eta b puntuetan zeinu berbera du.')
else
    disp('n      a          b          p_n          f(p_n)      errorea')
    disp('_____')
```

```

for n=1:nmax
    p_n=(a+b)/2;
    err=(b-a)/2;
    fp_n=f(p_n);
    fprintf('%3i    %11.6f %11.6f %11.6f %11.6f %11.6f\n',n,a,b,p_n,fp_n,err)
    if fp_n==0
        fprintf('p=%11.6f soluzio zehatz bat aurkitu dugu',p_n)
        break
    end
    if err<tol
        fprintf('p=%11.6f soluzio hurbildu bat aurkitu dugu',p_n)
        break
    end
    if i==nmax
        fprintf('Ez dugu lortu soluzioa %i iterazioetan',nmax)
        break
    end
    if fa*fp_n<0
        b=p_n;
    else
        a=p_n;
    end
end
end
end

```

### 4.2.3. *Regula falsi* metodoa

Metodo honen beste izen batzuk dira posizio faltsuaren eta interpolazio linealaren metodoak. Bisekzio-metodoan bezala,  $[a, b]$  tartean  $f$  jarraitua da,  $f(a) \cdot f(b) < 0$  eta, beraz,  $f(x) = 0$ -k erro bat dauka tarte horretan. Bisekzio-metodoan, hurrengo urratsa egiteko,  $[a, b]$  tartearen erdigunea erabiltzen da. Izan bedi  $(a, f(a))$ ,  $(b, f(b))$  puntuetatik igarotzen den  $y = f(x)$ -ren  $E$  zuzen ebakitzailea. Oraingo metodoan hurbilpena hobetzen da,  $E$  zuzen horretako  $(\bar{p}, 0)$  puntua erabiliz; hots,  $E$ -ren ebakitze-puntua abszisarekin. Orain  $\bar{p}$  aurkituko dugu.  $E$ -ren  $m$  malda kalkulatzeko, bi adierazpen hauek ditugu:

$$m = \frac{f(b) - f(a)}{b - a} \quad \text{eta} \quad m = \frac{0 - f(b)}{\bar{p} - b},$$

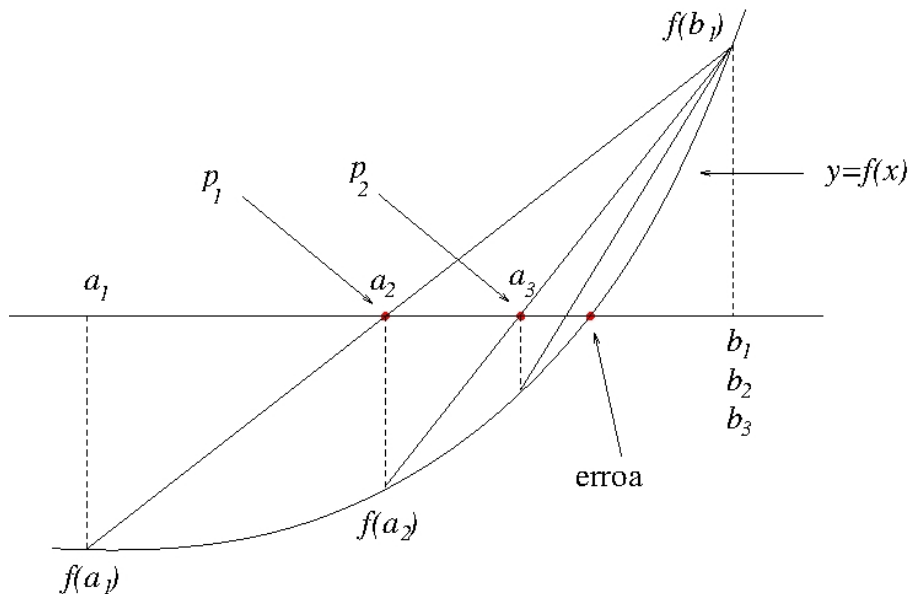
biak berdinduz,

$$\frac{f(b) - f(a)}{b - a} = \frac{0 - f(b)}{\bar{p} - b}$$

eta  $\bar{p}$  bakanduz formula hau ateratzen dugu:

$$\bar{p} = b - \frac{(b - a)}{f(b) - f(a)} f(b). \quad (4.3)$$





4.4. irudia. Regula falsi metodoa.

Hurrengo iteraziorako, metodo honen aukerak bisekzioaren berdinak dira, alegia:

- $f(a) \cdot f(\bar{p}) < 0$  bada,  $[a, \bar{p}]$  tartean erro bat dago.
- $f(\bar{p}) \cdot f(b) < 0$  bada,  $[\bar{p}, b]$  tartean erro bat dago.
- $f(\bar{p}) = 0$  bada,  $\bar{p}$  da  $f(x)$ -ren erro bat.

Aurreko (4.3) formula erabil dezakegu, aukera horiekin batera  $\{[a_n, b_n]\}$  tarteen segida sortzeko, horietako bakoitzak erro bat edukiz.

### Regula falsi metodoari buruzko oharrak

- Iterazio bakoitzean, hau izango da erroaren hurbilpena:

$$p_n = b_n - \frac{(b_n - a_n)}{f(b_n) - f(a_n)} f(b_n) = \frac{f(b_n)a_n - f(a_n)b_n}{f(b_n) - f(a_n)}$$

eta egiazta daiteke, goiko baldintzetan,  $\{p_n\}$ -k jotzen duela funtzioaren  $p$  erro batera.

- Iterazioak geldituko dira  $\varepsilon_n = |f(p_n)|$  erroa erabiltzaileak finkaturiko tolerantzia bat baino txikiago denean, edo iterazio kopurua ezarritako maximora heltzen denean.
- Sarritan, funtzioa gorako edo beherako ahurra da  $[a, b]$  tartean; orduan, tartearen mutur bat finko geratzen da iterazio guztietan, beste muturrak errorantz jotzen duen bitartean. Hots, zenbakizko soluzioa errorantz doa, alde batetik bakarrik.

## 4.3. Metodo irekiak

### 4.3.1. Puntu finkoaren metodoa

Metodo irekiek formula bat erabiltzen dute erro bat iragartzeko. Horrelako formula garatu dezakegu *puntu finkoaren iterazioa* erabiliz,  $f(x) = 0$  ekuazioa berrordenatuz ekuazioaren ezker aldean  $x$  izateko, alegia:

$$x = F(x). \quad (4.4)$$

Transformazio hori lor dezakegu manipulazio algebraikoaren bidez, edo  $x$  batuz jatorrizko ekuazioaren bi aldeetan. Orduan, (4.4) adierazpenaren bidez iragar daiteke  $x$ -rentzat balio berri bat,  $x$ -ren balio zahar batetik. Hots,

$$p_n = F(p_{n-1}). \quad (4.5)$$

### Konbergentziaren analisia

**4.2. teorema.** *Demagun  $F \in C[a, b]$  dela.*

(i)  $y = F(x)$  funtzioaren irudiak  $y \in [a, b]$  betetzen badu,  $x \in [a, b]$  puntu bakoitzeko, orduan,  $F$ -k puntu finko bat dauka  $[a, b]$  tartean.

(ii) Demagun, gainera,  $F'(x)$  definituta dagoela  $(a, b)$  tartean, eta  $|F'(x)| < 1$  dugula,  $x \in (a, b)$  denean; orduan,  $F$ -k puntu finko bakar bat dauka  $[a, b]$  tartean.

(i)-ren frogapena.  $F(a) = a$  edo  $F(b) = b$  bada, ondorioa egia da. Demagun, orduan,  $F(a) \in (a, b]$  eta  $F(b) \in [a, b)$  egiaztatzen direla.  $f(x) = x - F(x)$  funtzioak propietate hau dauka:

$$f(a) = a - F(a) < 0 \text{ eta } f(b) = b - F(b) > 0.$$

Bolzanoren teorema aplikatuz, badago  $p \in (a, b)$  puntu bat non  $f(p) = 0$  baita. Beraz,  $p = F(p)$  eta  $p$  puntua  $F(x)$ -ren puntu finkoa da.

(ii)-ren frogapena. Puntu finko bakar bat dagoela frogatu behar dugu. Aldiz, demagun bi puntu finko daudela,  $p_1$  eta  $p_2$  ( $p_1 \neq p_2$ ); orduan, batez besteko balioaren teorema (Lagranjeren teorema) erabiliz, badago  $d \in (a, b)$  puntu bat hau betetzen duena:

$$F'(d) = \frac{F(p_2) - F(p_1)}{p_2 - p_1}.$$

Baina bi puntu horiek puntu finkoak direnez, hau bete behar:

$$F'(d) = \frac{p_2 - p_1}{p_2 - p_1} = 1,$$

eta hori  $|F'(x)| < 1$  teoremako baldintzaren kontraesana da. Beraz, ezin dira egon bi puntu finko desberdin.  $\square$

**4.5. adibidea.** *Frogatu  $F(x) = \cos(x)$  funtzioak puntu finko bakar bat daukala  $[0, 1]$  tartean.*

*Ebazpena.* Bistakoa da  $F \in C[0, 1]$  dela. Gainera,  $F(x) = \cos(x)$  funtzioa tarte horretan beherakorra denez, haren irudiak  $F([0, 1]) = [\cos(1), 1] \subset [0, 1]$  betetzen du. Orduan, aurreko teoremaren (i) atala betetzen du. Alegia,  $F$ -k puntu finko bat du  $[0, 1]$  tartean. Azkenik,  $x \in (0, 1)$  bada,  $|F'(x)| = |-\sin(x)| = \sin(x) \leq \sin(1) \approx 0.8415 < 1$ . Beraz, aurreko teoremaren (ii) baldintza betetzen da, eta, ondorioz, bakarra da  $F$ -ren puntu finkoa  $[0, 1]$  tartean.  $\square$

**4.3. teorema. (Puntu finkoaren teorema).** *Demagun (a)  $F, F' \in C[a, b]$ , (b)  $K > 0$ , (c)  $p_0 \in (a, b)$  eta (d)  $F(x) \in [a, b]$ ,  $x \in [a, b]$  guztietarako. Orduan, badago  $F$ -ren  $p \in [a, b]$  puntu finko bat.*

(i) *Baldin  $|F'(x)| \leq K < 1$ ,  $x \in [a, b]$  guztietarako, orduan,  $p$  da  $F$ -ren puntu finko bakarra  $[a, b]$  tartean, eta  $p_n = F(p_{n-1})$  iterazioak  $p$  puntura jotzen du. Kasu horretan,  $p$  puntu finko erakargarria dela esaten da. Gainera,  $e_n = |p_n - p|$  erroreak ez badira zero,  $e_{n+1}/e_n \rightarrow F'(p)$ ,  $n \rightarrow \infty$  denean.*

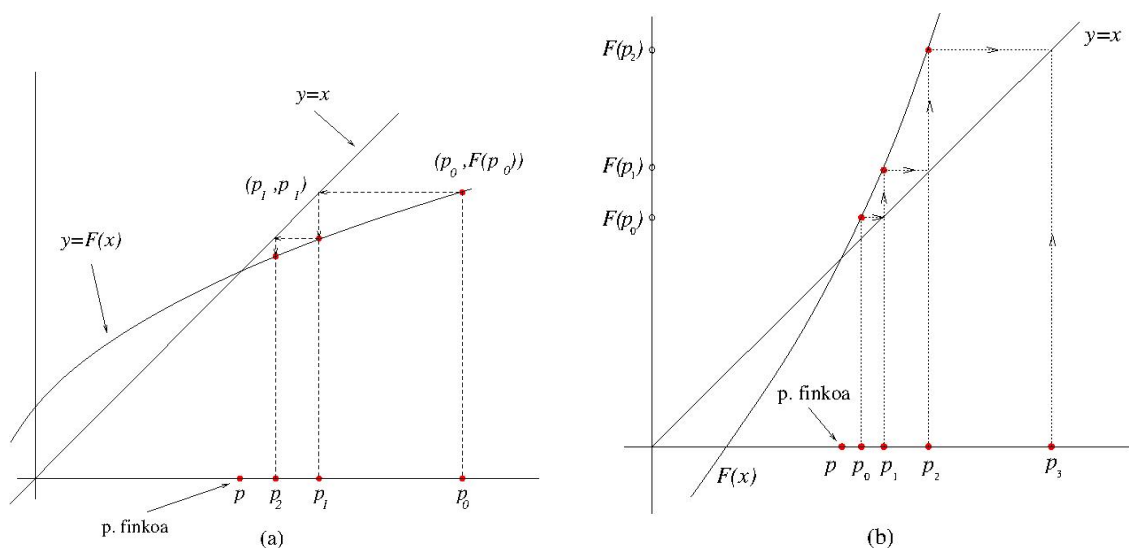
(ii) *Baldin  $|F'(x)| > 1$  eta  $p_0 \neq p$  badira, orduan,  $p_n = F(p_{n-1})$  iterazioak ez du jotzen  $p$  puntura. Kasu horretan,  $p$  puntu finko aldaragarria dela esaten da, eta iterazioak dibergentzia lokala du.*

(i)-ren frogapena. 4.2. teoremako (i) atalak eta oraingo teoremaren (a) eta (d) hipotesiek bermatzen dute  $F$ -k  $[a, b]$  tartean puntu finko bat daukala. Bestalde,  $|F'(x)| < 1$  izateak eta 4.2. teoremako (ii) atalak inplikatzeko dute  $F$ -ren puntu finko bakarra dela  $[a, b]$  tartean. Bestalde, indukzioz ikusiko dugun bezala, (c) eta (d) baldintzetatik ateratzen da  $\{p_n\}_{n=0}^{\infty}$  segidako puntu guztiak  $[a, b]$  tartean daudela. Jarraian,  $p_0 \in [a, b]$  puntutik hasiz, batez besteko balioaren teorema aplikatuko dugu,  $p_0$  eta  $p$ -ren arteko  $c_0 \in (a, b)$  existitzen dela ateratzeko, hau egiaztatzen duena:

$$\begin{aligned} |p - p_1| &= |F(p) - F(p_0)| = |F'(c_0)(p - p_0)| = |F'(c_0)| \cdot |p - p_0| \\ &\leq K|p - p_0| < |p - p_0|, \end{aligned} \quad (4.6)$$

Beraz,  $p_1$   $p$ -tik hurbilago dago  $p_0$  baino. Horrela arrazoituz, oro har, hau izango dugu:

$$\begin{aligned} |p - p_n| &= |F(p) - F(p_{n-1})| = |F'(c_{n-1})(p - p_{n-1})| = |F'(c_{n-1})| \cdot |p - p_{n-1}| \\ &\leq K|p - p_{n-1}| < |p - p_{n-1}|. \end{aligned} \quad (4.7)$$



**4.5. irudia.** (a) P. finko erakargarria ( $|F'(x)| < 1$ ).  
 (b) P. finko alderagarria ( $|F'(x)| > 1$ ).

Frogapena bukatzeko, hau betetzen dela ikusi behar dugu:

$$\lim_{n \rightarrow \infty} |p - p_n| = 0.$$

Lehenik, indukzioz (4.6)-(4.7) erlazioak erabiliz, hau erraz lortzen da:

$$|p - p_n| \leq K^n |p - p_0|. \quad (4.8)$$

Gainera,  $0 < K < 1$  denez,

$$0 \leq \lim_{n \rightarrow \infty} |p - p_n| \leq \lim_{n \rightarrow \infty} K^n |p - p_0| = 0.$$

Ondorioz,  $\lim_{n \rightarrow \infty} |p - p_n| = 0$ , eta, orduan,  $\lim_{n \rightarrow \infty} p_n = p$ .

Azkenik,  $e_n = |p_n - p|$  erroreak ez badira zero, deribatuaren definizioaz zera betetzen da:

$$\frac{e_{n+1}}{e_n} = \frac{|p_{n+1} - p|}{|p_n - p|} = \frac{|F(p_n) - F(p)|}{|p_n - p|} \rightarrow |F'(p)|.$$

(ii) frogapena ariketa gisa geratzen da.  $\square$

Kontuan izan 4.3. teoremaren (i) ataleko baldintzak betetzen direnean (4.8) betetzen dela, eta, ondorioz, puntu finkoaren iteraziorako erroreaken borne hau dugula:

$$|p - p_n| \leq K^n |p - p_0|, \quad n \geq 1 \text{ guztietarako.}$$

$K$  horren balio hurbildu on bat  $F'(\tilde{p})$  izan daiteke,  $\tilde{p}$  puntua  $p$  puntu finkotik nahiko hurbila bada ( $p$  bera sartuta, ezagutzen bada).

**4.1. Korolaria.** Demagun  $F, F' \in C(p - r_0, p + r_0)$ , non  $F(p) = p$ . Orduan, hau betetzen bada:

$$|F'(p)| < 1, \quad (4.9)$$

$r$  positibo bat existitzen da,  $r < r_0$ , non  $p_0 \in [p - r, p + r]$  betetzen bada,  $p_n = F(p_{n-1})$  segidako gai guztiak  $[p - r, p + r]$  tartean baitaude eta  $\{p_n\} \rightarrow p$ .

Aurreko teoremako (i) atala kontuan hartuz froga daiteke (ariketa gisa geratzen da).  $\square$

**4.6. adibidea.** Puntu finkoaren metodoa erabiliz, kalkulatu  $\cos x - x = 0$  ekuazioaren soluzio hurbildu bat, errore-tolerantzia = 0.0001 eta  $p_0 = \pi/4 = 0.785398$  hartuz. Gehienez, egin hamar iterazio. (Errorea  $|p - p_n|$  da).

*Ebazpena.* Kasu honetan,  $f(x) = \cos x - x = 0$  ekuazioa dugu, eta  $[0, \pi/2]$  tartean  $p$  erro bakar bat dauka, zeren eta  $F(x) = \cos x$  funtzioak 4.2. teoremaren baldintzak betetzen baititu ( $|F'(x)| = |-\sin(x)| < 1$  da  $(0, \pi/2)$  tartean). Gainera, 4.3. teoremaren (i) ataleko baldintzak betetzen direnez,  $p$  puntu finko hori erakargarria da. Ondorioz, metodo hori konbergentea da. Taula honetan  $|p - p_n|$  errorea  $|p_{n+1} - p_n| = |F(p_n) - p_n| = |f(p_n)|$  balioaz hurbiltzen da.

$n$	$p_n$	$F(p_n)$	$ f(p_n) $
0	0.785398	0.707107	0.078291
1	0.707107	0.760244	0.053137
2	0.760244	0.724668	0.035576
3	0.724668	0.748720	0.024052
4	0.748720	0.732561	0.016159
5	0.732561	0.743464	0.010903
6	0.743464	0.736128	0.007336
7	0.736128	0.741074	0.004945
8	0.741074	0.737744	0.003330
9	0.737744	0.739988	0.002244

**4.2. taula.** Puntu finkoaren iterazioak.

Hona hemen puntu finkoaren iterazioa:

$$p_n = F(p_{n-1}) = \cos(p_{n-1}), \quad n \geq 1.$$

Beraz, hamar iterazio egin ondoren, ez dugu lortu nahi genuen soluzio hurbildua. Iterazio gehiago eginez gero, eskatutako soluzio hurbildura helduko da. Zenbat iteraziotan gehienez?  $\square$

### 4.3.2. Konbergentziaren ordena

Segida baten konbergentzia neurtzeko erabiltzen den kontzeptu bat da. Askotan, segida hori algoritmo batek sortutakoa da. Demagun  $\{x_n\} \rightarrow x^*$  segida konbergentea. Orduan, hau betetzen da:

$$\lim_{k \rightarrow \infty} |x_k - x^*| = 0.$$

**4.1. Definizioa.** *Konbergentzia **Q-lineala** dela esango dugu  $r \in [0, 1)$  konstante bat existitzen bada, non hau betetzen baita:*

$$\frac{|x_{k+1} - x^*|}{|x_k - x^*|} \leq r, \quad k \text{ nahiko handi baterako.}$$

Horrek esan nahi du iterazio bakoitzean  $x^*$ -ra arteko distantzia gutxienez txikitzen dela faktore konstante batean.

Baldintza hori betetzen da hau badugu:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} < 1.$$

Adibidez,  $\{1 + (0.5)^k\} \rightarrow 1$  segidak konbergentzia Q-lineala du. Horren frogapena ariketa gisa geratzen da.

**4.2. Definizioa.** *Konbergentzia **Q-superlineala** dela esango dugu  $\{r_k\} \rightarrow 0$  segida bat existitzen bada, non hau betetzen baita:*

$$\frac{|x_{k+1} - x^*|}{|x_k - x^*|} \leq r_k.$$

Hots,  $\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = 0$  denean.

Adibidez,  $\{1 + k^{-k}\} \rightarrow 1$  segidak konbergentzia Q-superlineala du. Horren frogapena ariketa gisa geratzen da.

**4.3. Definizioa.** *Konbergentzia **Q-koadratiko**a dela esango dugu  $r \geq 0$  konstante bat existitzen bada, non hau betetzen baita:*

$$\frac{|x_{k+1} - x^*|}{|x_k - x^*|^2} \leq r, \quad k \text{ nahiko handi baterako.}$$

Kasu honetan,  $r$  konstantea ez da derrigorrez 1 baino txikiagoa izan behar.

Baldintza hori betetzen da hau badugu:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^2} < \infty,$$

hau da, limite hori finitua bada.

Adibidez,  $\{1 + (0.5)^{2^k}\} \rightarrow 1$  segidak konbergentzia Q-koadratikoa du. Horren frogapena ariketa gisa geratzen da.

**4.4. Definizioa.** Oro har, konbergentzia **Q-konbergentzia ordena  $p$  gutxienez dela esango dugu**  $r \geq 0$  konstante bat existitzen bada, non hau betetzen baita:

$$\frac{|x_{k+1} - x^*|}{|x_k - x^*|^p} \leq r, \quad k \text{ nahiko handi baterako.}$$

Konbergentzia koadratikoa bada, superlineala da (frogapena ariketa gisa geratzen da).

Oro har, «Q» ez da idazten (Q «quotient» hitzetik dator). Puntu finkoaren metodoaren konbergentzia (konbergente denean) lineala da (ikus (4.7) adierazpena).

### 4.3.3. Newton-Raphson-en metodoa

Metodo hau (Newtonen metodoa ere esaten zaio) oso azkarra da, baina badu arazo bat: hasierako puntuak errotik nahiko hurbil egon behar du.

Oinarri analitikoa Taylorren serie-garapenean datza. Izan bedi  $f \in C^2[a, b]$ , non  $p$  erroa baita eta  $p_0 \in [a, b]$  haren hurbilketa on bat. Orduan,  $|p_0 - p|$  txikia eta  $f(p_0) \neq 0$  izango dira. Funtzio horren  $p_0$ -ren ingurunekeo Taylorren garapena hau da:

$$f(x) = f(p_0) + (x - p_0)f'(p_0) + \frac{1}{2}(x - p_0)^2 f''(c),$$

non  $c$  balioa  $p_0$  eta  $x$ -ren artean baitago. Beraz,  $f(p) = 0$  bete behar denez, hau lortuko dugu:

$$f(p_0) + (p - p_0)f'(p_0) + \frac{1}{2}(p - p_0)^2 f''(c) = 0 \Rightarrow f(p_0) + (p - p_0)f'(p_0) \approx 0.$$

Baina, azken emaitza lortzeko,  $(p - p_0)^2 f''(c)/2$  batugaia ezabatu dugu, baztergarria dela suposatzen baita; izan ere,  $|p - p_0|$  tartea nahiko txikia dela jo dugu. Arrazoi horretan datza, hain zuzen ere, metodoaren baldintza.

Beraz,  $f(p_0) + (p - p_0)f'(p_0) \approx 0$  dela joz eta  $p$  askatuz,  $p \approx p_0 - \frac{f(p_0)}{f'(p_0)} = p_1$  hurbilketa lortzen dugu. Iterazio hori errepikatuz, metodoaren iterazio-formula hau lortuko dugu:

$$p_n = p_{n-1} - \frac{f(p_{n-1})}{f'(p_{n-1})}. \quad (4.10)$$

Argi dago  $\{p_n\}$  segida ondo definituta izateko,  $f'(p_n) \neq 0$  bete behar dela  $n$  azpiindize guztietarako.

Newtonen iterazioa puntu finkoaren iterazioa da, *iterazio-funtzio* honetarako:

$$F(x) = x - \frac{f(x)}{f'(x)}.$$

Kontuan izan  $p$  errorako  $p = F(p)$  dela. Kasu honetan,  $f'(p) \neq 0$  hartuz, zera dugu:

$$F'(x) = 1 - \frac{(f'(x))^2 - f(x)f''(x)}{(f'(x))^2} = \frac{f(x)f''(x)}{(f'(x))^2}.$$

Hipotesiaren arabera,  $f(p) = 0$  denez,  $F'(p) = 0$  dugu. Beraz,  $[a, b]$  tartean  $F'(x)$  jarraitua denez eta  $F'(p) = 0$ , badago  $p$ -ren  $(p - \delta, p + \delta)$  ingurune bat non  $|F'(x)| < 1$  den. Ondorioz, 4.3. teoremako (i) atalaren arabera, iterazio hori konbergentea da  $p_0$  ingurune horretako puntu bat bada. Baina, zein da haren konbergentzia-ordena? Jarraian erantzungo dugu.

**4.4. teorema. (Newtonen teorema.)** *Izan bitez  $f \in C^2[a, b]$  eta  $f(x) = 0$  ekuazioko  $p \in [a, b]$  erroa. Demagun  $\rho > 0$  baterako  $|f'(x)| > \rho$  betetzen dela  $x \in [a, b]$  guztietarako. Orduan,  $\eta > 0$  bat existitzen da, non:  $|p_0 - p| < \eta$  betetzen bada, (4.10) adierazpenak  $p_0$ -tik sortutako  $\{p_n\}$  segidak  $p$  puntura jotzen baitu konbergentzia koadratikoarekin.*

*Frogapena.*  $n = 1$ -erako hau dugu:

$$\begin{aligned} p_1 - p &= \left( p_0 - \frac{f(p_0)}{f'(p_0)} \right) - p = p_0 - p - \frac{f(p_0) - f(p)}{f'(p_0)} \\ &= \frac{1}{f'(p_0)} [f(p) - f(p_0) - f'(p_0)(p - p_0)]. \end{aligned}$$

Makoen arteko gaia Taylorren garapenaren  $O(|p - p_0|^2)$  hondarra da (ikus 3.4.3. azpiatala); hots, badago  $M > 0$  zenbaki bat hau betetzen duena:

$$|f(p) - f(p_0) - f'(p_0)(p - p_0)| \leq M|p - p_0|^2.$$

$M$  hori nahiko handia hartuko dugu  $|f''(x)/2| < M$  bete dadin  $[a, b]$  tartean ere ( $f''(x)$  jarraitua denez, bornatua da  $[a, b]$  tartean). Ondorioz,

$$|p_1 - p| \leq \frac{M}{|f'(p_0)|} |p_0 - p|^2 \leq \frac{M}{\rho} |p_0 - p|^2$$



$n = 2$  denean, hau dugu:

$$|p_2 - p| \leq \frac{M}{|f'(p_1)|} |p_1 - p|^2 \leq \frac{M}{\rho} |p_1 - p|^2.$$

Horrela,  $r = M/\rho$  bada,  $n = 1, 2, \dots$  guztietarako emaitza honetara heltzen gara:

$$|p_n - p| \leq r |p_{n-1} - p|^2.$$

Hori dela eta, ondorio hau lortuko dugu:

$$\frac{|p_n - p|}{|p_{n-1} - p|^2} \leq r. \quad (4.11)$$

Beraz,  $n = 1, 2, \dots$  zera dugu.

$$r |p_{n-1} - p| < 1 \Rightarrow |p_n - p| < |p_{n-1} - p|. \quad (4.12)$$

Izan bedi  $\eta \in \mathbb{R}$ , non  $\eta \leq 1/r$ , orduan nahikoa da  $p_0 \in (p - \eta, p + \eta)$  hartzea  $r |p_0 - p| < 1$  betetzeko. Aurrekoak (4.12)-ren bitartez  $r |p_1 - p| < 1$  inplikatzeko du, eta (4.12) berriro erabiliz,  $|p_2 - p| < |p_1 - p|$  dugu, eta horrela elkarren segidan. Hots:

$$|p_n - p| < |p_{n-1} - p| < \dots < |p_1 - p| < |p_0 - p|, \quad \forall n \in \mathbb{N}, n > 0.$$

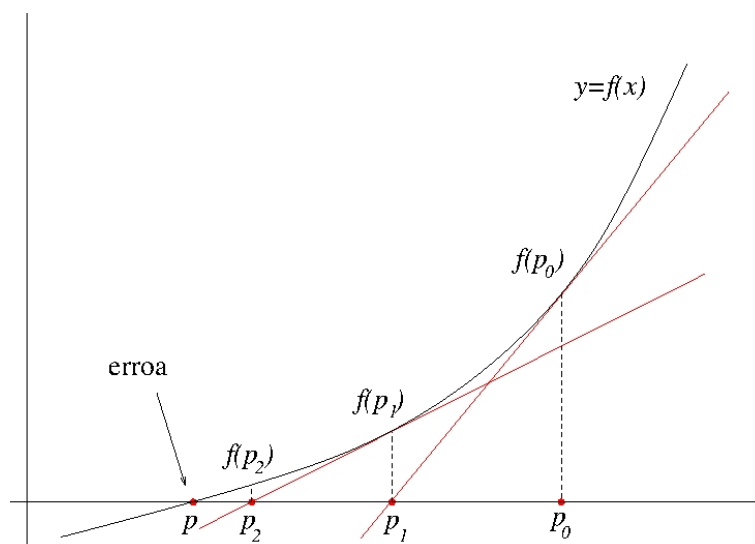
Gainera,  $K = r |p_0 - p|$  definituz  $r |p_n - p| < K$  da,  $n > 0$  azpiindize guztietarako, eta zera dugu:

$$|p_n - p| < K |p_{n-1} - p| < K^2 |p_{n-2} - p| < \dots < K^n |p_0 - p|, \quad K < 1.$$

Ondorioz,  $p_0$  hasierako puntuak  $|p_0 - p| < \eta$  betetzen badu, zera dugu:

$$\lim_{n \rightarrow \infty} |p_n - p| = 0$$

eta  $\{p_n\} \rightarrow p$  konbergentzia koadratikoarekin, ikus (4.11).  $\square$



4.6. irudia. Newtonen metodoa.

Geometrikoki, metodoak hau egiten du:  $OX$  ardatza ebakitzen duen  $y = f(x)$  kurbaren orde,  $(p_0, f(p_0))$  puntutik pasatzen den kurbaren zuzen ukitzaila erabiltzen du. Gero, ebakidura  $OX$ -rekin ( $y = 0$ -rekin) kalkulatzen da, eta, hala,  $p_1$  puntua lortzen da. Teknika hori errepikatuz,  $(p_n, f(p_n))$  puntutik pasatzen den  $f(x)$ -ren ukitzailaren ekuazioa lortzen da, hots,  $y - f(p_n) = f'(p_n)(x - p_n)$ . Orain,  $OX$  ardatzarekin ( $y = 0$ -rekin) ebakiz,  $x$  aska daiteke. Ondorioz,  $x = p_n - \frac{f(p_n)}{f'(p_n)}$  lortzen da, eta  $p_{n+1} = x$ . Begira ezazu 4.6. irudia.

#### 4.1. algoritmoa. Newtonen algoritmoa.

**0 urratsa.** SARRERA. Sartu:  $f(x)$  funtzioa,  $f'(x)$  funtzioa,  $p_0$  hasierako hurbilpena,  $\varepsilon_{max}$  errorearen tolerantzia,  $ze_{max}$  eskatutako zehaztasun erlatiboa eta  $i_{max}$  iterazioen kopuru maximoa.

**1 urratsa.** Jarri  $i = 0$  eta definitu:

$$\begin{aligned}fp_0 &= f(p_0); \\dfp_0 &= f'(p_0); \\ \varepsilon &= |fp_0|; \\ ze &= ze_{max} + 1;\end{aligned}$$

**2 urratsa.**  $i \leq i_{max}$  eta  $\varepsilon > \varepsilon_{max}$  eta  $ze > ze_{max}$  diren bitartean, egin hau:

(a) Kalkulatu  $p$  honela:

$$p = p_0 - fp_0/dfp_0;$$

(b) Egin hau:

$$\begin{aligned}p_{zaharra} &= p_0; \\ p_0 &= p; \\ fp_0 &= f(p_0); \\ dfp_0 &= f'(p_0);\end{aligned}$$

(c) Egin  $i = i + 1$ ;

(d) Kalkulatu zehaztasun erlatiboa:  $ze = |p_0 - p_{zaharra}|/|p_0|$ ;

(e) Kalkulatu errorea:  $\varepsilon = |fp_0|$ ;

**3 urratsa.** IRTEERA. Emaizak:  $p$  erroa,  $i$  erabilitako iterazioen kopurua,  $ze$  zehaztasun erlatiboa eta  $\varepsilon$  errorea.

**4.7. adibidea.** Newtonen metodoa erabiliz, kalkulatu  $x = \cos x$  ekuazioaren soluzio hurbil-du bat, errore-tolerantzia 0.00001 hartuz.

*Ebazpena.* Izan bedi  $f(x) = \cos x - x$ . Orain,  $f(x)$ -ren erro bat aurkitu behar dugu. Baina, lehendabizi, soluzio bat bakartu behar dugu tarte batean. Horretarako, alde zurretik ikusitako (ii) propietatea (Bolzanoren teorema) erabiliko dugu. Orduan,  $f(\pi/2) = -\pi/2 < 0$  eta  $f(0) = 1 > 0$  direnez, funtzioak gutxienez erro bat dauka  $(0, \pi/2)$  tartean (grafikoki ikus daiteke erro bakarra dela). Hona hemen Newtonen iterazioa kasu horretan:

$$p_n = p_{n-1} - \frac{\cos(p_{n-1}) - p_{n-1}}{-\sin(p_{n-1}) - 1}, \quad n \geq 1$$

$p_0$  geometriaren arabera hautatuko dugu. Kasu honetan,  $p_0 = \pi/4 = 0.785398$  egokia dela ikus daiteke.

$n$	$p_n$	$ f(p_n) $
0	0.785398	0.078291
1	0.739542	0.000765
2	0.739085	$2.2 \cdot 10^{-7}$

**4.3. taula.** Newtonen iterazioak.

Beraz, bi iterazio egin ondoren,  $x = 0.739085$  dugu ekuazioaren soluzio hurbildua.  $\square$

**4.8. adibidea.** Newtonen metodoa erabiliz eta  $p_0 = -3$  hartuz, kalkulatu  $e^x = \sin x$  ekuazioaren soluzio hurbildu bat errore-tolerantzia  $0.0001$  hartuz.

*Ebazpena.* Hemen  $f(x) = e^x - \sin x$  dugu eta

$$p_n = p_{n-1} - \frac{e^{p_{n-1}} - \sin(p_{n-1})}{e^{p_{n-1}} - \cos(p_{n-1})}, \quad n \geq 1$$

$n$	$p_n$	$ f(p_n) $
0	-3.0000	0.19090
1	-3.1836	0.00056
2	-3.1831	$3.8 \cdot 10^{-5}$

**4.4. taula.** Newtonen iterazioak.

Beraz, bi iterazio egin ondoren,  $x = -3.1831$  dugu ekuazioaren soluzio hurbildua.  $\square$

### Newtonen metodoari buruzko oharrak

- Ikusi dugun bezala, metodo hau konbergente izateko, hasierako puntuak  $p$  soluziotik nahiko hurbil egon behar du; horrelako metodo bati *metodo lokala* deritzogu.
- Oro har, metodoak konbergentzia arazoak izaten ditu  $p$  soluzioaren ingurune batean  $f'(x) \approx 0$  denean. Ariketa gisa, aurkitu  $x^{10} - 1 = 0$  ekuazioaren erro positibo bat  $p_0 = 0.5$  hartuz. Zer gertatzen da?
- Erroa anizkoitza denean, konbergentzia lineala da. Demagun  $f(x) = (x - p)^m G(x)$  funtzioa dugula eta  $G(p) \neq 0$ ; orduan,  $|p_n - p|/|p_{n-1} - p| \approx 1 - \frac{1}{m}$  da  $n$  nahiko handi baterako. Beraz, anizkoitzasuna handitzen den heinean, konbergentzia motelago bihurtzen da. Horrelako kasuetan, oso egokia da  $f(x)$  funtzioaren ordez  $g(x) = f(x)/f'(x)$  funtzioa hartzea; horrek *Newtonen metodo orokortua* ondorioztatzen du. Adibidez,  $f(x) = x^3 - 3x + 2$  funtzioak erro bikoitza du  $x = 1$  puntuan. Aztertu gertatzen dena Newtonen metodoa erabiltzen badugu erroa kalkulatzeko  $p_0 = 1.2$  hartuz. Geroago, erabili Newtonen metodo orokortua.
- Hasierako balioa,  $p_0$ , errotik nahiko urrun badago, gerta daiteke metodoa konbergente ez izatea. Adibidez,  $f(x) = xe^{-x}$  kasuan, bistan da erroa  $p = 0$  dela. Aztertu gertatzen dena  $p_0 = 2$  hartzen badugu.
- Beste gertaera berezi bat *periodikotasuna* da; hori gertatzen da  $\{p_n\}$  segidaren gaiak periodikoki errepikatzen (edo ia errepikatzen) hasten direnean. Adibidez, aztertu  $f(x) = x^3 - x - 3$  eta  $p_0 = 0$  hartzen dugunean.
- $|F'(x)| > 1$  betetzen bada  $p$  erroa daukan tartean, orduan, oszilazio dibergente bat gerta daiteke. Adibidez,  $f(x) = \arctan(x)$  funtzioa  $p_0 = 1.45$  hartuz. Aztertu kasu hori.

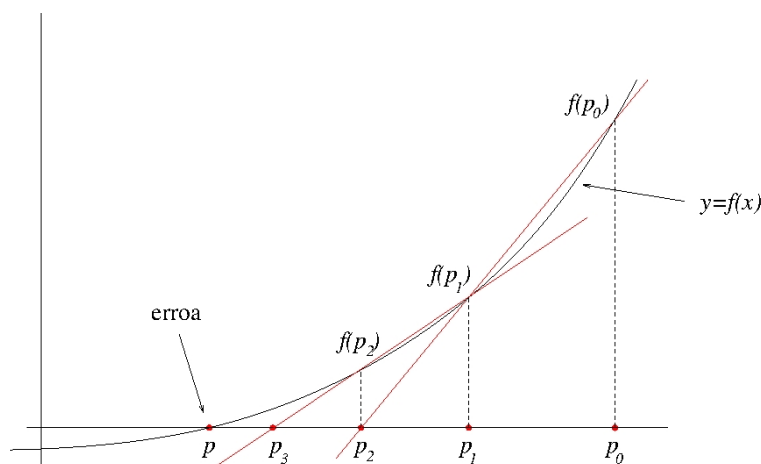
#### 4.3.4. Ebakitzaileren metodoa

Batzuetan, ezin da  $f'(x)$  deribatua kalkulatu iterazioaren formulatan ordezkatzeko; orduan, zenbakizko metodoak erabiliz malda zehaztu dezakegu. Hots, Newtonen metodoaren ordez *ebakitzaileren metodoa* erabiliko dugu. Kasu horietan, deribatua hurbil dezakegu adierazpen honen bitartez:

$$f'(p_k) \approx \frac{f(p_k) - f(p_{k-1})}{p_k - p_{k-1}},$$

orduan, Newtonen ekuazioan ordezkatuz  $k = n - 1$ -erako,

$$p_n = p_{n-1} - \frac{f(p_{n-1})(p_{n-1} - p_{n-2})}{f(p_{n-1}) - f(p_{n-2})}. \quad (4.13)$$



#### 4.7. irudia. Ebakitzaileren metodoa.

Adierazpen horretan ikusten den bezala, hasierako bi puntu behar ditugu. Hau da,  $p_0$  eta  $p_1$  puntuek errotik nahiko hurbil egon behar dute (Newtonen metodoan bezala);  $p_2$  izango da lortuko duguna 1. iterazioan.

Newtonen metodoan, beste aukera bat da deribatua hurbiltzea oraingo iterazioa zertxobait aldatuz, alegia:

$$f'(p_k) \approx \frac{f(p_k + \delta) - f(p_k)}{\delta},$$

non  $\delta > 0$  balio txiki bat baita, eta, orduan, *ebakitzaileren metodo aldatua* izango dugu:

$$p_n = p_{n-1} - \frac{\delta \cdot f(p_{n-1})}{f(p_{n-1} + \delta) - f(p_{n-1})}. \quad (4.14)$$

#### Ebakitzaileren metodoari buruzko oharrak

- Metodo hau Newtonen metodoaren deribatu gabeko hurbilpen bat denez (batez ere (4.14) iterazioa erabiltzen badugu), konbergentziaren ordena  $\alpha$  da, non:

$$\alpha = \frac{1 + \sqrt{5}}{2} \approx 1.618$$

*urrezko ratioa* baita. Bereziki, konbergentzia superlineala da, baina ez da heltzen koardatikoa izatera. Konbergentziaren ordena 1.618 da baldin  $f \in C^2$  bada,  $p$  erro bakun (hots, anizkoiztasuna=1) bada baten ingurune batean, eta  $p_0, p_1$  hasierako puntuak  $p$ -tik nahiko hurbil daudenean.

- Baina,  $\bar{p} \in [p_0, p_1]$  puntu baterako  $f'(\bar{p}) = 0$  bada, algoritmoa baliteke konbergentzia ez izatea.

## 4.2. algoritmoa. Ebakitzaileren algoritmoa.

**0 urratsa.** SARRERA. Sartu:  $f(x)$  funtzioa,  $p_0$  eta  $p_1$  hasierako hurbilpenak,  $\varepsilon_{max}$  erro-  
rearen tolerantzia,  $ze_{max}$  eskatutako zehaztasun erlatiboa, eta  $i_{max}$  iterazioen kopuru  
maximoa.

**1 urratsa.** Jarri  $i = 1$  eta definitu:

$$\begin{aligned}q_0 &= f(p_0); \\q_1 &= f(p_1); \\ \varepsilon &= |q_1|; \\ ze &= |p_1 - p_0|/|p_1|;\end{aligned}$$

**2 urratsa.**  $i \leq i_{max}$  eta  $\varepsilon > \varepsilon_{max}$  eta  $ze > ze_{max}$  diren bitartean, egin hau:

(a) Kalkulatu  $p$  honela:

$$p = p_1 - q_1(p_1 - p_0)/(q_1 - q_0);$$

(b) Egin hau:

$$\begin{aligned}p_0 &= p_1; \\ p_1 &= p; \\ q_0 &= q_1; \\ q_1 &= f(p_1);\end{aligned}$$

(c) Egin  $i = i + 1$ ;

(d) Kalkulatu zehaztasun erlatiboa:  $ze = |p_1 - p_0|/|p_1|$ ;

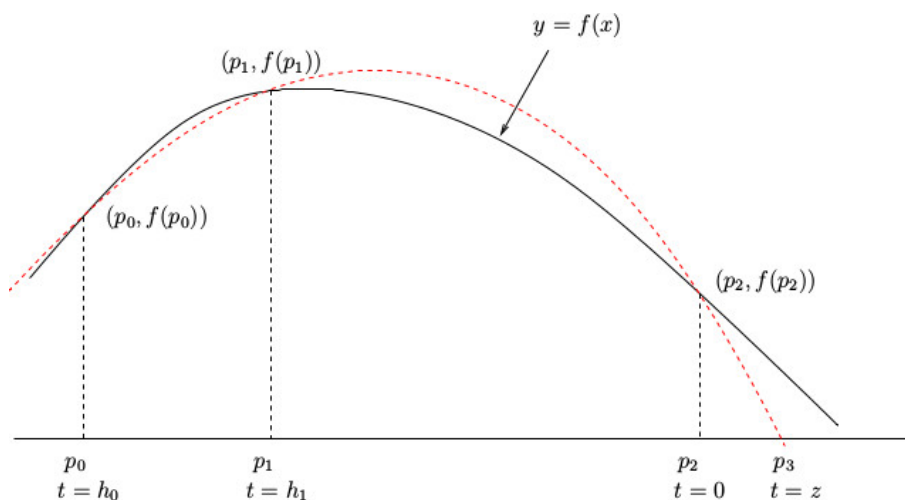
(e) Kalkulatu errorea:  $\varepsilon = |q_1|$ ;

**3 urratsa.** IRTEERA. Emaitzak:  $p$  erroa,  $i$  erabilitako iterazioen kopurua,  $ze$  zehaztasun  
erlatiboa eta  $\varepsilon$  errorea.

### 4.3.5. Muller-en metodoa

Ebakitzaileren metodoak bi puntu erabiltzen ditu puntu berri bat lortzeko. Jarraian azal-  
duko dugun metodoari *Mullerren metodoa* deritzogu.

Demagun iterazioen hasieran  $(p_0, f(p_0))$ ,  $(p_1, f(p_1))$  eta  $(p_2, f(p_2))$  hiru puntuak ditugula.  
Orduan, hiru puntu horietatik pasatzen den parabola bakarra eraikitzen dugu, eta, abszisa-  
ardatzarekin ebakitze-puntua kalkulatu, hurrengo  $p_3$  puntua izango dugu. Jo dezagun  $p_2$   
erroaren hurbilpen hoberena dela, eta demagun hiru puntuetatik igarotzen den 4.8. irudiko  
parabola dela.



4.8. irudia. Mullerren metodoa.

Izan bitez  $t = x - p_2$  aldagai-aldaketa, eta  $h_0 = p_0 - p_2$  eta  $h_1 = p_1 - p_2$  diferentziak. Izan bedi  $y = at^2 + bt + c$  polinomio koadratikoa. Hiru puntu horietatik igarotzen den parabolaren ekuazioko  $a$ ,  $b$  eta  $c$  koefizienteak kalkulatuko ditugu. Puntu bakoitzeko, hau dugu:

$$\begin{aligned} p_0 - \text{rako} \quad t = h_0 \text{ da eta hau bete behar:} \quad & ah_0^2 + bh_0 + c = f(p_0), \\ p_1 - \text{erako} \quad t = h_1 \text{ da eta hau bete behar:} \quad & ah_1^2 + bh_1 + c = f(p_1), \\ p_2 - \text{rako} \quad t = 0 \text{ da eta hau bete behar:} \quad & a0^2 + b0 + c = f(p_2). \end{aligned}$$

Hau da sistema horren soluzioa:

$$c = f(p_2),$$

$$a = \frac{e_0 h_1 - e_1 h_0}{h_1 h_0^2 - h_0 h_1^2},$$

$$b = \frac{e_1 h_0^2 - e_0 h_1^2}{h_1 h_0^2 - h_0 h_1^2},$$

non  $e_0 = f(p_0) - c$  eta  $e_1 = f(p_1) - c$ . Orain, parabola horren  $t = z_1, z_2$  erroak formula hau erabiliz lortzen dira:

$$z = \frac{-2c}{b \pm \sqrt{b^2 - 4ac}}.$$

Dakigunez, formula hori da bigarren mailako ekuazioen erroak aurkitzeko dugun ohiko formularen baliokidea; baina kasu honetan hobe da hura erabiltzea, jada ezagutzen baitugu  $c = f(p_2) \approx 0$  dela eta, ondorioz,  $b \approx \sqrt{b^2 - 4ac}$ .

Metodoaren egonkortasuna ziurtatzeko, balio absolutu txikiena duen erroa aukeratu dugu. Beraz,  $b > 0$  bada, zeinu positiboa aukeratu dugu erro karraturako. Aldiz,  $b < 0$  bada, zeinu negatiboa aukeratu dugu erro karraturako. Puntu berria hau izango da:

$$p_3 = p_2 + z.$$

Orduan,  $p_0$  eta  $p_1$  puntu berriak  $\{p_0, p_1, p_2\}$  puntu zaharren artean aukeratuko ditugu, eta  $p_3$ -tik hurbilenak izango dira;  $p_2$  berria  $p_3$  izango da.

Prozesu hau errepikatuz, iterazio-segida bat lortzen da, eta horrek, baldintza egokietan,  $f(x)$  funtzioaren errora joko du. Mullerren metodoak arazo larriak izan ditzake, ez badira desberdinak  $f(p_1)$ ,  $f(p_2)$  eta  $f(p_3)$ . Hala ere, errotik nahiko hurbil badago, oso azkarra da; haren konbergentzia-ordena 1.84 da. Ebakitzailearena 1.62 da eta Newtonen metodoarena 2 (erro bakunen kasurako).

Gerta daiteke sortutako parabolak  $x$  ardatza ez ukitzea; orduan, erro konplexuak izango ditugu, non alde irudikaria txikia baita alde errearekin konparatuta. Kasu horretan, alde irudikaria kentzen da eta prozesuarekin jarraitzen da.

#### 4.3.6. Alderantzizko interpolazio koadratikoa

Metodo honetan,  $x$ -rekiko koadratikoa izan beharrea,  $y$ -rekiko izango da. Hots,  $(a, f(a))$ ,  $(b, f(b))$  eta  $(c, f(c))$  hiru puntuak ezagutuz, alderantzizko funtzioarekin arituko gara; lortuko dugun polinomio koadratikoa  $P(y)$  motakoa izango da, eta interpolazio-baldintza hauek erabiliz zehaztuko dugu:

$$a = P(f(a)), \quad b = P(f(b)), \quad c = P(f(c)).$$

Parabola honek beti ebakitzen du  $x$  ardatza; hots,  $y = 0$ . Beraz,  $x = P(0)$  da hurrengo iterazioa. Metodo honi IQI (Inverse Quadratic Interpolation) izen laburtua emango diogu. MATLABen kodean; hau izan liteke:

```
function x=iqi(f,a,b,c)
    while abs(c-b) > eps*abs(b);
        x = polyinterp([f(a),f(b),f(c)], [a,b,c],0)
        a = b;
        b = c;
        c = x;
    end
end
```

Metodo horren arazoa zera da: interpolazio-polinomioak  $f(a)$ ,  $f(b)$  eta  $f(c)$  desberdinak izatea eskatzen du. Baina, guk ez dugu berme hori. Adibidez, kalkulatu nahi badugu  $\sqrt{2}$  balioa  $f(x) = x^2 - 2$  funtzioa erabiliz, eta  $a = -2$ ,  $b = 0$  eta  $c = 2$ , orduan,  $f(a) = f(c)$  eta lehenengo urratsa ez dago definituta. Aldiz,  $a = -2.001$ ,  $b = 0$  eta  $c = 1.999$  balioekin hasten bagara, hurrengo iterazioa  $x = 500$  izango da.



### 4.3.7. Zeroin algoritmoa

Metodo honetan, bisekzioaren konbergentziarekiko fidagarritasuna, ebakitzaille-metodoaren konbergentzia-azkartasuna eta IQI metodoa konbinatzen dira. Dekker-ek eta Brent-ek garatu zuten algoritmo honen lehenengo bertsioa (ikus [5, 8]); geroko bertsioek asko hobetu dute algoritmo hori. Hona hemen algoritmoaren laburpen bat:

#### 4.3. algoritmoa. Zeroin algoritmoa.

$a$  eta  $b$  emanda, puntu horietarako  $f(a) \cdot f(b) < 0$  betetzen da; beraz,  $f$  jarraitua bada  $[a, b]$ -n, erroa tartearen barnean dago.

- Erabili ebakitzaillearen iterazioa,  $a$  eta  $b$ -ren arteko  $c$  lortzeko.
- Gero, errepikatu honako urrats hauek,  $|b - a| < \varepsilon|b|$  edo  $f(b) = 0$  bete arte:

**1 urratsa.** Ordenatu  $a$ ,  $b$  eta  $c$  honela:

- ▷  $f(a)$ -k eta  $f(b)$ -k zeinu desberdinak dituzte.
- ▷  $|f(b)| \leq |f(a)|$ .
- ▷  $c$  da  $b$ -ren aurreko balioa; beraz,  $c = a$  izan liteke.

**2 urratsa.**  $c \neq a$  bada, IQI iterazioa kalkulatu dugu.

**3 urratsa.**  $c = a$ , ebakitzaillearen iterazioa kalkulatu dugu.

**4 urratsa.** IQIren edo ebakitzaillearen emaitza  $[a, b]$ -n badago, hartu egingo dugu; ez bada horrela, bisekzioaren metodoa erabiliko dugu hurrengo iterazioa lortzeko.

Algoritmo honen inplementazioa MATLABeko `fzero` funtzioa da, eta eraginkortasun handia du ekuazio ez-linealen erroen kalkuluan. Kontuan izan ez duela deribaturik erabiltzen. Adibidez,  $f(x) = x - 0.5 \sin(x) - 0.7$  funtzioaren erroaren kalkulurako,  $x_0 = 1$  hartuz, `fzero('x-0.5*sin(x)-0.7', 1)` instrukzioak 1.1580 ematen du eta  $f(1.1580) = -1.1102 \times 10^{-16}$ .

Honela ere idatz daiteke: `fzero('x-0.5*sin(x)-0.7', [1,2])`,  $x_0$ -ren ordeztan  $[1,2]$  tartea jarritz, baina bertan erro bat izan behar du, bestela errore-mezu bat emango du.

Badago `roots` funtzioa ere; `r=roots([1,-6,5])` idazten badugu,  $p(x) = x^2 - 6x + 5 = 0$  ekuazioaren erroak emango dizkigu, honela:  $\mathbf{r} = (1, 5)^T$ ; hots,  $\mathbf{r}$  bektoreak  $p(x)$  polinomioaren erroak gordeko ditu.

## 4.4. Algoritmoak gelditzeko irizpideak

Hauek izaten dira algoritmo horiek gelditzeko irizpideak:

1. Errorrea,  $\varepsilon$ , aldez aurretik finkatutako  $\varepsilon_{max} > 0$  errorearen tolerantzia baino txikiagoa izatea, alegia:

$$\varepsilon = |f(p_n)| < \varepsilon_{max}.$$

2. Ondoz ondoko bi iterazioen arteko diferentzia (zehaztasun absolutua) nahiko txikia izatea; hots, zehaztasun absolutua aldez aurretik finkatutako  $za > 0$  tolerantzia baino txikiagoa izatea, alegia:

$$|p_n - p_{n-1}| < za.$$

3. Bi iterazioen arteko zehaztasun erlatiboa nahiko txikia izatea; hots, diferentzia erlatibo hori aldez aurretik finkatutako  $ze > 0$  zehaztasun erlatiboa baino txikiagoa izatea, alegia:

$$\frac{|p_n - p_{n-1}|}{|p_n|} < ze.$$

Hau da  $p_n \approx 0$ -ren arazoa gainditzeko erabil dezakegun beste zehaztasun mota bat:

$$\frac{|p_n - p_{n-1}|}{1 + |p_n|} < ze.$$

Kontuan izan  $p_n = 0$  bada, zatidura hori zehaztasun absolutua dela.

4. Iterazioen kopuru maximo bat,  $n_{max}$ , finkatzea. Alegia, hau bada:

$$n > n_{max},$$

algoritmoa gelditu egingo da.

## 4.5. Problemak

### Eskuz ebazteko problemak:

1.  $e^x - 3x = 0$  ekuazioak erro bat dauka  $p \approx 0.61906129$  puntuan. Erabili bisekzio-metodoaren sei iterazio, erro hori aurkitzeko,  $[0, 1]$  tartetik hasiz. Zenbat iterazio behar dira erroaren hurbilpena balioztatzeko lau zifra esanguratsurekin (hots, 0.6190 lortzeko)?
2. Bistakoa da  $(x - 0.4)(x - 0.6) = x^2 - x + 0.24 = 0$  ekuazioak  $x = 0.4$  eta  $x = 0.6$  erroak dituela. Kontuan izan  $[0, 1]$  tarteko muturrak ez direla onak bisekzioaren metodoa hasteko; zergatik? Funtzioaren marrazketa erabiliz, aurkitu erroen tarte egokiak bisekzio-metodoak erro bakoitzera jotzeko. Bestalde,  $[0.5, 1]$  tartetik bisekzio-metodoaz bilaketa hasten badugu, zein da  $|p - p_5|$  zehaztasunaren bornea bost iterazio igaro ondoren? Zein da benetako zehaztasuna, bost iterazio igaro ondoren?
3. Erabili bisekzioaren metodoa, ekuazio hauen erro positibo txikiena aurkitzeko:
  - (a)  $e^x - x - 2 = 0$ ,
  - (b)  $x^2 - e^x = 0$ ,
  - (c)  $\sin(x) - 2 \cos(x) = 0$ ,
  - (d)  $x^3 - x^2 - 2x + 1 = 0$ ,
  - (e)  $2e^{-x} - \sin x = 0$ ,
  - (f)  $3x^3 + 4x^2 - 8x - 1 = 0$ .

Kasu bakoitzean, aurkitu tarte egoki bat eta, gero, kalkulatu erroaren hurbilpena % 0.5 zehaztasun erlatiboarekin. Lehen erabilitako tartearekin, kalkula ezazu *regula falsi* metodoaz erroaren hurbilpena % 0.5 zehaztasun erlatiboarekin. Zein metodok du konbergentzia azkarrena?

4. Izan bedi  $f(x) = e^x - x^2 + 3x - 2$  funtzioa.
  - (a) Aurkitu, konbergentzia irizpidea erabiliz, iterazioen kopuru maximoa  $f(x) = 0$  ekuazioa bisekzio-metodoaz ebazteko,  $a = 0$  eta  $b = 1$  hartuz,  $\tau = 10^{-3}$  zehaztasun absolutuarekin.
  - (b) Bisekzio-metodoa erabiliz, aurkitu  $f(x) = 0$  ekuazioaren  $p \in [0, 1]$  erroaren lehenengo hurbilpena, non  $\frac{|p_n - p_{n-1}|}{|p_n|} < 10^{-3}$  edo  $|f(p_n)| < 10^{-3}$  betetzen baita.
  - (c) *Regula falsi* metodoa erabiliz, aurkitu  $f(x) = 0$  ekuazioaren  $p \in [0, 1]$  erroaren lehenengo hurbilpena, non  $\frac{|p_n - p_{n-1}|}{|p_n|} < 10^{-3}$  edo  $|f(p_n)| < 10^{-3}$  betetzen baita.
  - (d) Zein metodok du konbergentzia azkarrena?
5. Izan bedi  $f(x) = x^4 + 2x^2 - x - 3$  funtzioa.

- (a) Aurkitu, konbergentzia irizpidea erabiliz, iterazioen kopuru maximoa  $f(x) = 0$  ekuazioa bisekzio metodoaz ebazteko,  $a = 1$  eta  $b = 2$  hartuz,  $\tau = 10^{-4}$  zehaztasun absolutuarekin.
- (b) Bisekzio metodoa erabiliz, aurkitu  $f(x) = 0$  ekuazioaren  $p \in [1, 2]$  erroaren lehenengo hurbilpena, non  $\frac{|p_n - p_{n-1}|}{|p_n|} < 10^{-2}$  edo  $|f(p_n)| < 5.0 \times 10^{-2}$  betetzen baita.
- (c) *Regula falsi* metodoa erabiliz, aurkitu  $f(x) = 0$  ekuazioaren  $p \in [1, 2]$  erroaren lehenengo hurbilpena, non  $\frac{|p_n - p_{n-1}|}{|p_n|} < 10^{-2}$  edo  $|f(p_n)| < 5.0 \times 10^{-2}$  betetzen baita.
- (d) Zein metodok du konbergentzia azkarrena?
6. Zer gertatuko da bisekzio-metodoa erabiltzen badugu  $f(x) = 1/(x - 2)$  funtzioarekin tarte hauetan: (a)  $[3, 7]$  eta (b)  $[1, 7]$ .
7. Emandako tartetan, aztertu funtzio hauek puntu finko bakar bat duten ala ez:
- (a)  $F(x) = 1 - x^2/4$ ,  $[0, 1]$ .
- (b)  $F(x) = 2^{-x}$ ,  $[0, 1]$ .
- (c)  $F(x) = 1/x$ ,  $[0.5, 5.2]$ .
- Aintzat hartu 4.5. adibidea.
8. Aztertu puntu finkoko iterazioaren konbergentzia  $F(x) = -4 + 4x - \frac{1}{2}x^2$  denean.
- (a) Ebatzi analitikoki  $F(x) = x$ , eta frogatu  $p = 2$  eta  $p = 4$  puntu finkoak direla.
- (b) Hartu  $p_0 = 1.9$  eta kalkula itzazu  $p_1, p_2, p_3$ , eta balio horiei dagozkien zehaztasun absolutuak eta erlatiboak.
- (c) Hartu  $p_0 = 3.8$  eta kalkula itzazu  $p_1, p_2, p_3$ , eta balio horiei dagozkien zehaztasun absolutuak eta erlatiboak.
- (d) 4.3. teorema erabiliz, zer ondoriozta daiteke?
9. Izan bedi  $F(x) = x^2 + x - 4$ . Orduan,  $x = F(x)$  ekuazioaren erroak kalkulatzeko puntu finkoko iterazioa erabil dezakegu? Zergatik?
10. Puntu finkoaren metodoan, zergatik da abantaila  $F'(p) \approx 0$  izatea?
11. Demagun  $F, F' \in C(p - r_0, p + r_0)$ ,  $r_0 > 0$ , non  $F(p) = p$ , eta ez dela  $n$  existitzen  $F(p_n) = p$  betetzen duenik. Froga ezazu hau betetzen bada:

$$|F'(p)| > 1,$$

$p$  puntu finko aldaragarria dela. Hots,  $p_n$  iterazioa  $p$ -tik nahiko hurbil badago,  $|p_{n+1} - p| > |p_n - p|$  dugula.

12. Demagun 4.1. korolararioaren hipotesiak betetzen direla,  $F'(p) = 0$ , eta  $F''(p)$  existitzen dela. Froga ezazu  $r$  positibo bat existitzen dela,  $r < r_0$ , non  $p_0 \in [p - r, p + r]$  betetzen bada,  $p_n = F(p_{n-1})$  segidako gai guztiak  $[p - r, p + r]$  tartean baitaude eta  $\{p_n\} \rightarrow p$  konbergentzia koadratikoarekin. (*Iradokizuna*: erabili  $F$ -ren bigarren mailako Taylorren polinomioa  $p$ -ren ingurunean  $F(p_n)$  hurbiltzeko).
13. Izan bedi  $e^x - 3x^2 = 0$  ekuazioa; aurkitu  $x = -0.5$  eta  $x = 4$  puntuetatik gertuko erroak Newton-Raphson metodoaz, sei digituko zehaztasun erlatiboarekin (hots,  $ze = 10^{-6}$  izanik).
14. Izan bedi  $x^2 - x + 2 = 0$  ekuazioa.
- Zehaztu Newton-Raphsonen  $p_n = F(p_{n-1})$  iterazio-funtzioa.
  - Hasi  $p_0 = -1.5$  puntutik, eta aurkitu  $p_1$ ,  $p_2$  eta  $p_3$ .
  - Zein da  $|f(p_3)|$  errorea  $p_3$  puntuan?
  - Zergatik gertatzen da hori?
15. Izan bedi  $x^2 - x - 3 = 0$  ekuazioa.
- Zehaztu Newton-Raphsonen  $p_n = F(p_{n-1})$  iterazio-funtzioa.
  - Hasi  $p_0 = 1.6$  puntutik, eta aurkitu  $p_1$ ,  $p_2$ ,  $p_3$  eta  $p_4$ .
  - Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan?
16. Izan bedi  $(x - 2)^4 = 0$  ekuazioa.
- Zehaztu Newton-Raphsonen formulako  $p_n = F(p_{n-1})$  iterazio-funtzioa.
  - Hasi  $p_0 = 2.1$  puntutik, eta aurkitu  $p_1$ ,  $p_2$ ,  $p_3$  eta  $p_4$ .
  - Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan?
  - Aztertu  $\{p_n\}$  segidaren konbergentzia-ordena.
17. Izan bedi  $x^3 - 3x - 2 = 0$  ekuazioa.
- Zehaztu Newton-Raphsonen formulako  $p_n = F(p_{n-1})$  iterazio-funtzioa.
  - Hasi  $p_0 = 2.1$  puntutik, eta aurkitu  $p_1$ ,  $p_2$ ,  $p_3$  eta  $p_4$ .
  - Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan?
  - Aztertu  $\{p_n\}$  segidaren konbergentzia-ordena.
18. Izan bedi  $xe^{-x} = 0$  ekuazioa.
- Zehaztu Newton-Raphsonen formulako  $p_n = F(p_{n-1})$  iterazio-funtzioa.
  - Hasi  $p_0 = 0.2$  puntutik, eta aurkitu  $p_1$ ,  $p_2$ ,  $p_3$  eta  $p_4$ . Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan? Zein da  $\lim_{n \rightarrow \infty} p_n$ ? Zein da zehaztasun absolutua  $p_4$  puntuan?
  - Hasi  $p_0 = 20$  puntutik, eta aurkitu  $p_1$ ,  $p_2$ ,  $p_3$  eta  $p_4$ . Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan? Zein da  $\lim_{n \rightarrow \infty} p_n$ ? Zein da zehaztasun absolutua  $p_4$  puntuan?

(d) Azaldu zer gertatzen den (b) eta (c) kasuetan.

19. Ebakitzailaren metodoa erabiliz, aurki itzazu ekuazio hauen  $p_2$  eta  $p_3$  hurbilpenak, emandako hastapen-puntuetatik hasiz:

(a)  $x^2 - 2x - 1 = 0$ , non  $p_0 = 2.6$  eta  $p_1 = 2.5$ .

(b)  $x^2 - x - 3 = 0$ , non  $p_0 = 1.7$  eta  $p_1 = 1.67$ .

(c)  $x^2 - 2x - 1 = 0$ , non  $p_0 = -1.5$  eta  $p_1 = -1.52$ .

20. Erabil dezakegu Newton-Raphsonen metodoa  $x^2 - 14x + 50 = 0$  ebazteko? Zergatik?

21. Erabil dezakegu Newton-Raphsonen metodoa  $x^{1/3} = 0$  ebazteko? Zergatik?

22. Erabil dezakegu Newton-Raphsonen metodoa  $(x - 3)^{1/2} = 0$  ebazteko,  $p_0 = 4$  hartuz? Zergatik?

23. Erabili Mullerren metodoa  $x^3 - x - 2 = 0$  ekuazioaren erro bat hurbiltzeko. Hasi  $p_0 = 1.0$ ,  $p_1 = 1.2$  eta  $p_2 = 1.4$  puntuetatik, eta kalkula itzazu  $p_3$  eta  $p_4$ . Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan?

24. Erabili Mullerren metodoa  $4x^2 - e^x = 0$  ekuazioaren erro bat hurbiltzeko. Hasi  $p_0 = 4.0$ ,  $p_1 = 4.1$  eta  $p_2 = 4.2$  puntuetatik, eta kalkula itzazu  $p_3$  eta  $p_4$ . Zein da  $|f(p_4)|$  errorea  $p_4$  puntuan?

25. Izan bedi  $xe^{0.5x} + 1.2x - 5 = 0$  ekuazioa.  $f(x) = xe^{0.5x} + 1.2x - 5$  funtzioa marrazten badugu,  $x = 1$  eta  $x = 2$  artean erro bat dagoela ikusten dugu. Ekuazio hori  $x = F(x)$  era desberdinetan berridatz dezakegu. Honako hiru aukera hauetatik, zein da zure aburuz egokiena? Eman arrazoia funtzioen deribatuak aztertuz:

(a)  $x = \frac{5 - xe^{0.5x}}{1.2}$ ;

(b)  $x = \frac{5}{e^{0.5x} + 1.2}$ ;

(c)  $x = \frac{5 - 1.2x}{e^{0.5x}}$ .

### MATLABez ebazteko problemak:

26. Metodo grafikoa erabiliz, bakartu tartetean funtzio hauen erroak:

(a)  $f(x) = \sin(x) + 0.8 \cos(x)$ .

(b)  $f(x) = x^2 - 4x + 3.5 - \ln(x)$ .

(c)  $f(x) = (x - 2.1)^2 - 7x \cos(x)$ .

Egiaztatu analitikoki tarte bakoitzean erro bat dagoela (hots, Bolzanoren teorema erabiliz).

27. Inplementatu bisekzio-metodoa honelako MATLABeko funtzio bat garatuz:

(a) Lehenengo lerroa honela izan behar:

```
function [p,z,e,i]=bisekzioa(f,a,b,ztol,etol,imax)
```

(b) Sarrerak: **f**= funtzioa, **a,b**= funtzio horrek erro bat duen tarte bat, **ztol**= zehaztasunaren tolerantzia, **etol**= errorearen tolerantzia, eta **imax**= iterazio kopuru maximoa.

(c) Irteerak: **p**= erroaren azken hurbilpena tolerantzia horiekin, **i**= erabilitako iterazio kopurua, **z**= zehaztasun absolutua, eta **e**= erroa azken iterazioan.

(d) MATLABeko funtzio horrek, datu horiek emateaz gain, ondo eraikitako (formatu egokiak erabiliz) taula bat sortu behar du zutabe desberdinak erabiliz era honetan:

$i$	$a_i$	$b_i$	$p_i$	$z_i$	$e_i$
-----	-------	-------	-------	-------	-------

$i$  iterazio bakoitzerako  $a_i$  eta  $b_i$  zutabe horietan agertu behar dira erroa daukan tartearen muturrak,  $p_i = (a_i + b_i)/2$  hurbilpen berria,  $z_i = |p_i - p_{i-1}|$  zehaztasun absolutua eta  $e_i = |f(p_i)|$  erroa.

(e) Egiaztatu kode hori ondo dabilela, zuk aukeratuz ekuazio egoki bat non soluzio zehatza aurrez ezagutzen duzun eta bisekzio-metodoaren urrats batzuk eginda daukazun. Aukera itzazu tarte egoki bat, iterazio kopuru maximoa, eta zehaztasunaren eta errorearen tolerantziak.

28. Bisekzio-metodoaren kodea erabiliz, aurkitu 26. problemako funtzioen erroak.

(i) Kalkulatu erroa  $10^{-4}$  errore tolerantziarekin,  $10^{-4}$  zehaztasun erlatiboarekin eta gehienez 100 iterazio erabiliz.

(ii) Kalkula ezazu, arkatzez, kasu bakoitzean zenbat iterazio beharko dituen gehienez metodoak, zehaztasun absolutua  $10^{-6}$  izateko. Orain, MATLABen bidez, kalkulatu metodoak erabiltzen duen benetako iterazio kopurua zehaztasun hori lortzeko. Horretarako, hartu **etol** oso txikia (esate baterako,  $10^{-10}$ ).

29. Zure kodean zati batzuk aldatuz, lortu **regula\_falsi.m** fitxategi-kode bat; metodo honetan erroaren hurbilpen berria ez da lortzen  $p_i = (a_i + b_i)/2$  eginez, honela baizik:

$$p_i = b_i - f(b_i) \frac{b_i - a_i}{f(b_i) - f(a_i)}.$$

(i) Aurkitu 26. problemako funtzioen erroak  $M$  fitxategi berri hau erabiliz eta  $10^{-4}$  errore tolerantziarekin,  $10^{-4}$  zehaztasun erlatiboarekin eta gehienez 100 iterazio erabiliz.

(ii) Zein du konbergentzia azkarrena? Bisekzioarena edo regula falsirena?

(iii) Orain, MATLABen bidez, kalkulatu metodoak erabiltzen duen benetako iterazio kopurua zehaztasun absolutua  $10^{-6}$  izateko. Horretarako, hartu **etol** oso txikia (esate baterako,  $10^{-10}$ ). Konparatu 28. problemako (ii) atalean lortutako emaitzarekin. Ondorioren bat atera dezakezu?

30. `newton.m` fitxategia lortzeko, inplementa ezazu Newtonen metodoa (4.1. algoritmoa) honelako MATLABeko funtzio bat garatuz:

(a) Lehenengo lerroa honela izan behar:

```
function [p,ze,e,i]=newton(f,df,p0,zetol,etol,imax)
```

(b) Sarrerak: `f`= funtzioa, `df`= funtzioaren deribatua, `p0`= hasierako puntua, `zetol`= zehaztasun erlatiboaren tolerantzia, `etol`= errorearen tolerantzia eta `imax`= iterazio kopuru maximoa.

(c) Irteerak: `p`= erroaren hurbilpen hoberena tolerantzia horiekin, `i`= iterazio kopurua, `ze`= zehaztasun erlatiboa eta `e`= errorea.

(d) MATLABeko funtzio horrek, datu horiek emateaz gain, ondo eraikitako (formatu egokiak erabiliz) taula bat sortu behar du zutabe desberdinak erabiliz era honetan:

$i$	$p_i$	$ze_i$	$e_i$
-----	-------	--------	-------

$i$  iterazio bakoitzerako,  $p_i$  hurbilpenak,  $ze_i = |p_i - p_{i-1}|/|p_i|$  zehaztasun erlatiboa eta  $e_i = |f(p_i)|$  errorea.

(e) Egiaztatu kode hori ondo dabilela, zuk aukeratuz ekuazio egoki bat non soluzio zehatza aurrez ezagutzen duzun eta Newtonen metodoaren urrats batzuk egin da daukazun. Aukera itzazu zuk hasierako puntu egoki bat, iterazio kopuru maximoa, eta zehaztasunaren eta errorearen tolerantziak.

31. Aurreko ariketa egin ondoren, ebatzi ekuazio hauek hasierako puntuetatik abiatuz eta Newton-en metodoa erabiliz:

(a)  $f(x) = 2^{-x} - x = 0$ , non  $p_0 = 1$ .

(b)  $f(x) = x^2 - 4x + 3.5 - \ln(x) = 0$ , non  $p_0 = 1$ .

(c)  $f(x) = (x - 2.1)^2 - 7x \cos(x) = 0$ , non  $p_0 = 1.5$ .

(i) Kalkulu hauek egiteko, erabili lehenengo ariketako `newton.m` fitxategia.

(ii) Kalkulatu erroa ekuazio bakoitzean  $10^{-6}$  errore-tolerantziarekin,  $10^{-6}$  zehaztasun erlatiboarekin, eta gehienez 100 iterazio erabiliz.

32. Newtonen metodoari dagokion inplementazioa adibide moduan hartuz, inplementa ezazu ebakitzailaren metodoa; ikus ezazu 4.2. algoritmoaren sasikodea. Ebatzi:

(a)  $f(x) = 2^{-x} - x = 0$ ,  $p_0 = 1, p_1 = 1.5$ .

(b)  $f(x) = x^2 - 4x - 3.5 - \log(x) = 0$ ,  $p_0 = 1, p_1 = 1.5$ .

(c)  $f(x) = (x - 2.1)^2 - 7x \cos(x) = 0$ ,  $p_0 = 1.5, p_1 = 2$ .

Konparatu emaitza horiek MATLABeko `fzero.m` funtzioaz lortutakoekin.

33. (a) Garatu sasikode bat puntu finkoaren metodorako,  $x = F(x)$  ekuazioa ebazteko gai dena.



- (b) Inplementatu sasikode hori MATLABeko `fpuntufinko.m` fitxategi berri bat bezala.
- (c) Ebatzi  $x = F(x) = 2^{-x}$  ekuazioa,  $p_0 = 1$  puntutik hasiz. Zer gertatzen da? Zergatik?
34. Garatu 4.3.5. ataleko Mullerren metodoaren sasikode bat, eta, gero, inplementatu MATLABeko `M` fitxategi batean, non sartutako balioak `f,p0,p1,p2,etol,zetol,imax` baitira, eta ateratzen diren balioak `p,err,ze,i` (ikus ezazu 28. ariketako (b) eta (c) ataletan parametro horien esanahia). Kontuan hartu  $\sqrt{b^2 - 4ac}$  zenbaki konplexua izan daitekeela, eta zein erro gorde behar dugun.
35. Erabili Mullerren metodoa  $p_0 = 1.5$ ,  $p_1 = 1.4$  eta  $p_2 = 1.3$  hasierako balioak  $f(x) = 1 + 2x - \tan(x)$  funtzioaren erro baten hurbilpena lortzeko, `etol`= $10^{-8}$  errorearen tolerantziarekin, `zetol`= $10^{-12}$  zehaztasun erlatiboarekin, eta `imax`=50 iterazio kopuru maximoarekin.
36. Konparatu Newtonen, ebakitzaillearen eta Mullerren metodoen konbergentzia aurreko problemaren funtziorako, Newtonen metodorako  $p_0 = 1.5$  hartuz, eta ebakitzaillearen metodorako  $p_0 = 1.5$  eta  $p_1 = 1.4$ . Hartu aurreko problemako balio berdinak `etol`, `zetol` eta `imax` parametroetarako.



## 5. kapitulua

# Sistema linealak: metodo zuzenak eta iteratiboak

Bi (edo hiru) ekuazio eta bi (edo hiru) ezezagun dituen sistema bat ebaztea eskuz egin dezakegu, ordezkapenaz edo beste metodo bat erabiliz (esate baterako, Cramer-en metodoa). Sistema bat horrela ebaztea, praktikan, ezinezko bihurtzen da ekuazioen eta ezezagunen kopurua handiagoa denean.

### 5.1. Sistema linealen ebazpena

Notazio matrizialarekin, honela idatz daiteke ekuazio linealen sistema bat:

$$\mathbf{Ax} = \mathbf{b}.$$

Askotan, ekuazioen eta ezezagunen kopurua berdina eta handia da,  $n$  ordenako  $\mathbf{A}$  matrize karratua ezaguna da, eta  $n$  dimentsioko  $\mathbf{b}$  zutabe-bektorea ere bai, eta  $\mathbf{x}$  ezezagun zutabe-bektorea  $n$  dimentsiokoa da.

Jakin badakigu  $\mathbf{Ax} = \mathbf{b}$ -ren soluzioa  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$  idatz daitekeela, non  $\mathbf{A}^{-1}$  matrizea  $\mathbf{A}$ -ren alderantzizkoa baita. Hala ere, konputazio praktikako problema gehienetan, ez da beharrezkoa, ezta gomendagarri ere,  $\mathbf{A}^{-1}$  kalkulatzeko. Adibide erakusgarri moduan, ekuazio bateko eta ezezagun bateko ekuazio hau hartuko dugu:

$$7x = 21.$$

Sistema hori ebazteko modu hoberena zatiketa da:

$$x = \frac{21}{7} = 3.$$

Alderantzizko matrizea erabiltzeak honetara eramaten gaitu:

$$x = 7^{-1} \times 21 = 0.142857 \times 21 = 2.99997.$$

Alderantzizkoak aritmetika gehiago behar du (zatiketa bat eta biderketa bat, zatiketa bat bakarrik izan beharrean) eta emaitzaren zehaztasuna txikiagoa da. Antzeko zerbait gertatzen da sistema handiagoetan. Ondorioz, ekuazio-sistemen ebazpen zuzenean zentratuko gara, alderantzizkoaren kalkulua egin beharrean.

MATLABeko atzeranzko barra “\” erabiliko dugu  $\mathbf{AX} = \mathbf{B}$  sistema ebazteko, non  $\mathbf{A}$ -ren lerroen kopurua eta  $\mathbf{B}$ -rena berdinak baitira. Orduan, sistema horren soluzioa  $\mathbf{X}=\mathbf{A}\backslash\mathbf{B}$  da eta horrek  $\mathbf{X} = \mathbf{A}^{-1}\mathbf{B}$  ematen du. Alegia, *ezker zatiketa* da.

MATLABeko aurreranzko barra “/” erabiliko dugu  $\mathbf{XA} = \mathbf{B}$  sistema ebazteko, non  $\mathbf{A}$ -ren zutabeen kopurua eta  $\mathbf{B}$ -rena berdinak baitira. Orduan, sistema horren soluzioa  $\mathbf{X}=\mathbf{B}/\mathbf{A}$  da eta horrek  $\mathbf{X} = \mathbf{BA}^{-1}$  ematen du. Alegia, *eskuin zatiketa* da.

Notazio hori erabiltzen da  $\mathbf{A}$  karratua ez denean ere; hots, ekuazioen kopurua eta ezezagunen kopurua desberdinak izan arren. Oraingoz, matrize karratuko sistemekin arituko gara.

## 5.2. 3x3 adibide bat

Sistema lineal baten ebazte-algoritmo bat erakutsiko dugu. Izan bedi  $\mathbf{Ax} = \mathbf{b}$  ekuazio-sistema hau:

$$\begin{bmatrix} 10 & -7 & 0 \\ -3 & 2 & 6 \\ 5 & -1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 4 \\ 6 \end{bmatrix}.$$

Adierazpen hori ekuazio hauen sistemari dagokio:

$$\begin{aligned} E_1 : & \quad 10x_1 - 7x_2 = 7, \\ E_2 : & \quad -3x_1 + 2x_2 + 6x_3 = 4, \\ E_3 : & \quad 5x_1 - x_2 + 5x_3 = 6. \end{aligned}$$

Algoritmoaren lehenengo urratsak lehenengo ekuazioa erabiltzen du beste ekuazioetako  $x_1$  ezabatzeke. Hori lortzeko,  $E_2 + 0.3 \cdot E_1$  eta  $E_3 - 0.5 \cdot E_1$  egiten da.  $E_1$  ekuazioko  $x_1$ -en 10 koefizienteari *pibot* deritzogu, eta *biderkatzaile* deritze beste ekuazioetako  $x_1$ -en koefizienteak 10 pibotaz zatituz lortutako -0.3 eta 0.5 kopuruei. Lehenengo urratsak honela aldatzen ditu ekuazioak:

$$\begin{bmatrix} 10 & -7 & 0 \\ 0 & -0.1 & 6 \\ 0 & 2.5 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 6.1 \\ 2.5 \end{bmatrix}.$$

Ohartu hau betetzen dela:

$$\begin{bmatrix} 1 & 0 & 0 \\ -0.3 & 1 & 0 \\ 0.5 & 0 & 1 \end{bmatrix} \begin{bmatrix} 10 & -7 & 0 \\ 0 & -0.1 & 6 \\ 0 & 2.5 & 5 \end{bmatrix} = \begin{bmatrix} 10 & -7 & 0 \\ -3 & 2 & 6 \\ 5 & -1 & 5 \end{bmatrix}.$$

Bigarren urratsean, bigarren ekuazioa erabil dezakegu hirugarren ekuazioko  $x_2$  ezabatzeko. Baina, bigarren pibota, bigarren ekuazioko  $x_2$ -ren koefizientea,  $-0.1$  da, eta hori hurrengo koefizientea baino txikiagoa da. Ondorioz, azken bi ekuazioak trukutzen dira. Horri *pibotatzea* deritzogu. Adibide honetan, egia esan, hori ez da beharrezkoa, ezen ez baitago biribiltze-errererik; baina, oro har, erabakigarria da. Beraz, pibotatu eta gero, zera dugu:

$$\begin{bmatrix} 10 & -7 & 0 \\ 0 & 2.5 & 5 \\ 0 & -0.1 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 2.5 \\ 6.1 \end{bmatrix}.$$

Orain, bigarren pibota 2.5 da, eta erabil dezakegu hirugarren ekuazioko  $x_2$  ezabatzeko. Hori lortuko dugu hirugarren ekuazioari 0.04 bider bigarrena batuz (hots,  $-0.04$  biderkatzailea da), hau lortuz:

$$\begin{bmatrix} 10 & -7 & 0 \\ 0 & 2.5 & 5 \\ 0 & 0 & 6.2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 2.5 \\ 6.2 \end{bmatrix}.$$

Orain, azken ekuazioa hau da:

$$6.2x_3 = 6.2$$

eta, ondorioz,  $x_3 = 1$ . Balio hori bigarren ekuazioan ordezkatu dezakegu:

$$2.5x_2 + (5) \cdot (1) = 2.5$$

eta  $x_2$  askatuz,  $x_2 = -1$  dugu. Azkenik, lehenengo ekuazioan ordezkatzeko ditugu  $x_2$ -ren eta  $x_3$ -ren balioak:

$$10x_1 + (-7) \cdot (-1) = 7$$

eta  $x_1$  askatuz,  $x_1 = 0$  dugu. Emaitza hau da:

$$\mathbf{x} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}.$$

Emaitza hau erraz egiazta daiteke jatorrizko ekuazioak erabiliz:

$$\begin{bmatrix} 10 & -7 & 0 \\ -3 & 2 & 6 \\ 5 & -1 & 5 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 7 \\ 4 \\ 6 \end{bmatrix}.$$

Ohartu  $\mathbf{L}_1\mathbf{L}_2 = \mathbf{L}$  betetzen dela, non

$$\mathbf{L}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ -0.3 & 0 & 1 \end{bmatrix}, \quad \mathbf{L}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -0.04 & 1 \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ -0.3 & -0.04 & 1 \end{bmatrix}.$$

Bestalde,  $\mathbf{A}$  matrizean bigarren eta hirugarren lerroak trukatu (permutatu) ditugu; hori  $\mathbf{P}$  permutazio-matrize batez aurrebiderkatuz lortzen da, alegia:

$$\mathbf{PA} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 10 & -7 & 0 \\ -3 & 2 & 6 \\ 5 & -1 & 5 \end{bmatrix} = \begin{bmatrix} 10 & -7 & 0 \\ 5 & -1 & 5 \\ -3 & 2 & 6 \end{bmatrix}.$$

Algoritmo osoa era trinko batean adieraz daiteke, notazio matriziala erabiliz. Adibide honetarako hau dugu:

$$\mathbf{LU} = \mathbf{PA},$$

non

$$\mathbf{U} = \begin{bmatrix} 10 & -7 & 0 \\ 0 & 2.5 & 5 \\ 0 & 0 & 6.2 \end{bmatrix}.$$

Hau da,  $\mathbf{L}$  matrizeak ezabatze-prozesuko biderkatzaileak gordetzen ditu bateko diagonal azpian,  $\mathbf{U}$  matrizea azken koefiziente-matrizea da, eta  $\mathbf{P}$  matrizeak pibotatze-prozesua deskribatzen du.

### 5.3. Permutazio- eta triangelu-matrizeak

*Permutazio-matrize* bat identitate-matrize bat da, lerroak eta zutabeak trukatuta dauzkana. Hain zuzen ere, 1 bakarria dauka lerro eta zutabe bakoitzean, eta zero dira beste gaiak. Adibidez,

$$\mathbf{P} = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

$\mathbf{A}$  matrizea ezkerretik biderkatzen badugu  $\mathbf{P}$  permutazio-matrize batez (hots,  $\mathbf{PA}$ ),  $\mathbf{A}$  matrizearen lerroak permutatzen dira. Aldiz, eskuinetik biderkatzen badugu (hots,  $\mathbf{AP}$ ), orduan,  $\mathbf{A}$  matrizearen zutabeak permutatzen ditu.

MATLABek badauka  $\mathbf{p}$  permutazio-bektore bat matrize baten lerroak edo zutabeak berordenatzeko. Goiko  $\mathbf{P}$  matrizearen kasuan,  $\mathbf{p}$  hau da:

$$\mathbf{p}=[4 \ 1 \ 3 \ 2]$$

Orduan,  $\mathbf{P}*\mathbf{A}$  eta  $\mathbf{A}(\mathbf{p}, :)$  berdinak dira. Lortutako matrizean 1. lerroa  $\mathbf{A}$ -ren 4.a izango da, 2. lerroa  $\mathbf{A}$ -ren 1.a, 3. lerroa  $\mathbf{A}$ -ren 3.a, eta 4. lerroa  $\mathbf{A}$ -ren 2.a. Halaber,  $\mathbf{A}*\mathbf{P}$  eta  $\mathbf{A}(:, \mathbf{p})$  berdinak dira, biek  $\mathbf{A}$ -ren zutabeen permutazio berdina sortzen dute.

$\mathbf{Px} = \mathbf{b}$  ekuazio-sistematarako soluzioa kalkulatzeko, hau bakarrik egin behar dugu:

$$\mathbf{x} = \mathbf{P}^T \mathbf{b},$$

zeren  $\mathbf{P}^{-1} = \mathbf{P}^T$  baita; hots,  $\mathbf{P}$  ortogonala da.

Matrize goi-triangeluar batean, zeroak dira diagonal nagusiaren azpiko gai guztiak. Matrize behe-triangeluar batean, zeroak dira diagonal nagusiaren goiko gai guztiak. Adibidez, aurreko atalean,  $\mathbf{L}$  matrizea behe-triangeluarra da eta  $\mathbf{U}$  matrizea goi-triangeluarra.

Sistema bateko matrizea triangeluarra bada, erraz ebazten da. Izan bedi  $\mathbf{U}\mathbf{x} = \mathbf{b}$  sistema goi-triangeluar hau:

$$\begin{aligned} u_{11}x_1 + u_{12}x_2 + \dots + u_{1,n-1}x_{n-1} + u_{1n}x_n &= b_1 \\ u_{22}x_2 + \dots + u_{2,n-1}x_{n-1} + u_{2n}x_n &= b_2 \\ \dots &\dots \\ u_{n-1,n-1}x_{n-1} + u_{n-1,n}x_n &= b_{n-1} \\ u_{nn}x_n &= b_n. \end{aligned} \tag{5.1}$$

Beraz, soluzioa lortzeko, atzeranzko ordezkatzeko-prozesu honi jarraitu behar diogu (hots, behetik gora):

$$\begin{aligned} x_n &= \frac{b_n}{u_{nn}} \\ x_{n-1} &= \frac{b_{n-1} - u_{n-1,n}x_n}{u_{n-1,n-1}}. \end{aligned} \tag{5.2}$$

Eta, honela, ondoz ondo jardunez, zera dugu:

$$x_i = \frac{b_i - u_{in}x_n - u_{i,n-1}x_{n-1} - \dots - u_{i,i+1}x_{i+1}}{u_{ii}} = \frac{b_i - \sum_{j=i+1}^n u_{ij}x_j}{u_{ii}}, \tag{5.3}$$

non  $i = n - 1, n - 2, \dots, 3, 2, 1$ .

Eragiketa horiek, MATLABen bitartez, honela egin ditzakegu:

```
x = zeros(n,1);
for k = n:-1:1
    x(k) = b(k)/U(k,k);
    i = (1:k-1)';
    b(i) = b(i)-x(k)*U(i,k);
end
```

Bestalde, demagun  $\mathbf{L}\mathbf{x} = \mathbf{b}$  sistema behe-triangeluar hau dugula, eta  $\mathbf{L}$ -ren diagonaleko gaiak batak direla:

$$\begin{aligned} x_1 &= b_1 \\ l_{21}x_1 + x_2 &= b_2 \\ \dots &\dots \\ l_{n1}x_1 + l_{n2}x_2 + \dots + l_{n,n-1}x_{n-1} + x_n &= b_n. \end{aligned} \tag{5.4}$$

Soluzioa lortzeko, *aurreranzko ordezkatzeko-prozesu* honi jarraitu behar diogu (hots, goitik behera):

$$\begin{aligned}x_1 &= b_1 \\x_2 &= b_2 - l_{2,1}x_1.\end{aligned}\tag{5.5}$$

Eta, honela, ondoz ondo jardunez, zera dugu:

$$x_i = b_i - l_{i,1}x_1 - l_{i,2}x_2 - \dots - l_{i,i-1}x_{i-1} = b_i - \sum_{j=1}^{i-1} l_{ij}x_j,\tag{5.6}$$

non  $i = 2, 3, \dots, n - 1, n$ .

**Ariketa moduan**, sortu MATLABen kode bat eragiketa horiek egiteko.

## 5.4. Pibotatzearen beharra

$U$  matrizearen diagonaleko gaiei *pibotak* deritzegu. Goiko adibidean, pibotak dira 10, 2.5 eta 6.2. Biderkatzaileen kalkuluak eta atzeranzko ordezkapenak pibotekin zatitu behar dira. Ondorioz, piboten bat zero bada, ezin dugu burutu algoritmoa. Intuizioz, badakigu ez dela ideia ona kalkulua aurrera eramatea piboten bat ia zero bada. Hori erakusteko, pixka bat aldatuko dugu gure goiko adibidea:

$$\begin{bmatrix} 10 & -7 & 0 \\ -3 & 2.099 & 6 \\ 5 & -1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 3.901 \\ 6 \end{bmatrix}.$$

Matrizeko (2, 2) gaia aldatu egin da 2.000 izatetik 2.099 izatera, eta berdintzaren eskuineko aldea aldatu egin dugu emaitza berbera izateko moduan,  $(0, -1, 1)^T$ . Demagun bost zifra esanguratsutako puntu higikor hamartarra duen makina batean kalkulatu dugula soluzioa.

Ezabapenaren lehenengo urratsak hau ematen du:

$$\begin{bmatrix} 10 & -7 & 0 \\ 0 & -0.001 & 6 \\ 0 & 2.5 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 6.001 \\ 2.5 \end{bmatrix}.$$

Orain, (2, 2) gaia nahiko txikia da matrizeko beste gaiekin konparatzen badugu. Hala ere, trukerik gabe beteko dugu ezabapena. Hurrengo urratsean, hirugarren ekuazioari  $2.5 \cdot 10^3$  bider bigarrena batuko diogu, honela:

$$(5 + (2.5 \cdot 10^3) \cdot 6)x_3 = 2.5 + (2.5 \cdot 10^3) \cdot 6.001.$$

Berdintzaren eskuin aldeko gaien  $6.001 \cdot 2.5 \cdot 10^3 = 1.50025 \cdot 10^4$  dugu. Emaitza  $1.50025 \cdot 10^4$  da, eta hori ezin da adierazi zehatz-mehatz gure puntu higikorraren zenbaki-sistema



hipotetikoan. Hura  $1.5003 \cdot 10^4$ ra biribildu beharko dugu. Gero, emaitza hori 2.5 zenbakiari gehituko diogu  $1.50055 \cdot 10^4$  lortuz, eta berriro biribilduz  $1.5006 \cdot 10^4$  lortzen da. Ondorioz, gure makina hipotetikoan, azken ekuazioa hau bihurtzen da:

$$1.5005 \cdot 10^4 x_3 = 1.5006 \cdot 10^4.$$

Beraz, honekin hasten da atzeranzko ordezkapena:

$$x_3 = \frac{1.5006 \cdot 10^4}{1.5005 \cdot 10^4} = 1.0001 \quad \text{biribilduz.}$$

Emaitza zehatza  $x_3 = 1$  denez, erroreak ez dirudi serioegia. Zoritxarrez,  $x_2$  ekuazio honetatik aurkitu behar dugu:

$$-0.001x_2 + 6 \cdot (1.0001) = 6.001,$$

eta horrek hau ematen du:

$$x_2 = \frac{4 \cdot 10^{-4}}{-1.0 \cdot 10^{-3}} = -0.4.$$

Azkenik,  $x_1$  lehenengo ekuazioaz honela zehazten da:

$$10x_1 + (-7) \cdot (-0.4) = 7,$$

hau lortuz:

$$x_1 = 0.42.$$

Alegia,  $(0, -1, 1)^T$  lortu beharrean  $(0.42, -0.4, 1.0001)^T$  lortu dugu.

Non dago errorea? Ez dago «biribiltze-errorearen pilaketa» sortuta milaka eragiketa aritmetiko eginez. Matrizea ez da ia singularra. Zailtasuna ezabapeneko bigarren urratsean pibot txiki bat hartzetik dator. Ondorioz, biderkatzailea  $2.5 \cdot 10^3$  da, eta azken ekuazioak dauzkan koefizienteak ia  $10^3$  bider jatorrizko problemaren koefizienteak dira.

**Ariketa moduan**, aztertu zer gertatuko litzatekeen bigarren urratsean trukatu bage-nitu bigarren eta hirugarren ekuazioak.

Biderkatzaile guztiak balio absolutuan 1 edo txikiagoak badira, kalkulatuako soluzioa zuzena dela froga daiteke. Biderkatzaileen balio absolutuak 1 baino handiagoak ez izateko, *pibotatze partziala* izendatutako prozesua erabil dezakegu.

### 5.4.1. Pibotatze partziala

Ezabatze-prozesuan  $i$ -garren lerroan bagaude,  $i$ -garren zutabeen geratzen diren gaien artean (hots,  $a_{ji}$ , non  $j = i, i + 1, \dots, n$  baita) balio absolutu handieneko gaia bilatu behar da. Gai hori  $p$ -garrena bada,  $i$ -garren eta  $p$ -garren lerroak trukatu ditugu. Alegia, hau badugu:

$$\max_{j \geq i} \{|a_{ji}|\} = |a_{pi}|,$$

$i$ -garren eta  $p$ -garren lerroak trukatu ditugu. Truke berdina egiten dira berdintzaren eskuineko  $b$  bektorean. Hots,  $b_p$  eta  $b_i$  ere trukatu dira.

## 5.5. $LU$ faktORIZAZIOA

Oro har, ezabatze gaussiarrek bi etapa ditu: *aurreranzko ebazpena* eta *atzeranzko ebazpena*. Aurreranzko ebazpenak  $n - 1$  urrats ditu, aurreko atalean ikusi dugun bezala.  $i$ -garren urratsean,  $i$ -garren ezezaguna ezabatzeko  $i$ -garren ekuazioaren multiploak kentzen dizkiegu gainerako ekuazioei. Baldin  $x_i$ -ren koefizientea «txikia» bada, gomendagarria da (5.2)-(5.3) prozesua egin baino lehen ekuazioak trukatzea. Ezabapen-urratsak berdintzaren eskuinaldeko gaiei batera aplikatu diezazkiekegu, edo trukeak eta biderkatzaileak gorde ditzakegu eta geroago aplikatu eskuinaldeko gaiei. Azken hori ikusiko dugu, hain zuzen ere, jarraian. Azkenik, sistemaren atzeranzko ebazpena (5.5)-(5.6) adierazpenak erabiliz lortuko dugu.

Izan bedi  $\mathbf{P}_i$ ,  $i = 1, \dots, n - 1$ , ezabapenaren  $i$ -garren urratsean erabilitako permutazio-matrizea. Izan bedi  $\mathbf{M}_i$   $i$ -garren urratsean biderkatzaileen negatiboak diagonalaren azpian sartuz lortutako matrize behe-triangeluar unitate bat; horrelako matrizeei *ezabapen-matrizeak* deritzegu. Izan bedi  $\mathbf{U}$  ezabapenaren  $n - 1$  urratsak egin ondoren lortutako azken matrize goi-triangeluarra. Prozesu osoa ekuazio batez deskriba daiteke:

$$\mathbf{M}_{n-1}\mathbf{P}_{n-1} \dots \mathbf{M}_2\mathbf{P}_2\mathbf{M}_1\mathbf{P}_1\mathbf{A} = \mathbf{U}. \quad (5.7)$$

Jarraian ikusiko dugu  $\mathbf{A}$  matrizea ez-singularra bada, beti  $\mathbf{P}$  permutazio-matrize bat aurki dezakegula  $\mathbf{PA} = \mathbf{LU}$  idatzi ahal izateko moduan.

Jo dezagun  $\mathbf{A} \in \mathbb{R}^{3 \times 3}$  matrizea eta bi permutazio-matrize,  $\mathbf{P}_1, \mathbf{P}_2$ , eta bi ezabapen-matrize,  $\mathbf{M}_1, \mathbf{M}_2$ , eraiki ditugula hau gertatzeko moduan:

$$\mathbf{M}_2\mathbf{P}_2\mathbf{M}_1\mathbf{P}_1\mathbf{A} = \mathbf{U}.$$

$\mathbf{P}_i$  ortogonalakenez, badakigu  $\mathbf{P}_i^T\mathbf{P}_i = \mathbb{1}$ . Kasu honetan, gainera,  $\mathbf{P}_i$  hauek simetrikoak direnez,  $\mathbf{P}_i^T = \mathbf{P}_i$  dugu. Beraz,  $\mathbf{P}_2\mathbf{P}_2 = \mathbb{1}$  dugu. Orduan,  $\mathbf{P}_2\mathbf{P}_2$  adierazpena  $\mathbf{M}_1$  eta  $\mathbf{P}_1$  artean sar dezakegu:

$$\mathbf{M}_2\mathbf{P}_2\mathbf{M}_1\mathbf{P}_2\mathbf{P}_2\mathbf{P}_1\mathbf{A} = \mathbf{M}_2\tilde{\mathbf{M}}_1\mathbf{P}_2\mathbf{P}_1\mathbf{A} = \mathbf{U},$$

non  $\tilde{\mathbf{M}}_1$  berrordenatutako ezabapen-matrize bat baita; alegia,  $\mathbf{M}_1$ , baina simetrikoki trukaturako lerroekin eta zutabeekin:

$$\tilde{\mathbf{M}}_1 = \mathbf{P}_2\mathbf{M}_1\mathbf{P}_2.$$

Adibidez, hau badugu:

$$\mathbf{P}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{M}_1 = \begin{bmatrix} 1 & 0 & 0 \\ -m_{21} & 1 & 0 \\ -m_{31} & 0 & 1 \end{bmatrix},$$

zera lortuko dugu:

$$\tilde{\mathbf{M}}_1 = \mathbf{P}_2\mathbf{M}_1\mathbf{P}_2 = \begin{bmatrix} 1 & 0 & 0 \\ -m_{31} & 1 & 0 \\ -m_{21} & 0 & 1 \end{bmatrix}.$$

Demagun orain  $\mathbf{A} \in \mathbb{R}^{4 \times 4}$  matrizea eta bi permutazio-matrize,  $\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3$ , eta bi ezabapen-matrize,  $\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3$ , eraiki ditugula hau gertatzeko moduan:

$$\mathbf{M}_3 \mathbf{P}_3 \mathbf{M}_2 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A} = \mathbf{U}.$$

Orduan, honela eraldatuko dugu, ezkerretik eskuinera:

$$\mathbf{M}_3 \mathbf{P}_3 \mathbf{M}_2 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A} = \mathbf{M}_3 (\mathbf{P}_3 \mathbf{M}_2 \mathbf{P}_3) \mathbf{P}_3 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A} = \mathbf{M}_3 \widetilde{\mathbf{M}}_2 \mathbf{P}_3 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A},$$

non  $\widetilde{\mathbf{M}}_2 = \mathbf{P}_3 \mathbf{M}_2 \mathbf{P}_3$  definitzen baita. Eta gero, hau eginez:

$$\mathbf{M}_3 \widetilde{\mathbf{M}}_2 \mathbf{P}_3 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A} = \mathbf{M}_3 \widetilde{\mathbf{M}}_2 (\mathbf{P}_3 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_2 \mathbf{P}_3) \mathbf{P}_3 \mathbf{P}_2 \mathbf{P}_1 \mathbf{A} = \mathbf{M}_3 \widetilde{\mathbf{M}}_2 \widetilde{\mathbf{M}}_1 \mathbf{P}_3 \mathbf{P}_2 \mathbf{P}_1 \mathbf{A} = \mathbf{U},$$

non  $\widetilde{\mathbf{M}}_1 = \mathbf{P}_3 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_2 \mathbf{P}_3$ .

Horrela ere aplika diezaiokegu kasu orokorrean (5.7) adierazpenari, hots:

$$\mathbf{M}_{n-1} \widetilde{\mathbf{M}}_{n-2} \dots \widetilde{\mathbf{M}}_1 \mathbf{P}_{n-1} \dots \mathbf{P}_1 \mathbf{A} = \mathbf{U}, \quad (5.8)$$

non

$$\widetilde{\mathbf{M}}_k = \mathbf{P}_{n-1} \dots \mathbf{P}_{k+1} \mathbf{M}_k \mathbf{P}_{k+1} \dots \mathbf{P}_{n-1}, \quad k = 1, \dots, n-2. \quad (5.9)$$

Orain, (5.8) adierazpenetik hau ondoriozta dezakegu:

$$\mathbf{L}_1 \mathbf{L}_2 \dots \mathbf{L}_{n-1} \mathbf{U} = \mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1 \mathbf{A}.$$

non  $\mathbf{L}_k = \widetilde{\mathbf{M}}_k^{-1}$  baita, eta hori kalkulatzen da  $\widetilde{\mathbf{M}}_k$  matrizeko diagonalaren azpiko biderkatzaileen zeinuak aldatuz, eta gero permutazioak aplikatuz (5.9) berdintzak adierazten duen bezala. Beraz, hau badugu:

$$\begin{aligned} \mathbf{L} &= \mathbf{L}_1 \mathbf{L}_2 \dots \mathbf{L}_{n-1} \\ \mathbf{P} &= \mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1, \end{aligned}$$

orduan

$$\mathbf{LU} = \mathbf{PA}.$$

Ondorioz,  $\mathbf{L}$  ezabapenean erabilitako biderkatzaile guztiak gordetzen ditu, eta  $\mathbf{P}$  permutazio-matrizeak lerro-truke guztiak gordetzen ditu.

Aurreko atalaren adibiderako, hau dugu:

$$\mathbf{A} = \begin{bmatrix} 10 & -7 & 0 \\ -3 & 2 & 6 \\ 5 & -1 & 5 \end{bmatrix},$$

Hauek dira ezabapenean definitutako matrizeak:

$$\mathbf{P}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{M}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0.3 & 1 & 0 \\ -0.5 & 0 & 1 \end{bmatrix}$$

$$\mathbf{P}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{M}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0.04 & 1 \end{bmatrix}.$$

Beraz,

$$\widetilde{\mathbf{M}}_1 = \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_2 = \begin{bmatrix} 1 & 0 & 0 \\ -0.5 & 1 & 0 \\ 0.3 & 0 & 1 \end{bmatrix}.$$

Hauek dira horiei dagozkien  $\mathbf{L}$  matrizeak:

$$\mathbf{L}_1 = \widetilde{\mathbf{M}}_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ -0.3 & 0 & 1 \end{bmatrix}, \quad \mathbf{L}_2 = \mathbf{M}_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -0.04 & 1 \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ -0.3 & -0.04 & 1 \end{bmatrix}.$$

$\mathbf{LU} = \mathbf{PA}$  erlazioari  $\mathbf{A}$ -ren  $\mathbf{LU}$  faktORIZAZIOA (edo *deskonposizio trianguluarra*) deritzo. Egia esan,  $\mathbf{LU}$  faktORIZAZIOA ezabatze gaussiarra da, notazio matritzialarekin adierazita.

FaktORIZAZIO horrekin, ekuazio-sistema orokor baterako hau dugu:

$$\mathbf{Ax} = \mathbf{b} \quad \iff \quad \mathbf{PAx} = \mathbf{Pb}, \quad (5.10)$$

beraz,  $\mathbf{LUx} = \mathbf{Pb}$  sistema sistema trianguluarren bikote hau bihurtzen da:

$$\begin{cases} \mathbf{Ly} = \mathbf{Pb} \\ \mathbf{Ux} = \mathbf{y}. \end{cases} \quad (5.11)$$

**5.1. teorema.**  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizeak  $\mathbf{LU}$  faktORIZAZIO bat du, baldin  $\det(\mathbf{A}(1:k, 1:k)) \neq 0$  bada,  $k = 1 : n - 1$  guztietarako. Baldin  $\mathbf{LU}$  faktORIZAZIOA existitzen bada eta  $\mathbf{A}$  ez bada singularra, orduan  $\mathbf{LU}$  faktORIZAZIOA bakarra da eta  $\det(\mathbf{A}) = u_{11} \cdot \dots \cdot u_{nn}$  dugu.

*Frogantza.* Demagun  $k - 1$  urratsetan  $\mathbf{LU}$  metodoa erabili dugula eta  $\mathbf{A}^{(k-1)}$  lortu dugula. Kontuan izan  $a_{kk}^{(k-1)}$  gaia  $k$ . pibota dela. Ezabatze gaussiarrengatik

$$\det(\mathbf{A}(1:k, 1:k)) = a_{11}^{(k-1)} \cdot \dots \cdot a_{kk}^{(k-1)}$$

izango dugu. Beraz,  $\mathbf{A}(1:k, 1:k)$  ez bada singularra,  $a_{kk}^{(k-1)}$  ez da zero izango, hots  $\mathbf{A}$ -ren  $\mathbf{LU}$  faktORIZAZIO bat existitzen da.

Bestalde, demagun bi faktORIZAZIO daudela,  $\mathbf{A}$  matrize ez-singularrerako,  $\mathbf{A} = \mathbf{L}_1 \mathbf{U}_1$  eta  $\mathbf{A} = \mathbf{L}_2 \mathbf{U}_2$ , orduan  $\mathbf{L}_2^{-1} \mathbf{L}_1 = \mathbf{U}_2 \mathbf{U}_1^{-1}$  dugu.  $\mathbf{L}_2^{-1} \mathbf{L}_1$  unitate behe-trianguluarra denez eta  $\mathbf{U}_2 \mathbf{U}_1^{-1}$  goi-trianguluarra, biak identitatea izan behar dira berdinak izateko. Beraz,  $\mathbf{L}_1 = \mathbf{L}_2$  eta  $\mathbf{U}_1 = \mathbf{U}_2$ .

Azkenik,  $\mathbf{A} = \mathbf{LU}$  denez, orduan,

$$\det(\mathbf{A}) = \det(\mathbf{LU}) = \det(\mathbf{L}) \det(\mathbf{U}) = \det(\mathbf{U}) = u_{11} \cdot \dots \cdot u_{nn}. \quad \square$$

**5.2. teorema.**  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizerako  $\det(\mathbf{A}(1:k, 1:k)) \neq 0$ ,  $k = 1 : n$  guztietarako, badugu, orduan  $\mathbf{L}$  eta  $\mathbf{M}$  matrize behe-trianguluar bakarrik eta  $\mathbf{D}$  matrize diagonal bakar bat existitzen dira, non  $\mathbf{A} = \mathbf{LDM}^T$  betetzen baita.

*Frogantza.* Aurreko teoremaren emaitzak erabiliz  $\mathbf{A} = \mathbf{LU}$  deskonposizio bakarra existitzen da, eta  $d_i = u_{ii}$ ,  $i = 1 : n$ ,  $\mathbf{D}$ -ren diagonaleko gaiak dira.  $\mathbf{A}$  singularra ez denez,  $\mathbf{D}$  ere ez-singularra da eta  $\mathbf{M}^T = \mathbf{D}^{-1}\mathbf{U}$  unitate-matrize goi-triangeluarra da. Ondorioz,  $\mathbf{A} = \mathbf{LU} = \mathbf{LD}(\mathbf{D}^{-1}\mathbf{U}) = \mathbf{LDM}^T$  eta bakarra da.  $\square$

**5.3. teorema.** *Baldin  $\mathbf{A}$  matrize simetriko ez-singular baten deskonposizioa  $\mathbf{A} = \mathbf{LDM}^T$  bada,  $\mathbf{M} = \mathbf{L}$  da.*

*Frogantza:*

$$\mathbf{A}(\mathbf{M}^{-1})^T = \mathbf{LD} \Rightarrow \mathbf{M}^{-1}\mathbf{A}(\mathbf{M}^{-1})^T = \mathbf{M}^{-1}\mathbf{LD}.$$

Azken matrizea simetrikoa eta behe-triangeluarra da, eta, beraz, diagonal.  $\mathbf{D}$  ez-singularra denez,  $\mathbf{M}^{-1}\mathbf{L}$  unitate-matrize behe-triangeluarra da eta, simetrikoa izateagatik,  $\mathbf{M}^{-1}\mathbf{L} = \mathbf{1}$ .  $\square$

### 5.5.1. Pibotatze baztergarria

Matrize batzuetarako, ez da beharrezkoa pibotatzea.

**5.1. Definizioa.**  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizea *hertsiki diagonal menperatzailea* dela esango dugu, hau betetzen bada:

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|, \quad i = 1, \dots, n. \quad (5.12)$$

[15] liburuan honako teorema hau frogatzen da:

**5.4. teorema.**  $\mathbf{A}^T$  diagonal menperatzailea bada,  $\mathbf{A}$  matrizeak LU faktORIZAZIO bat du eta  $|l_{ij}| \leq 1$ .

Alegia,  $\mathbf{A}^T$  diagonal menperatzailea bada, aurretiko  $\mathbf{LU} = \mathbf{PA}$  faktORIZAZIOA eginez gero,  $\mathbf{P} = \mathbf{1}$  izango da.

## 5.6. Matematikako problema baten baldintza

Hitz gutxitan, problema baten *baldintza* neurri bat da problemako datuen aldaketekiko soluzio zehatzaren sentikortasuna adierazteko. Ideia hori kuantifikatzeko, demagun problema hori datuen  $d$  multzo batek definitua dela. Izan bedi  $s(d)$  problemaren *soluzio zehatza*  $d$

datu horietarako.  $d$  datuetan aldaketa txikiek  $s(d)$ -ko aldaketa txikitara eramaten badute,  $d$  datuetarako problema *ondo baldintzatua* dela esango dugu. Aldiz,  $d$ -ren aldaketa txikiek  $s(d)$ -ko aldaketa handietara eramaten badute,  $d$  datuetarako problema *txarto baldintzatua* dela esango dugu. Demagun  $d_1$  eta  $d_2$  gerta litezkeen bi datu multzo direla. Problemaren *baldintza* (edo *baldintzazko zenbakia*) honelako ratioen maximoa da,  $\|d_1 - d_2\|$  txikia denean:

$$\frac{\|s(d_1) - s(d_2)\|}{\|d_1 - d_2\|}. \quad (5.13)$$

Argi izan behar dugu problema baten baldintza *propietate matematikoa* dela, eta askea dela kalkulu- eta biribiltze-errorearekiko.

Jarraian, problema baten baldintzaren erakusketa erraz bat ikusiko dugu. Izan bedi polinomio honen erroak kalkulatzeko problema:

$$(x - 1)^4 = 0, \quad (5.14)$$

haren lau erroak 1 dira, hain zuzen. Demagun aldaketa txiki bat (esate baterako,  $10^{-8}$ ) egina dela (5.14) berdintzaren eskuinaldean; orain, hau da ebatzi behar dugun ekuazioa:

$$(x - 1)^4 = 10^{-8}, \quad (5.15)$$

Erro zehatza  $1 + 10^{-2}$  da. Beraz, zera dugu kasu honetan:

$$\frac{\|s(d_1) - s(d_2)\|}{\|d_1 - d_2\|} = \frac{|1 - (1 + 10^{-2})|}{|0 - 10^{-8}|} = \frac{10^{-2}}{10^{-8}} = 10^6.$$

Ikusi dugun bezala, datuen aldaketa txiki batek,  $10^{-8}$ , soluzioaren aldaketa handi bat sortzen du,  $10^{-2}$ , zeren hori  $10^6$  bider handiagoa baita datuen aldaketa baino. Alegia, ratio hori askoz handiagoa da 1 baino. Beraz, problema hori *txarto baldintzatua* da. Kontuan izan propietate horrek ez duela zerikusirik erabilitako kalkulu-metodoarekin.

## 5.7. Matrizeen normak

*Matrize-norma* bat, bektoreena bezala,  $\|\cdot\|$  notazioaz adierazten da, eta, bektore-normak bezala, hiru propietate hauek betetzen ditu:

- (i)  $\|\mathbf{A}\| > 0, \quad \forall \mathbf{A} \neq \mathbf{0}$ ;
- (ii)  $\|c\mathbf{A}\| = |c| \cdot \|\mathbf{A}\|, \quad c$  eskalar guztietarako;
- (iii)  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ .

Matrize-norma baterako, erabilgarria da laugarren propietate hau ere, *trinkotasun* izendatua:

$$(iv) \|\mathbf{AB}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|.$$

Jarraian, edozein bektore-normari elkarturiko matrize-norma definituko dugu. Intuitiboki,  $\mathbf{A}$  matrizeak norma «handia» izan beharko luke,  $\mathbf{Ax}$  bektorearen norma handia balitz  $\mathbf{x}$  bektorearen normarekiko.

**5.2. Definizioa.** *Izan bitez  $\mathbf{A}$  matrize bat eta  $\|\cdot\|$  bektore-norma bat.  $\|\mathbf{A}\|$  eragindako matrize-norma honela definitzen da:*

$$\|\mathbf{A}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|}. \quad (5.16)$$

$\mathbf{x}$ -ren tamainaren menpe ez izateko, hau da  $\|\mathbf{A}\|$ -ren beste definizio baliokide bat:

$$\|\mathbf{A}\| = \max_{\|\mathbf{u}\|=1} \|\mathbf{Au}\|.$$

$\|\mathbf{A}\|$ -ren behe-borne bat lor dezakegu  $\mathbf{x}$  bektore ezagun baterako,  $\|\mathbf{Ax}\|$ -ren zatidura  $\|\mathbf{x}\|$ -rekin kalkulatu, hots:

$$\|\mathbf{A}\| \geq \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|}. \quad (5.17)$$

Hauek dira bektoreen bat-, bi- eta infinitu-normek eragindako matrize-normak:

- $\|\mathbf{A}\|_1 = \max_j \|\mathbf{A}_{:,j}\|_1$  (zutabe guztien bat-normetako maximoa)
- $\|\mathbf{A}\|_2 = \sigma_1(A)$  (balio singular handiena).
- $\mathbf{A} \in \mathbb{R}^{n \times n}$  simetrikoa bada,  $\|\mathbf{A}\|_2 = \max_{1 \leq i \leq n} |\lambda_i|$  ( $\lambda_i$ ,  $1 \leq i \leq n$ ,  $\mathbf{A}$ -ren autobalioak dira).
- $\|\mathbf{A}\|_\infty = \max_i \|\mathbf{A}_{i,:}\|_1$  (lerro guztien bat-normetako maximoa)

**5.3. Definizioa.**  $\|\cdot\|_b$  bektore-norma bat eta  $\|\cdot\|_m$  matrize-norma bat **bateragarriak** direla esaten da,  $\mathbf{A}$  eta  $\mathbf{x}$  guztietarako hau betetzen bada:

$$\|\mathbf{Ax}\|_b \leq \|\mathbf{A}\|_m \|\mathbf{x}\|_b.$$

Bektore-norma bat eta berak eragindako matrize-norma beti dira bateragarriak; ikus (5.17).

Bektore-norma batek ez eragindako matrize-norma garrantzitsu bat *Frobeniusen* norma da; hori honela definitzen da  $\mathbf{A} \in \mathbb{R}^{m \times n}$  matrize baterako:

$$\|\mathbf{A}\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{1/2},$$

eta froga daiteke hau egiaztatzen duela:

$$\|\mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}^T \mathbf{A}),$$

$\text{tr}(\mathbf{B})$ -k adierazten du  $\mathbf{B}$ -ren traza, diagonaleko gaien batura.

Frobeniusen norma bateragarria da bektore-norma euklidearrarekin; hots,  $\mathbf{A}$  eta  $\mathbf{x}$  guztietarako,

$$\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2.$$

Bektore edo matrize baterako norma desberdinen balioak oro har desberdinak izan arren, haiek zentzu batean «baliokidetzat» har daitezke. Edozein bi bektore-normetarako, esate baterako  $\|\cdot\|$  eta  $\|\cdot\|'$ , badaude  $c_b$  eta  $c_g$  konstanteak, bektorearen tamainaren menpekoak bakarrik, non  $\mathbf{x}$  guztietarako hau betetzen baita:

$$c_b \|\mathbf{x}\| \leq \|\mathbf{x}\|' \leq c_g \|\mathbf{x}\|. \quad (5.18)$$

Antzeko emaitza bat egiaztatzen da bi matrize-normetarako.

### Matrize-normen propietate batzuk

Edozein  $\mathbf{A} \in \mathbb{R}^{m \times n}$  matritzetarako propietate hauek betetzen dira:

- $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{n} \|\mathbf{A}\|_2$ .
- $\frac{1}{\sqrt{n}} \|\mathbf{A}\|_\infty \leq \|\mathbf{A}\|_2 \leq \sqrt{m} \|\mathbf{A}\|_\infty$ .
- $\frac{1}{\sqrt{m}} \|\mathbf{A}\|_1 \leq \|\mathbf{A}\|_2 \leq \sqrt{n} \|\mathbf{A}\|_1$ .

Propietate horien arabera, bat-, bi- eta infinitu-normak eta Frobeniusen norma balioki-deak dira. Beraz, matrize-norma baten balioa ezagutzen badugu, beste norma batena borna dezakegu propietate horiek erabiliz. Ondorioz, matrize baten norma estimatzeko orduan, bat- edo infinitu-norma aukera ditzakegu, zeren haietarako kalkuluak merkeagoak baitira.

$\|\mathbf{X}\|_p$  norma kalkulatzeko MATLABeko `norm(X,p)` funtzioa erabil daiteke,  $\mathbf{X}$  bektore bat edo matrize bat izanik.



### 5.7.1. Bi-norma eta espektro-erradioa

Izan bedi  $\mathbf{B} \in \mathbb{R}^{n \times n}$  matrize karratua.  $\mathbf{B}$  matrizearen *espektro-erradioa* honela definitzen da:

$$\rho(\mathbf{B}) = \max\{|\lambda|, \lambda \text{ da } \mathbf{B} \text{ -ren autobalore bat}\}.$$

Beraz,  $\rho(\mathbf{B}) = |\bar{\lambda}|$  eta  $\bar{\lambda}$  autobaloreari dagokion autobektorea  $\bar{\mathbf{x}}$  bada, non  $\|\bar{\mathbf{x}}\| = 1$ , hau betetzen da:

$$\rho(\mathbf{B}) = |\bar{\lambda}| = |\bar{\lambda}| \cdot \|\bar{\mathbf{x}}\| = \|\bar{\lambda}\bar{\mathbf{x}}\| = \|\mathbf{B}\bar{\mathbf{x}}\| \leq \|\mathbf{B}\|.$$

Ondorioz, espektro-erradioak beti behe-bornatzen du edozein norma induzitu. Kontuan izan matrize karratu orokor baten norma ez dela espektro-erradioa. Sarritan, espektro-erradioa norma baterako behe-borne on bat izaten da. Are gehiago, edozein  $\mathbf{A}$  matritzearako hau froga daiteke:

$$\|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^T \mathbf{A})} = \sigma_1(\mathbf{A}).$$

Eta  $\mathbf{B}$  matrize simetrikoa bada,  $\|\mathbf{B}\|_2 = \rho(\mathbf{B})$ .

## 5.8. Sistema lineal baten baldintzazko zenbakia

Izan bedi  $\mathbf{A}$  matrize ez-singular bat; hots,  $\mathbf{A}^{-1}$  existitzen da eta bakarra da. Demagun sistema lineal hau:

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

bere soluzio zehatza (eta bakarra)  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$  da. Demagun berdintzaren eskuinaldeko gaia  $\mathbf{b}$  izatetik  $\mathbf{b} + \delta\mathbf{b}$  izatera aldatzen dela (hots,  $\mathbf{b}$  perturbatzen dugula),  $\mathbf{A}$  berdina izanik. Izan bedi  $\mathbf{x} + \delta\mathbf{x}_b$  aldatutako problemaren soluzio zehatza, alegia:

$$\mathbf{A}(\mathbf{x} + \delta\mathbf{x}_b) = \mathbf{b} + \delta\mathbf{b}. \quad (5.19)$$

Atal honetan, bektore edo matrize baten aurrean « $\delta$ » izateak beren dimentsio bereko aldaketa txiki (perturbazio) bat esan nahi du; esate baterako,  $\delta\mathbf{b}$ ,  $\mathbf{b}$ -ren aldaketa bat da. Bestalde,  $\delta\mathbf{x}_b$ -ren  $b$  azpiindizeak esan nahi du  $\mathbf{x}$ -ren aldaketa hau  $\mathbf{b}$ -ren aldaketak eragindakoa dela.

$\mathbf{A}\mathbf{x} = \mathbf{b}$  dugunez, (5.19) erlazioak  $\mathbf{A}\delta\mathbf{x}_b = \delta\mathbf{b}$  inplikatzan du, eta ondorioz:

$$\delta\mathbf{x}_b = \mathbf{A}^{-1}\delta\mathbf{b}.$$

Edozein eragindako norma erabiliz,  $\|\delta\mathbf{x}_b\|$ -ren borne bat lortuko dugu, zeren (5.17) desberdintzaz hau baitugu:

$$\|\delta\mathbf{x}_b\| \leq \|\mathbf{A}^{-1}\| \cdot \|\delta\mathbf{b}\|. \quad (5.20)$$

Perturbazio erlatiboa bornatzeko, kontuan izan  $\mathbf{Ax} = \mathbf{b}$  berdintzatik eta (5.17) -etik hau ondorioztatzen dela:

$$\|\mathbf{b}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\| \Rightarrow \frac{1}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \cdot \frac{1}{\|\mathbf{b}\|}$$

eta adierazpen hori atalez atal biderkatuz (5.20) desberdintzarekin, emaitza garrantzitsu hau lortzen dugu:

$$\frac{\|\delta\mathbf{x}_b\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\| \cdot \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}. \quad (5.21)$$

Kontuan izan desberdintzaren eskuinaldeko kantitatea soluzio zehatzeko perturbazio erlatiboaren balio maximo posiblea dela. Kasu batzuetan, nahiz eta goi-bornea oso handia izan,  $\mathbf{A}$ ,  $\mathbf{b}$  eta  $\delta\mathbf{b}$  berezi batzuetarako bakarrik betetzen da (5.21) adierazpeneko berdintza.

Jarraian,  $\mathbf{A}$  matrizea aldatuko dugu pixka bat (hots, perturbatuko dugu) eta  $\mathbf{b}$  finko eutsiko dugu; orduan, hau dugu:

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}_A) = \mathbf{b}, \quad (5.22)$$

non  $\delta\mathbf{x}_A$ -ren  $\mathbf{A}$  azpiindizeak esan nahi baitu  $\mathbf{x}$ -ren aldaketa  $\mathbf{A}$ -ren aldaketak eragindakoa dela. Demagun  $\|\delta\mathbf{A}\|$  nahiko txikia dela  $\mathbf{A} + \delta\mathbf{A}$  ez-singularra gordetzeko moduan. Orduan, (5.22) -ren eragiketak eginez,

$$\mathbf{Ax} + \mathbf{A}(\delta\mathbf{x}_A) + \delta\mathbf{A}(\mathbf{x} + \delta\mathbf{x}_A) = \mathbf{b} \Rightarrow \mathbf{A}(\delta\mathbf{x}_A) = -\delta\mathbf{A}(\mathbf{x} + \delta\mathbf{x}_A)$$

eta hortik,  $\mathbf{A}$  ez-singularra denez, zera ondorioztatzen da:

$$\delta\mathbf{x}_A = -\mathbf{A}^{-1}\delta\mathbf{A}(\mathbf{x} + \delta\mathbf{x}_A).$$

Orain, (5.17) erabiliz, honako hau ateratzen da:

$$\|\delta\mathbf{x}_A\| \leq \|\mathbf{A}^{-1}\| \cdot \|\delta\mathbf{A}(\mathbf{x} + \delta\mathbf{x}_A)\| \leq \|\mathbf{A}^{-1}\| \cdot \|\delta\mathbf{A}\| \cdot \|\mathbf{x} + \delta\mathbf{x}_A\|$$

eta hortik:

$$\frac{\|\delta\mathbf{x}_A\|}{\|\mathbf{x} + \delta\mathbf{x}_A\|} \leq \|\mathbf{A}^{-1}\| \cdot \|\delta\mathbf{A}\|,$$

eta, azkenik, aldaketa erlatiboaren bornapen hau lortzen da soluzio zehatzean:

$$\frac{\|\delta\mathbf{x}_A\|}{\|\mathbf{x} + \delta\mathbf{x}_A\|} \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\| \cdot \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|}, \quad (5.23)$$

eta berdintza gertatzen da  $\delta\mathbf{A}$ -ren eta  $\mathbf{b}$ -ren balio berezi batzuetarako.

Ikus dezakegunez,  $\|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|$  kantitatea agertzen da (5.21) eta (5.23) bornapenetan, eta horrek erakusten du sistema lineal bateko datuen aldaketa batek eragindako aldaketa maximo posiblea soluzio zehatzaren gainean. Gogoan (5.13) izanik,  $\mathbf{A}$  matrize ez-singular baten *baldintzazko zenbakia* (edo *baldintza* bakarrik) honela definitzen dugu ( $\mathbf{Ax} = \mathbf{b}$  sistemaren ebazpenarekiko):

$$\kappa(\mathbf{A}) = \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\|. \quad (5.24)$$

Edozein eragindako normatarako, identitate-matrizearen norma 1 da. Izan ere,  $\mathbb{1} = \mathbf{A}^{-1}\mathbf{A}$  eta  $\|\mathbf{A}^{-1}\mathbf{A}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\|$ , orduan  $\kappa(\mathbf{A}) \geq 1$ . Beraz, matrize *ondo baldintzatu* baten baldintzazko zenbakia unitatearen ordenakoa da, eta matrize *txarto baldintzatu* baten baldintzazko zenbakia unitatea baino askoz handiagoa da. Nahiz eta  $\kappa(\mathbf{A})$  balioa aldatu bere kalkuluan erabilitako normaren arabera, balio horiek konparagarriak dira normen baliokidetasunagatik.

Bestalde, sistema linealaren datuen perturbazio baterako soluzio zehatzaren aldaketa erlatiboa ezagutzen badugu, baldintzazko zenbakiaren *behe-borne* bat lor dezakegu. Berrordenatuz (5.21) eta (5.23) desberdintzak eta (5.24) kontuan hartuz, hau lortzen da:

$$\kappa(\mathbf{A}) \geq \frac{\|\delta\mathbf{x}_b\|/\|\mathbf{x}\|}{\|\delta\mathbf{b}\|/\|\mathbf{b}\|} \quad \text{eta} \quad \kappa(\mathbf{A}) \geq \frac{\|\delta\mathbf{x}_A\|/(\|\mathbf{x} + \delta\mathbf{x}_A\|)}{\|\delta\mathbf{A}\|/\|\mathbf{A}\|}. \quad (5.25)$$

**5.1. adibidea.** *Izan bedi sistema lineal hau:*

$$\begin{aligned} 0.550x_1 + 0.423x_2 &= 0.127 \\ 0.484x_1 + 0.372x_2 &= 0.112. \end{aligned}$$

*Sistema horretarako, hau dugu:*

$$\mathbf{Ax} = \mathbf{b} \quad \text{non} \quad \mathbf{A} = \begin{bmatrix} 0.550 & 0.423 \\ 0.484 & 0.372 \end{bmatrix} \quad \text{eta} \quad \mathbf{b} = \begin{bmatrix} 0.127 \\ 0.112 \end{bmatrix}.$$

*Ikusiko dugu sistema horren  $\mathbf{A}$  matrizea txarto baldintzatu dela.*

*Ebazpena.* Jo dezagun  $\mathbf{b}$  bektorea honela perturbatzen dugula:

$$\tilde{\mathbf{b}} = \mathbf{b} + \delta\mathbf{b} = \begin{bmatrix} 0.127 \\ 0.112 \end{bmatrix} + \begin{bmatrix} 0.00007 \\ 0.00028 \end{bmatrix} = \begin{bmatrix} 0.12707 \\ 0.11228 \end{bmatrix}.$$

$\mathbf{Ax} = \mathbf{b}$  sistemaren soluzio zehatza  $\mathbf{x} = (1, -1)^T$  da, baina  $\mathbf{A}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$  sistemaren soluzio zehatza  $\mathbf{x} = (1.7, -1.91)^T$  da, hots  $\delta\mathbf{x}_b = (0.7, -0.91)^T$ . Infinitu-norma erabiltzen badugu,  $\mathbf{b}$ -ren eta  $\mathbf{x}$ -ren perturbazio erlatiboak kalkulatzeko, zera dugu:

$$\frac{\|\delta\mathbf{b}\|_\infty}{\|\mathbf{b}\|_\infty} = \frac{0.00028}{0.127} \approx 2.2 \cdot 10^{-3} \quad \text{eta} \quad \frac{\|\delta\mathbf{x}_b\|_\infty}{\|\mathbf{x}\|_\infty} = 0.91.$$

Beraz,  $\delta\mathbf{b}$  horretarako, soluzioan egon den aldaketa erlatiboa 413 bider baino gehiago izan da; horrek erakusten du  $\mathbf{A}$ -ren baldintzazko zenbakia gutxienez 413 dela; hots,  $\kappa_\infty(\mathbf{A}) \geq 413$ , ikus (5.25).

$\mathbf{A}$ -ren baldintzapen txarra ikus dezakegu  $a_{21}$  gaia pixka bat aldatuz ere (0.001 kantitatean), hau izanik:

$$\mathbf{A} + \delta\mathbf{A} = \begin{bmatrix} 0.550 & 0.423 \\ 0.483 & 0.372 \end{bmatrix}.$$

Orain,  $(\mathbf{A} + \delta\mathbf{A})\hat{\mathbf{x}} = \mathbf{b}$  sistemaren soluzio zehatza  $\hat{\mathbf{x}} = (-0.4536, 0.8900)^T$ , zeinak  $\|\delta\mathbf{x}_A\|_\infty = \|\hat{\mathbf{x}} - \mathbf{x}\|_\infty = 1.89$  ematen baitu. Ondorioz, zera dugu:

$$\frac{\|\delta\mathbf{x}_A\|_\infty / \|\hat{\mathbf{x}}\|_\infty}{\|\delta\mathbf{A}\|_\infty / \|\mathbf{A}\|_\infty} = \frac{1.89/0.89}{0.001/0.973} \approx 2066.$$

Kasu honetan, soluzio zehatzaren aldaketa erlatiboa 2066 bider  $\mathbf{A}$ -ren aldaketa erlatiboa da; hots,  $\kappa_\infty(\mathbf{A}) \geq 2066$ , ikus (5.25) .

Benetan,  $\mathbf{A}$ -ren alderantzizko zehatza (5 digitutara biribildua) hau da:

$$\mathbf{A}^{-1} = \begin{bmatrix} -2818.2 & 3204.5 \\ 3666.7 & -4166.7 \end{bmatrix},$$

eta (infinitu-norma erabiliz)  $\|\mathbf{A}\|_\infty = 0.973$  eta  $\|\mathbf{A}^{-1}\|_\infty = 7833.4$  dugunez, hau da baldintzazko zenbakia:

$$\kappa_\infty(\mathbf{A}) = 7833.4 \cdot 0.973 \approx 7622.$$

Beraz, benetako baldintzazko zenbakia askoz handiagoa da  $\delta\mathbf{b}$ -ren eta  $\delta\mathbf{A}$ -ren bitartez aurkitutako 413 eta 2066, hurrenez hurren,  $\kappa_\infty(\mathbf{A})$ -ren behe-borneak baino.  $\square$

$\mathbf{A}$  matrize bat txarto baldintzatua da tamaina bereko bektore desberdinei aplikatuz gero, biderkadura-bektoreen tamainak oso desberdinak direnean. Alegia,  $\mathbf{x}$  bektore unitario baterako  $\|\mathbf{Ax}\|$  handia bada eta beste bektore unitario baterako txikia bada, orduan,  $\mathbf{A}$  txarto baldintzatua da. Baina,  $\|\mathbf{Ax}\|/\|\mathbf{x}\|$  handia bada  $\mathbf{x}$  guztietarako (esate baterako,  $\mathbf{A} = \text{diag}(10^8, 10^8)$  kasuan), orduan,  $\mathbf{A}$  ez da txarto baldintzatua. Aldiz, jo ditzagun matrize eta bektore hauek:

$$\hat{\mathbf{A}} = \begin{bmatrix} 10^4 & \\ & 10^{-4} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \text{eta} \quad \hat{\mathbf{x}} = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

non  $\|\mathbf{x}\| = \|\hat{\mathbf{x}}\| = 1$ . Orduan, hau dugu:

$$\hat{\mathbf{A}}\mathbf{x} = \begin{bmatrix} 10^4 \\ 0 \end{bmatrix} \quad \text{eta} \quad \hat{\mathbf{A}}\hat{\mathbf{x}} = \begin{bmatrix} 0 \\ 10^{-4} \end{bmatrix},$$

zera gertatuz:  $\|\hat{\mathbf{A}}\mathbf{x}\| = 10^4$  (handia) eta  $\|\hat{\mathbf{A}}\hat{\mathbf{x}}\| = 10^{-4}$  (txikia). Alegia,  $\hat{\mathbf{A}}$  matrizeak txarto baldintzatua izan behar du. Egia esan, erraza da ikustea aldaketa handienak  $\mathbf{x}$  eta  $\hat{\mathbf{x}}$  bektore horietarako lortzen dituela, eta  $\kappa(\hat{\mathbf{A}}) = 10^8$  dela.

$\mathbf{A}$ -ren baldintza eta bere singularitasuna oso kontzeptu erlazionatuak dira. Informalki, txarto baldintzatutako matrize bat «ia singularra» da. Baina, nahiz eta matrize singular baten determinantea zero izan, determinantea ia zero duen matrize baten baldintza ez da nahitaez txarra izan behar. Adibidez,  $n \times n$  dimentsioko  $\mathbf{A} = \text{diag}(10^{-10})$  matrizearen determinantea  $10^{-10n}$  da, baina bere baldintza perfektua da,  $\kappa(\mathbf{A}) = 1$  da.

Bestalde, matrize singular batek, gutxienez, autobalio nulu bat dauka. Gainera,  $\mathbf{A}$  simetrikoa denean, bere baldintzazko zenbakia da autobalio handienak txikienarekiko duen

zatidura; alegia:

$$\kappa_2(\mathbf{A}) = \frac{\max_{1 \leq i \leq n} |\lambda_i|}{\min_{1 \leq i \leq n} |\lambda_i|}.$$

Ondorioz, matrize simetriko txarto baldintzatu batek tamaina txikiko autobalio bat izan behar du bere normarekiko (izan ere,  $\|\mathbf{A}\|_2 = \max_{1 \leq i \leq n} |\lambda_i|$ ). Oro har, irizpide hori ez da zuzena matrizea simetrikoa ez denean.

$\mathbf{P}$  permutazio-matrizea bada,  $\mathbf{P}\mathbf{x}$ -ren osagaiak  $\mathbf{x}$ -ren osagaien berrantolaketa direnez,  $\|\mathbf{P}\mathbf{x}\| = \|\mathbf{x}\|$  betetzen da  $\mathbf{x}$  guztietarako, eta beraz:

$$\kappa(\mathbf{P}) = 1.$$

$\mathbf{A}$  matrizea  $c \neq 0$  eskalar batez biderkatzen badugu,  $\|c\mathbf{A}\| = |c| \cdot \|\mathbf{A}\|$  eta  $\|(c\mathbf{A})^{-1}\| = |1/c| \cdot \|\mathbf{A}^{-1}\|$  dira, eta, ondorioz:

$$\kappa(c\mathbf{A}) = \|(c\mathbf{A})^{-1}\| \cdot \|c\mathbf{A}\| = \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\| = \kappa(\mathbf{A}).$$

$\mathbf{D}$  matrize diagonalada bada, orduan:

$$\kappa(\mathbf{D}) = \frac{\max |d_{ii}|}{\min |d_{ii}|}.$$

MATLABek badauzka funtzio batzuk baldintzazko zenbakia kalkulatzeko:

- $\kappa_2(\mathbf{A})$  (bi-norma erabiliz) kalkulatu du `cond(A)` edo `cond(A,2)` funtzioez. Balio singularren kalkuluan erabiltzen du `svd(A)` funtzioa.
- $\kappa_1(\mathbf{A})$  (bat-norma erabiliz) kalkulatu du `cond(A,1)` funtzioaz. Horrek `inv(A)` funtzioa erabiltzen du. Kalkuluak eragiketa gutxiago behar ditu `cond(A,2)` baino.
- $\kappa_\infty(\mathbf{A})$  (infinitu-norma erabiliz) kalkulatu du `cond(A,inf)` funtzioaz. Horrek `inv(A)` funtzioa erabiltzen du eta `cond(A,1)`-en lan berdina behar du.

## 5.9. Cholesky-ren faktORIZAZIOA

**5.5. teorema.** *Izan bedi  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrize simetriko bat. Orduan,  $\mathbf{A}$ -k  $n$  autobalio errealak ditu,  $\lambda_1, \dots, \lambda_n$ , eta horiei elkartutako  $\mathbf{v}_1, \dots, \mathbf{v}_n$  autobektore unitarioek  $\mathbb{R}^n$ -rako oinarri ortonormal bat osatzen dute (hots,  $\mathbf{v}_i^T \mathbf{v}_j = 1$ ,  $i = j$  bada, eta  $\mathbf{v}_i^T \mathbf{v}_j = 0$ ,  $i \neq j$  bada).*

**5.6. teorema.** *Izan bedi  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrize simetriko bat. Orduan,  $\mathbf{A}$  definitu positiboa ( $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \ \forall \mathbf{x} \neq \mathbf{0}$ ) da, baldin eta soilik baldin autobalio guztiak positiboak badira.*

*Frogantza.* Izan bitez  $\lambda_1, \dots, \lambda_n$  autobalioak eta horiei dagozkien  $\mathbf{v}_1, \dots, \mathbf{v}_n$  autobektore ortonormalak. Demagun  $j$  baterako  $\lambda_j \leq 0$ . Orduan:

$$\mathbf{v}_j^T \mathbf{A} \mathbf{v}_j = \mathbf{v}_j^T (\lambda_j \mathbf{v}_j) = \lambda_j \mathbf{v}_j^T \mathbf{v}_j = \lambda_j \leq 0,$$

horrek erakusten du  $\mathbf{A}$  ez dela definitu positiboa. Beraz,  $\mathbf{A}$  definitu positiboa bada, autobalio guztiak positiboak dira.

Orain, frogatuko dugu  $\mathbf{A}$ -ren autobalio guztiak positiboak badira,  $\mathbf{A}$  definitu positiboa dela.

Demagun  $\lambda_i > 0$ ,  $i = 1, \dots, n$ . Orduan,  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$   $\mathbb{R}^n$ -rako oinarri bat osatzen dutenez, hau betetzen da edozein  $\mathbf{v} \in \mathbb{R}^n$  bektore ez-nulu baterako:

$$\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{v}_i \Rightarrow \mathbf{A} \mathbf{v} = \mathbf{A} \sum_{i=1}^n \alpha_i \mathbf{v}_i = \sum_{i=1}^n \alpha_i \mathbf{A} \mathbf{v}_i = \sum_{i=1}^n \alpha_i \lambda_i \mathbf{v}_i$$

eta  $\alpha_j$  bat gutxienez ez da zero. Orduan,  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  oinarria ortonormala izateagatik, hau dugu:

$$\mathbf{v}^T \mathbf{A} \mathbf{v} = \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \mathbf{v}_i)^T (\lambda_j \alpha_j \mathbf{v}_j) = \sum_{i=1}^n \sum_{j=1}^n \lambda_j \alpha_i \alpha_j \mathbf{v}_i^T \mathbf{v}_j = \sum_{j=1}^n \lambda_j \alpha_j^2 > 0.$$

Beraz,  $\mathbf{A}$  definitu positiboa da.  $\square$

### 5.7. teorema (Gerschgorin).

Izan bedi  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrize simetriko bat  $\lambda_1, \dots, \lambda_n$  autobalioekin. Orduan, hau dugu:

$$\min_{1 \leq i \leq n} \lambda_i \geq \min_{1 \leq i \leq n} \left\{ a_{ii} - \sum_{j=1, j \neq i}^n |a_{ij}| \right\}$$

$$\max_{1 \leq i \leq n} \lambda_i \leq \max_{1 \leq i \leq n} \left\{ a_{ii} + \sum_{j=1, j \neq i}^n |a_{ij}| \right\}.$$

**5.1. korolaria.**  $\mathbf{A}$  matrizea definitu positiboa da,  $i = 1, \dots, n$  guztietarako hau betetzen bada:

$$a_{ii} - \sum_{j=1, j \neq i}^n |a_{ij}| > 0.$$

*Frogantza.* (Kontuan izan  $\mathbf{A}$  hertsiki diagonal menperatzailea dela). Gerschgorin-en teorema-ren lehenengo desberdintzagatik, autobalio txikiena zero baino handiagoa denez,  $\mathbf{A}$  definitu positiboa da.  $\square$

Hala ere, gerta daiteke  $\mathbf{A}$  definitu positiboa izatea eta ez hertsiki diagonal menperatzailea.

$\mathbf{A}$  simetrikoa bada, bere autobalio txikienaren hurbilpen gisa balio hau hartu ohi da:

$$d = \min_{1 \leq i \leq n} \left\{ a_{ii} - \sum_{j=1, j \neq i}^n |a_{ij}| \right\}.$$

Ariketa gisa, frogatu  $\mathbf{A}$  matrizea simetrikoa eta definitu positiboa bada, diagonaleko gai guztiak positiboak direla, eta matrizearen tamaina handieneko gaia diagonalean dagoela.

**5.8. teorema.**  $\mathbf{A} \in \mathbb{R}^{n \times n}$  definitu positiboa bada eta  $\mathbf{X} \in \mathbb{R}^{n \times k}$  matrizeak  $k$  heina badu,  $\mathbf{B} = \mathbf{X}^T \mathbf{A} \mathbf{X} \in \mathbb{R}^{k \times k}$  ere definitu positiboa da.

*Frogantza.* Demagun  $\mathbf{v} \in \mathbb{R}^k$  bektoreak hau betetzen duela:

$$0 \geq \mathbf{v}^T \mathbf{B} \mathbf{v} = \mathbf{v}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{v} = (\mathbf{X} \mathbf{v})^T \mathbf{A} (\mathbf{X} \mathbf{v}),$$

baina,  $\mathbf{A}$  definitu positiboa denez,  $\mathbf{X} \mathbf{v} = \mathbf{0}$  izan behar da; gainera,  $\mathbf{X}$ -ren  $k$  zutabeak linealki askeak direnez,  $\mathbf{v} = \mathbf{0}$  izan behar da. Horrek esan nahi du aurreko desberdintza ezin dela bete eta  $\mathbf{B}$  definitu positiboa dela.  $\square$

$\mathbf{A}$  simetrikoa eta definitu positiboa bada, 5.3. teoremagatik,  $\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^T$  da, non  $\mathbf{L}$  unitate-matrize behe-triangeluarra eta  $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$  matrize diagonal baitira. Beraz,  $\mathbf{L}^{-1} \mathbf{A} (\mathbf{L}^{-1})^T = \mathbf{D}$  dugu, eta aurreko teoremagatik  $\mathbf{A}$  definitu positiboa bada,  $\mathbf{D}$  ere bai; eta  $\mathbf{D}$  diagonal denez,  $d_i > 0$ ,  $i = 1, \dots, n$ .

**5.9. teorema (Choleskyren faktORIZAZIOA).**  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizea simetrikoa eta definitu positiboa bada,  $\mathbf{R} \in \mathbb{R}^{n \times n}$  matrize goi-triangeluar bakar bat existitzen da diagonaleko gai positiboekin eta  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  betetzen duena.

*Frogantza.*  $\mathbf{L}$  unitate-matrize behe-triangeluar bakar bat eta  $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$  matrize diagonal bakar bat daude  $\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^T$  betetzen dutenak (ikus 5.3. teorema).  $d_k$  gaiak positiboak direnez,  $\mathbf{R}^T = \mathbf{L} \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_n})$  matrizea erreal eta behe-triangeluarra da gai diagonal positiboekin. Gainera,  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  betetzen du. Bakartasuna  $\mathbf{L} \mathbf{D} \mathbf{L}^T$  faktORIZAZIOAREN bakartasunak inplikutzen du.  $\square$

$\mathbf{A}$  simetrikoa ( $\mathbf{A} = \mathbf{A}^T$ ) eta definitu positiboa denean,  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  faktORIZAZIOA (edo deskonposizioa) kalkula dezakegu, non  $\mathbf{R}$  matrize goi-triangeluar bat baita, hots:

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{12} & a_{22} & & a_{2n} \\ \vdots & & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{nn} \end{bmatrix} = \begin{bmatrix} r_{11} & & & \\ r_{12} & r_{22} & & \\ \vdots & & \ddots & \\ r_{1n} & r_{2n} & \cdots & r_{nn} \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ & r_{22} & & r_{2n} \\ & & \ddots & \vdots \\ & & & r_{nn} \end{bmatrix}.$$

Hori dela eta, zera dugu:

- $\mathbf{A}$ -ren  $a_{11}$  gaitik 1. lerroari jarraituz, ekuazio hauek ditugu:

$$\begin{aligned} r_{11}^2 &= a_{11} \\ r_{11}r_{12} &= a_{12} \\ &\vdots \\ r_{11}r_{1n} &= a_{1n}. \end{aligned}$$

- $\mathbf{A}$ -ren  $a_{22}$  gaitik 2. lerroari jarraituz, ekuazio hauek ditugu:

$$\begin{aligned} r_{12}^2 + r_{22}^2 &= a_{22} \\ r_{12}r_{13} + r_{22}r_{23} &= a_{23} \\ &\vdots \\ r_{12}r_{1n} + r_{22}r_{2n} &= a_{2n}. \end{aligned}$$

- Horrela jarraituz,  $i = 1, \dots, n$  guztietarako,  $\mathbf{A}$ -ren  $a_{ii}$  gaitik  $i$ . lerroko ekuazio hauek ditugu:

$$\begin{aligned} r_{1i}^2 + r_{2i}^2 + \dots + r_{ii}^2 &= a_{ii} \\ r_{1i}r_{1,i+1} + r_{2i}r_{2,i+1} + \dots + r_{ii}r_{i,i+1} &= a_{i,i+1} \\ &\vdots \\ r_{1i}r_{1n} + r_{2i}r_{2n} + \dots + r_{ii}r_{in} &= a_{in}. \end{aligned}$$

- Ondorioz,  $i = 1, \dots, n$  guztietarako,  $\mathbf{R}$ -ren koefizienteak honela kalkula ditzakegu:

$$r_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} r_{ki}^2}, \quad i = 1, \dots, n, \quad (5.26)$$

$$r_{ij} = \frac{a_{ij} - \sum_{k=1}^{i-1} r_{ki}r_{kj}}{r_{ii}}, \quad j = i + 1, \dots, n. \quad (5.27)$$

**5.2. adibidea.** Kalkula ezazu matrize simetriko eta definitu positibo honen Choleskyren faktORIZAZIOA:

$$\mathbf{A} = \begin{bmatrix} 6 & 15 & 55 \\ 15 & 55 & 225 \\ 55 & 225 & 979 \end{bmatrix}.$$

*Ebazpena.* Lehenengo lerroko ( $i = 1$ ) (5.26) erabiliz, zera dugu:

$$r_{11} = \sqrt{a_{11}} = \sqrt{6} = 2.44949.$$

Orduan, (5.27) erabiliz koefiziente hauek lortzeko:

$$r_{12} = \frac{a_{12}}{r_{11}} = \frac{15}{2.44949} = 6.123724$$

$$r_{13} = \frac{a_{13}}{r_{11}} = \frac{55}{2.44949} = 22.45366.$$



Hurrengo lerroko ( $i = 2$ ), zera dugu:

$$r_{22} = \sqrt{a_{22} - r_{12}^2} = \sqrt{55 - (6.123724)^2} = 4.1833$$

$$r_{23} = \frac{a_{23} - r_{12}r_{13}}{r_{22}} = \frac{225 - 6.123724 \cdot 22.45366}{4.1833} = 20.9165.$$

Hirugarren lerroko ( $i = 3$ ), hau lortuko dugu:

$$r_{33} = \sqrt{a_{33} - r_{13}^2 - r_{23}^2} = \sqrt{979 - (22.45366)^2 - (20.9165)^2} = 6.110101.$$

Ondorioz, Choleskyren faktORIZAZIOAK zera ematen du:

$$\begin{bmatrix} 2.44949 & 6.123724 & 22.45366 \\ & 4.1833 & 20.9165 \\ & & 6.110101 \end{bmatrix}.$$

FaktORIZAZIO horren baliotasuna ontzat emango dugu,  $\mathbf{R}^T \mathbf{R} = \mathbf{A}$  berdintza betetzen dela egiaztatuz.  $\square$

$\mathbf{R}^T \mathbf{R} = \mathbf{A}$  faktORIZAZIO hori lortu ondoren,  $\mathbf{A}\mathbf{x} = \mathbf{b}$  sistema ebatz dezakegu  $LU$  faktORIZAZIOAREKIN egiten den bezala; alegia, sistema triangeluar hauek ebatziz:

$$\begin{aligned} \mathbf{R}^T \mathbf{y} &= \mathbf{b} \\ \mathbf{R}\mathbf{x} &= \mathbf{y}. \end{aligned}$$

$\mathbf{A}$  matrize simetriko bat definitu positiboa bada, ez da beharrezkoa lerroen trukaketa; aldiz,  $LU$  metodoak pibotaze partziala behar du zenbakizko egonkorra izateko. Choleskyren metodoan nahiz eta  $\mathbf{A}$  txarto baldintzatua izan, ez dugu pibotaziorik egin behar, zeren  $i = 1, \dots, n$  guztietarako hau betetzen baita:

$$r_{1i}^2 + r_{2i}^2 + \dots + r_{ii}^2 = a_{ii} \quad (5.28)$$

eta, ondorioz,  $\mathbf{R}$ -ko gai guztiak honela bornatuta daude:

$$|r_{ji}| \leq \sqrt{a_{ii}}, \quad j = 1, \dots, i.$$

Metodo hau oso erabilia da optimizazioan, bere zenbakizko egonkortasunagatik. MATLABen `chol(A)` funtzioak  $\mathbf{A}$  matrizearen Choleskyren faktorea ematen digu.

**5.10. teorema.** *Izan bedi  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrize simetrikoa eta definitu positiboa eta  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$  bere Choleskyren deskonposizioa. Orduan:*

$$\kappa_2(\mathbf{A}) \geq \frac{\max_i r_{ii}^2}{\min_i r_{ii}^2}.$$

Gainera,  $\mathbf{A} = \mathbf{LDL}^T$  Choleskyren deskonposizioa erabiltzen badugu,  $\mathbf{L}$  unitate-matrize behe-triangeluarra eta  $\mathbf{D} = \text{diag}(d_i)$  matrize diagonalak izanik, hau betetzen da:

$$\kappa_2(\mathbf{A}) \geq \frac{\max_i d_i}{\min_i d_i}.$$

Alegia,  $\mathbf{D}$ -ren gai diagonal maximoaren eta minimoaren zatidurak  $\mathbf{A}$  matrizearen baldintzaren behe-borne bat ematen digu, eta, bide batez, hurbilpen bat.

Frogantza:

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{u}\|_2=1} \|\mathbf{A}\mathbf{u}\|_2 \geq \max_i \|\mathbf{A}\mathbf{e}_i\|_2 = \max_i \sqrt{a_{1i}^2 + \dots + a_{ni}^2} \geq \max_i a_{ii} \geq \max_i r_{ii}^2$$

azken desberdintza (5.28) berdintzaren ondorioa da.

Izan bedi  $\mathbf{B} = \mathbf{A}^{-1}$ , orduan  $\mathbf{B} = (\mathbf{R}^T \mathbf{R})^{-1} = \mathbf{R}^{-1} (\mathbf{R}^T)^{-1} = \mathbf{R}^{-1} (\mathbf{R}^{-1})^T$ . Beraz,

$$\|\mathbf{B}\|_2 \geq \max_i b_{ii} \geq \max_i \frac{1}{r_{ii}^2} = \frac{1}{\min_i r_{ii}^2}.$$

Ondorioz,

$$\kappa_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{B}\|_2 \geq \frac{\max_i r_{ii}^2}{\min_i r_{ii}^2}.$$

Bestalde,  $\mathbf{A} = \mathbf{LDL}^T = \mathbf{R}^T \mathbf{R}$  denez,  $d_i = r_{ii}^2$  eta horrek frogantza hau bukatzen du.  $\square$

## 5.10. Metodo iteratiboak

$LU$  eta Choleskyren faktORIZAZIOA *metodo zuzenak* dira; izan ere, ez balira egongo biribiltze-erroreak, ikusi dugun bezala, urratsen kopuru finitu batean soluzio zehatzera helduko lirateke. Baina, sistemaren matrizea eskasa bada (hots, gaien ehuneko handi bat zeroak direnean), Gaussen metodoak zeroak hondatzen ditu (hots, ez-zero bihurtzen ditu). Edonola ere, matrizearen egitura hau eskasa denean, badaude estrategia egokiak esplotatzeko; ikus [10, 13]. MATLABeko `nnz(A)` funtzioak  $\mathbf{A}$  matrizean zero ez diren gaien kopurua zenbatzen du.

*Metodo iteratiboak* egokiak dira matrizearen egitura aprobetxatuz eragiketa gutxiago egiteko eta ordenagailuaren memoria gutxiago erabiltzeko, eta hori oso interesgarria da sistema oso handia denean, deribatu partzialtako ekuazioen sistemetan bezala.

### 5.10.1. Jacobi-ren iterazioa

5.3. adibidea. *Ebatzi sistema hau:*

$$\begin{aligned} 4x - y + z &= 7 \\ 4x - 8y + z &= -21 \\ -2x + y + 5z &= 15. \end{aligned}$$

*Ebazpena.* Ekuazio horiek honela idatz ditzakegu:

$$\begin{aligned} x &= \frac{7 + y - z}{4} \\ y &= \frac{21 + 4x + z}{8} \\ z &= \frac{15 + 2x - y}{5} \end{aligned}$$

eta, hortik, prozesu iteratibo hau ateratzen dugu:

$$\begin{aligned} x^{(k+1)} &= \frac{7 + y^{(k)} - z^{(k)}}{4} \\ y^{(k+1)} &= \frac{21 + 4x^{(k)} + z^{(k)}}{8} \\ z^{(k+1)} &= \frac{15 + 2x^{(k)} - y^{(k)}}{5}. \end{aligned} \tag{5.29}$$

$k$	$x^{(k)}$	$y^{(k)}$	$z^{(k)}$
0	1.0	2.0	2.0
1	1.75	3.375	3.0
2	1.84375	3.875	3.025
3	1.9625	3.925	2.9625
4	1.99062500	3.97656250	3.00000000
5	1.99414063	3.99531250	3.00093750
...	...	...	...
15	1.99999993	3.99999985	2.99999993
...	...	...	...
19	2.00000000	4.00000000	3.00000000

**5.1. taula.** Jacobiren iterazioaren konbergentzia 5.3. adibidean.

Hasierako puntua  $(x^{(0)}, y^{(0)}, z^{(0)})^T = (1, 2, 2)^T$  bada, ikusiko dugu iterazio horrek  $(2, 4, 3)^T$  soluziora jotzen duela.

$x^{(0)} = 1$ ,  $y^{(0)} = 2$  eta  $z^{(0)} = 2$  ordezkatuz (5.29) eskuineko ataletan, hau lortzen dugu:

$$\begin{aligned}x^{(1)} &= \frac{7 + 2 - 2}{4} = 1.75 \\y^{(1)} &= \frac{21 + 4 \cdot 1 + 2}{8} = 3.375 \\z^{(1)} &= \frac{15 + 2 \cdot 1 - 2}{5} = 3.00.\end{aligned}$$

Prozesu hori jarraituz, 5.10.1. taula lortzen da. Taula horrek adierazten du (5.29) iterazioak  $(2, 4, 3)^T$  soluziora jotzen duela.  $\square$

Prozesu horri *Jacobiren iterazioaren metodoa* deritzogu.

Izan bedi  $\mathbf{Ax} = \mathbf{b}$  sistema eta  $\mathbf{A}$  matrizearen deskonposizio hau:

$$\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U},$$

non

$$\mathbf{L} = \begin{bmatrix} 0 & 0 & \dots & \dots & 0 \\ a_{21} & 0 & \dots & & 0 \\ a_{31} & a_{32} & \ddots & & 0 \\ \vdots & & & 0 & 0 \\ a_{n1} & a_{n2} & \dots & a_{n,n-1} & 0 \end{bmatrix}$$

$$\mathbf{D} = \text{diag}(a_{11}, \dots, a_{nn}) \tag{5.30}$$

$$\mathbf{U} = \begin{bmatrix} 0 & a_{12} & \dots & \dots & a_{1n} \\ 0 & 0 & \dots & \dots & \vdots \\ 0 & 0 & \ddots & & a_{n-2,n} \\ \vdots & & & \ddots & a_{n-1,n} \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}.$$

Orduan:

$$\mathbf{Ax} = (\mathbf{L} + \mathbf{D} + \mathbf{U})\mathbf{x} = \mathbf{b}$$

eta ondorioz,

$$\mathbf{Dx} = \mathbf{b} - (\mathbf{L} + \mathbf{U})\mathbf{x}.$$

Hortaz, honela idatz dezakegu Jacobiren  $k$ -garren iterazioa:

$$\mathbf{Dx}^{(k+1)} = \mathbf{b} - (\mathbf{L} + \mathbf{U})\mathbf{x}^{(k)}. \tag{5.31}$$

Beraz, matrizeak ez dira aldatzen eta, eskasak badira (hots, zero asko badituzte), eragiketara gutxi egin behar ditugu. Orain, (5.31) iterazioa garatuz  $x_i$  bakoitzeko,  $i = 1, \dots, n$  guztietarako, zera dugu:

$$a_{ii}x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right).$$

Azkenik,  $i = 1, \dots, n$  guztietarako, hau egingo da Jacobiren iterazioan:

$$x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right) / a_{ii}. \quad (5.32)$$

### 5.10.2. Gauss-Seidel-en iterazioa

Jacobiren metodoaren konbergentzia azeleratzeko asmoz, beste metodo hau dugu. Aurreko 5.3. adibideko  $\{x^{(k)}\}$ ,  $\{y^{(k)}\}$ ,  $\{z^{(k)}\}$  segidek 2ra, 4ra eta 3ra jotzen dute, hurrenez hurren. Izan ere,  $x^{(k+1)}$  seguraski  $x^{(k)}$  baino bere limitearen hurbilpen hobea denez, arrazoizkoa izango litzateke  $y^{(k+1)}$  kalkulatzeko  $x^{(k)}$ -ren ordez  $x^{(k+1)}$  erabiltzea. Hortaz, hobe izango litzateke  $z^{(k+1)}$  kalkuluan  $y^{(k)}$ -ren ordez  $y^{(k+1)}$  erabiltzea ere. Hori kontuan hartuz, honela geratzen da iterazioa:

$$\begin{aligned} x^{(k+1)} &= \frac{7 + y^{(k)} - z^{(k)}}{4} \\ y^{(k+1)} &= \frac{21 + 4x^{(k+1)} + z^{(k)}}{8} \\ z^{(k+1)} &= \frac{15 + 2x^{(k+1)} - y^{(k+1)}}{5}. \end{aligned} \quad (5.33)$$

Metodo horri *Gauss-Seidelen iterazioaren metodoa* deritzogu. Eta aldaketa horiek (5.32) ekuazioari eramanez, hau lortzen dugu  $i = 1, \dots, n$  guztietarako:

$$x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right) / a_{ii}. \quad (5.34)$$

Era matritzialean adierazita, honela da:

$$D\mathbf{x}^{(k+1)} = \mathbf{b} - L\mathbf{x}^{(k+1)} - U\mathbf{x}^{(k)} \quad (5.35)$$

edo  $(D + L)\mathbf{x}^{(k+1)} = \mathbf{b} - U\mathbf{x}^{(k)}$ .

Orain 5.3. adibidearen sistema ebartziko dugu Gauss-Seidelen metodoa erabiliz (hots, (5.33) ); orduan,  $y^{(0)} = 2$  eta  $z^{(0)} = 2$  ordezkatuz, hau dugu:

$$x^{(1)} = \frac{7 + 2 - 2}{4} = 1.75$$

$k$	$x^{(k)}$	$y^{(k)}$	$z^{(k)}$
0	1.0	2.0	2.0
1	1.75	3.75	2.95
2	1.95	3.96875	2.98625
3	1.995625	3.999609375	2.99903125
...	...	...	...
8	1.99999983	3.99999988	2.99999996
9	1.99999998	3.99999999	3.00000000
10	2.00000000	4.00000000	3.00000000

**5.2. taula.** Gauss-Seidelen iterazioaren konbergentzia 5.3. adibidean.

eta  $x^{(1)} = 1.75$  eta  $z^{(0)} = 2$  ordezkatur, zera lortzen da:

$$y^{(1)} = \frac{21 + 4 \cdot 1.75 + 2}{8} = 3.75.$$

Azkenik,  $x^{(1)} = 1.75$  eta  $y^{(1)} = 3.75$  ordezkatur, hau lortzen dugu:

$$z^{(1)} = \frac{15 + 2 \cdot 1.75 - 3.75}{5} = 2.95.$$

Emaitza hau aurrekoa baino hurbilago dago benetako soluziotik,  $(x, y, z)^T = (2, 4, 3)^T$ .

Prozesu horri jarraituz, 5.2. taula lortzen da.

### 5.10.3. Metodo egonkorren konbergentzia

Jacobiren eta Gauss-Seidelen metodoak puntu finkoko metodo bezala ikus ditzakegu (*erlaxazio-metodoak* ere esaten zaie).

Izan bedi  $\mathbf{A} = \mathbf{M} - \mathbf{N}$  gure matrizearen erdibitze bat; hots,  $\mathbf{N} = \mathbf{M} - \mathbf{A}$  dugu. Beraz,  $\mathbf{Ax} = (\mathbf{M} - \mathbf{N})\mathbf{x} = \mathbf{b}$ , eta  $\mathbf{Mx} = \mathbf{Nx} + \mathbf{b}$  idatz dezakegu. Horrek puntu finkoko iterazio honetara eramaten gaitu:

$$\begin{aligned}
 \mathbf{x}_{k+1} &= \mathbf{M}^{-1}\mathbf{Nx}_k + \mathbf{M}^{-1}\mathbf{b} \\
 &= \mathbf{M}^{-1}(\mathbf{M} - \mathbf{A})\mathbf{x}_k + \mathbf{M}^{-1}\mathbf{b} \\
 &= \mathbf{x}_k - \mathbf{M}^{-1}\mathbf{Ax}_k + \mathbf{M}^{-1}\mathbf{b} \\
 &= \mathbf{x}_k + \mathbf{M}^{-1}(\mathbf{b} - \mathbf{Ax}_k)
 \end{aligned} \tag{5.36}$$

eta hori  $\mathbf{x}_{k+1} = \mathbf{g}(\mathbf{x}_k)$  puntu finkoko iterazioa da, non  $\mathbf{g}(\mathbf{x}) = \mathbf{x} + \mathbf{M}^{-1}(\mathbf{b} - \mathbf{Ax})$ . Metodo horri egonkorra deritzogu; funtzio hori definitzen duten  $\mathbf{A}$ ,  $\mathbf{M}$  eta  $\mathbf{b}$  bektoreak konstante baitira, ez daude iterazioaren menpe.

Jarraian ikusiko dugu nola  $\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k$  erroreak jokatzen duen adierazpen orokor horretan, eta nola errorea erlazionatuta dagoen  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$  hondarrarekin.

Errorearen adierazpena  $\mathbf{A}$ -rekin biderkatuz zera dugu:

$$\mathbf{A}\mathbf{e}_k = \mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k = \mathbf{r}_k \Rightarrow \underline{\mathbf{e}_k = \mathbf{A}^{-1}\mathbf{r}_k}$$

eta ondorioz

$$\mathbf{x} = \mathbf{x}_k + \mathbf{e}_k = \mathbf{x}_k + \mathbf{A}^{-1}\mathbf{r}_k,$$

baina kalkulu horren zailtasuna eta  $\mathbf{A}\mathbf{x} = \mathbf{b}$  sistema ebaztekoa berdina da. Hori dela eta,  $\mathbf{A}\mathbf{e}_k = \mathbf{r}_k$  sistema ebatzi beharrean,  $\mathbf{e}_k$  errorea hurbilduko dugu  $\mathbf{M}\mathbf{p}_k = \mathbf{r}_k$  sistema ebatziz, non  $\mathbf{p}_k \approx \mathbf{e}_k$ . Beraz, gure iterazioa aurkitzeko  $k$  bakoitzerako  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$  hondarra kalkulatu dugu, eta gero hau kalkulatu dugu:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{M}^{-1}\mathbf{r}_k,$$

hots, goiko metodo finkoko iterazioa.

Orain, ikus dezagun Jacobiren eta Gauss-Seidelen metodoak mota horretakoak direla.

Baldin  $\mathbf{M} = \mathbf{D}$  hartzen badugu, zera dugu:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k + \mathbf{D}^{-1}(\mathbf{b} - \mathbf{A}\mathbf{x}_k) \\ &= \mathbf{D}^{-1}\mathbf{D}\mathbf{x}_k + \mathbf{D}^{-1}(\mathbf{b} - \mathbf{A}\mathbf{x}_k) \\ &= \mathbf{D}^{-1}(\mathbf{b} - (\mathbf{A} - \mathbf{D})\mathbf{x}_k) \\ &= \mathbf{D}^{-1}(\mathbf{b} - (\mathbf{L} + \mathbf{U})\mathbf{x}_k) \end{aligned}$$

eta hori Jacobiren iterazioa da; ikus (5.31).

Aldiz,  $\mathbf{M} = \mathbf{L} + \mathbf{D}$  hartzen badugu, zera dugu:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k + (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{b} - \mathbf{A}\mathbf{x}_k) \\ &= (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{L} + \mathbf{D})\mathbf{x}_k + (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{b} - \mathbf{A}\mathbf{x}_k) \\ &= (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{b} - (\mathbf{A} - (\mathbf{L} + \mathbf{D}))\mathbf{x}_k) \\ &= (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{b} - \mathbf{U}\mathbf{x}_k) \end{aligned}$$

eta hori Gauss-Seidelen iterazioa da; ikus (5.35).

## Konbergentziaren analisia

Gogora dezagun  $\mathbf{x}_k$  iterazio bakoitzean  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$  hondarra ezaguna dela; aldiz,  $\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k = \mathbf{A}^{-1}\mathbf{r}_k$  errorea ezezaguna da. Jakin nahi dugu noiz gertatzen den  $\mathbf{e}_k \rightarrow \mathbf{0}$ ,  $k \rightarrow \infty$  denean. (5.36) erabiliz eta adierazpen hauek kenduz:

$$\begin{aligned} \mathbf{x}_k &= \mathbf{x}_{k-1} - \mathbf{M}^{-1}\mathbf{A}\mathbf{x}_{k-1} + \mathbf{M}^{-1}\mathbf{b} = (\mathbb{1} - \mathbf{M}^{-1}\mathbf{A})\mathbf{x}_{k-1} + \mathbf{M}^{-1}\mathbf{b} \\ \mathbf{x} &= \mathbf{x} - \mathbf{M}^{-1}\mathbf{A}\mathbf{x} + \mathbf{M}^{-1}\mathbf{b} = (\mathbb{1} - \mathbf{M}^{-1}\mathbf{A})\mathbf{x} + \mathbf{M}^{-1}\mathbf{b}, \end{aligned}$$

$\mathbf{x} - \mathbf{x}_k = (\mathbb{1} - \mathbf{M}^{-1}\mathbf{A})(\mathbf{x} - \mathbf{x}_{k-1})$  lortu dugu. Orain,  $\mathbf{T} = \mathbb{1} - \mathbf{M}^{-1}\mathbf{A}$  izendatuz, zera dugu:

$$\mathbf{x} - \mathbf{x}_k = \mathbf{T}(\mathbf{x} - \mathbf{x}_{k-1}),$$

eta, ondorioz,  $\mathbf{e}_k = \mathbf{T}\mathbf{e}_{k-1} = \mathbf{T}(\mathbf{T}\mathbf{e}_{k-2}) = \dots = \mathbf{T}^k\mathbf{e}_0$  betetzen da.  $\mathbf{T}$  matrizeari *iterazio-matrizea* deritzogu.  $\mathbf{e}_0$  hasierako errorea denez, konbergentzia (hots,  $\mathbf{e}_k \rightarrow \mathbf{0}$ ) izango dugu baldin eta soilik baldin  $\mathbf{T}^k \rightarrow \mathbf{0}$ . Jakina, hori gertatzen da eragindako matrize-norma baterako hau betetzen bada:

$$\|\mathbf{T}\| < 1,$$

zeren

$$\|\mathbf{e}_k\| = \|\mathbf{T}\mathbf{T}\dots\mathbf{T}\mathbf{e}_0\| \leq \|\mathbf{T}\|\|\mathbf{T}\|\dots\|\mathbf{T}\|\|\mathbf{e}_0\| = \|\mathbf{T}\|^k\|\mathbf{e}_0\|.$$

Puntu finkoaren metodoan dugun konbergentzia-baldintzaren antzekoa da. Aurrekoa kontuan hartuz, teorema hau ondorioztatzen da.

### 5.11. teorema: metodo egonkorren konbergentzia.

$\mathbf{Ax} = \mathbf{b}$  sistema linealerako, izan bedi metodo iteratibo hau:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{M}^{-1}\mathbf{r}_k, \quad k = 0, 1, \dots,$$

non  $\mathbf{r}_k = \mathbf{b} - \mathbf{Ax}_k$  hondar-bektorea eta  $\mathbf{T} = \mathbb{1} - \mathbf{M}^{-1}\mathbf{A}$  iterazio-matrizea definitzen baititugu.

Orduan, metodoa konbergentea da baldin eta soilik baldin iterazio-matrizearen espektro-erradioak hau betetzen badu:

$$\rho(\mathbf{T}) < 1.$$

Zenbat eta  $\rho(\mathbf{T})$  txikiago den, hainbat eta azkarrago konbergitzen da.

*Frogantza.* Lehendabizi,  $\rho(\mathbf{T}) < 1 \Rightarrow \mathbf{e}_k \rightarrow 0$  betetzen dela frogatuko dugu.

Izan bitez  $\mathbf{v}_1, \dots, \mathbf{v}_n$  bektore linealki independenteak  $\mathbf{T}$  matrizearen autobektoreak. Orduan, hau idatz dezakegu:

$$\mathbf{e}_0 = \sum_{i=1}^n \alpha_i \mathbf{v}_i.$$

Izan bitez  $\mathbf{T}$ -ren  $\mathbf{v}_i$  autobektoreei dagozkien  $\lambda_i$  autobalioak. Baldin  $\rho(\mathbf{T}) < 1$  bada, zera dugu:

$$\mathbf{T}\mathbf{v}_i = \lambda_i \mathbf{v}_i, \quad |\lambda_i| < 1, \quad i = 1, \dots, n,$$

ondorioz,

$$\mathbf{e}_k = \mathbf{T}^k \mathbf{e}_0 = \sum_{i=1}^n \alpha_i \mathbf{T}^k \mathbf{v}_i = \sum_{i=1}^n \alpha_i \lambda_i^k \mathbf{v}_i \rightarrow \mathbf{0} \quad (k \rightarrow \infty).$$

Orain  $\mathbf{e}_k \rightarrow 0 \Rightarrow \rho(\mathbf{T}) < 1$  frogatuko dugu absurdora eramanez. Baldin  $\rho(\mathbf{T}) \geq 1$  bada, badago  $i$  bat non  $|\lambda_i| \geq 1$  betetzen baita. Izan bedi  $\mathbf{v}_i$  dagokion autobektorea. Hortaz,



$\mathbf{x}_0 = \mathbf{x} - \mathbf{v}_i$  hartzen badugu,  $\mathbf{e}_k = \mathbf{T}^k \mathbf{e}_0 = \mathbf{T}^k(\mathbf{x} - \mathbf{x}_0) = \mathbf{T}^k \mathbf{v}_i = \lambda_i^k \mathbf{v}_i \not\rightarrow \mathbf{0}$  da eta ez da konbergentea.  $\square$

Jakina,  $\|\mathbf{T}\| < 1$  bada,  $\rho(\mathbf{T}) < 1$  da (ikus 5.7.1. atala). Hala ere, espektro-erradioaren gaineko baldintzak gehiago esaten digu  $\mathbf{T}$ -ren normaren baldintzak baino. Zeren eta  $\rho(\mathbf{T}) < 1$  baldintza beharrezkoa eta nahikoa baita, eta, aldiz,  $\|\mathbf{T}\| < 1$  baldintza nahikoa bakarrik da.

### Zenbat iterazio behar dugu errorea 10 bider txikiagoa egiteko?

Beste hitz batzuetan,  $0.1\|\mathbf{e}_0\| = \|\mathbf{e}_k\| \approx \rho(\mathbf{T})^k \|\mathbf{e}_0\|$  gertatzeko, zein izan behar du  $k$  iterazio kopuruak?

Bi ataletan  $\log_{10}$  hartuz eta  $k$  bakanduz, hau lortzen da:

$$k \approx -\frac{1}{\log_{10} \rho(\mathbf{T})}.$$

*Konbergentzia-ratioa* honela definitzen da:  $\text{ratioa} = -\log_{10} \rho(\mathbf{T})$ . Orduan  $k \approx 1/\text{ratioa}$ .

Ondorioz, zenbat eta txikiago den espektro-erradioa, hainbat eta handiagoa da ratioa eta, beraz, iterazio gutxiago behar izango dira errore-txikitze maila berdina lortzeko.

### Hertsiki diagonal menperatzailea

Izan bedi  $\mathbf{Ax} = \mathbf{b}$  sistema. Aurreko metodo iteratiboak konbergenteak izateko, nahikoa da (5.12) baldintza betetzea, hots, hertsiki diagonal menperatzailea izatea. Jarraian frogatuko dugu.

#### 5.12. teorema.

*Izan bedi  $\mathbf{Ax} = \mathbf{b}$  sistema. Jacobiren eta Gauss-Seidelen metodo iteratiboak konbergenteak izateko, nahikoa da  $\mathbf{A}$  hertsiki diagonal menperatzailea izatea.*

*Frogantza.*

Jacobirena. Metodo honetan hau dugu:

$$\mathbf{T} = \mathbb{1} - \mathbf{D}^{-1}\mathbf{A} = \mathbf{D}^{-1}\mathbf{D} - \mathbf{D}^{-1}\mathbf{A} = \mathbf{D}^{-1}(\mathbf{D} - \mathbf{A}) = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U}),$$

eta  $\rho(\mathbf{T}) \leq \|\mathbf{T}\|$  denez, zera lortzen da:

$$\rho(\mathbf{T}) \leq \|\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\|_{\infty} = \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1,$$

Izan ere,  $\mathbf{A}$  hertsiki diagonal menperatzailea izateagatik azken desberdintza betetzen da. Beraz,  $\rho(\mathbf{T}) < 1$  denez, metodo hau konbergentea da.

Gauss-Seidelena. Metodo honetan hau dugu:

$$\begin{aligned}\mathbf{T} &= \mathbf{1} - (\mathbf{L} + \mathbf{D})^{-1}\mathbf{A} = (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{L} + \mathbf{D}) - (\mathbf{L} + \mathbf{D})^{-1}\mathbf{A} \\ &= (\mathbf{L} + \mathbf{D})^{-1}((\mathbf{L} + \mathbf{D}) - \mathbf{A}) = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}.\end{aligned}$$

Izan bedi

$$\gamma = \max_{1 \leq i \leq n} \frac{\sum_{j>i} |a_{ij}|}{|a_{ii}| - \sum_{j<i} |a_{ij}|}.$$

Orain,  $\mathbf{A}$  hertsiki diagonal menperatzailea denez,  $i$  guztietarako hau dugu:

$$|a_{ii}| - \sum_{j<i} |a_{ij}| > \sum_{j>i} |a_{ij}|,$$

ondorioz  $\gamma < 1$ . Jarraian,  $\|\mathbf{T}\|_\infty \leq \gamma$  frogatuko dugu.

Demagun  $\|\mathbf{x}\|_\infty = 1$ . Izan bedi  $\mathbf{y} = \mathbf{T}\mathbf{x}$ , hots,  $(\mathbf{D} + \mathbf{L})\mathbf{y} = -\mathbf{U}\mathbf{x}$ . Sistema horren  $i$ . ekuazioa,  $i = 1, \dots, n$ , hau da:

$$a_{ii}y_i + \sum_{j<i} a_{ij}y_j = -\sum_{j>i} a_{ij}x_j. \quad (5.37)$$

Izan bedi  $i$  non  $\|\mathbf{y}\|_\infty = |y_i|$ . Orduan (5.37)-en ezker aldearen balio absolutua adierazpen honek behetik bornatzen du ( $|a| - |b| \leq |a + b|$  baita):

$$\begin{aligned}|a_{ii}y_i| - \left| \sum_{j<i} a_{ij}y_j \right| &\geq |a_{ii}| |y_i| - \sum_{j<i} |a_{ij}| |y_j| \\ &\geq |a_{ii}| |y_i| - \sum_{j<i} |a_{ij}| |y_i| \\ &= (|a_{ii}| - \sum_{j<i} |a_{ij}|) \|\mathbf{y}\|_\infty.\end{aligned}$$

Era antzeko batean, (5.37)-en eskuin aldearen balio absolutua adierazpen honek goitik bornatzen du:

$$\sum_{j>i} |a_{ij}| |x_j| \leq \sum_{j>i} |a_{ij}| \|\mathbf{x}\|_\infty = \sum_{j>i} |a_{ij}|.$$

Aurreko desberdintzak konbinatuz, zera lortzen da:

$$(|a_{ii}| - \sum_{j<i} |a_{ij}|) \|\mathbf{y}\|_\infty \leq \sum_{j>i} |a_{ij}|,$$

eta orduan:

$$\|\mathbf{T}\|_\infty = \max_{\|\mathbf{x}\|=1} \|\mathbf{T}\mathbf{x}\|_\infty = \max_{\|\mathbf{y}\|=1} \|\mathbf{y}\|_\infty \leq \max_{1 \leq i \leq n} \frac{\sum_{j>i} |a_{ij}|}{|a_{ii}| - \sum_{j<i} |a_{ij}|} = \gamma < 1. \quad \square$$

Kontuan izan baldintza hori nahikoa dela, ez beharrezkoa. Beraz, hertsiki diagonal menperatzailea ez izan arren, metodo iteratibo horiek konbergenteak izan daitezke  $\rho(\mathbf{T}) < 1$  betetzen bada.

## 5.11. Problemak

Eskuz ebazteko problemak:

1. Faktorizatu matrize hauek  $LU$  deskonposizioaren bidez eta pibotatze partziala erabiliz:

$$a) \begin{bmatrix} 1 & 3 & 2 \\ 1 & 5 & 3 \\ 2 & 4 & -6 \end{bmatrix} \quad b) \begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{bmatrix}.$$

$$c) \begin{bmatrix} -5 & 2 & -1 \\ 1 & 0 & 3 \\ 3 & 1 & 6 \end{bmatrix} \quad d) \begin{bmatrix} 2 & -1 & 1 \\ 3 & 3 & 9 \\ 3 & 3 & 5 \end{bmatrix}.$$

$$e) \begin{bmatrix} 1 & 1 & 0 & 4 \\ 2 & -1 & 5 & 0 \\ 5 & 2 & 1 & 2 \\ -3 & 0 & 2 & 6 \end{bmatrix} \quad f) \begin{bmatrix} 4 & 8 & 4 & 0 \\ 1 & 5 & 4 & -3 \\ 1 & 4 & 7 & 2 \\ 1 & 3 & 0 & -2 \end{bmatrix}.$$

Zein dira matrize horien bat-norma eta infinitu-norma?

Aurkitu matrize horien baldintzazko zenbakiak bi norma horietarako (hots,  $\kappa_1$  eta  $\kappa_\infty$ ).

2. Ebatzi sistema lineal hauek  $LU$  deskonposizioaren bidez eta pibotatze partziala erabiliz:

$$a) \begin{aligned} x_1 + 3x_2 + 2x_3 &= 5 \\ x_1 + 5x_2 + 3x_3 &= 10 \\ 2x_1 + 4x_2 - 6x_3 &= -4. \end{aligned} \quad b) \begin{aligned} 2x_1 + x_2 + x_3 &= 3 \\ 4x_1 + 3x_2 + 3x_3 + x_4 &= 7 \\ 8x_1 + 7x_2 + 9x_3 + 5x_4 &= 17 \\ 6x_1 + 7x_2 + 9x_3 + 8x_4 &= 15. \end{aligned}$$

$$c) \begin{aligned} -5x_1 + 2x_2 - x_3 &= 2 \\ x_1 + 3x_3 &= 2 \\ 3x_1 + x_2 + 6x_3 &= 2. \end{aligned} \quad d) \begin{aligned} 2x_1 - x_2 + x_3 &= -1 \\ 3x_1 + 3x_2 + 9x_3 &= 0 \\ 3x_1 + 3x_2 + 5x_3 &= 4. \end{aligned}$$

$$e) \begin{aligned} x_1 + x_2 + 4x_4 &= 5 \\ 2x_1 - 1x_2 + 5x_3 &= -6 \\ 5x_1 + 2x_2 + x_3 + 2x_4 &= 3 \\ -3x_1 + 2x_3 + 6x_4 &= 4. \end{aligned} \quad f) \begin{aligned} 4x_1 + 8x_2 + 4x_3 &= 8 \\ x_1 + 5x_2 + 4x_3 - 3x_4 &= -4 \\ x_1 + 4x_2 + 7x_3 + 2x_4 &= 10 \\ x_1 + 3x_2 - 2x_4 &= -4. \end{aligned}$$

Kalkula ezazu sistema guztietarako  $\kappa(A)_\infty$ -ren behe-borne bat.

3. Izan bedi  $\mathbf{U}\mathbf{x} = \mathbf{b}$  sistema, non  $\mathbf{U} \in \mathbb{R}^{n \times n}$  matrize goi-triangeluarra eta  $\mathbf{b} \in \mathbb{R}^{n \times 1}$  baitira. Zenbat eragiketa aritmetiko behar ditugu sistema hori ebazteko?

4. (a) Izan bedi  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizea. Zenbat biderketa/zatiketa behar ditu pibotatzerik gabeko  $LU$  deskonposizioak? (*Iradokizuna*:  $1^2 + 2^2 + \dots + m^2 = m(m+1)(2m+1)/6$  berdintza erabili, zeina indukzioz frogatu baitaiteke).
- (b) Zenbat biderketa/zatiketa behar ditu aurreranzko ordezkapenak? Eta atzeranzko ordezkapenak?
- (c) Sistema lineal bat pibotatzerik gabeko  $LU$  metodoaz ebazteko, zenbat biderketa/zatiketa egin behar ditugu?
5. (a) Izan bedi  $\mathbf{U}$  matrize ez-singular goi-triangeluarra. Frogatu  $\mathbf{U}^{-1}$ -en diagonaleko gaiak  $\mathbf{U}$ -ren diagonaleko gaien erreziprokoak direla.
- (b) Aurreko atalaren emaitza erabiliz, frogatu hau betetzen dela:

$$\|\mathbf{U}\|_{\infty} \geq \max_i |u_{ii}| \quad \text{eta} \quad \|\mathbf{U}^{-1}\|_{\infty} \geq \frac{1}{\min_i |u_{ii}|}.$$

Infinitu-normako bi desberdintza horiek hau inplikatzeko dute:

$$\kappa_{\infty}(\mathbf{U}) \geq \frac{\max_i |u_{ii}|}{\min_i |u_{ii}|}.$$

Desberdintza horren eskuin aldeko kantitatea sarri erabiltzen dugu matrize goi-triangeluar baten baldintzazko zenbakiaren hurbilpena bezala.

6. Faktorizatu matrize simetriko hauek,  $\mathbf{R}^T \mathbf{R}$  eran, Choleskyren metodoa erabiliz:

$$a) \begin{bmatrix} 4 & -2 & 2 \\ -2 & 2 & -1 \\ 2 & -1 & 10 \end{bmatrix}. \quad b) \begin{bmatrix} 4 & 2 & -2 \\ 2 & 10 & 2 \\ -2 & 2 & 3 \end{bmatrix}.$$

$$c) \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}. \quad d) \begin{bmatrix} 4 & 1 & 1 & 1 \\ 1 & 3 & -1 & 1 \\ 1 & -1 & 2 & 0 \\ 1 & 1 & 0 & 2 \end{bmatrix}.$$

$$e) \begin{bmatrix} 6 & 2 & 1 & -1 \\ 2 & 4 & 1 & 0 \\ 1 & 1 & 4 & -1 \\ -1 & 0 & -1 & 3 \end{bmatrix}. \quad f) \begin{bmatrix} 4 & 1 & -1 & 0 \\ 1 & 3 & -1 & 0 \\ -1 & -1 & 5 & 2 \\ 0 & 0 & 2 & 4 \end{bmatrix}.$$

Zein kasutan da matrizea hertsiki diagonal menperatzaile? Hori gertatzen denean, nolakoa da matrizea? Kalkulatu matrize horietarako baldintzazko zenbakiaren beheborne bat.

7. Aurkitu  $a$ -ren eta  $b$ -ren balio guztiak, matrize hau simetriko definitu positiboa izateko:

$$\begin{bmatrix} a & 1 & 1+b \\ 1 & a & 1 \\ 1-b^2 & 1 & a \end{bmatrix}.$$

8. Sistema hauetarako, aurkitu Jacobiren eta Gauss-Seidelen lehenengo hiru iterazioak,  $\mathbf{x}^{(0)} = \mathbf{0}$  erabiliz:

$$\begin{array}{ll}
 a) & \begin{array}{l} -x_1 + 3x_2 = 1 \\ 6x_1 - 2x_2 = 2. \end{array} \\
 b) & \begin{array}{l} 5x_1 - x_2 + x_3 = 10 \\ 2x_1 + 8x_2 - x_3 = 11 \\ -x_1 + x_2 + 4x_3 = 3. \end{array} \\
 c) & \begin{array}{l} 3x_1 - 0.1x_2 - 0.2x_3 = 7.85 \\ 0.1x_1 + 7x_2 - 0.3x_3 = -19.3 \\ 0.3x_1 - 0.2x_2 + 10x_3 = 71.4. \end{array} \\
 d) & \begin{array}{l} 2x_1 - x_2 + x_3 = -1 \\ 3x_1 + 3x_2 + 9x_3 = 0 \\ 3x_1 + 3x_2 + 5x_3 = 4. \end{array} \\
 e) & \begin{array}{l} 2x_1 - x_2 + 10x_3 = -11 \\ 3x_2 - x_3 + 8x_4 = -11 \\ 10x_1 - x_2 + 2x_3 = 6 \\ -x_1 + 11x_2 - x_3 + 3x_4 = 25. \end{array} \\
 f) & \begin{array}{l} 4x_1 - 2x_2 = 0 \\ -2x_1 + 5x_2 - x_3 = 2 \\ -x_2 + 4x_3 + 2x_4 = 3 \\ 2x_3 + 3x_4 = -2. \end{array}
 \end{array}$$

Beste galdera batzuk:

- i) Zein da metodo bakoitzeko zehaztasun erlatiboa hirugarren iterazioan? Metodo horiek konbergenteak badira, zein baliotara konbergitzen dira? Egiaztatu hori dela sistemaren soluzioa sisteman ordezkatzuz.
- ii) Aztertu sistema bakoitzean iterazioekin hasi baino lehen jakin dezakegun Jacobiren metodoa konbergentea den ala dibergentea. Zergatik?
- iii) Ez bada konbergente, nola bihur dezakegu konbergente (ahal bada)? Kasu horretan, ebatzi berriro sistema hori hasierako puntu berdina erabiliz.
- iv) Zein metodok egiten du konbergentzia azkarren (egiten badu)? Justifikatu erantzuna konbergentziari buruzko teoria erabiliz.
- v) Zein da problema bakoitzean Jacobiren metodoaren konbergentzia-ratioa?
- vi) Zenbat iterazio gehienez behar ditu Jacobiren metodoak hasierako errorea  $10^4$  bider txikiago egiteko? Erantzuna konbergentziaren teoreman oinarritu behar da.

$$\left( \text{Oharra: } k. \text{ iterazioko zehaztasun erlatiboa} = \frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_{\infty}}{\|\mathbf{x}^{(k)}\|_{\infty}} \right).$$

9. Izan bedi matrize hau:

$$\mathbf{A} = \begin{bmatrix} 1 & a & a \\ a & 1 & a \\ a & a & 1 \end{bmatrix}.$$

Aurkitu  $a$ -ren balioak  $\mathbf{A}$  simetriko definitu positiboa izan dadin, baina Jacobiren iterazioa ez konbergitzeko.

10. Froga ezazu  $\mathbf{A}$   $2 \times 2$ -matrize simetriko definitu positiboa bada, Jacobiren metodoa konbergentea dela edozein hasierako baliotarako.

**MATLABez ebazteko problemak:**

11. Izan bedi matrize hau:

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & -1 & 3 & 5 \\ 0 & 0 & 2 & 5 \\ -2 & -6 & -3 & 1 \end{bmatrix}$$

eta izan bitez matrize-norma hauek:

- $\|\mathbf{A}\|_1 = \max_j \|\mathbf{A}_{:,j}\|_1$  (zutabe guztien bat-normetako maximoa)
  - $\|\mathbf{A}\|_\infty = \max_i \|\mathbf{A}_{i,:}\|_1$  (lerro guztien bat-normetako maximoa)
- (a) Aztertu **eskuz**  $\mathbf{A}$  matrize honen singularitasuna, eta kalkulatu bat-norma eta infinitu-norma.
  - (b) Kalkulatu  $\mathbf{A}^{-1}$  MATLABeko `inv(A)` funtzioa erabiliz, eta kalkulatu matrize horren bat-norma eta infinitu-norma.
  - (c) Kalkulatu **eskuz**  $\mathbf{A}$  matrizearen baldintzazko zenbakia, hau da:  $\kappa(\mathbf{A}) = \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\|$  bi norma horietarako.
  - (d) Egiaztatu lortutako emaitzak MATLABeko `cond(A, 1)` eta `cond(A, inf)` funtzioak erabiliz.
  - (e) Eranskinetik kopiatu kodea MATLABeko `lup.m` funtzio-fitxategia sortzeko, eta kalkulatu  $\mathbf{A}$  matrizearen  $\mathbf{L}$ ,  $\mathbf{U}$  eta  $\mathbf{P}$  matrizeak funtzio hori erabiliz.
  - (f) Kalkulatutako matrize horiek erabiliz (ikus (5.10) - (5.11) adierazpenak), ebatzi **eskuz** sistema hau:

$$\begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & -1 & 3 & 5 \\ 0 & 0 & 2 & 5 \\ -2 & -6 & -3 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}.$$

12. Izan bedi  $\mathbf{L}\mathbf{y} = \mathbf{P}\mathbf{b}$  sistema. Sortu MATLABeko funtzio bat `aurrerantz.m` izenekoa, adierazpen hauek kontuan hartuz:

$$x_1 = b_1$$

$$x_i = b_i - \sum_{j=1}^{i-1} l_{ij}x_j, \quad i = 2, 3, \dots, n-1, n,$$

ikus (5.5) - (5.6) adierazpenak.

Egiaztatu aurreko problemako  $\mathbf{L}\mathbf{y} = \mathbf{P}\mathbf{b}$  sistema ondo ebatzi duzula funtzio hori erabiliz.

13. Izan bedi  $\mathbf{U}\mathbf{x} = \mathbf{y}$  sistema. Sortu MATLABeko funtzio bat `atzerantz.m` izenekoa, adierazpen hauek kontuan hartuz:

$$x_n = \frac{b_n}{u_{nn}}$$

$$x_i = \frac{b_i - \sum_{j=i+1}^n u_{ij}x_j}{u_{ii}}, \quad i = n-1, n-2, \dots, 3, 2, 1,$$

ikus (5.2) - (5.3) adierazpenak.

Egiaztatu aurreko problemako  $\mathbf{U}\mathbf{x} = \mathbf{y}$  sistema ondo ebatzi duzula funtzio hori erabiliz.

14. Sortu M fitxategi bat,  $\mathbf{A}\mathbf{x} = \mathbf{b}$  sistema  $LU$  faktORIZAZIOAZ ebazteko gai izan dadin. Programa horrek sistema hori ebazteko ondoz ondoko lan hauek egin behar ditu:

- (a) Programa horrek goiburu hau izan behar du:  
`function x=ebatzi(A,b).`
- (b) Egiaztatu  $\mathbf{A}$  ez dela singularra, hots  $\det(\mathbf{A}) \neq 0$ . Kalkulu hori egiteko, MATLABek `det` funtzioa du. Baldin singularra bada, mezu bat eman ez bukatuko da.
- (c)  $\mathbf{LU} = \mathbf{PA}$  faktORIZAZIOA egin behar du. Lan hori `lup` funtzioak egin dezake.
- (d)  $\mathbf{Pb}$  kalkulatu behar du. Lan hori egiteko modu errazena `b(p)` idaztea da (MATLABi esker), non `p` permutazio-bektorea baita (`lup` funtzioak kalkulatzen du).
- (e)  $\mathbf{LU}\mathbf{x} = \mathbf{Pb}$  sistema ebatzi behar du, urrats hauei jarraituz:  
 (L) Ebatzi  $\mathbf{Ly} = \mathbf{Pb}$ . Lan hori `aurrerantz` funtzioak egin dezake.  
 (U) Ebatzi  $\mathbf{U}\mathbf{x} = \mathbf{y}$ . Lan hori `atzerantz` funtzioak egin dezake.

Egiaztatu `ebatzi.m` funtzioa ondo dabilela aurreko sistema ebatziz (ez ahaztu M fitxategiak izen hori izan behar duela).

15. Izan bedi sistema hau:

$$\begin{array}{rcccc} 4x_1 & -2x_2 & & & = & 0 \\ -2x_1 & +5x_2 & -x_3 & & = & 2 \\ & -x_2 & +4x_3 & +2x_4 & = & 3 \\ & & 2x_3 & +3x_4 & = & -2. \end{array}$$

- (a) Garatu Jacobiren metodorako M fitxategi bat, `jacobi.m` izenekoa. Programa horrek gai izan behar du  $\mathbf{A}\mathbf{x} = \mathbf{b}$  sistema lineala ebazteko, non  $\mathbf{A} \in \mathbb{R}^{n \times n}$  hertsiki diagonal menperatzailea baita. Programa horrek goiburu hau izan behar du:

`function x=jacobi(A,b,x0,ze,imax),`

non  $\mathbf{x}$  soluzioaren hurbilpen bat baita,  $\mathbf{x}_0$  hasierako hurbilpen bat, eta  $ze$  eta  $imax$  eskatutako zehaztasun erlatiboa eta iterazioen kopuru maximoa baitira, hurrenez hurren. Programaren exekuzioa bukatuko da,  $k$ -garren iterazioan lortutako zehaztasun erlatiboak hau betetzen duenean:

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_{\infty}}{\|\mathbf{x}^{(k)}\|_{\infty}} \leq ze,$$

edo  $k = imax$  denean.

- (a) Asmatu zure kabuz sistema bat diagonal menperatzailekoa, eta ebatzi egindako programa erabiliz. Egiaztatu emaitza.
  - (b) Emandako sistema Jacobiren metodoaz ebatzi gabe, esan konbergituko den ala ez. Eta 11. problemako (f) atalekoa? Arrazoitu erantzunak.
  - (c) Gutxi gorabehera, zenbat iterazio erabiliko ditu hasierako errorea mila bider txikiagoa egiteko?
  - (d) Ebatzi emandako sistema kode hori erabiliz, eta egiaztatu aurreko analisia betetzen dela.
- (b) Garatu Gauss-Seidelen metodorako MATLABeko M fitxategi bat `gaussseidel.m` izeneko. Programa horrek gai izan behar du  $\mathbf{Ax} = \mathbf{b}$  sistema lineala ebazteko, non  $\mathbf{A} \in \mathbb{R}^{n \times n}$  hertsiki diagonal menperatzailea baita. Programa horrek goiburu hau izan behar du:

```
function x=gaussseidel(A,b,x0,ze,imax),
```

non  $\mathbf{x}$  soluzioaren hurbilpen bat baita,  $\mathbf{x}_0$  hasierako hurbilpen bat, eta  $ze$  eta  $imax$  eskatutako zehaztasun erlatiboa eta iterazioen kopuru maximoa baitira, hurrenez hurren. Programaren exekuzioa bukatuko da  $k$ -garren iterazioan lortutako zehaztasun erlatiboak hau betetzen duenean:

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_{\infty}}{\|\mathbf{x}^{(k)}\|_{\infty}} \leq ze,$$

edo  $k = imax$  denean.

- (a) Asmatu zure kabuz sistema bat diagonal menperatzailekoa, eta ebatzi egindako programa erabiliz. Egiaztatu emaitza.
- (b) Emandako sistema Gauss-Seidelen metodoaz ebatzi gabe, esan konbergituko den ala ez. Eta 11. problemako (f) atalekoa? Arrazoitu erantzunak.
- (c) Gutxi gorabehera zenbat iterazio erabiliko ditu hasierako errorea mila bider txikiagoa egiteko?
- (d) Ebatzi emandako sistema kode hori erabiliz, eta egiaztatu aurreko analisia betetzen dela.



## ERANSKINA

Kode hau  $LU = PA$  faktORIZAZIOA kalkulatzeko erabil daiteke:

```
function [L,U,p] = lup(A)
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% LUP  Triangelu faktORIZAZIOA L*U = PA bete dadin.
%   SARRERAK:
%   A deskonposatu nahi dugun matrizea,
%   EMAITZAK:
%   L matrize behe-triangeluarra,
%   U matrize goi-triangeluarra eta
%   p=permutazio-bektorea.
%
%   (P=permutazio-matrizea eman dezake goiburuan p P-rekin trukaturaz.)
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
[n,n] = size(A);
p = (1:n)';

for k = 1:n-1

    % Pibotatze partziala egiten du.
    % Aurkitu diagonaleko (k,k) gaitik behera dagoen tamaina
    % handieneko gaia (pibota) eta posizioa (m) azpibektore horretan.
    [pibota,m] = max(abs(A(k:n,k)));

    % pibotaren lerroa A matrizean:
    m = m+k-1;

    % Jauzi ezabapena zutabea 0 bada.
    if (A(m,k) ~= 0)

        % Trukatu m eta k lerroak, eta eguneratu permutazio-bektorea.
        if (m ~= k)
            A([k,m], :) = A([m,k], :);
            p([k,m]) = p([m,k]);
        end

        % Biderkatzaileak kalkulatzeko ditugu, eta A matrizearen diagonal
        % azpian gordetzen ditugu. Gero L matrizeko diagonal azpian jarriko
        % ditugu.
        i = k+1:n;
        A(i,k) = A(i,k)/A(k,k);

        % Eguneratu matrizearen gainerakoa. Gero, U matrize
        % goi-triangeluarrean gordeko dugu.
        j = k+1:n;
        A(i,j) = A(i,j) - A(i,k)*A(k,j);
    end
end

% L, U eta P matrizeak idazten dira;
```

```
% tril(A,-1): A-ren diagonal azpikoa gordetzen du, bestea zero da.  
% triu(A): A-ren diagonal eta bere goiko aldea gordetzen ditu, bestea  
%          zero da.  
L = tril(A,-1) + eye(n,n);  
U = triu(A);  
P=zeros(n,n);  
for i=1:n  
    P(i,p(i))=1;  
end
```

## 6. kapitulua

# $QR$ faktORIZAZIOA eta minimo karratu linealak

Kapitulu honetan problema hau ebatziko dugu:

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2,$$

non  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^m$  eta  $m \geq n$ . Demagun  $\mathbf{A}$  zutabe hein betekoa dela. Kontuan hartu behar da kasu gaindeterminatuan,  $m > n$  denean, ez dela  $\mathbf{x}$  egoten  $\mathbf{Ax} = \mathbf{b}$  sistema zehazki bete dezan, nahiz eta biribiltze-errorerik ez egon.

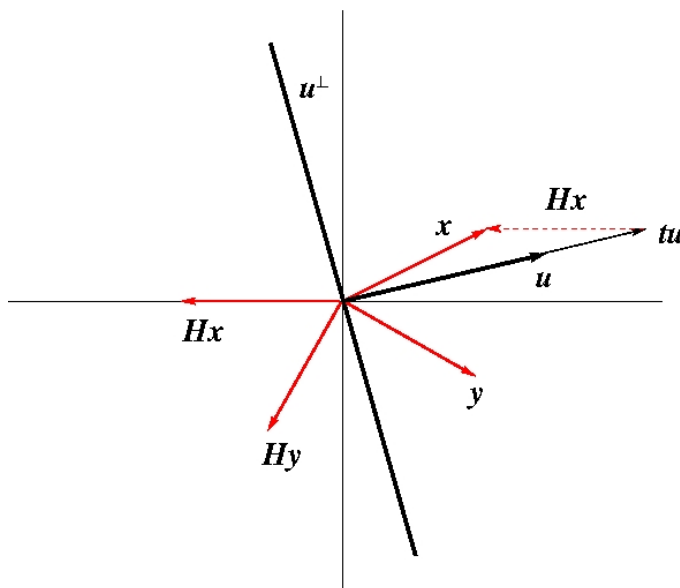
Minimo karratuen problema askotan agertzen da mundu errealeko aplikazioetan, bereziki datu-doitzea behar denean. Problema horri aurre egiteko, oso garrantzitsuak dira transformazio ortogonalak, eta horrelako tresnak azaltzen hasiko gara.

### 6.1. Householder-en islapenak

Householderren islapenak matrize-transformazioak dira, eta zenbakizko algoritmo moldagarri eta eraginkorrenetariko batzuen oinarria dira. Oso teknika ezaguna da matrize ortogonalen segida bat eraikitzeke,  $\mathbf{A}$  matrize baten diagonalaren pean dauden gaiak zero bihurtzeko moduan. Beraz,  $\mathbf{A}$  karratua bada, Householderren transformazioen bidez, matrize triangeluar bat ere lor dezakegu.

**6.1. Definizioa.** *Edozein  $\mathbf{u} \neq \mathbf{0}$  bektore bati dagokion **Householderren islapena** (edo **Householderren transformazioa** edo **Householderren matrizea**) itxura honetako matrizea da:*

$$\mathbf{H} = \mathbb{1} - \rho \mathbf{u} \mathbf{u}^T, \quad \text{non} \quad \rho = \frac{2}{\|\mathbf{u}\|_2^2}. \quad (6.1)$$



6.1. irudia. Householderren islapena.

$\mathbf{u}$  bektoreari **Householderren bektorea** deritzogu.

$\mathbf{H}$  matrizeak simetrikoak ( $\mathbf{H} = \mathbf{H}^T$ ) eta ortogonalak ( $\mathbf{H}^{-1} = \mathbf{H}^T$ ) dira (egiaztatu). Kontuan izan  $\mathbf{H}$  matrizea  $\mathbf{u}$  bektorearen menpe bakarrik dagoela.

Praktikan,  $\mathbf{H}$  ez da inoiz eratzten. Izan ere,  $\mathbf{H}$ -ren aplikazioa  $\mathbf{x}$  bektore baten gainean egiten dugunean, hau dugu:

$$\mathbf{H}\mathbf{x} = (\mathbb{1} - \rho\mathbf{u}\mathbf{u}^T)\mathbf{x} = \mathbf{x} - \rho\mathbf{u}(\mathbf{u}^T\mathbf{x})$$

eta, ondorioz, honela kalkulatzen da  $\mathbf{H}\mathbf{x}$ :

$$\begin{aligned} t &= \rho\mathbf{u}^T\mathbf{x} \\ \mathbf{H}\mathbf{x} &= \mathbf{x} - t\mathbf{u}. \end{aligned} \tag{6.2}$$

Geometrikoki,  $\mathbf{x}$  bektorea  $\mathbf{u}$  gainean proiektatzen da, eta, gero,  $\mathbf{x}$ -tik bi bider proiektzio hori kentzen da, zeren:

$$t\mathbf{u} = \rho(\mathbf{u}^T\mathbf{x})\mathbf{u} = 2 \left( \frac{\mathbf{u}^T\mathbf{x}}{\|\mathbf{u}\|_2} \right) \frac{\mathbf{u}}{\|\mathbf{u}\|_2}.$$

6.1. irudiak erakusten ditu  $\mathbf{u}$  bektore bat eta bere azpiespazio ortogonalak (irudian, lerrozuzen bat):  $\mathbf{u}^\perp$ . Hark erakusten ditu  $\mathbf{x}$  eta  $\mathbf{y}$  eta beren irudiak,  $\mathbf{H}\mathbf{x}$  eta  $\mathbf{H}\mathbf{y}$ .  $\mathbf{H}$  matrizeak edozein bektore bere islapen bihurtzen du  $\mathbf{u}^\perp$  lerroarekiko. Edozein  $\mathbf{x}$  bektoretarako, bektore honek ematen digu  $\mathbf{x}$ -ren eta  $\mathbf{H}\mathbf{x}$ -ren arteko erdigunea:

$$\mathbf{x} - \frac{t}{2}\mathbf{u}$$

eta  $\mathbf{u}^\perp$  lerroan (azpiespazioan) dago (aurrekoaren frogapena ariketa gisa geratzen da). Espazioak bi dimentsio baino gehiago dituenean, azpiespazio hori  $\mathbf{u}$  bektorearekiko plano perpendikularra (edo hiperplano ortogonal) izango da.

Irudiak erakusten du, baita ere, zer gertatzen den  $\mathbf{u}$  bektoreak  $\mathbf{x}$ -k eta ardatz batek osatutako angelua erdibitzen duenean. Orduan,  $\mathbf{H}\mathbf{x}$  ardatz horren gainean dago. Alegia,  $\mathbf{H}\mathbf{x}$ -ren osagai guztiak zero dira, bat izan ezik. Gainera,  $\mathbf{H}$  ortogonal denez, bektorearen luzera gordetzen du, zeren:

$$\|\mathbf{H}\mathbf{x}\|_2^2 = (\mathbf{H}\mathbf{x})^T(\mathbf{H}\mathbf{x}) = \mathbf{x}^T\mathbf{H}^T\mathbf{H}\mathbf{x} = \mathbf{x}^T\mathbf{x} = \|\mathbf{x}\|_2^2 \quad \Rightarrow \quad \|\mathbf{H}\mathbf{x}\|_2 = \|\mathbf{x}\|_2. \quad (6.3)$$

Ondorioz,  $\mathbf{H}\mathbf{x}$  bektorearen zero ez den osagai bakarraren balioa  $\pm\|\mathbf{x}\|_2$  da.

Orduan,  $\mathbf{x}$  bektore baterako  $\mathbf{H}\mathbf{x}$ -ren  $k$ -garren osagaia izan ezik beste guztiak zero bihurtzeko,  $\mathbf{x}$  bektorearen luzera gordez, honelako  $\mathbf{H}$  Householderren matrizea eraiki behar dugu:

$$\sigma = \pm\|\mathbf{x}\|_2, \quad (6.4)$$

$$\mathbf{u} = \mathbf{x} + \sigma\mathbf{e}_k, \quad (6.5)$$

$$\rho = 2/\|\mathbf{u}\|_2^2 = 1/(\sigma u_k), \quad (6.6)$$

$$\mathbf{H} = \mathbb{1} - \rho\mathbf{u}\mathbf{u}^T, \quad (6.7)$$

non  $\mathbf{e}_k$  oinarri kanonikoaren  $k$ -garren bektorea baita. Biribiltze-errorearen aurrean eta (6.5) berdintza kontuan hartuz, hobe da  $\sigma$ -ren zeinua  $x_k$ -renaren berdina hartzea; hau da,  $\sigma = \text{zeinu}(\mathbf{x}_k)\|\mathbf{x}\|_2$ .

(6.6) berdintzan, ohartu hau betetzen dela:

$$\begin{aligned} \|\mathbf{u}\|_2^2 &= \mathbf{u}^T\mathbf{u} = (\mathbf{x} + \sigma\mathbf{e}_k)^T(\mathbf{x} + \sigma\mathbf{e}_k) \\ &= \mathbf{x}^T\mathbf{x} + \sigma^2\mathbf{e}_k^T\mathbf{e}_k + \sigma\mathbf{x}^T\mathbf{e}_k + \sigma\mathbf{e}_k^T\mathbf{x} \\ &= \|\mathbf{x}\|_2^2 + \sigma^2 + 2\sigma x_k \\ &= 2\sigma^2 + 2\sigma x_k \quad \text{dugu (6.4) bidez} \\ &= 2\sigma(x_k + \sigma) \\ &= 2\sigma u_k \quad \text{dugu (6.5) erabiliz,} \end{aligned}$$

non  $u_k$  zenbakia  $\mathbf{u}$  bektorearen  $k$ -garren osagaia baita.

Azkenik, (6.2) eta (6.5) berdintzak kontuan hartuz, hau lortzen da:

$$\begin{aligned} \mathbf{H}\mathbf{x} &= \mathbf{x} - \rho(\mathbf{u}^T\mathbf{x})\mathbf{u} \\ &= \mathbf{x} - 2\frac{\mathbf{x}^T\mathbf{x} + \sigma\mathbf{x}^T\mathbf{e}_k}{\|\mathbf{u}\|_2^2}\mathbf{u} \\ &= \mathbf{x} - 2\frac{\sigma^2 + \sigma x_k}{2\sigma u_k}\mathbf{u} \\ &= \mathbf{x} - \mathbf{u} \\ &= -\sigma\mathbf{e}_k. \end{aligned} \quad (6.8)$$

Hots,  $\mathbf{H}$  transformazioak  $\mathbf{x}$  garraiatzen du  $k$ . ardatzera. Zer gertatzen da  $\sigma = -\text{zeinu}(\mathbf{x}_k)\|\mathbf{x}\|_2$  hartzen bada?

## 6.2. QR faktORIZAZIOA

### 6.2.1. A matrize karratu ez-singularra

Demagun  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizea ez dela singularra; orduan, aurreko atalaren arabera,  $n - 1$  Householderren matrizeak eraiki ditzakegu hau bete dadin:

$$\mathbf{H}_{n-1} \dots \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \mathbf{R}, \quad (6.9)$$

non  $\mathbf{R} \in \mathbb{R}^{n \times n}$  matrizea goi-triangeluarra baita.

Prozesu honen lehenengo urratsa  $\mathbf{H}_1$  matrizea eraikitzea da,  $\mathbf{A}$ -ren lehenengo zutabea,  $\mathbf{a}_1$ ,  $\mathbf{e}_1$ -en multiplo bihurtzeko ( $\mathbf{e}_1$  oinarri kanonikoaren lehenengo bektorea da), zutabearen norma euklidearra gordez. Alegia, 2. gaitik  $n$ . gaira dauden gaiak zero bihurtuz. Orduan, (6.8) kontuan hartuz, zera lortzen da:

$$\mathbf{H}_1 \mathbf{a}_1 = (\mathbf{1} - \rho_1 \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{a}_1 = -\sigma_1 \mathbf{e}_1 = \begin{bmatrix} r_{11} \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad (6.10)$$

non  $r_{11} = -\sigma_1$  eta  $\sigma_1 = \text{zeinu}(a_{11}) \|\mathbf{a}_1\|_2$  (biribiltze-erroreak saihesteko). Aurreko atalagatik, badakigu  $\mathbf{u}_1 = \mathbf{a}_1 + \sigma_1 \mathbf{e}_1$  hartu behar dugula. Beraz, Householderren lehenengo bektorea hau izango da:

$$\mathbf{u}_1 = \begin{bmatrix} a_{11} + \sigma_1 \\ a_{21} \\ \vdots \\ a_{n1} \end{bmatrix}. \quad (6.11)$$

Bektore horrekin,  $\mathbf{A}$  matrizean  $\mathbf{H}_1$  aplikatuz, zera dugu:

$$\mathbf{A}^{(2)} = \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} r_{11} & a_{12}^{(2)} & \dots & a_{1n}^{(2)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix}, \quad (6.12)$$

non  $\mathbf{A}$  matrizearen gai guztiak aldatu baitira. Hori ez da gertatzen ezabatze gaussiarrean; metodo horretan, lehenengo lerroa ez da aldatzen. Gainerako matrizea,  $\tilde{\mathbf{A}}_2$ , lehenengo lerroa eta lehenengo zutabea kenduz geratzen den  $(n - 1) \times (n - 1)$  matrizea da.

**6.1. adibidea.** Egin dezagun Householderren ezabatzearen lehenengo urratsa matrize honetarako:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 3 & 1 \\ 3 & -1 & -1 \end{bmatrix}.$$

*Ebazpena.* (6.4) eta (6.10) aplikatuz, hau dugu:

$$\|\mathbf{a}_1\|_2 = \sqrt{14} = 3.742, \quad \sigma_1 = \text{zeinu}(+1)3.742 = +3.742, \quad r_{11} = -3.742$$

eta (6.11) -ren bitartez

$$\mathbf{u}_1 = \begin{bmatrix} 1 + 3.742 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 4.742 \\ 2 \\ 3 \end{bmatrix}.$$

$\mathbf{H}_1$  eta  $\mathbf{A}^{(2)}$  kalkulatzeko, (6.6) - (6.7) eta (6.12) erabiliz, hurrenez hurren, emaitza hauek ditugu:

$$\rho_1 = 1/(3.742 \cdot 4.742) = 0.05637$$

eta

$$\mathbf{H}_1 = \mathbf{I} - \rho_1 \mathbf{u}_1 \mathbf{u}_1^T = \begin{bmatrix} -0.2673 & -0.5345 & -0.8018 \\ -0.5345 & 0.7745 & -0.3382 \\ -0.8018 & -0.3382 & 0.4927 \end{bmatrix},$$

$$\mathbf{A}^{(2)} = \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} -3.742 & -1.069 & -0.2673 \\ 0 & 2.127 & 0.04368 \\ 0 & -2.309 & -2.434 \end{bmatrix},$$

non zenbaki guztiak lau zifra esanguratsutara biribildu baititugu.  $\square$

Bigarren Householderren transformazioaren eraikuntzan, helburua da  $\tilde{\mathbf{A}}_2$ -ren lehenengo zutabea egokiro eraldatzea,  $\mathbf{A}^{(2)}$ -ren lehenengo lerroa eta lehenengo zutabea aldatu barik. Hori lortzeko, nahikoa da  $\mathbf{u}_2$ -ren lehenengo gaia zero hartzea; ikus (6.2) -ren bigarren berdintza. Aukera horrekin;  $\mathbf{H}_2$ -ren aplikazioak bektore orokor bati ez dio aldatzen lehenengo osagaia, eta  $\mathbf{H}_2$ -ren aplikazioak berdin uzten dio  $\mathbf{e}_1$ -en edozein multiplori.

Aurreko adibidera itzuliz:

$$\tilde{\mathbf{A}}_2 = \begin{bmatrix} 2.127 & 0.04368 \\ -2.309 & -2.434 \end{bmatrix}.$$

Orain,  $\tilde{\mathbf{a}}_2$  lehenengo zutabea da eta

$$\|\tilde{\mathbf{a}}_2\|_2 = 3.139, \quad \sigma_2 = \text{zeinu}(+2.127)3.139 = +3.139, \quad r_{22} = -3.139,$$

gainera, (6.5) erabiliz:

$$\mathbf{u}_2 = \begin{bmatrix} 0 \\ 2.127 + 3.139 \\ -2.309 \end{bmatrix} = \begin{bmatrix} 0 \\ 5.266 \\ -2.309 \end{bmatrix}.$$

$\mathbf{H}_2$  eta  $\mathbf{A}^{(3)}$  kalkulatzeko, (6.6) - (6.7) eta (6.12) erabiliz, hurrenez hurren, emaitza hauek ditugu:

$$\rho_2 = 1/(3.139 \cdot 5.266) = 0.06050$$

eta

$$\mathbf{H}_2 = \mathbb{1} - \rho_2 \mathbf{u}_2 \mathbf{u}_2^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -0.678 & 0.7357 \\ 0 & 0.7357 & 0.6775 \end{bmatrix}, \quad \mathbf{A}^{(3)} = \mathbf{H}_2 \mathbf{A}^{(2)} = \begin{bmatrix} -3.742 & -1.069 & -0.2673 \\ 0 & -3.139 & -1.820 \\ 0 & 0 & -1.617 \end{bmatrix}.$$

Ondorioz, hau lortu dugu:

$$\mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \mathbf{R},$$

non  $\mathbf{R} = \mathbf{A}^{(3)} \in \mathbb{R}^{3 \times 3}$  matrize goi-triangeluarra baita.  $\square$

Baldin  $\mathbf{A} \in \mathbb{R}^{n \times n}$  ez bada singularra, Householderren ezabapenaren  $n - 1$  urrats egin ditzakegu, eta  $\tilde{\mathbf{A}}_i$  gainerako matrizearen  $\tilde{\mathbf{a}}_i$  lehenengo zutabea ez da zero izango  $i$ -garren urrats bakoitzean,  $i = 1, \dots, n - 1$  ( $\tilde{\mathbf{A}}_1 = \mathbf{A}$  da); eta azken zutabea ez ditugu zeroak egin behar. Orduan, hau dugu:

$$\mathbf{H}_{n-1} \dots \mathbf{H}_1 \mathbf{A} = \mathbf{R},$$

non  $\mathbf{R} \in \mathbb{R}^{n \times n}$  matrize goi-triangeluarra baita. Izan bedi  $\mathbf{Q}^T \in \mathbb{R}^{n \times n}$  matrize ortogonal hau:

$$\mathbf{Q}^T = \mathbf{H}_{n-1} \dots \mathbf{H}_1. \quad (6.13)$$

Beraz,

$$\mathbf{Q} = \mathbf{H}_1 \dots \mathbf{H}_{n-1}. \quad (6.14)$$

Honako adierazpen bakoitzari  $\mathbf{A}$ -ren  $QR$  faktORIZAZIO deritzo:

$$\mathbf{Q}^T \mathbf{A} = \mathbf{R} \quad \text{edo} \quad \mathbf{A} = \mathbf{Q} \mathbf{R}. \quad (6.15)$$

Praktikan, ez dugu  $\mathbf{Q}$  kalkulatu behar sistema bat ebazteko, zeren  $\mathbf{Q}^T \mathbf{v}$  matrize-bektore biderketetan bakarrik agertzen baita. Bektore hori kalkula daiteke (6.13) adierazpeneko  $\mathbf{H}_{n-1}, \dots, \mathbf{H}_1$  Householderren banako transformazioak aplikatuz; matrize horiek ere ez dira esplizituki kalkulatu behar.

$k$ -garren Householderren matrizea adierazteko,  $\mathbf{u}_k$  Householderren bektorearen  $n - (k - 1)$  osagai gorde behar ditugu, eta  $\rho_k$  balioa. Nahiz eta  $\rho_k$  kalkula dezakegun  $\mathbf{u}_k$  erabiliz, bere kalkulurako eragiketak ez errepikatzeko, gordetzen dugu.

FaktORIZAZIOA	Biderketak/Zatiketak	Batuketak/Kenketak
$\mathbf{A} = \mathbf{PLU}$	$n^3/3$	$n^3/3$
$\mathbf{A} = \mathbf{QR}$	$2n^3/3$	$2n^3/3$
$\mathbf{A} = \mathbf{R}^T \mathbf{R}$	$n^3/6$	$n^3/6$

**6.1. taula.** Matrize-faktORIZAZIOEN kostu aritmetikoa.

6.2.1. taularen arabera,  $QR$  faktORIZAZIOAK  $LU$  faktORIZAZIOAREN eragiketen bikoitza erabiltzen du, hau da, lan bikoitza egin behar du matrize bat faktORIZATZECO. Hala ere, askotan,  $QR$  faktORIZAZIOA aukeratzen da bere zenbakizko egonkortasunagatik.  $\mathbf{Q}$  ortogonal denez  $\kappa_2(\mathbf{R}) = \kappa_2(\mathbf{A})$ ; hau da,  $QR$  faktORIZAZIOAK ez du okerragotzen  $\mathbf{A}$ -ren baldintza. Frogatu hori ariketa gisa.



### 6.2.2. Sistema lineal karratu determinatu baten ebazpena

Behin  $\mathbf{A}$ -ren  $QR$  faktORIZAZIOA ezagutu eta gero,  $\mathbf{Ax} = \mathbf{b}$  kalkula daiteke, hau kontuan hartuz:

$$\mathbf{Ax} = \mathbf{QRx} = \mathbf{b} \quad \Rightarrow \quad \mathbf{Rx} = \mathbf{Q}^T \mathbf{b}. \quad (6.16)$$

Beraz, lehenengo  $\mathbf{Q}^T \mathbf{b}$  kalkulatu da, eta, gero, (6.16) sistema ebazten dugu.

Aurreko 6.1. adibideari jarraituz, sistema hau dugu:

$$\begin{bmatrix} 1 & 1 & 2 \\ 2 & 3 & 1 \\ 3 & -1 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 3 \\ 2 \\ 6 \end{bmatrix},$$

orduan:

$$\mathbf{R} = \begin{bmatrix} -3.742 & -1.069 & -0.2673 \\ 0 & -3.139 & -1.820 \\ 0 & 0 & -1.617 \end{bmatrix}, \quad \mathbf{Q}^T \mathbf{b} = \begin{bmatrix} -6.682 \\ 1.320 \\ -1.617 \end{bmatrix}$$

(lau zifra esanguratsutara biribildua). Atzeranzko ebazpena eginez, (6.16) sisteman,  $\mathbf{x} = (2, -1, 1)^T$  emaitza lortzen da ( $\mathbf{R}$  eta  $\mathbf{Q}^T \mathbf{b}$ -ren bertsio zehatzekin!).  $\square$

## 6.3. Givens-en biraketak

Householderren islapenak arras erabilgarriak dira bektore batean zero asko sartzeko. Hala ere, badaude kalkulu batzuk non zeroak era selektiboago batean sartu behar baititugu; orduan *Givens-en biraketak* aukeratuko ditugu, eta haiek dira honelako identitaterako bi-heineko zuzenketak:

$$G(i, k, \theta) = \begin{bmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & c & \dots & s & \dots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \dots & -s & \dots & c & \dots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix} \begin{matrix} i \\ k \end{matrix},$$

$i \qquad k$

non  $\theta$  angelu baterako  $c = \cos(\theta)$  eta  $s = \sin(\theta)$ . Givens-en biraketak direnez, matrize ortogonalak dira (egiaztatu).  $\mathbf{x}$  bektore bat  $G(i, k, \theta)^T$  matrizeaz aurrebiderkatzea ( $i, k$ ) koordinatu-planoan,  $\theta$  radianeko erlojuaren kontrako noranzko biraketa bat ematea da. Hots,  $\mathbf{y} = G(i, k, \theta)^T \mathbf{x}$  bada, orduan:

$$y_j = \begin{cases} cx_i - sx_k, & j = i \\ sx_i + cx_k, & j = k \\ x_j, & j \neq i, k. \end{cases}$$

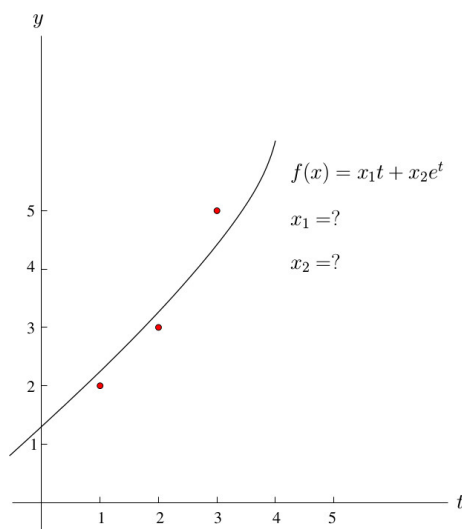
Formula horietatik  $y_k = 0$  izatera behartu dezakegu balio hauek hartuz:

$$c = \frac{x_i}{\sqrt{x_i^2 + x_k^2}}, \quad s = \frac{-x_k}{\sqrt{x_i^2 + x_k^2}}.$$

**6.2. adibidea.**  $\mathbf{x} = [1, 2, 3, 4]^T$ ,  $\cos(\theta) = 1/\sqrt{5}$  eta  $\sin(\theta) = -2/\sqrt{5}$  badira, orduan  $\mathbf{y} = G(2, 4, \theta)\mathbf{x} = [1, \sqrt{20}, 3, 0]^T$ .

## 6.4. Minimo karratu linealak: sistema gaindeterminatuak

Demagun  $m$  datu-puntu dauzkagula,  $(t_i, y_i)$ ,  $f(\mathbf{x}, t)$  funtzio batekin, lineala dena  $x_1, x_2, \dots, x_n$  parametro askeekiko. Adibidez, demagun  $(1, 2)$ ,  $(2, 3)$  eta  $(3, 5)$  hiru bikoteei  $f(\mathbf{x}, t) = x_1 t + x_2 e^t$  funtzioa egokitu nahi diegula (6.2. irudia).



**6.2. irudia.** Minimo karratu linealen adibidea.

Doikuntza zehatz bat lortzeko, hots:

$$\begin{aligned} f(\mathbf{x}, t_1) &= y_1 \\ f(\mathbf{x}, t_2) &= y_2 \\ f(\mathbf{x}, t_3) &= y_3, \end{aligned}$$

$\mathbf{Ax} = \mathbf{b}$  sistema bete behar du, non:

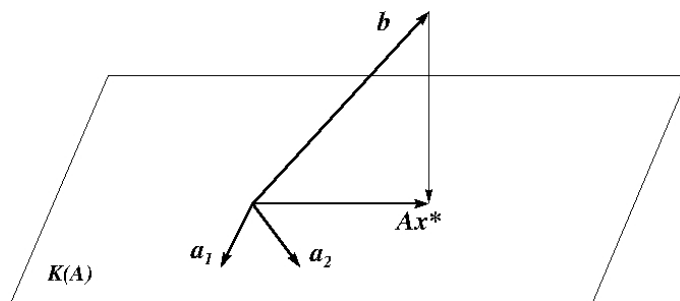
$$\mathbf{A} = \begin{bmatrix} 1 & e \\ 2 & e^2 \\ 3 & e^3 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \text{eta} \quad \mathbf{b} = \begin{bmatrix} 2 \\ 3 \\ 5 \end{bmatrix}.$$

Sistema lineal hau *gaindeterminatua*enez ( $m > n$  zentzuan), ezin dugu ebatzi ezagutzen ditugun metodoez. Izan ere, guk aukeratuko dugu  $\mathbf{x}$  bektorea  $\mathbf{Ax} - \mathbf{b}$  hondar-bektorearen neurriren bat minimizatzeke. Izan bedi bi-norma aukeratutako neurria. Beraz, problema hau ebatzi behar dugu:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{b}\|_2, \quad (6.17)$$

non  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ ,  $\mathbf{b} \in \mathbb{R}^m$ , gure adibidean  $m = 3$  eta  $n = 2$ .

Adibidearen soluzioa kalkulatzeko eta « $(2, 3, 5)^T$  bektoretik  $(1, 2, 3)^T$  eta  $(e, e^2, e^3)^T$  bektoreen konbinazio lineal hurbilena aurkitzea (bi-norman)», problema baliokideak dira. Orokorrean, (6.17) problema ebatzea honela uler daiteke: «Aurkitu  $\mathbf{b}$ -tik  $\mathbf{A}$  zutabeen konbinazio lineal hurbilena bi-norman». Geometrikoki, horrek esan nahi du  $\mathbf{b}$  bektoretik  $\mathbf{A}$ -ren zutabeek sortutako  $K(\mathbf{A})$  azpiespazioko  $n$  dimentsioko bektore hurbilena (norma euklidearrean) aurkitzea; ikus 6.3. irudia. Izan bitez  $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^n$  bektoreak  $\mathbf{A}$  matrizeko  $m$  zutabeak.



6.3. irudia. Minimo karratu linealen soluzioa.

Badakigu  $\mathbf{b}$ -tik  $K(\mathbf{A})$  azpiespazioko bektore hurbilena  $\mathbf{Ax}^* \in K(\mathbf{A})$  izango dela; izan ere,  $\mathbf{Ax}^* - \mathbf{b}$  perpendikularra da  $K(\mathbf{A})$  azpiespazioarekin. Beraz,  $x^*$ -k hau bete behar du:

$$\mathbf{a}_i^T (\mathbf{Ax}^* - \mathbf{b}) = 0, \quad i = 1, \dots, n,$$

edo, baliokideki,

$$\mathbf{A}^T (\mathbf{Ax}^* - \mathbf{b}) = \mathbf{0}$$

Bistan dago  $\mathbf{Ax}^*$  bakarra dela, eta  $\mathbf{A}$ -ren zutabeak linealki askeak badira,  $\mathbf{x}^*$  bektorea ere bakarra dela. Ondorioz,  $\mathbf{x}^*$  ekuazio-sistema honen soluzioa da:

$$(\mathbf{A}^T \mathbf{A}) \mathbf{x} = \mathbf{A}^T \mathbf{b},$$

bateragarria dena.

Informazio hori hurrengo teoreman finkatzen da.

**6.1. teorema.** Izan bitez  $m \geq n > 0$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ . Orduan, (6.17) minimo karratuen problemaren soluzioa  $\{\mathbf{x}^* \mid \mathbf{A}^T (\mathbf{Ax}^* - \mathbf{b}) = \mathbf{0}\}$  bektoreen multzoa da.

Baldin  $\mathbf{A}$ -ren zutabeak linealki askeak badira,  $\mathbf{x}^*$  soluzio bakarra da,  $\mathbf{A}^T \mathbf{A}$  ez da singularra eta  $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ .

*Frogantza.* Defini dezagun  $f(\mathbf{x})$  funtzio hau:

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{Ax} - \mathbf{b}\|_2^2 = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) \\ &= (\mathbf{Ax})^T (\mathbf{Ax}) - (\mathbf{Ax})^T \mathbf{b} - \mathbf{b}^T (\mathbf{Ax}) + \mathbf{b}^T \mathbf{b} \\ &= \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - 2\mathbf{x}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b}. \end{aligned}$$

Orduan,  $\min \|\mathbf{Ax} - \mathbf{b}\|_2 = \min f(\mathbf{x})$  eta soluzioak  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  bete behar du; alegia:

$$\nabla f(\mathbf{x}^*) = 2\mathbf{A}^T \mathbf{Ax}^* - 2\mathbf{A}^T \mathbf{b} = \mathbf{0},$$

ondorioz,  $\mathbf{x}^*$  soluzioak ekuazio sistema hau bete behar du:

$$\mathbf{A}^T (\mathbf{Ax} - \mathbf{b}) = \mathbf{0}. \quad (6.18)$$

eta  $\mathbf{A}^T \mathbf{A}$ , behintzat, erdidefinitu positiboa denez,  $f(\mathbf{x})$  funtzio konbexua da, eta, ondorioz,  $\mathbf{x}^*$  minimizatzaile bat da. Beraz, teoremaren lehenengo atala frogatuta dago. Bigarrenean,  $\mathbf{A}$ -ren zutabeak linealki askeak direnez,  $\mathbf{A}^T \mathbf{A}$  ez da singularra, eta orduan  $\mathbf{x}^*$  minimizatzaile bakarra da.  $\square$

Berrordenatuz (6.18), ekuazio hauek lortzen dira:

$$(\mathbf{A}^T \mathbf{A})\mathbf{x} = \mathbf{A}^T \mathbf{b}, \quad (6.19)$$

*ekuazio normalak* deritzenak. Ekuazio normalen matrizea,  $\mathbf{A}^T \mathbf{A}$ , edozein  $\mathbf{A}$ -tarako matrize simetrikoa eta erdidefinitu positiboa da, eta definitu positiboa da baldin eta soilik baldin  $\mathbf{A}$  matrizearen zutabeak linealki askeak badira, hots,  $\mathbf{A}$  matrizea *zutabe hein betekoa* bada (horen frogapena ariketa gisa geratzen da). Bestalde, ekuazio normalak beti dira bateragarriak,  $\mathbf{A}^T \mathbf{A}$  singularra izan arren.

Ekuazio normalek garrantzi praktiko handia dute, haiek bide zuzen bat ematen baitute minimo karratuen soluzioa kalkulatzeko.  $\mathbf{A}^T \mathbf{A}$  definitu positiboa denean ( $\mathbf{A}$  hein betekoa denean), ekuazio normalek soluzio bakar bat dute eta  $\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$  da.  $\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \in \mathbb{R}^{n \times m}$  matrizea ezaguna da  $\mathbf{A}$  matrizearen *sasialderantzizkoa* izenez.  $\mathbf{A}$  hein betekoa denean, minimo karratuen problema ebatz dezakegu algoritmo honekin.

**6.1. algoritmoa. Hein beteko minimo karratu linealen problemaren ebazpena, ekuazio normalen bitartez.**

**0 urratsa.** SARRERA. Sartu:  $\mathbf{A}$  eta  $\mathbf{b}$ .

**1 urratsa.** Eratu ekuazio normalen matrizea:  $\mathbf{A}^T \mathbf{A}$  eta  $\mathbf{A}^T \mathbf{b}$  bektorea.

**2 urratsa.** Kalkula ezazu Choleskyren faktORIZAZIOA:  $\mathbf{A}^T \mathbf{A} = \mathbf{R}^T \mathbf{R}$ , non  $\mathbf{R}$  goi-triangeluarra baita.

**3 urratsa.** Ebatzi  $\mathbf{R}^T \mathbf{y} = \mathbf{A}^T \mathbf{b}$  aurreranzko ordezkapenaz; gero, ebatzi  $\mathbf{R} \mathbf{x} = \mathbf{y}$  atzeranzko ordezkapenaz.

**4 urratsa.** IRTEERA. Emaitzak:  $\mathbf{x}$ .

**6.3. adibidea.** Ebatzi  $\mathbf{A} \mathbf{x} = \mathbf{b}$  sistema gaindeterminatua, hau badugu:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{eta} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ -5 \end{bmatrix}.$$

*Ebazpena.*  $\mathbf{A}$  matrizearen zutabeak bektore askeak direnez,  $\mathbf{A}$  hein betekoa da.

$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$  eta  $\mathbf{A}^T \mathbf{b} = \begin{bmatrix} 1 \\ -4 \end{bmatrix}$  eta  $(\mathbf{A}^T \mathbf{A}) \mathbf{x}^* = \mathbf{A}^T \mathbf{b}$  sistemaren soluzioa  $\mathbf{x}^* = (2, -3)^T$  da.

Hau da soluzio horri dagokion *hondar-bektorea*:

$$\mathbf{A} \mathbf{x}^* - \mathbf{b} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ -5 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \\ 2 \end{bmatrix}.$$

Beraz, hau da sistema horren hondar-bektorearen norma euklidearra:  $\min_{\mathbf{x}} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 = \|\mathbf{A} \mathbf{x}^* - \mathbf{b}\|_2 = \sqrt{12} = 2\sqrt{3} \approx 3.464$  (ikus 6.3. irudia).  $\square$

Ariketa moduan, ebatzi ekuazio normalen bidez atal honen hasierako  $f(\mathbf{x}, t) = x_1 t + x_2 e^t$  funtzioaren parametroen kalkulua, funtzioa (1, 2), (2, 3) eta (3, 5) bikoteei egokitzeko.

Nahiz eta  $\mathbf{A}$  hein betekoa izan eta, beraz, ekuazio normalek  $\mathbf{x}^*$  soluzio bakar bat izan, metodo hori erabiltzea ez da beti egokiena minimo karratuen problema ebazteko. Hori da  $\mathbf{A}^T \mathbf{A}$  matrizearen baldintza  $\kappa_2(\mathbf{A}^T \mathbf{A}) = (\kappa_2(\mathbf{A}))^2$  betetzeagatik (egiaztatu ariketa gisa). Adibidez,  $\kappa_2(\mathbf{A}) = 10^3$  baldintza ez da oso txarra, baina  $\kappa_2(\mathbf{A}^T \mathbf{A}) = 10^6$  askoz txarragoa da eta zenbakizko egonkortasunaren arazo larriak sor ditzake. Beraz, metodo horren baldintza  $\mathbf{A}$  matrizearena baino okerragoa da. Aldiz, guk  $\mathbf{A}$ -ren  $QR$  deskonposizioa erabil dezakegu, honela:

$$\begin{array}{c} n \\ \boxed{\mathbf{A}} \\ m \end{array} = \begin{array}{c} m \\ \boxed{\mathbf{Q}} \\ m \end{array} \begin{array}{c} n \\ \left. \begin{array}{c} \boxed{\begin{array}{c} \mathbf{0} \\ \mathbf{0} \end{array}} \\ \mathbf{R} \end{array} \right\} \\ m \end{array} \quad \mathbf{R}_u \quad (6.20)$$

( $\mathbf{Q} \in \mathbb{R}^{m \times m}$  ortogonala,  $\mathbf{R} \in \mathbb{R}^{m \times n}$  goi-triangeluarra), faktORIZAZIO hori eginez Householderren islapenen bitartez,  $\mathbf{A}$  matrizea karratua eta ez-singularra den kasuan bezala. Deskonposizio horrek ere ematen du  $\mathbf{A}$ -ren zutabeen ortonormalizazio bat zenbaki egonkorra dela. Honako teorema honek erakusten du nola erabili (6.20) faktORIZAZIOA (6.17) problema ebazteko.

**6.2. teorema.** *Izan bitez  $m \geq n > 0$ ,  $\mathbf{b} \in \mathbb{R}^m$  bektorea eta  $\mathbf{A} \in \mathbb{R}^{m \times n}$  zutabe hein beteko matrizea. Orduan, (6.20) erako  $\mathbf{A} = \mathbf{QR}$  deskonposizio bat existitzen da, non  $\mathbf{Q} \in \mathbb{R}^{m \times m}$  matrize ortogonala baita,  $\mathbf{R} \in \mathbb{R}^{m \times n}$  goi-triangeluarra eta  $\mathbf{R}_u$  (hots,  $\mathbf{R}$ -ren lehenengo  $n$  erroak) matrize goi-triangeluarra eta ez-singularra.*

*Gainera, hau da (6.17) minimo karratuen problemaren soluzio bakarra:*

$$\mathbf{x}^* = \mathbf{R}_u^{-1}(\mathbf{Q}^T \mathbf{b})_u, \quad \text{non} \quad (\mathbf{Q}^T \mathbf{b})_u^T = ((\mathbf{Q}^T \mathbf{b})_1, \dots, (\mathbf{Q}^T \mathbf{b})_n),$$

*eta hondarraren luzeraren karratua hau da:*

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \sum_{i=n+1}^m (\mathbf{Q}^T \mathbf{b})_i^2.$$

*Frogantza.* Householderren transformazioetatik ondorioztatzen da  $\mathbf{A}$ -ren  $QR$  deskonposizioaren existentzia, eta  $\mathbf{A}$ -ren zutabe hein betekoa izatetik  $\mathbf{R}_u$ -ren ez-singularitatea ( $\mathbf{Q}$  ez baita singularra). Orain,  $\mathbf{Q}^T$  ortogonala denez, bektore batez biderkatzean ez dio aldatzen bere norma euklidearra ( $\|\mathbf{Q}^T \mathbf{x}\|_2 = \|\mathbf{x}\|_2$ ) eta, orduan, zera dugu:

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \|\mathbf{QRx} - \mathbf{b}\|_2 = \|\mathbf{Q}^T(\mathbf{QRx} - \mathbf{b})\|_2 = \|\mathbf{Rx} - \mathbf{Q}^T \mathbf{b}\|_2,$$

hori dela eta, honela berridatz dezakegu (6.17) problema:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Rx} - \mathbf{Q}^T \mathbf{b}\|_2.$$

Orduan,  $(\mathbf{Q}^T \mathbf{b})_l^T = ((\mathbf{Q}^T \mathbf{b})_{n+1}, \dots, (\mathbf{Q}^T \mathbf{b})_m)$  bada, hau dugu:

$$\|\mathbf{Rx} - \mathbf{Q}^T \mathbf{b}\|_2^2 = \|\mathbf{R}_u \mathbf{x} - (\mathbf{Q}^T \mathbf{b})_u\|_2^2 + \|(\mathbf{Q}^T \mathbf{b})_l\|_2^2$$

eta hori minimizatzen da  $\mathbf{x} = \mathbf{R}_u^{-1}(\mathbf{Q}^T \mathbf{b})_u$  denean (lehenengo batugaia zero egiten baita), gainera  $\|(\mathbf{Q}^T \mathbf{b})_l\|_2^2 = \sum_{i=n+1}^m (\mathbf{Q}^T \mathbf{b})_i^2$  dugu.  $\square$

**6.2. algoritmoa.** Hein beteko minimo karratu linealen problemaren ebazpena,  $QR$  faktORIZAZIOAREN bitartez.

**0 urratsa.** SARRERA. Sartu:  $\mathbf{A}$  eta  $\mathbf{b}$ .

- 1 urratsa.** Kalkula ezazu  $\mathbf{A}$ -ren  $QR$  faktORIZAZIOA Householderren transformazioak erabiliz,  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , non  $\mathbf{Q}$  ortogonala baita eta  $\mathbf{R}$  goi-triangeluarra.
- 2 urratsa.** Eratu  $\tilde{\mathbf{b}} = \mathbf{Q}^T \mathbf{b}$  eta  $\tilde{\mathbf{b}}_u$  (hots,  $\tilde{\mathbf{b}}$ -ren lehenengo  $n$  osagaiak).
- 3 urratsa.** Ebatzi  $\mathbf{R}_u \mathbf{x} = \tilde{\mathbf{b}}_u$  atzeranzko ordezkapenez, non  $\mathbf{R}_u$   $\mathbf{R}$ -ren lehenengo  $n$  lerroek osatutako azpimatrizen karratu goi-triangeluarra baita.
- 4 urratsa.** IRTEERA. Eraitza:  $\mathbf{x}$ .

Alegia,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  matrizea zutabe hein betekoa ( $m > n$ ) eta  $\mathbf{b} \in \mathbb{R}^m$  badira, eta  $\mathbf{H}_n \dots \mathbf{H}_1 \mathbf{A} = \mathbf{R}$  bada,  $\mathbf{H}_n \dots \mathbf{H}_1 \mathbf{b}$  kalkulatu behar dugu. Orduan,  $\mathbf{Q}^T [\mathbf{A} \ \mathbf{b}]$  guztia batera kalkulatu dugu,  $\mathbf{Q}^T = \mathbf{H}_n \dots \mathbf{H}_1$  izanik. Gainera, gogoratu  $\mathbf{H}\mathbf{x}$  kalkulatzeko ez dugula  $\mathbf{H}$  aurkitu behar.

**6.4. adibidea.** Adibide honetan,  $QR$  metodoaren bitartez ebatziko da aurreko 6.3. adibidearen  $\mathbf{A}\mathbf{x} = \mathbf{b}$  sistema gaindeterminatua.

*Ebazpena.* Gogora dezagun honako hau:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{eta} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ -5 \end{bmatrix}.$$

Izan bitez  $\mathbf{a}_1 = [1 \ 1 \ 0]^T$  eta  $\mathbf{a}_2 = [1 \ 0 \ 1]^T$   $\mathbf{A}$  matrizearen zutabe-bektoreak. Orduan, lau digitu esanguratsutara biribilduz, hau dugu:

$$\sigma_1 = \text{zeinu}(a_{11}) \|\mathbf{a}_1\|_2 = +\sqrt{2} = 1.414, \quad \mathbf{u}_1 = \begin{bmatrix} 1 + 1.414 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.414 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{eta} \quad \rho_1 = \frac{1}{\sigma_1 \cdot [\mathbf{u}_1]_1} = \frac{1}{1.414 \cdot 2.414} = 0.2929.$$

Beraz, (6.2) adierazpenak erabiliz, hau lortzen dugu:

$$\mathbf{H}_1 \mathbf{a}_1 = \mathbf{a}_1 - \rho_1 (\mathbf{u}_1^T \mathbf{a}_1) \mathbf{u}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} - 0.2929 \left( [2.414 \ 1 \ 0] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \right) \begin{bmatrix} 2.414 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -1.414 \\ 0 \\ 0 \end{bmatrix},$$

espero genuen bezala, ikus (6.10) eta (6.11) adierazpenak (hots, praktikan kalkulu horiek ez ditugu egin behar). Orain  $\mathbf{H}_1 \mathbf{a}_2$  kalkulatu dugu:

$$\mathbf{H}_1 \mathbf{a}_2 = \mathbf{a}_2 - \rho_1 (\mathbf{u}_1^T \mathbf{a}_2) \mathbf{u}_1 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} - 0.2929 \left( [2.414 \ 1 \ 0] \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right) \begin{bmatrix} 2.414 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.7071 \\ -0.7071 \\ 1 \end{bmatrix}.$$

Hori dela eta, hau dugu:

$$\mathbf{A}^{(2)} = \begin{bmatrix} -1.414 & -0.7071 \\ 0 & -0.7071 \\ 0 & 1 \end{bmatrix}$$

( $\mathbf{A}^{(1)} = \mathbf{A}$  baita).

Ikusi den bezala, ez dugu  $\mathbf{H}_1$  kalkulatzeko beharrik. Hala ere, kalkulatu egingo dugu:

$$\mathbf{H}_1 = \mathbb{1} - \rho_1 \mathbf{u}_1 \mathbf{u}_1^T = \begin{bmatrix} -0.7071 & -0.7071 & 0 \\ -0.7071 & 0.7071 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

(ariketa gisa, egiaztatu  $\mathbf{H}_1 \mathbf{A} = \mathbf{A}^{(2)}$  betetzen dela).

Gainera,  $\mathbf{Q}^T \mathbf{b}$  kalkulatzeko,  $\mathbf{b}^{(2)} = \mathbf{H}_1 \mathbf{b}$  ( $\mathbf{b}^{(1)} = \mathbf{b}$  baita) honela aurkituko dugu:

$$\mathbf{H}_1 \mathbf{b} = \mathbf{b} - \rho_1 (\mathbf{u}_1^T \mathbf{b}) \mathbf{u}_1 = \begin{bmatrix} 1 \\ 0 \\ -5 \end{bmatrix} - 0.2929 \left( \begin{bmatrix} 2.414 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -5 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.7071 \\ -0.7071 \\ -5 \end{bmatrix}.$$

Orain,  $\mathbf{A}^{(2)}$  matrizeari lehenengo zutabea eta lehenengo lerroa ezabatuz, azpimatrizatze hau geratzen da:

$$\tilde{\mathbf{a}}_2 = \tilde{\mathbf{A}}_2 = \begin{bmatrix} -0.7071 \\ 1 \end{bmatrix}.$$

Jarraian,  $\tilde{\mathbf{a}}_2$  bektoreko lehenengo gaiaren azpiko gaia zero bihurtzeko, hau egingo dugu:

$$\sigma_2 = \text{zeinu}(-0.7071) \|\tilde{\mathbf{a}}_2\|_2 = -1.225,$$

eta, beraz, hau da  $\mathbf{A}^{(2)}$ -ren bigarren zutaberako erabili behar dugun  $\mathbf{u}_2$  Householderren bektorea:

$$\mathbf{u}_2 = \begin{bmatrix} 0 \\ -0.7071 - 1.225 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1.932 \\ 1 \end{bmatrix}, \quad \rho_2 = \frac{1}{\sigma_2 \cdot [\mathbf{u}_2]_2} = \frac{1}{(-1.225)(-1.932)} = 0.4226;$$

kontuan izan  $\mathbf{u}_2$ -ren lehenengo gaia zero dela,  $\mathbf{A}^{(2)}$ -ren lehenengo zutabea eta lehenengo lerroa ez aldatzeko.

Beraz, (6.2) adierazpenak erabiliz, hau lortzen dugu:

$$\begin{aligned} \mathbf{H}_2 \mathbf{a}_2^{(2)} &= \mathbf{a}_2^{(2)} - \rho_2 (\mathbf{u}_2^T \mathbf{a}_2^{(2)}) \mathbf{u}_2 \\ &= \begin{bmatrix} -0.7071 \\ -0.7071 \\ 1 \end{bmatrix} - 0.4226 \left( \begin{bmatrix} 0 & -1.932 & 1 \end{bmatrix} \begin{bmatrix} -0.7071 \\ -0.7071 \\ 1 \end{bmatrix} \right) \begin{bmatrix} 0 \\ -1.932 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} -0.7071 \\ 1.225 \\ 0 \end{bmatrix}, \end{aligned}$$



espero genuen bezala (hots, praktikan kalkulu horiek ez ditugu egin behar).

Hori dela eta, hau dugu:

$$\mathbf{A}^{(3)} = \begin{bmatrix} -1.414 & -0.7071 \\ 0 & 1.225 \\ 0 & 0 \end{bmatrix}.$$

Ikusi den bezala, ez dugu  $\mathbf{H}_2$  kalkulatzeko beharrik. Hala ere, kalkulatu egingo dugu:

$$\mathbf{H}_2 = \mathbb{1} - \rho_2 \mathbf{u}_2 \mathbf{u}_2^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -0.5780 & 0.8165 \\ 0 & 0.8165 & 0.5774 \end{bmatrix},$$

(ariketa gisa, egiaztatu  $\mathbf{H}_2 \mathbf{A}^{(2)} = \mathbf{A}^{(3)}$  betetzen dela).

Gainera,  $\mathbf{Q}^T \mathbf{b}$ -ren kalkulua bukatzeko,  $\mathbf{b}^{(3)} = \mathbf{H}_2 \mathbf{b}^{(2)}$  (hots,  $\mathbf{b}^{(3)} = \mathbf{Q}^T \mathbf{b} = \mathbf{H}_2 \mathbf{H}_1 \mathbf{b}$ ) honela aurkituko dugu:

$$\begin{aligned} \mathbf{H}_2 \mathbf{b}^{(2)} &= \mathbf{b}^{(2)} - \rho_2 (\mathbf{u}_2^T \mathbf{b}^{(2)}) \mathbf{u}_2 \\ &= \begin{bmatrix} -0.7071 \\ -0.7071 \\ -5 \end{bmatrix} - 0.4226 \left( \begin{bmatrix} 0 & -1.932 & 1 \end{bmatrix} \begin{bmatrix} -0.7071 \\ -0.7071 \\ 5 \end{bmatrix} \right) \begin{bmatrix} 0 \\ -1.932 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} -0.7071 \\ -3.674 \\ -3.464 \end{bmatrix}. \end{aligned}$$

Azkenik,  $\mathbf{R} = \mathbf{A}^{(3)}$  eta  $\mathbf{Q}^T \mathbf{b} = \mathbf{b}^{(3)}$  direnez,  $\mathbf{R}_u \mathbf{x} = (\mathbf{Q}^T \mathbf{b})_u$  sistema ebatziko dugu; alegia:

$$\mathbf{R}_u \mathbf{x} = \begin{bmatrix} -1.414 & -0.7071 \\ 0 & 1.225 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -0.7071 \\ -3.674 \end{bmatrix} = (\mathbf{Q}^T \mathbf{b})_u.$$

Sistema horren soluzio bakarra  $x_1 = 2$  eta  $x_2 = -3$  da (biribilduz).

Bestalde,  $|(\mathbf{Q}^T \mathbf{b})_l| = 3.464$  da hondar-bektorearen norma euklidearra (kasu honetan, eskalar bat denez, balio absolutua da). Konparatu emaitza horiek 6.3. adibidean lortutako emaitzekin.  $\square$

## 6.5. QR faktORIZAZIOAREN PROPIETATEAK

Aurreko algoritmoak frogatzen du QR faktORIZAZIOA existitzen dela. Izan bedi  $K(\mathbf{A})$   $\mathbf{A}$  matrizeko zutabeen konbinazio linealen multzoa; hots, zutabe horiek sortzen duten azpiespazio bektoriala. Jarraian ikusiko dugu zer erlazio dagoen  $\mathbf{Q}$ -ren zutabeen eta  $K(\mathbf{A})$  eta  $K(\mathbf{A})^\perp$  azpiespazioen artean.

**6.3. teorema.**  $\mathbf{A} = \mathbf{QR}$  deskonposizioa  $\mathbf{A} \in \mathbb{R}^{m \times n}$ -ren zutabe hein beteko QR faktORIZAZIOA bada, eta  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ ,  $\mathbf{Q} = [\mathbf{Q}_u, \mathbf{Q}_l]$  eta  $\mathbf{Q}_u = [\mathbf{q}_1, \dots, \mathbf{q}_n]$ ,  $\mathbf{Q}_l = [\mathbf{q}_{n+1}, \dots, \mathbf{q}_m]$  zutabekako partiketak badira, orduan, hau betetzen da:

$$\begin{aligned} K(\mathbf{A}) &= K(\mathbf{Q}_u), \\ K(\mathbf{A})^\perp &= K(\mathbf{Q}_l) \end{aligned}$$

eta  $\mathbf{A} = \mathbf{Q}_u \mathbf{R}_u$ , non  $\mathbf{R}_u \in \mathbb{R}^{n \times n}$   $\mathbf{R}$  matrizearen lehenengo  $n$  lerroetako azpimatrizare goi-triangeluarra baita.

Frogantza. 6.20 adierazpenean  $\mathbf{Q} = [\mathbf{Q}_u, \mathbf{Q}_l]$  partiketa sartuta, hau dugu:

$$\begin{array}{c} n \\ m \end{array} \mathbf{A} = \begin{array}{c} n \quad m-n \\ m \end{array} \left[ \begin{array}{c|c} \mathbf{Q}_u & \mathbf{Q}_l \end{array} \right] = \begin{array}{c} n \\ m \end{array} \left[ \begin{array}{c|c} \mathbf{R}_u & \mathbf{0} \\ \hline \mathbf{0} & \end{array} \right] \mathbf{R} \quad (6.21)$$

Beraz, hau betetzen da:

$$\mathbf{A} = \mathbf{Q}_u \mathbf{R}_u + \mathbf{Q}_l \mathbf{0} = \mathbf{Q}_u \mathbf{R}_u.$$

Alegia,  $k = 1, \dots, n$  guztietarako, hau betetzen da:

$$\mathbf{a}_k = \sum_{i=1}^k r_{ik} \mathbf{q}_i \in K(\mathbf{Q}_u),$$

$(r_{1k}, \dots, r_{kk}, 0, \dots, 0)^T$  bektorea  $\mathbf{R}_u$ -ren  $k$ -garren zutabea izanik. Ondorioz,  $K(\mathbf{A}) \subset K(\mathbf{Q}_u)$ . Azkenik,  $\mathbf{A}$ -ren eta  $\mathbf{Q}_u$ -ren heinak berdinak direnez,  $K(\mathbf{A}) = K(\mathbf{Q}_u)$ .

Bestalde,

$$\mathbf{A} = \mathbf{QR} \Rightarrow \mathbf{Q}^T \mathbf{A} = \mathbf{R} \Rightarrow \begin{bmatrix} \mathbf{Q}_u^T \\ \mathbf{Q}_l^T \end{bmatrix} \mathbf{A} = \begin{bmatrix} \mathbf{R}_u \\ \mathbf{0} \end{bmatrix} \Rightarrow \mathbf{Q}_l^T \mathbf{A} = \mathbf{0}.$$

Beraz,  $K(\mathbf{Q}_l) \subset K(\mathbf{A})^\perp$ . Gainera,  $K(\mathbf{A})^\perp$ -ren dimentsioa eta  $\mathbf{Q}_l$ -ren heina berdinak dira  $(m - n)$ . Ondorioz,  $K(\mathbf{Q}_l) = K(\mathbf{A})^\perp$ .  $\square$

Aurreko 6.4. adibidean,  $\mathbf{H}_1$  eta  $\mathbf{H}_2$  kalkulatu ditugu. Orduan, hau dugu (lau digitu esanguratsutara biribilduz):

$$\begin{aligned} \mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 &= \begin{bmatrix} -0.7071 & -0.7071 & 0 \\ -0.7071 & 0.7071 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & -0.5780 & 0.8165 \\ 0 & 0.8165 & 0.5774 \end{bmatrix} \\ &= \begin{bmatrix} -0.7071 & 0.4087 & -0.5773 \\ -0.7071 & -0.4087 & 0.5773 \\ 0 & 0.8165 & 0.5774 \end{bmatrix}. \end{aligned}$$

Beraz, zera dugu:

$$\mathbf{Q}_u = \begin{bmatrix} -0.7071 & 0.4087 \\ -0.7071 & -0.4087 \\ 0 & 0.8165 \end{bmatrix} \quad \text{eta} \quad \mathbf{Q}_l = \begin{bmatrix} -0.5773 \\ 0.5773 \\ 0.5774 \end{bmatrix}.$$

Egiaztatu adibide horretarako  $\mathbf{Q}_u \mathbf{R}_u = \mathbf{A}$  eta  $\mathbf{Q}_l^T \mathbf{A} = \mathbf{0}$  betetzen direla, 6.3. teoremak dioen bezala.

**6.4. teorema.** *Demagun  $\mathbf{A} \in \mathbb{R}^{m \times n}$  matrizea zutabe hein betekoa dela.  $QR$  faktORIZAZIO «garbia»  $\mathbf{A} = \mathbf{Q}_u \mathbf{R}_u$  bakarra da, non  $\mathbf{Q}_u \in \mathbb{R}^{m \times n}$  matrizeak zutabe ortogonalak baititu eta  $\mathbf{R}_u$  goi-triangeluarra eta diagonaleko gai positiboduna baita. Gainera,  $\mathbf{R}_u$  matrizea  $\mathbf{A}^T \mathbf{A}$  matrizearen Choleskyren faktorea da.*

*Frogantza.*  $\mathbf{A}^T \mathbf{A} = (\mathbf{Q}_u \mathbf{R}_u)^T (\mathbf{Q}_u \mathbf{R}_u) = \mathbf{R}_u^T \mathbf{Q}_u^T \mathbf{Q}_u \mathbf{R}_u = \mathbf{R}_u^T \mathbf{R}_u$  denez gero,  $\mathbf{R}_u$  matrizea  $\mathbf{A}^T \mathbf{A}$  matrizearen Choleskyren faktorea da. Faktore hori bakarra da 5.9. teoremagatik eta,  $\mathbf{Q}_u = \mathbf{A} \mathbf{R}_u^{-1}$  denez,  $\mathbf{Q}_u$  ere bakarra da. Horrela da,  $\mathbf{Q}^T \mathbf{A} = \mathbf{R}$  eraikitzean Householderren matrizeak aukeratzten direnean  $\mathbf{R}$ -ren diagonaleko gaiak positiboak izan daitezzen.  $\square$

## 6.6. Hein urriko $QR$ faktORIZAZIOA

Baldin  $\mathbf{A}$ -ren heina urria bada,  $QR$  faktORIZAZIOAK ez du ematen oinarri bat  $K(\mathbf{A})$  azpiespaziorako. Problema hori zuzendu dezakegu  $\mathbf{A}$ -ren bertsio permutatu baten  $QR$  faktORIZAZIOA kalkulatz; hots,  $\mathbf{A} \mathbf{P} = \mathbf{Q} \mathbf{R}$ , non  $\mathbf{P}$  permutazio-matrize bat baita.

Orokorki,  $\mathbf{A}$ -ren heina ezagutzen ez dugunez, zutabeen trukeak egin beharko ditugu. Matrize horren heina urria bada, trukerik gabe bat-batean bukatuko litzateke Householderren ezabapena. Baina, kasu horretan, bi gauza gerta daitezke: bata da beste zutabeak zero izatea, orduan, bukatu da; bestea da beste zutabeen artean baten bat zero ez izatea eta, kasu horregatik, orokorrean egiten dugu zutabeen permutazioa. Adibidez,  $\mathbf{A}$  matrize honek hiru zutabe ditu, eta hau da bere  $QR$  faktORIZAZIOA:

$$\mathbf{A} = [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \mathbf{a}_3] = [\mathbf{q}_1 \quad \mathbf{q}_2 \quad \mathbf{q}_3] \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix},$$

orduan,  $\text{hein}(\mathbf{A}) = 2$  da, baina  $K(\mathbf{A})$  ez da  $K([\mathbf{q}_1 \quad \mathbf{q}_2])$ , ezta  $K([\mathbf{q}_1 \quad \mathbf{q}_3])$  edo  $K([\mathbf{q}_2 \quad \mathbf{q}_3])$  ere.

Zorionez, Householderren  $QR$  faktORIZAZIOA erraz alda dezakegu  $K(\mathbf{A})$ -rako oinarri orto-normal bat lortzeko. Pibotatzearen estrategia izaten da norma handiena duen zutabea pibot gisa hartzekoa.

Demagun hau dugula  $k$ -garren urratsean:

$$(\mathbf{H}_{k-1} \dots \mathbf{H}_1) \mathbf{A} (\mathbf{P}_1 \dots \mathbf{P}_{k-1}) = \mathbf{R}^{(k-1)} = \begin{bmatrix} \mathbf{R}_{11}^{(k-1)} & \mathbf{R}_{12}^{(k-1)} \\ \mathbf{0} & \mathbf{R}_{22}^{(k-1)} \end{bmatrix},$$

non  $\mathbf{R}_{11}^{(k-1)}$  ez-singularra eta goi-triangeluarra baita. Demagun  $\mathbf{R}_{22}^{(k-1)}$ -ren zutabekako partiketa hau:

$$\mathbf{R}_{22}^{(k-1)} = [\mathbf{z}_k^{(k-1)}, \dots, \mathbf{z}_n^{(k-1)}]$$

eta  $p$  dela (non  $k \leq p \leq n$ ) hau betetzen duen azpiindize txikiena:

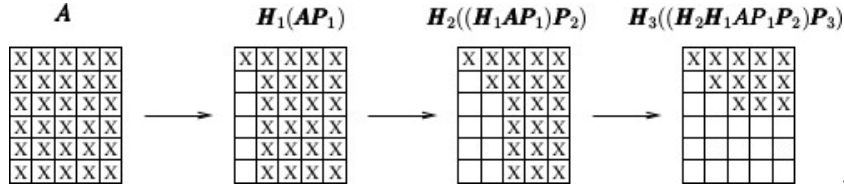
$$\|\mathbf{z}_p^{(k-1)}\|_2 = \max \{ \|\mathbf{z}_k^{(k-1)}\|_2, \dots, \|\mathbf{z}_n^{(k-1)}\|_2 \}.$$

Kontuan izan, hein( $\mathbf{A}$ ) =  $k - 1$  bada, norma handiena zero izango dela eta faktORIZAZIOA bukatu dugula. Bestela, izan bedi  $\mathbf{P}_k$  permutazioa,  $p$  eta  $k$  zutabeak trukatzan dituen, eta aurkitzen dugula  $\mathbf{H}_k$  Householder transformazio egokia,  $\mathbf{R}^{(k)} = \mathbf{H}_k \mathbf{R}^{(k-1)} \mathbf{P}_k$ , non  $\mathbf{R}^{(k)} = \mathbf{R}^{(k)}(k+1 : m, k) = \mathbf{0}$  (hots,  $(k, k)$  gaiaren azpiko gaiak zero dira).

$\mathbf{A} \in \mathbb{R}^{m \times n}$ -ren heina  $r$  bada,  $r$  zutabe-permutazio beharko ditugu. Beraz,  $r$  urrats horiek egin ondoren, hau izango dugu:

$$(\mathbf{H}_r \dots \mathbf{H}_1) \mathbf{A} (\mathbf{P}_1 \dots \mathbf{P}_r) = \tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix},$$

non  $\mathbf{R}_{11}$  ez-singularra eta goi-triangeluarra baita,  $\mathbf{R}_{11}$ -ren zutabe kopuruak  $\mathbf{A}$ -ren heina ematen digu. Adibide honetan ikus daitezke  $n = 5$ ,  $m = 6$  eta  $r = 3$  kasuari dagozkion urratsak:



Orain,  $\mathbf{Q}^T = \mathbf{H}_r \dots \mathbf{H}_1$  eta  $\mathbf{P} = \mathbf{P}_1 \dots \mathbf{P}_r$  definitzen baditugu, honela geratuko da aurreko adierazpena:

$$\mathbf{Q}^T \mathbf{A} \mathbf{P} = \tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{matrix} r \\ m-r \\ r \\ n-r \end{matrix} \quad (6.22)$$

edo, baliokideki,  $\mathbf{A} \mathbf{P} = \mathbf{Q} \tilde{\mathbf{R}}$ . Ondorioz, minimo karratuen problema ebazteko, hau dugu ( $\mathbf{P}$  eta  $\mathbf{Q}^T$  ortogonalak direla kontuan hartuz):

$$\|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2 = \|\mathbf{Q}^T (\mathbf{A} \mathbf{x} - \mathbf{b})\|_2^2 = \|(\mathbf{Q}^T \mathbf{A} \mathbf{P})(\mathbf{P}^T \mathbf{x}) - \mathbf{Q}^T \mathbf{b}\|_2^2. \quad (6.23)$$

Izan bitez

$$\mathbf{P}^T \mathbf{x} = \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} \begin{matrix} r \\ n-r \end{matrix} \quad \text{eta} \quad \mathbf{Q}^T \mathbf{b} = \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} \begin{matrix} r \\ m-r \end{matrix},$$

orduan, (6.22) eta (6.23) ekuazioen bidez hau lortzen da:

$$\begin{aligned} \|\mathbf{Ax} - \mathbf{b}\|_2^2 &= \left\| \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} - \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} \right\|_2^2 = \left\| \begin{bmatrix} \mathbf{R}_{11}\mathbf{y} + \mathbf{R}_{12}\mathbf{z} - \mathbf{c} \\ -\mathbf{d} \end{bmatrix} \right\|_2^2 \\ &= \|\mathbf{R}_{11}\mathbf{y} - (\mathbf{c} - \mathbf{R}_{12}\mathbf{z})\|_2^2 + \|\mathbf{d}\|_2^2. \end{aligned} \quad (6.24)$$

Hori dela eta, minimo karratuen soluzioak  $\mathbf{R}_{11}\mathbf{y} = (\mathbf{c} - \mathbf{R}_{12}\mathbf{z})$  bete behar du. Hots, edozein  $\mathbf{z}$ -rako  $\mathbf{y} = \mathbf{R}_{11}^{-1}(\mathbf{c} - \mathbf{R}_{12}\mathbf{z})$  soluzio minimizatzaile bat izango dugu, eta, ondorioz:

$$\mathbf{x} = \mathbf{P} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = \mathbf{P} \begin{bmatrix} \mathbf{R}_{11}^{-1}(\mathbf{c} - \mathbf{R}_{12}\mathbf{z}) \\ \mathbf{z} \end{bmatrix}.$$

Baldin  $\mathbf{z} = \mathbf{0}$  hartzen badugu,  $\mathbf{x}_B$  oinarri-soluzio hau lortuko dugu:

$$\mathbf{x}_B = \mathbf{P} \begin{bmatrix} \mathbf{R}_{11}^{-1}\mathbf{c} \\ \mathbf{0} \end{bmatrix}.$$

Izen hori  $\mathbf{z} = \mathbf{0}$  hartzetik eta, praktikan,  $\mathbf{A}$ -ren hein beteko azpimatrizen baterako soluzioa izatetik dator.

### 6.3. algoritmoa. Hein urriko minimo karratu linealen problemaren ebazpena, QR faktORIZAZIOAREN BITARTEZ.

**0 urratsa.** SARRERA. Sartu:  $\mathbf{A}$  eta  $\mathbf{b}$ .

**1 urratsa.** Kalkula ezazu  $\mathbf{A}$ -ren QR faktORIZAZIOA Householderren transformazioak erabiliz,

$\mathbf{AP} = \mathbf{Q}\tilde{\mathbf{R}}$ , non  $\mathbf{Q}$  ortogonal baia eta  $\tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ , non  $\mathbf{R}_{11}$  goi-triangeluarra baia eta  $r = \text{hein}(\tilde{\mathbf{R}})$ .

**2 urratsa.** Kalkulatu  $\begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \mathbf{Q}^T\mathbf{b}$ , non  $\mathbf{c}$  bektorea  $\mathbf{Q}^T\mathbf{b}$ -ren lehenengo  $r$  osagaiek osatzen baitute.

**3 urratsa.** Ebatzi  $\mathbf{R}_{11}\mathbf{y} = (\mathbf{c} - \mathbf{R}_{12}\mathbf{z})$ . (Baldin  $\mathbf{z} = \mathbf{0}$  hartzen badugu,  $\mathbf{R}_{11}\mathbf{y}_B = \mathbf{c}$  ebatzen dugu.)

**4 urratsa.** Aurkitu  $\mathbf{x} = \mathbf{P} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix}$ . ( $\mathbf{z} = \mathbf{0}$  eta  $\mathbf{y}_B$  baditugu,  $\mathbf{x}_B = \mathbf{P} \begin{bmatrix} \mathbf{y}_B \\ \mathbf{0} \end{bmatrix}$ ).

**5 urratsa.** IRTEERA. Emaitza:  $\mathbf{x}$ . (Baldin  $\mathbf{z} = \mathbf{0}$  hartu badugu,  $\mathbf{x}_B$  oinarri-soluzioa izango dugu.)

Normalean, norma euklidear minimoa (luzera minimoa) duen  $\mathbf{x}_{LM}$  soluzioarekin geratzen gara. Alegia,

$$\|\mathbf{x}_{LM}\|_2 = \min_{\mathbf{z} \in \mathbb{R}^{n-r}} \left\| \mathbf{x}_B - \mathbf{P} \begin{bmatrix} \mathbf{R}_{11}^{-1}\mathbf{R}_{12} \\ -\mathbf{1}_{n-r} \end{bmatrix} \mathbf{z} \right\|_2. \quad (6.25)$$

Froga daiteke (ikus [11]) hau betetzen dela:

$$1 \leq \frac{\|\mathbf{x}_B\|_2}{\|\mathbf{x}_{LM}\|_2} \leq \sqrt{1 + \|\mathbf{R}_{11}^{-1}\mathbf{R}_{12}\|_2^2}.$$

**6.5. adibidea.** *Izan bitez*

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 2 \\ 7 & 6 & 10 \\ 4 & 4 & 6 \\ 1 & 0 & 1 \end{bmatrix} \quad \text{eta} \quad \mathbf{b} = \begin{bmatrix} 6 \\ 6 \\ 8 \\ 3 \end{bmatrix}.$$

Aurkitu  $\min \|\mathbf{Ax} - \mathbf{b}\|_2$  problemaren oinarri-soluzioa. Aurkitu  $\mathbf{x}_{LM}$  soluzioa ere (hots, luzera minimoa duena).

*Ebazpena.* Hirugarren zutabearen norma euklidearra handiena denez, 1. eta 3. zutabeak trukatu ditugu permutazio-matrize honen bitartez:

$$\mathbf{P} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix};$$

hots,

$$\mathbf{AP} = \begin{bmatrix} 2 & 2 & 1 \\ 10 & 6 & 7 \\ 6 & 4 & 4 \\ 1 & 0 & 1 \end{bmatrix}.$$

Householderren transformazioak erabiliz,  $\mathbf{AP} = \mathbf{Q}\tilde{\mathbf{R}}$  deskonposizio hau kalkulatu da:

$$\mathbf{Q} = \begin{bmatrix} -0.1684 & 0.7241 & -0.4453 & 0.4991 \\ -0.8422 & -0.2322 & -0.3881 & -0.2936 \\ -0.5053 & 0.2459 & 0.8049 & 0.1908 \\ -0.0842 & -0.6011 & -0.0571 & 0.7927 \end{bmatrix}, \quad \tilde{\mathbf{R}} = \begin{bmatrix} -11.87 & -7.411 & -8.169 \\ 0 & 1.038 & -0.5191 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix};$$

alegia,  $\mathbf{A}$ -ren heina  $r = 2$  da, eta

$$\mathbf{R}_{11} = \begin{bmatrix} -11.87 & -7.411 \\ 0 & 1.038 \end{bmatrix}, \quad \mathbf{R}_{12} = \begin{bmatrix} -8.169 \\ -0.5191 \end{bmatrix}, \quad \mathbf{Q}^T\mathbf{b} = \begin{bmatrix} -10.36 \\ 3.115 \\ 1.267 \\ 5.138 \end{bmatrix},$$

$\mathbf{Q}^T\mathbf{b}$  bektorearen lehenengo  $r = 2$  osagaiek  $\mathbf{c} = [-10.36 \ 3.115]^T$  bektorea osatzen dute, eta azken  $m - r = 4 - 2 = 2$  osagaiek  $\mathbf{d} = [1.267 \ 5.138]^T$  bektorea. Orain  $z = 0$  hartuz,  $\mathbf{R}_{11}\mathbf{y}_B = \mathbf{c}$  dugu, eta sistema hori ebatziz hau lortzen dugu:

$$\mathbf{y}_B = \begin{bmatrix} -1 \\ 3 \end{bmatrix} \Rightarrow \mathbf{x}_B = \mathbf{P} \begin{bmatrix} \mathbf{y}_B \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \\ -1 \end{bmatrix}.$$

Beraz,  $\min \|\mathbf{Ax} - \mathbf{b}\|_2 = \|\mathbf{Ax}_B - \mathbf{b}\|_2 = \|\mathbf{d}\|_2 = \sqrt{1.267^2 + 5.138^2} = 5.291$  hondar optimoa (minimoa) da.

Orain,  $\mathbf{x}_{LM}$  kalkulatu dugu, (6.25) erabiliz. Lehenik,  $\mathbf{w} = \mathbf{R}_{11}^{-1}\mathbf{R}_{12}$  aurkitu dugu; hots,  $\mathbf{R}_{11}\mathbf{w} = \mathbf{R}_{12}$  ebatzi behar dugu:

$$\begin{bmatrix} -11.87 & -7.411 \\ 0 & 1.038 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} -8.169 \\ -0.5191 \end{bmatrix} \Rightarrow \mathbf{w} = \begin{bmatrix} 1 \\ -0.5 \end{bmatrix}.$$

Gero, bektore hau kalkulatu dugu:

$$\mathbf{v} = \mathbf{x}_B - \mathbf{P} \begin{bmatrix} \mathbf{R}_{11}^{-1}\mathbf{R}_{12} \\ -\mathbf{1}_{n-r} \end{bmatrix} \mathbf{z} = \begin{bmatrix} 0 \\ 3 \\ -1 \end{bmatrix} - \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ -0.5 \\ -1 \end{bmatrix} \mathbf{z} = \begin{bmatrix} 0 \\ 3 \\ -1 \end{bmatrix} - \begin{bmatrix} -1 \\ -0.5 \\ 1 \end{bmatrix} \mathbf{z} = \begin{bmatrix} z \\ 3 + 0.5z \\ -1 - z \end{bmatrix}.$$

Azkenik (6.25) kontuan hartuz, hauxe dugu:

$$\|\mathbf{x}_{LM}\|_2 = \min_z \|\mathbf{v}\|_2 \Rightarrow \min\{z^2 + (3+0.5z)^2 + (-1-z)^2\} = \min\{2.25z^2 + 5z + 10\} \Rightarrow z = -1.111$$

Beraz,

$$\mathbf{x}_{LM} = \begin{bmatrix} -1.111 \\ 3 + 0.5 \cdot (-1.111) \\ -1 - (-1.111) \end{bmatrix} = \begin{bmatrix} -1.111 \\ 2.444 \\ 0.111 \end{bmatrix},$$

non  $\|\mathbf{x}_{LM}\|_2 = 2.687$  baita ( $\|\mathbf{x}_B\|_2 = 3.162$  da) .  $\square$

## 6.7. Deskonposizio ortogonal osoa

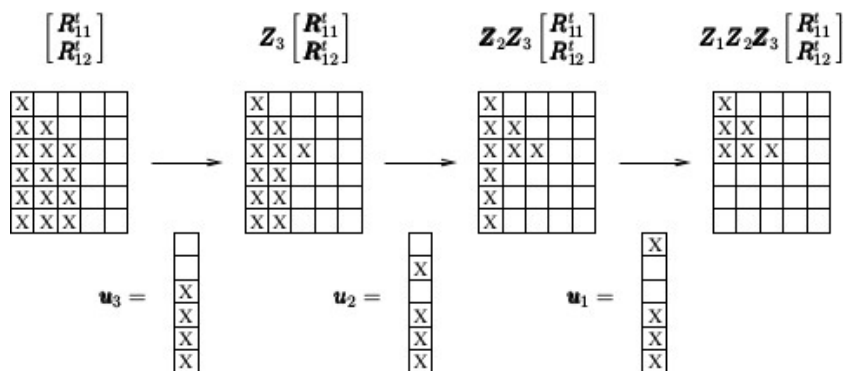
Hein urriko  $QR$  faktORIZAZIOAN  $\tilde{\mathbf{R}}$  matrize hau lortu dugu (ikus (6.22)):

$$\mathbf{Q}^T \mathbf{A} \mathbf{P} = \tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{matrix} r \\ m-r \end{matrix} \\ \begin{matrix} r & n-r \end{matrix}$$

Gainera,  $\mathbf{Z}_i$  Householderren matrize bereziak erabiliz, hau lor dezakegu ( $\mathbf{T}$  definituz):

$$\mathbf{T}^T = \mathbf{Z}_1 \dots \mathbf{Z}_r \begin{bmatrix} \mathbf{R}_{11}^T & \mathbf{0} \\ \mathbf{R}_{12}^T & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{T}_{11}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{matrix} r \\ n-r \end{matrix},$$

$\mathbf{T}_{11}^T$  behe-triangeluarra erdietsiz. Hori lortzeko, honela jokatzen da (ikus [12]-ko 193-194 orrialdeak):



Hots,  $T^T = Z_1 \dots Z_r (Q^T A P)^T = Z_1 \dots Z_r P^T A^T Q$ . Ondorioz:

$$T = Q^T A P Z_r \dots Z_1 = Q^T A Z = \begin{bmatrix} T_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{matrix} r \\ m-r \\ r & n-r \end{matrix}, \quad (6.26)$$

non  $Z = P Z_r \dots Z_1$  ere ortogonala baita. Deskonposizio horri *deskonposizio ortogonal osoa* deritzogu.

Ondorioz, minimo karratuen probleman hau betetzen da:

$$\begin{aligned} \|A\mathbf{x} - \mathbf{b}\|_2^2 &= \|Q^T(A\mathbf{x} - \mathbf{b})\|_2^2 = \|(Q^T A Z)Z^T \mathbf{x} - Q^T \mathbf{b}\|_2^2 = \left\| T \begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} \right\|_2^2 = \\ &= \left\| \begin{bmatrix} T_{11} \mathbf{w} - \mathbf{c} \\ -\mathbf{d} \end{bmatrix} \right\|_2^2 = \|T_{11} \mathbf{w} - \mathbf{c}\|_2^2 + \|\mathbf{d}\|_2^2, \end{aligned}$$

non

$$Z^T \mathbf{x} = \begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \begin{matrix} r \\ n-r \end{matrix} \quad \text{eta} \quad Q^T \mathbf{b} = \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} \begin{matrix} r \\ m-r \end{matrix}.$$

Bistan denez, baldin  $\mathbf{x}$ -k minimo karratuak minimizatu behar baditu, orduan  $\mathbf{w} = T_{11}^{-1} \mathbf{c}$  bete behar du. Gainera,  $\mathbf{x}$ -ren norma euklidearra minimoa izateko  $\mathbf{y}$  zero izan behar, eta, horrela bada, hau dugu:

$$\mathbf{x}_{LM} = Z \begin{bmatrix} \mathbf{w} \\ \mathbf{0} \end{bmatrix} = Z \begin{bmatrix} T_{11}^{-1} \mathbf{c} \\ \mathbf{0} \end{bmatrix}.$$

## 6.8. Balio singularretako deskonposizioa

Edozein  $m \times n$   $A$  matrizea honela idatz dezakegu:

$$A = U \Sigma V^T,$$





**6.5. teorema.** *Izan bedi  $\mathbf{A} \in \mathbb{R}^{m \times n}$  eta  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  bere balio singularretako deskonposizioa. Izan bitez  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  balio singular ez-nuluak eta  $\sigma_{r+1} = \dots = \sigma_p = 0$ , non  $p = \min\{m, n\}$ . Orduan baieztapen hauek betetzen dira:*

1.  $\mathbf{A}$ -ren heina  $r$  da.
2.  $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$  bektoreek  $K(\mathbf{A})$ -ren oinarri ortonormal bat osatzen dute.
3.  $\{\mathbf{u}_{r+1}, \mathbf{u}_2, \dots, \mathbf{u}_m\}$  bektoreek  $K(\mathbf{A})$  azpiespazioaren  $\mathbb{R}^m$ -ko  $K(\mathbf{A})^\perp$  azpiespazio osagarri ortogonaleko oinarri ortonormal bat osatzen dute.
4.  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$  bektoreek  $K(\mathbf{A}^T)$ -ren oinarri ortonormal bat osatzen dute.
5.  $\{\mathbf{v}_{r+1}, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  bektoreek  $K(\mathbf{A}^T)$  azpiespazioaren  $\mathbb{R}^n$ -ko  $K(\mathbf{A}^T)^\perp$  azpiespazio osagarri ortogonaleko oinarri ortonormal bat osatzen dute.

### 6.8.1. Moore-Penrose-ren sasiialderantzizko matrizea

$\mathbf{A}$  matrize ez-nulu baten alderantzizkoaren orokortze klasiko bat *Moore-Penroseren sasiialderantzizkoa* da,  $\mathbf{A}^+$ .  $\mathbf{A}$  karratua eta ez-singularra denean,  $\mathbf{A}^+ = \mathbf{A}^{-1}$  dugu.

Izan bedi  $\mathbf{A}$   $m \times n$  matrize bat eta  $p = \min\{m, n\}$ . Demagun  $\mathbf{A}$  matrizearen balio singularretako deskonposizioa  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  dela. Orduan,  $\text{hein}(\mathbf{A}) = r$  bada,  $r \leq p$  izango da. Beraz, zera izango dugu:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0 \quad \text{eta} \quad \sigma_{r+1} = \dots = \sigma_p = 0.$$

$\mathbf{D} \in \mathbb{R}^{m \times n}$  matrize diagonal baten sasiialderantzizkoa honela definitzen dugu:  $\mathbf{D}^+$   $n \times m$  matrize diagonal da, non diagonaleko  $d_i^+$  gaiak honela definitzen baitira:

$$d_i^+ = \begin{cases} 1/d_i, & d_i \neq 0 \text{ bada,} \\ 0, & \text{bestela.} \end{cases}$$

Orduan,  $\mathbf{A}$  matrize orokor baterako  $\mathbf{A}^+$  honela definitzen da:

$$\mathbf{A}^+ = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^+ = (\mathbf{V}^T)^+ \mathbf{\Sigma}^+ \mathbf{U}^+.$$

Eta  $\mathbf{U}$  eta  $\mathbf{V}^T$  matrize ortogonalak direnez, haien sasiialderantzizko matrizeak alderantzizko arruntak dira, beraz:  $(\mathbf{V}^T)^+ = \mathbf{V}$  eta  $\mathbf{U}^+ = \mathbf{U}^T$ . Ondorioz,  $\mathbf{A}$  matrize orokor baten sasiialderantzizkoak hau betetzen du:

$$\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T.$$

Baina, ohartu  $\mathbf{\Sigma}^+$ -ren gai ez-nuluak ez daudela ordenatuta handienetik txikienera  $\mathbf{\Sigma}$ -n geratzen den bezala.

**Propietate batzuk**

- 1)  $(\mathbf{A}^+)^+ = \mathbf{A}$ .
- 2)  $(\mathbf{A}^T)^+ = (\mathbf{A}^+)^T$ .
- 3)  $(\mathbf{A}\mathbf{A}^+)^T = \mathbf{A}\mathbf{A}^+$ .
- 4)  $(\mathbf{A}^+\mathbf{A})^T = \mathbf{A}^+\mathbf{A}$ .
- 5)  $\mathbf{A}\mathbf{A}^+\mathbf{v} = \mathbf{v}, \forall \mathbf{v} \in K(\mathbf{A})$ .
- 6)  $\mathbf{A}^+ = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$ ,  $\mathbf{A}$ -ren zutabeak linealki askeak badira ( $\mathbf{A}$  ezkerraldeetik alderantzikagarria da,  $\mathbf{A}^+\mathbf{A} = \mathbb{1}$ ).
- 7)  $\mathbf{A}^+ = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}$ ,  $\mathbf{A}$ -ren lerroak linealki askeak badira ( $\mathbf{A}$  eskuinaldeetik alderantzikagarria da,  $\mathbf{A}\mathbf{A}^+ = \mathbb{1}$ ).
- 8)  $\mathbf{A}^+\mathbf{z} = \mathbf{0}, \forall \mathbf{z} \in K(\mathbf{A})^\perp$ .

**6.8.2. Baldintzazko zenbaki orokortua**

$\mathbf{A}$  matrizea hein betekoa bada, *baldintzazko zenbakia* honela definitzen da:

$$\kappa(\mathbf{A}) = \|\mathbf{A}^+\| \|\mathbf{A}\|.$$

Definizio hori zabaldu daiteke  $\mathbf{A}$  matrize orokor baterako, eta, orduan, *baldintzazko zenbaki orokortua* izango dugu.

Jarraian, definizio orokortu hori aplikatuko diogu norma euklidearrerako baldintzazko zenbakiari. Dakigunez, matrize bat matrize ortogonal batez biderkatzen badugu, ez da aldatzen norma euklidearra. Beraz,  $\|\mathbf{A}\|_2 = \|\boldsymbol{\Sigma}\|_2 = \sigma_1$  dugu,  $\mathbf{A}$ -ren balio singular handiena. Orain  $\|\mathbf{A}^+\|_2$  kalkulatu dugu.

$\mathbf{A}$ -ren heina  $r$  bada,  $\sigma_r$  da balio singular ez-nulu txikiena. Gainera,  $\boldsymbol{\Sigma}^+$  diagonaleko gaiak  $\mathbf{A}$ -ren balio singular ez-nuluaren erreziprokoak direnez,  $\mathbf{A}^+ = \mathbf{V}\boldsymbol{\Sigma}^+\mathbf{U}^T$  berdintzak erakusten du  $\mathbf{A}^+$ -ren balio singular handiena  $1/\sigma_r$  dela. Beraz,  $\|\mathbf{A}^+\|_2 = 1/\sigma_r$  eta  $\mathbf{A}^+$ -ren bi-normarekiko baldintzazko zenbaki orokortua hau da:

$$\kappa_2(\mathbf{A}) = \|\mathbf{A}^+\|_2 \|\mathbf{A}\|_2 = \frac{\sigma_1}{\sigma_r}.$$

Beste kasuetan bezala, baldintzak erakusten du sistema linealen sentzibilitatea datuen aldatetarekiko. Baina,  $\mathbf{A}$ -ren zutabeen menpekotasun linealari buruzko informazioa ere ematen du. Adibidez, matrize hauetarako:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0.0001 \end{bmatrix}$$



## 6.9. Problemak

### Eskuz ebazteko problemak:

1. (i) Izan bitez sistema hauek:

$$\begin{array}{ll}
 a) & \begin{array}{l} x_1 + 3x_2 + 2x_3 = 5 \\ x_1 + 5x_2 + 3x_3 = 10 \\ 2x_1 + 4x_2 - 6x_3 = -4. \end{array} \\
 b) & \begin{array}{l} 2x_1 + x_2 + x_3 = 3 \\ 4x_1 + 3x_2 + 3x_3 + x_4 = 7 \\ 8x_1 + 7x_2 + 9x_3 + 5x_4 = 17 \\ 6x_1 + 7x_2 + 9x_3 + 8x_4 = 15. \end{array} \\
 c) & \begin{array}{l} -5x_1 + 2x_2 - x_3 = 2 \\ x_1 + 3x_3 = 2 \\ 3x_1 + x_2 + 6x_3 = 2. \end{array} \\
 d) & \begin{array}{l} 2x_1 - x_2 + x_3 = -1 \\ 3x_1 + 3x_2 + 9x_3 = 0 \\ 3x_1 + 3x_2 + 5x_3 = 4. \end{array} \\
 e) & \begin{array}{l} x_1 + x_2 + 4x_4 = 5 \\ 2x_1 - x_2 + 5x_3 = -6 \\ 5x_1 + 2x_2 + x_3 + 2x_4 = 3 \\ -3x_1 + 2x_3 + 6x_4 = 4. \end{array} \\
 f) & \begin{array}{l} 4x_1 + 8x_2 + 4x_3 = 8 \\ x_1 + 5x_2 + 4x_3 - 3x_4 = -4 \\ x_1 + 4x_2 + 7x_3 + 2x_4 = 10 \\ x_1 + 3x_2 - 2x_4 = -4. \end{array}
 \end{array}$$

Aurkitu sistemei elkartutako matrizeen  $QR$  faktORIZAZIOA, erabilitako Householderren bektoreak emanez eta Householderren matrizeak kalkulatu barik. Erakutsi, urratsez urrats, egindako eragiketa nagusiak.

- (ii) Ebatzi sistemak, aurreko atalean aurkitutako  $QR$  faktORIZAZIOAK erabiliz.
2. (a) Izan bedi  $\mathbf{x} = [1, 12, 2, 5, 7]^T$  bektorea. Kalkula ezazu biraketa bat  $(4, 2)$  planoan  $x_2$  zero bihurtzeko. Zenbat radianetako biraketa eman behar du  $\mathbf{x}$  bektoreak plano horretan emaitza hori lortzeko?
- (b) Bektorea errenkada bat bada, esate baterako  $\mathbf{x} = [1, 12, 2, 5, 7]$ , nola egingo zenuke  $(4, 2)$  planoan  $x_2$  zero bihurtzeko?
3. Aplikatu Givens-en biraketak matrize honi,  $(1, 2)$  planoan lehenengo zutabeko bigarren lerroko gaia zero bihurtzeko:

$$\begin{bmatrix} 2 & 1 & 0 & 0 \\ 4 & 3 & 2 & 0 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & 1 & 2 \end{bmatrix}.$$

- (a) Horretarako,  $c$  eta  $s$  balioak honela kalkulatu ditugu (ikus 6.3. atala):

$$c = \cos(\theta) = \frac{x_i}{\sqrt{x_i^2 + x_k^2}} \quad \text{eta} \quad s = \sin(\theta) = \frac{-x_k}{\sqrt{x_i^2 + x_k^2}},$$

non  $(i, k)$  bikoteak esaten baitigu zein planotan egiten den  $\theta$  angeluko biraketa, bektore baten  $i$ . lerroko gaiaz  $k$ . lerroko gaia zero bihurtzeko.

- (b) Gero, kalkula ezazu dagokion  $\mathbf{G}(i, k, \theta)$  matrizea.
- (c) Kalkula ezazu  $\mathbf{G}(i, k, \theta)^T \mathbf{A}$ .
- (d) Diseinatu algoritmo bat  $\mathbf{M} \in \mathbb{R}^{m \times m}$  matrize orokor batekin  $\mathbf{G}(i, k, \theta)^T \mathbf{M}$  kalkulua egiteko.  $c$ -ren eta  $s$ -ren kalkulua algoritmoaren barnean egin daiteke.
- (e) Bigarren zutabeko 2. lerroko gaiatz 3. lerroko gaia zero egiteko, zein izango litzateke  $\mathbf{G}(i, k, \theta)$  Givensen matrizea?

4. Izan bitez  $\mathbf{Ax} = \mathbf{b}$  sistema gaindeterminatu hauek:

$$a) \quad \mathbf{A} = \begin{bmatrix} 1 & 1 \\ 2 & -3 \\ 0 & 0 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix}. \quad b) \quad \mathbf{A} = \begin{bmatrix} -1 & 1 \\ 2 & 1 \\ 1 & -2 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} 10 \\ 5 \\ 20 \end{bmatrix}.$$

$$c) \quad \mathbf{A} = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} 4 \\ 0 \\ 1 \\ 2 \end{bmatrix}. \quad d) \quad \mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}.$$

$$e) \quad \mathbf{A} = \begin{bmatrix} 1 & 1 \\ -1 & 3 \\ 1 & 2 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} -2 \\ 0 \\ 8 \end{bmatrix}. \quad f) \quad \mathbf{A} = \begin{bmatrix} 1 & -1 & 0 \\ 3 & -1 & 2 \\ -1 & 5 & 4 \\ 0 & 2 & 2 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} 1 \\ 5 \\ 3 \\ 4 \end{bmatrix}.$$

Ebatzi sistema horiek ekuazio normalak erabiliz. Zein da beren hondarraren norma euklidearra?

5. Ebatzi aurreko problemaren sistemak,  $QR$  faktORIZAZIOAREN bitartez. Kalkula ezazu hondarraren norma euklidearra.
6. Izan bedi  $\mathbf{Ax} = \mathbf{b}$  sistema gaindeterminatua, non

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0.001 \end{bmatrix} \text{ eta } \mathbf{b} = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix}.$$

- (a) Zein da ekuazio normalei dagokien matrizearen bi-normarako baldintzazko zenbakia?
- (b) Zein da  $\mathbf{A}$ -ren bi-normarako baldintzazko zenbaki orokortua? Zein dira  $QR$  faktORIZAZIOARAKO  $\mathbf{R}$  matrizea eta  $\mathbf{R}$ -ren bi-normarako baldintzazko zenbaki orokortua? Galdera hauek erantzuteko, kontuan hartu 6.8.2. atala.
- (c) Konparatu aurreko bi ataletan lortutako baldintzazko zenbakiak. Zer esan dezakegu?
- (d) Ebatzi sistema hori bi metodoekin, lau zifra esanguratsuz, eta konpara itzazu emaitzak.

7. Aurkitu  $f(x) = ax + b$  funtzio lineal hoberena, minimo karratuen zentzuan, datu multzo hauetarako:

$$a) \quad \frac{x}{y} \left| \begin{array}{cccc} -1 & 0 & 1 & 2 \\ 3 & 2 & 0 & 4 \end{array} \right. \quad b) \quad \frac{x}{y} \left| \begin{array}{cccc} -3 & -1 & 1 & 3 \\ 15 & 5 & 1 & 5 \end{array} \right.$$

Zein dira hondarrak? Marraztu  $a)$  eta  $b)$  kasuetan (bakoitza grafiko desberdin batean) funtzio linealak eta taulen puntuak.

8. Aurkitu  $f(x) = ax^2 + bx + c$  parabola hoberena, minimo karratuen zentzuan, datu multzo hauetarako:

$$a) \quad \frac{x}{y} \left| \begin{array}{cccc} -1 & 0 & 1 & 2 \\ 3 & 2 & 0 & 4 \end{array} \right. \quad b) \quad \frac{x}{y} \left| \begin{array}{cccc} -3 & -1 & 1 & 3 \\ 15 & 5 & 1 & 5 \end{array} \right.$$

Zein dira hondarrak? Marraztu  $a)$  eta  $b)$  kasuetan (bakoitza grafiko desberdin batean) funtzio koadratikoa eta taulen puntuak. Konparatu lortutako  $a)$  kasuko grafikoa aurreko problemari dagokion  $a)$  kasuko grafikoarekin. Egin berdina  $b)$  kasuetan. Zer ondorioztatzen da?

9. Izan bedi matrize hau:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 1 & 2.5 \\ 1 & 0 & 2 \end{bmatrix}.$$

- Aurkitu  $K(\mathbf{A})$  azpiespazioaren oinarri ortonormal bat,  $QR$ -faktORIZAZIOA erabiliz.
- Aurkitu  $K(\mathbf{A})$ -ren azpiespazio nuluen oinarri ortonormal bat,  $QR$ -faktORIZAZIOA erabiliz.
- Aurkitu  $\mathbf{Ax} = \mathbf{b}$  sistemaren  $\mathbf{x}_B$  oinarri-soluzioa, baldin  $\mathbf{b} = [1 \ 3 \ 5]^T$  bada.
- Aurkitu  $\mathbf{Ax} = \mathbf{b}$  sistemaren  $\mathbf{x}_{LM}$  soluzioa; hots, norma euklidear minimoa duen minimo karratuen problemaren soluzioa.

10. Izan bedi matrize hau:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 5 & 6 \\ 1 & 8 & 9 \\ 1 & 11 & 12 \end{bmatrix}.$$

- Hein urrikoa da? Arrazoitu erantzuna.
- Kalkula ezazu bere  $QR$ -faktORIZAZIOA, zutabeen permutazioak eginez, beharrezkoa bada.
- Aurkitu  $K(\mathbf{A})$  azpiespazioaren oinarri ortonormal bat,  $QR$ -faktORIZAZIOA erabiliz.
- Aurkitu  $K(\mathbf{A})$ -ren azpiespazio nuluen oinarri ortonormal bat,  $QR$ -faktORIZAZIOA erabiliz.
- Aurkitu  $\mathbf{Ax} = \mathbf{b}$  sistemaren  $\mathbf{x}_B$  oinarri-soluzioa, baldin  $\mathbf{b} = [2 \ 1 \ 4 \ 3]^T$  bada.
- Aurkitu  $\mathbf{Ax} = \mathbf{b}$  sistemaren  $\mathbf{x}_{LM}$  soluzioa; hots, norma euklidear minimoa duen minimo karratuen problemaren soluzioa.

11. Izan bedi  $\mathbf{A} \in \mathbb{R}^{m \times n}$  matrizea eta  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  bere balio singularretako deskonposizioa, non  $\mathbf{A}$ -ren heina  $r$  baita eta  $m \geq n \geq r$  ( $n \geq m \geq r$  ere izan liteke). Izan bedi  $\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T \in \mathbb{R}^{n \times m}$  non

$$\mathbf{\Sigma}^+ = \text{diag} \left\{ \frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r}, 0, \dots, 0 \right\}.$$

$\mathbf{A}^+$  matrizeari  $\mathbf{A}$ -ren Moore-Penroseren sasiialderantzizko matrize deritzo.

Froga ezazu  $\mathbf{X} = \mathbf{A}^+$ -ek Moore-Penroseren lau baldintza hauek betetzen dituela:

- (i)  $\mathbf{A}\mathbf{X}\mathbf{A} = \mathbf{A}$
- (ii)  $\mathbf{X}\mathbf{A}\mathbf{X} = \mathbf{X}$
- (iii)  $(\mathbf{A}\mathbf{X})^T = \mathbf{A}\mathbf{X}$
- (iv)  $(\mathbf{X}\mathbf{A})^T = \mathbf{X}\mathbf{A}$

Ikus 6.8.1. atala.

12. Izan bedi  $\mathbf{A} \in \mathbb{R}^{m \times n}$  matrizea eta  $\mathbf{P} \in \mathbb{R}^{m \times m}$  eta  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  matrize ortogonalak. Froga ezazu  $\|\mathbf{P}\mathbf{A}\mathbf{Q}\|_2 = \|\mathbf{A}\|_2$  betetzen dela.

13. Izan bedi  $\mathbf{A} \in \mathbb{R}^{m \times n}$   $r$  heineko matrize bat eta  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  balio singularretako deskonposizioa, non  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  balio singular ez-nuluak eta  $\sigma_{r+1} = \dots = \sigma_n = 0$  baitira.

- (a) Froga ezazu  $\|\mathbf{A}\|_2 = \|\mathbf{\Sigma}\|_2 = \sigma_1$ .
- (b) Froga ezazu  $r = n$  bada,  $\kappa_2(\mathbf{A}) = \sigma_1/\sigma_n$ .
- (c) Baldin balio singular ez-nuluak  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  badira, froga ezazu  $\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$ .  
 $\mathbf{A} = \widehat{\mathbf{U}} \widehat{\mathbf{\Sigma}} \widehat{\mathbf{V}}^T$  adierazpenari balio singularretako deskonposizio laburtu deritzogu, non

$$\begin{aligned} \widehat{\mathbf{\Sigma}} &= \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_r\} \\ \widehat{\mathbf{U}} &= [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r] \\ \widehat{\mathbf{V}} &= [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r] \end{aligned}$$

14. (a) Froga ezazu balio singularretako deskonposiziotik  $\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}$  atera dezakegula.
- (b) Aurreko berdintza erabiliz, froga ezazu  $\mathbf{A}^T \mathbf{A}$  eta  $\mathbf{A}\mathbf{A}^T$  autobalio berdinak dituztela eta balio singularren karratuak direla.
- (c) Froga ezazu  $\mathbf{A}$ -ren heina  $\mathbf{A}^T \mathbf{A}$ -ren heinaren berdina dela, eta hori dela balio singular ez-nuluaren kopuruaren berdina.

15. (a) Froga ezazu espazio bektorial baten oinarri ortonormal baten bektoreak matrize ortogonal batez biderkatzean lortutako bektore berriek espazio horretako beste oinarri ortonormal bat ematen dutela.



- (b) Espazio bektorial baten  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$  oinarri kanonikoa  $\mathbf{Q}$  matrize ortogonal batez biderkatzen badugu, zein da lortzen den oinarri ortonormal berria?
16. Izan bedi  $\mathbf{A} \in \mathbb{R}^{m \times n}$   $r$  heineko matrize bat eta  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  balio singularretako deskonposizioa. Izan bitez  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$  ezker bektore singularrak,  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$  eskuin bektore singularrak eta  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  balio singular ez-nuluak eta  $\sigma_{r+1} = \dots = \sigma_p = 0$ , non  $p = \min\{m, n\}$ .
- (a) Froga ezazu  $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r$   $\mathbb{R}^m$ -ko bektoreak ortogonalak direla.
- (b) Kalkula ezazu  $\mathbf{A}^T\mathbf{A}$ , eta hortik ondorioztatu nola kalkula ditzakezun  $\sigma_i$  balio singularrak eta  $\sigma_i > 0$  betetzen duten balioei dagozkien  $\mathbf{v}_i$  eskuin bektore singularrak.
- (c)  $\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}$  erabiliz, azaldu nola kalkula dezakezun  $\sigma_i > 0$  bakoitzari dagokion  $\mathbf{u}_i$  ezker bektore singularra.
- (d) Froga ezazu  $\mathbf{u}_1, \dots, \mathbf{u}_r$   $\mathbb{R}^m$ -ko bektoreek  $K(\mathbf{A})$  azpiespazioaren oinarri ortonormal bat osatzen dutela (kontuan hartu  $\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}$ ).
- (e) Froga ezazu  $\mathbf{v}_1, \dots, \mathbf{v}_r$   $\mathbb{R}^n$ -ko bektoreek  $K(\mathbf{A}^T)$  azpiespazioaren oinarri ortonormal bat osatzen dutela (kontuan hartu  $\mathbf{A}^T$ -ren balio singularretako deskonposiziotik  $\mathbf{A}^T\mathbf{U} = \mathbf{V}\mathbf{\Sigma}^T$  lortzen dela).
- (f) Froga ezazu  $\mathbf{V}$  matrizeko  $\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$   $\mathbb{R}^n$ -ko zutabe-bektoreek  $K(\mathbf{A}^T)^\perp$  azpiespazioaren oinarri ortonormal bat osatzen dutela (kontuan hartu  $\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}$ ).
- (g) Froga ezazu  $\mathbf{U}$  matrizeko  $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$   $\mathbb{R}^m$ -ko zutabe-bektoreek  $K(\mathbf{A})^\perp$  azpiespazioaren oinarri ortonormal bat osatzen dutela (kontuan hartu  $\mathbf{A}^T\mathbf{U} = \mathbf{V}\mathbf{\Sigma}^T$ ).
- (h)  $\mathbf{A}\mathbf{A}^T$  erabiliz, nola kalkula ditzakezu  $\sigma_i$  balio singularrak eta  $\sigma_i > 0$  betetzen duten balioei dagozkien  $\mathbf{u}_i$  ezker bektore singularrak? Gero, nola kalkulatu zenuke  $\sigma_i > 0$  bakoitzari dagokion  $\mathbf{v}_i$  ezker bektore singularra?

17. (a) Izan bedi matrize hau:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

16. problemaren emaitzak erabiliz, aurkitu  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  balio singularretako deskonposizioa, eta egiaztatu berdintza hori betetzen dela.

- (b)  $\mathbf{b} = [3, 2, 1]^T$  bada,  $\mathbf{A}$ -ren deskonposizio hori aprobeztatuz, ebatzi  $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$  minimo karratuen problema  $\mathbf{x}$  soluzioa norma minimokoa izanik.

18. Izan bedi matrize hau:

$$\mathbf{A} = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix}.$$

16. problemaren emaitzak erabiliz, aurkitu  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  balio singularretako deskonposizioa, eta egiaztatu berdintza hori betetzen dela.

**MATLABez programatzeko problemak:**

19. Sortu M fitxategi bat  $\mathbf{Ax} = \mathbf{b}$  sistema lineal bat ( $\mathbf{A}$  karratua eta ez-singularra denean) ebazteko gai izan dadin QR faktORIZAZIOAREN laguntzaz, baina  $\mathbf{Q}$  matrize ortogonalak eta Householderren matrizeak kalkulatu barik. Aldiz, gorde itzazu lortutako Householderren bektoreak  $\mathbf{U}$  matrize batean.

Programa horrek sistema hori ebazteko ondoz ondoko lan hauek egin behar ditu:

- Programa horrek goiburua hau izan behar du: `function [U,R,x]=qrkes(A,b)`
- Egiaztatu  $\mathbf{A}$  karratua eta ez-singularra dela.
- $\mathbf{R} = \mathbf{Q}^T \mathbf{A}$  eta  $\mathbf{c} = \mathbf{Q}^T \mathbf{b}$  kalkulatu behar ditu.
- $\mathbf{Rx} = \mathbf{c}$  sistema ebatzi behar du. Lan hori egiteko, `atzerantz` funtzioa erabili.

20. Izan bedi sistema hau:

$$\begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & -1 & 3 & 5 \\ 0 & 0 & 2 & 5 \\ -2 & -6 & -3 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}.$$

Egiaztatu `qrkes.m` funtzioa ondo dabilela aurreko sistema ebatziz. Egiaztatu emaitza.

21. Sortu M fitxategi bat  $\min_x \|\mathbf{Ax} - \mathbf{b}\|_2$  minimo karratuen problema ebazteko, non  $\mathbf{A} \in \mathbb{R}^{m \times n}$  eta  $m > n$ .  $\mathbf{A}$  hein betekoa dela suposatuko dugu. Bi metodo erabiliko ditugu.

- Ekuazio normalen metodoa.* Kasu honetan,  $\mathbf{A}^T \mathbf{A}$  matrizea deskonposatu Choleskyren faktORIZAZIOA erabiliz. Gero, `aurrerantz` eta `atzerantz` funtzioen bidez ebatzi sistema. Sortu funtzio berri hau:

```
function [h,x]=eknorm(A,b)
```

- QR metodoa.* Kasu honetan, aprobeitza dezakezu lehenengo ariketan egindako `qrkes` funtzioa. Baina, orain  $\mathbf{A}$  matrizea ez da karratua izan behar. Sortu funtzio berri hau:

```
function [h,x]=qrmknb(A,b)
```

Bi metodoetan  $h$  hondar bektorearen bi-norma da.

22. Aurreko ariketan lortutako bi metodoak erabiliz, ebatzi  $\min_x \|\mathbf{Ax} - \mathbf{b}\|_2$  minimo karratuen problema balio hauetarako:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{eta} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ -5 \end{bmatrix}.$$

Problema hau 6.3. eta 6.4. adibideetan agertzen da.

## 7. kapitulua

# Ekuzio ez-linealen sistemen ebazpena

Ikasgai honetan aldagai anitzeko problemekin arituko gara, eta gure helburua izango da ekuzio ez-linealen sistemak ebaztea, algoritmo lokalen bidez. Newtonen metodoarekin hasiko gara, eta haren implementazioa eta ezaugarriak aztertuko ditugu.

### 7.1. Newtonen metodoa

Izan bedi  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  funtzioa eta  $\mathbf{f} \in C^1$  (hots, lehenengo ordenako deribatu partzial guztiak funtzio jarraituak dira), non  $f_1, \dots, f_n$  bere osagaiak baitira. Ebatzi behar dugun ekuzioa zera da:

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}. \quad (7.1)$$

Aldagai bateko funtzioen kasurako bezala, hemen ere metodo iteratibo bat erabiliko dugu  $\mathbf{x}_0 \in \mathbb{R}^n$  puntu batetik abiatuz. Demagun  $k$ -garren iterazioan dagoela prozesu hori; orduan, puntu horren inguruan, hau da Taylorren lehenengo ordenako garapena:

$$\mathbf{M}_k(\mathbf{x}) = \mathbf{f}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k),$$

non  $\mathbf{J}(\mathbf{x}_k)$  matrizea  $\mathbf{f}$ -ren jacobiarra baita  $\mathbf{x}_k$  puntuan (hots,  $\mathbf{f}$ -ren lehenengo ordenako deribatua da puntu horretan); honela definitzen da:

$$\mathbf{J}(\mathbf{x}_k) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}_k) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}_k) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1}(\mathbf{x}_k) & \cdots & \frac{\partial f_n}{\partial x_n}(\mathbf{x}_k) \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x}_k)^T \\ \vdots \\ \nabla f_n(\mathbf{x}_k)^T \end{bmatrix}$$

Jarraian, (7.1) ekuzioan  $\mathbf{f}(\mathbf{x})$ -ren ordeztu  $\mathbf{x}_k$ -ren inguruko  $\mathbf{M}_k(\mathbf{x})$  hurbilpen lineala jarriko dugu; alegia:

$$\mathbf{f}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) = \mathbf{0}.$$

Ekuazio horren soluzioa  $\mathbf{x} = \mathbf{x}_k - \mathbf{J}(\mathbf{x}_k)^{-1}\mathbf{f}(\mathbf{x}_k)$  denez, hau da Newtonen metodoaren iterazioa:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{J}(\mathbf{x}_k)^{-1}\mathbf{f}(\mathbf{x}_k). \quad (7.2)$$

Newtonen  $\mathbf{x}_{k+1} - \mathbf{x}_k$  urratsak  $\mathbf{x}^* - \mathbf{x}_k$ -ren hurbilpen bat ematen digu.

Geometrikoki, aldagai bateko funtzioetarako, Newton-Raphsonen metodoan  $f(x)$ -ren hurbilpena  $x_k$  puntuko zuzen ukitzaila da. Ekuazio ez-linealen sistemetarako,  $f_i(\mathbf{x}) = 0$  ekuazio bakoitzeko  $f_i(\mathbf{x})$  hurbiltzen da  $\mathbf{x}_k$  puntuko (hiper)plano ukitzailaz, hots:

$$f_i(\mathbf{x}_k) + \nabla f_i(\mathbf{x}_k)^T(\mathbf{x} - \mathbf{x}_k) = 0$$

eta ekuazio linealen sistema horren soluzioa, (7.2) adierazpenak definitutakoa, plano ukitzaila horien guztien ebakitze-puntua da.

**7.1. adibidea.** *Izan bitez funtzio hauek:*

$$\begin{aligned} f_1(x, y) &= x^2 - 2x - y + 0.5 \\ f_2(x, y) &= x^2 + 4y^2 - 4. \end{aligned}$$

*Gure helburua ekuazio-sistema hau ebaztea da:*

$$\begin{aligned} f_1(x, y) &= 0 \\ f_2(x, y) &= 0. \end{aligned}$$

*Ebazpena.*  $f_1(x, y) = 0$  eta  $f_2(x, y) = 0$  ekuazioek  $XOY$  planoko bi kurba definitzen dituzte; beraz, sistema horren soluzioa bi kurben  $(\bar{x}, \bar{y})$  ebakitze-puntu bat da. Sistemako kurbak oso ezagunak dira:

$$\begin{aligned} x^2 - 2x - y + 0.5 = 0 & \quad \text{parabola bat da,} \\ x^2 + 4y^2 - 4 = 0 & \quad \text{elipse bat da.} \end{aligned}$$

Sistema horrek bi soluzio (bi ebakitze-puntu) ematen ditu,  $(-0.2, 1.0)$  eta  $(1.9, 0.3)$  puntuetatik hurbil daudenak.

Newtonen metodoan,  $(x_0, y_0) = (0, 1)$  puntutik abiatzen bagara, zera dugu:

$$\begin{aligned} \nabla f_1(x, y)^T &= (2x - 2, -1) & \Rightarrow & \quad \nabla f_1(0, 1)^T = (-2, -1) \\ \nabla f_2(x, y)^T &= (2x, 8y) & \Rightarrow & \quad \nabla f_2(0, 1)^T = (0, 8). \end{aligned}$$

Ondorioz, hauek dira  $(0, 1)$  puntuko hiperplano ukitzaila dagozkien ekuazioak:

$$\begin{aligned} f_1(0, 1) + \nabla f_1(x_0, y_0)^T \cdot (x - x_0, y - y_0) &= 0 & \Rightarrow & \quad -0.5 + (-2, -1) \cdot (x - 0, y - 1) = 0 \\ f_2(0, 1) + \nabla f_2(x_0, y_0)^T \cdot (x - x_0, y - y_0) &= 0 & \Rightarrow & \quad 0 + (0, 8) \cdot (x - 0, y - 1) = 0. \end{aligned}$$

Horrek sistema hau inplikatzeko du:

$$\begin{aligned} -2x - y + 0.5 &= 0 \\ y - 1 &= 0 \end{aligned}$$

eta ebatziz, hau lortzen dugu:

$$x = -0.25 \quad y = 1$$

eta (7.2) iterazioaz  $(x_1, y_1) = (-0.25, 1)$  dugu.

Prozesu hori jarraituz, lau iteraziotan  $(x_4, y_4) = (-0.2223147, 0.9938121)$  dugu, eta  $f_1(x_4, y_4) = 0.000241126$  eta  $f_2(x_4, y_4) = 0.000073786$ .

Hau izango da errorean tamaina:

$$\|(f_1(x_4, y_4), f_2(x_4, y_4))\|_\infty = \|(0.000241126, 0.000073786)\|_\infty = 0.000241126.$$

Jarritako errorearen tolerantzia maximoa 0.001 izan bada, soluzio hori onargarria da.  $\square$

### 7.1. algoritmoa. Newtonen metodoa, ekuazio ez-linealen sistemetarako.

**0. urratsa.** SARRERA. Sartu:  $\mathbf{f}(\mathbf{x})$ ,  $\mathbf{x}_0$  (hasierako puntua)  $ze > 0$  (zehaztasun erlatiboa),  $emax > 0$  (errorearen tolerantzia) eta  $imax$  (iterazioen kopuru maximoa). Jarri  $k = 0$

**1. urratsa.** Kalkulatu  $\mathbf{f}(\mathbf{x}_k)$  eta  $\mathbf{J}(\mathbf{x}_k)$ .

**2. urratsa.** Ebatzi sistema lineal hau:

$$\mathbf{J}(\mathbf{x}_k)\mathbf{p}_k = -\mathbf{f}(\mathbf{x}_k). \quad (7.3)$$

**3. urratsa.** Kalkulatu honako puntu hau:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$$

**4. urratsa.** Baldin  $\|\mathbf{p}_k\|/\|\mathbf{x}_{k+1}\| < ze$  edo  $\|\mathbf{f}(\mathbf{x}_{k+1})\| < emax$  edo  $k \geq imax$  bada, gelditu egiten da. Bestela,  $k = k + 1$  egiten da, eta 1. urratsera goaz.

**5. urratsa.** IRTEERA. Emaitza:  $\mathbf{x}_{k+1}$ .

Algoritmo horren 2. urratsean (7.3) Newtonen sistema ebatzi behar dugu. Urrats horretan,  $LU$  deskonposizioa edo  $QR$  faktORIZAZIOA erabil ditzakegu.

**7.2. adibidea.** Izan bedi  $\mathbf{f}(\mathbf{x}) = \begin{bmatrix} x + y - 3 \\ x^2 + y^2 - 9 \end{bmatrix}$ . Ebatzi  $\mathbf{f}(\mathbf{x}) = \mathbf{0}$  sistema.

*Ebazpena.* Sistema horrek  $(3, 0)^T$  eta  $(0, 3)^T$  erroak ditu. Izan bedi  $\mathbf{x}_0 = (1, 5)^T$ . Orduan, hauek dira Newtonen lehenengo bi iterazioak:

$$\mathbf{J}(\mathbf{x}_0)\mathbf{p}_0 = -\mathbf{f}(\mathbf{x}_0) \quad \Rightarrow \quad \begin{bmatrix} 1 & 1 \\ 2 & 10 \end{bmatrix} \mathbf{p}_0 = - \begin{bmatrix} 3 \\ 17 \end{bmatrix} \quad \Rightarrow \quad \mathbf{p}_0 = \begin{bmatrix} -13/8 \\ -11/8 \end{bmatrix},$$

$$\begin{aligned} \mathbf{x}_1 &= \mathbf{x}_0 + \mathbf{p}_0 = (-0.625, 3.625)^T, \\ \mathbf{J}(\mathbf{x}_1)\mathbf{p}_1 = -\mathbf{f}(\mathbf{x}_1) &\Rightarrow \begin{bmatrix} 1 & 1 \\ -5/4 & 29/4 \end{bmatrix} \mathbf{p}_1 = - \begin{bmatrix} 0 \\ 145/32 \end{bmatrix} \Rightarrow \mathbf{p}_1 = \begin{bmatrix} 145/272 \\ -145/272 \end{bmatrix} \\ \mathbf{x}_2 &= \mathbf{x}_1 + \mathbf{p}_1 = (-0.092, 3.092)^T. \end{aligned}$$

Ikus dezakegunez,  $(0, 3)^T$  soluziotik nahiko hurbil dago  $\mathbf{x}_2$ .  $\square$

### 7.1.1. Newtonen metodoaren konbergentzia lokala

**7.1. lema.** *Izan bedi  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  etengabe diferentziagarria  $\mathcal{D} \subset \mathbb{R}^n$  multzo konbexu irekian. Orduan, hau betetzen da  $\mathbf{x}, \mathbf{x} + \mathbf{p} \in \mathcal{D}$  guztietarako:*

$$\mathbf{f}(\mathbf{x} + \mathbf{p}) - \mathbf{f}(\mathbf{x}) = \int_0^1 \mathbf{J}(\mathbf{x} + t\mathbf{p})\mathbf{p} dt \equiv \int_{\mathbf{x}}^{\mathbf{x}+\mathbf{p}} \mathbf{f}'(\mathbf{z}) d\mathbf{z}. \quad (7.4)$$

**7.1. definizioa.** *Izan bedi  $\mathbf{G} : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ . Orduan,  $\mathbf{G}$  funtzio matritziala  $\mathbf{x}$  puntuko  $\gamma$ -Lipschitz jarraitua dela esango dugu,  $\mathcal{D} \subset \mathbb{R}^n$  multzo irekia,  $\mathbf{x} \in \mathcal{D}$ , eta  $\gamma$  konstante bat existitzen badira, non  $\mathbf{v} \in \mathcal{D}$  guztietarako hau betetzen baita:*

$$\|\mathbf{G}(\mathbf{v}) - \mathbf{G}(\mathbf{x})\| \leq \gamma \|\mathbf{v} - \mathbf{x}\|. \quad (7.5)$$

Hori gertatzen denean, honela adieraziko dugu:  $\mathbf{G} \in Lip_\gamma(\mathcal{D})$ .

**7.2. lema.** *Izan bedi  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  etengabe diferentziagarria  $\mathcal{D} \subset \mathbb{R}^n$  multzo konbexu irekian, eta izan bedi  $\mathbf{J}$   $\gamma$ -Lipschitz jarraitua  $\mathbf{x} \in \mathcal{D}$  puntuaren  $\mathcal{D}$  ingurunean, hots  $\mathbf{J} \in Lip_\gamma(\mathcal{D})$ , bektore-norma bat eta berak eragindako matrize-norma erabiliz. Orduan, hau betetzen da  $\mathbf{x} + \mathbf{p} \in \mathcal{D}$  guztietarako:*

$$\|\mathbf{f}(\mathbf{x} + \mathbf{p}) - \mathbf{f}(\mathbf{x}) - \mathbf{J}(\mathbf{x})\mathbf{p}\| \leq \frac{\gamma}{2} \|\mathbf{p}\|^2. \quad (7.6)$$

*Frogantza.* Aurreko lemaren arabera,

$$\begin{aligned} \mathbf{f}(\mathbf{x} + \mathbf{p}) - \mathbf{f}(\mathbf{x}) - \mathbf{J}(\mathbf{x})\mathbf{p} &= \int_0^1 \mathbf{J}(\mathbf{x} + t\mathbf{p})\mathbf{p} dt - \mathbf{J}(\mathbf{x})\mathbf{p} \\ &= \int_0^1 (\mathbf{J}(\mathbf{x} + t\mathbf{p}) - \mathbf{J}(\mathbf{x}))\mathbf{p} dt. \end{aligned}$$

Orain, eragindako matrize-normaren definizioa eta  $\mathbf{J}$ -ren  $\gamma$ -Lipschitz jarraitutasuna erabiliz  $\mathbf{x}$ -ren  $\mathcal{D}$  ingurunean, hau dugu:

$$\begin{aligned} \|\mathbf{f}(\mathbf{x} + \mathbf{p}) - \mathbf{f}(\mathbf{x}) - \mathbf{J}(\mathbf{x})\mathbf{p}\| &\leq \int_0^1 \|(\mathbf{J}(\mathbf{x} + t\mathbf{p}) - \mathbf{J}(\mathbf{x}))\| \cdot \|\mathbf{p}\| dt \\ &\leq \int_0^1 \gamma \|t\mathbf{p}\| \cdot \|\mathbf{p}\| dt \\ &= \gamma \|\mathbf{p}\|^2 \int_0^1 t dt \\ &= \frac{\gamma}{2} \|\mathbf{p}\|^2. \quad \square \end{aligned}$$

Jarraian, ekuazio ez-linealen sistemetarako, Newtonen metodoaren konbergentzia koadratiko lokala frogatuko dugu. Orain,  $\mathbf{x}$ -ren  $r$  erradioko ingurunea  $\mathcal{B}_r(\mathbf{x}) = \{\bar{\mathbf{x}} \in \mathbb{R}^n \mid \|\bar{\mathbf{x}} - \mathbf{x}\| < r\}$  dugu. Kontuan izan  $\mathbf{J} \in Lip_\gamma(\mathcal{B}_r(\mathbf{x}))$  adierazpenak zera esan nahi duela:  $\mathbf{J}$   $\gamma$ -Lipschitz jarraitua dela  $\mathbf{x} \in \mathcal{D}$  puntuaren  $\mathcal{B}_r(\mathbf{x})$  ingurunean.

**7.1. teorema.** *Izan bedi  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  etengabe diferentziagarria  $\mathcal{D} \subset \mathbb{R}^n$  multzo konbezu irekian. Demagun  $\mathbf{x}^* \in \mathbb{R}^n$  eta  $r, \beta > 0$  existitzen direla, non  $\mathcal{B}_r(\mathbf{x}^*) \subset \mathcal{D}$ ,  $\mathbf{f}(\mathbf{x}^*) = \mathbf{0}$ ,  $\mathbf{J}(\mathbf{x}^*)^{-1}$  existitzen baita,  $\|\mathbf{J}(\mathbf{x}^*)^{-1}\| \leq \beta$  izanik, eta  $\mathbf{J} \in Lip_\gamma(\mathcal{B}_r(\mathbf{x}^*))$ . Orduan,  $\varepsilon > 0$  dago, non  $\mathbf{x}_0 \in \mathcal{B}_\varepsilon(\mathbf{x}^*)$  guztietarako, honako adierazpen honek sortutako  $\{\mathbf{x}_k\}$  segida ondo definituta baitago eta  $\mathbf{x}^*$ -era jotzen baitu:*

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{J}(\mathbf{x}_k)^{-1} \mathbf{f}(\mathbf{x}_k), \quad k = 0, 1, \dots$$

eta, gainera,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \beta\gamma \|\mathbf{x}_k - \mathbf{x}^*\|^2, \quad k = 0, 1, \dots \quad (7.7)$$

*Frogantza.*  $\varepsilon$  aukeratuko dugu  $\mathbf{J}(\mathbf{x})$  ez-singularra izateko  $\mathbf{x} \in \mathcal{B}_\varepsilon(\mathbf{x}^*)$  guztietan. Jarraian, eredu lineal honek:

$$\mathbf{M}_k(x) = \mathbf{f}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k),$$

sortutako errorea  $O(\|\mathbf{x}_k - \mathbf{x}^*\|^2)$  denez gero, konbergentzia koadratikoa dela frogatuko dugu.

Izan bedi

$$\varepsilon = \min \left\{ r, \frac{1}{2\beta\gamma} \right\}. \quad (7.8)$$

$k$ -ren gaineko indukzioz (7.7) betetzen dela frogatuko dugu, eta hau ere bai:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \frac{1}{2} \|\mathbf{x}_k - \mathbf{x}^*\|$$

eta, beraz,

$$\mathbf{x}_{k+1} \in \mathcal{B}_\varepsilon(\mathbf{x}^*). \quad (7.9)$$

Lehendabizi,  $\mathbf{J}(\mathbf{x}_0)$  ez dela singularra ikusiko dugu. Erabiltzen baditugu  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq \varepsilon$ ,  $\mathbf{J}$ -ren Lipschitz jarraitutasuna  $\mathbf{x}^*$  puntuan eta (7.8), hau lortzen dugu:

$$\begin{aligned} \|\mathbf{J}(\mathbf{x}^*)^{-1}[\mathbf{J}(\mathbf{x}_0) - \mathbf{J}(\mathbf{x}^*)]\| &\leq \|\mathbf{J}(\mathbf{x}^*)^{-1}\| \cdot \|\mathbf{J}(\mathbf{x}_0) - \mathbf{J}(\mathbf{x}^*)\| \\ &\leq \beta\gamma \|\mathbf{x}_0 - \mathbf{x}^*\| \leq \beta \cdot \gamma \cdot \varepsilon \leq 1/2. \end{aligned}$$

Horregatik eta matrize-normaren jarraitutasunagatik,  $\mathbf{J}(\mathbf{x}_0)$  ez dela singularra ateratzen da eta alderantzizkoaren normak hau betetzen duela froga daiteke (ikus [9] 3.1.4. teorema):

$$\begin{aligned} \|\mathbf{J}(\mathbf{x}_0)^{-1}\| &\leq \frac{\|\mathbf{J}(\mathbf{x}^*)^{-1}\|}{1 - \|\mathbf{J}(\mathbf{x}^*)^{-1}[\mathbf{J}(\mathbf{x}_0) - \mathbf{J}(\mathbf{x}^*)]\|} \\ &\leq 2\|\mathbf{J}(\mathbf{x}^*)^{-1}\| \\ &\leq 2 \cdot \beta. \end{aligned} \quad (7.10)$$

Beraz,  $\mathbf{x}_1$  ondo definituta dago eta hau betetzen du:

$$\begin{aligned}\mathbf{x}_1 - \mathbf{x}^* &= \mathbf{x}_0 - \mathbf{x}^* - \mathbf{J}(\mathbf{x}_0)^{-1}\mathbf{f}(\mathbf{x}_0) \\ &= \mathbf{x}_0 - \mathbf{x}^* - \mathbf{J}(\mathbf{x}_0)^{-1}[\mathbf{f}(\mathbf{x}_0) - \mathbf{f}(\mathbf{x}^*)] \\ &= \mathbf{J}(\mathbf{x}_0)^{-1}[\mathbf{f}(\mathbf{x}^*) - \mathbf{f}(\mathbf{x}_0) - \mathbf{J}(\mathbf{x}_0)(\mathbf{x}^* - \mathbf{x}_0)].\end{aligned}$$

Kontuan izan parentesien arteko gaia  $\mathbf{f}(\mathbf{x}^*)$ -ren eta  $\mathbf{M}_0(\mathbf{x}^*)$  ereduaren arteko kendura dela. Beraz, 7.2. lema eta (7.10) erabiliz, hau dugu:

$$\begin{aligned}\|\mathbf{x}_1 - \mathbf{x}^*\| &\leq \|\mathbf{J}(\mathbf{x}_0)^{-1}\| \cdot \|\mathbf{f}(\mathbf{x}^*) - \mathbf{f}(\mathbf{x}_0) - \mathbf{J}(\mathbf{x}_0)(\mathbf{x}^* - \mathbf{x}_0)\| \\ &\leq 2 \cdot \beta \cdot \frac{\gamma}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &= \beta \cdot \gamma \|\mathbf{x}_0 - \mathbf{x}^*\|^2.\end{aligned}$$

Horrek (7.7) frogatzen du. Gainera,  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq 1/(2\beta\gamma)$  denez, hau lortzen dugu:

$$\|\mathbf{x}_1 - \mathbf{x}^*\| \leq \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|,$$

horrek (7.9) frogatzen du eta  $k = 0$  kasua osatzen du. Indukzioaren beste urratsak frogatzeko, antzeko era batean jokatzeko da.  $\square$

## 7.2. Newtonen metodoaren aldaketak, sistema ez-linealak ebazteko

Newtonen metodorako aldaera batzuk existitzen dira; metodoari elkartutako sistemaren definizioan eta ebazpenean desberdintzen dira. Horien guztien helburua da algoritmoaren fase hori murriztea, eta, horretarako, matrize jacobiarra hurbiltzen da era desberdinetan.

### 7.2.1. Diferentzia finituzko Newtonen metodoa

Jacobiarraren adierazpen analitikoa kalkulatu ezin denean, Newtonen metodoaren aldaera honek matrize jacobiarra ordezkatzeko bere diferentzia finituzko hurbilpenaz.

**7.2. teorema.** *Izan bitez  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  eta  $\mathbf{x}^*$  7.1. teoremaren hipotesiak betetzen dituztenak. Orduan,  $\varepsilon, h > 0$  existitzen dira, non baldin  $\{h_k\}$  segida erreala,  $0 < |h_k| \leq h$ , eta*



$\mathbf{x}_0 \in \mathcal{B}_\varepsilon(\mathbf{x}^*)$  badira, honako adierazpen honek sortutako  $\{\mathbf{x}_k\}$  segida ondo definituta dago eta  $\mathbf{x}^*$ -ra jotzen du linealki:

$$\mathbf{a}_j^{(k)} = \frac{\mathbf{f}(\mathbf{x}_k + h_k \mathbf{e}_j) - \mathbf{f}(\mathbf{x}_k)}{h_k}, \quad j = 1, \dots, n,$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{A}_k^{-1} \mathbf{f}(\mathbf{x}_k), \quad k = 0, 1, \dots,$$

( $\mathbf{a}_j^{(k)}$  bektorea  $\mathbf{A}_k$ -ren  $j$ -garren zutabea da).

Baldin hau betetzen bada:

$$\lim_{k \rightarrow \infty} h_k = 0,$$

konbergentzia superlineala da.

Baldin  $c_1$  konstante bat badago hau betetzen duena:

$$|h_k| \leq c_1 \|\mathbf{x}_k - \mathbf{x}^*\|_1,$$

edo, baliokideki,  $c_2$  konstante bat badago hau betetzen duena:

$$|h_k| \leq c_2 \|\mathbf{f}(\mathbf{x}_k)\|_1,$$

orduan, konbergentzia koadratikoa da.

Teorema honen frogantza 7.1. teoremakoaren antzekoa da, eta [9] erreferentzian aurki daiteke.

Edonola ere, iterazio bakoitzean  $\mathbf{J}(\mathbf{x})$ -ren hurbilpen hori egiteko,  $n^2 + n$  funtzio-balioztatzea egin behar dugu:  $n^2$ , jacobiarra kalkulatzeko, eta  $n$ , funtzioaren balioa kalkulatzeko. Sistema lineala ebazteko,  $O(n^3)$  eragiketa.

Honako hau da  $h$  parametroa aukeratzeko era arrazonagarri bat.  $\mathbf{f}(\mathbf{x})$   $t$  digitu zuzenekin kalkula daitekeenean, orduan,  $\mathbf{f}(\mathbf{x} + h\mathbf{e}_j)$  eta  $\mathbf{f}(\mathbf{x})$  desberdinak izan beharko lirateke azken  $t/2$  digitu horietan. Hain zuzen ere,  $\mathbf{f}(\mathbf{x})$ -ren kalkuluaren errore erlatiboa  $\eta$  balitz, hau lortu beharko genuke:

$$\frac{\|\mathbf{f}(\mathbf{x} + h\mathbf{e}_j) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{f}(\mathbf{x})\|} \leq \sqrt{\eta}, \quad j = 1, \dots, n. \quad (7.11)$$

Ez badugu informazio hoberik, (7.11) lortzeko era arrazonagarri bat  $x_j$  osagai bakoitza honela aldatzea da:

$$h_j = \sqrt{\eta} \cdot x_j, \quad (7.12)$$

eta gero,  $\mathbf{a}_j$  zutabe bakoitza honela kalkulatzeko:

$$\mathbf{a}_j = \frac{\mathbf{f}(\mathbf{x} + h_j \mathbf{e}_j) - \mathbf{f}(\mathbf{x})}{h_j}. \quad (7.13)$$

Gainera,  $\mathbf{f}(\mathbf{x})$  formula erraz batez ematen denean,  $\eta = \varepsilon_M$  (makinaren errorea) hartzea ere arrazonagarria da. Honako adibide honetan, Newtonen metodoan deribatuak diferentzia finituen bidez (d.d.f.) kalkulatu dira, hots (7.12)-(7.13) erabiliz  $\eta = \varepsilon_M$ -rekin, eta baita deribatu analitikoak (d.a.) erabiliz ere; emaitzak hasierako puntu berdinetik abiatuz konparatu dira. Ikus daitekeenez, emaitzak ia berdinak dira, eta hori da praktikan gertatu ohi dena. Hori dela eta, software pakete batzuek ez dute behar deribatu analitikorik; haiek beti erabiltzen dituzte diferentzia finituak.

**7.3. adibidea.** Izan bedi  $\mathbf{f}(x) = \begin{bmatrix} x_1^2 + x_2^2 - 2 \\ e^{x_1-1} + x_2^3 - 2 \end{bmatrix}$ . Ebatzi  $\mathbf{f}(\mathbf{x}) = 0$ ,  $\mathbf{x}_0 = [2, 3]^T$  hartuz. Dakigunez  $\mathbf{x}^* = [1, 1]^T$ .

*Ebazpena.* Kalkuluak 14 digitu esanguratsuko ordenagailu batean egin dira. Ondorioz,  $\eta = \varepsilon_M = 10^{-14}$  dugu; beraz,  $h_j = 10^{-7}|x_j|$ .

Newtonen metodoa d.a. erabiliz	$\mathbf{x}_k$	Newtonen metodoa d.d.f. erabiliz
$\hat{\mathbf{u}}[2, 3]^T$	$\mathbf{x}_0$	$[2, 3]^T$
$\begin{bmatrix} 0.57465515807608 \\ 2.1168965612826 \end{bmatrix}$	$\mathbf{x}_1$	$\begin{bmatrix} 0.57465515450268 \\ 2.1168966735234 \end{bmatrix}$
$\begin{bmatrix} 0.31178766389307 \\ 1.5241979559460 \end{bmatrix}$	$\mathbf{x}_2$	$\begin{bmatrix} 0.31178738552306 \\ 1.5241981016335 \end{bmatrix}$
$\begin{bmatrix} 1.4841388323960 \\ 1.1464779176945 \end{bmatrix}$	$\mathbf{x}_3$	$\begin{bmatrix} 1.4841386151178 \\ 1.1464781318492 \end{bmatrix}$
$\begin{bmatrix} 1.0592959013664 \\ 1.0348194625183 \end{bmatrix}$	$\mathbf{x}_4$	$\begin{bmatrix} 1.0592958450507 \\ 1.0348195092235 \end{bmatrix}$
$\begin{bmatrix} 1.0008031050945 \\ 1.0014625483617 \end{bmatrix}$	$\mathbf{x}_5$	$\begin{bmatrix} 1.0008031056081 \\ 1.0014625533494 \end{bmatrix}$
$\begin{bmatrix} 0.99999872187461 \\ 1.0000026672636 \end{bmatrix}$	$\mathbf{x}_6$	$\begin{bmatrix} 0.99999872173640 \\ 1.0000026674316 \end{bmatrix}$
$\begin{bmatrix} 0.9999999999548 \\ 1.0000000000089 \end{bmatrix}$	$\mathbf{x}_7$	$\begin{bmatrix} 0.9999999999535 \\ 1.0000000000091 \end{bmatrix}$
4	$[\mathbf{J}(\mathbf{x}_0)]_{11}$	4.0000003309615
6	$[\mathbf{J}(\mathbf{x}_0)]_{12}$	6.0000002122252
2.7182818284590	$[\mathbf{J}(\mathbf{x}_0)]_{21}$	2.7182824169358
27	$[\mathbf{J}(\mathbf{x}_0)]_{22}$	27.000002470838

Baina,  $\mathbf{f}(\mathbf{x})$  kalkulatzeko adierazpen luzeko kode bat edo prozedimendu iteratibo bat erabili behar dugunean,  $\eta \gg \varepsilon_M$  izan daiteke, eta diferentzia finituzko deribatuak ez dira

izango oso zehatzak; orduan hobe da ebakitzaileren metodo bat erabiltzea (geroago ikusiko duguna). Gainera, kasu horretan,  $\mathbf{A}$  kalkulatzeko iterazio bakoitzean  $\mathbf{f}(\mathbf{x})$ -ren  $n$  ebaluatze beharko dira, eta horren kostua handia da.

### 7.2.2. Newtonen metodo aldatua

Lehenengo aldaera matrize jacobiar berdina,  $\mathbf{J}(\mathbf{x}_0)$ , prozesu iteratibo osoan izatean datza, edo, iterazio kopuru finko batean zehar, gutxienez.

### 7.2.3. Jacobiren aldaera

Aldaera honetan, diagonal nagusiko gaien bidez hurbiltzen da matrize jacobiarra. Hots,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{D}_k^{-1} \mathbf{f}(\mathbf{x}_k),$$

non  $\mathbf{D}_k = \text{diag}\{d_{11}^{(k)}, \dots, d_{nn}^{(k)}\}$  eta  $d_{ii}^{(k)} = [\mathbf{J}(\mathbf{x}_k)]_{ii}$ ,  $i = 1, \dots, n$ . Ikus daitekeenez, sistema linealetarako Jacobiren metodoaren antzekoa da. Diagonalean ez dauden gaiak txikiak badira diagonaleko gaiekiko (hots, diagonal menperatzailea bada), interesgarria izango da metodo hau.

### 7.2.4. Gauss-Seidelen aldaera

Aldaera honetan, matrize jacobiarra bere diagonalarekin eta diagonalaren azpiko azpimatrizearekin hurbiltzen da (diagonalaren gaineko gaiak ez dira erabiltzen). Hots,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{L}_k^{-1} \mathbf{f}(\mathbf{x}_k),$$

non  $[\mathbf{L}_k]_{ij} = [\mathbf{J}(\mathbf{x}_k)]_{ij}$ ,  $i \geq j$ . Iterazio bakoitzean  $\mathbf{L}_k \mathbf{p}_k = -\mathbf{f}(\mathbf{x}_k)$  sistema ebatziko da aurrerantz eta, gero,  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$  egingo da.

## 7.3. Quasi-Newton metodoak

Diferentzia finituzko Newtonen metodoko iterazio bakoitzean,  $\mathbf{J}(\mathbf{x})$ -ren hurbilpena egiteko,  $n^2 + n$  funtzio-balioztatze egin behar ditugu ( $n^2$ , jacobiarra kalkulatzeko, eta  $n$ , funtzioaren balioa kalkulatzeko), eta sistema lineala ebatzeko  $O(n^3)$  eragiketa gehiago. Ahalegin konputazionalaren kopuru hori handiegia da,  $n$  txikia izan ezean.

Atal honetan ikusiko dugu sistema ez-linealak ebazteko ebakitzaileren metodoaren hedatze arrakastatsuen. Metodo honek iterazio bakoitzean  $n$  funtzio-balioztatze bakarrik behar izango ditu, eta  $O(n^2)$  eragiketa aritmetiko.

Gogora dezagun ebakitzaileren metodoan, Newtonen metodoaren  $f'(x_1)$ -en ordeztapen hau erabiltzen genuela:

$$\frac{f(x_1) - f(x_0)}{x_1 - x_0}.$$

Aldagai anitzeko  $\mathbf{f}(\mathbf{x})$  funtzioaren deribatua  $\mathbf{x}_1$  puntuan,  $\mathbf{J}(\mathbf{x}_1)$ ,  $\mathbf{A}_1$  matrizeaz ordezkatzeko da propietate honekin:

$$\mathbf{A}_1(\mathbf{x}_1 - \mathbf{x}_0) = \mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_0).$$

Berdintza horri *ebakitzaileren ekuazio* deritzogu. Gainera,  $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$  notazioa erabiliko dugu  $k$ -garren iterazioari dagokion urratsa adierazteko, eta  $\mathbf{y}_k = \mathbf{f}(\mathbf{x}_{k+1}) - \mathbf{f}(\mathbf{x}_k)$  notazioa urrats horri dagokion funtzioaren balioaren aldakuntza emateko. Beraz,  $k$ -garren iterazioan,  $\mathbf{A}_{k+1}$  matrizeak ebakitzaileren ekuazio hau bete behar du:

$$\mathbf{A}_{k+1}\mathbf{s}_k = \mathbf{y}_k. \quad (7.14)$$

(7.14) ekuazioa ez da nahikoa espezifikatzeko  $\mathbf{A}_{k+1}$  bakar bat  $n > 1$  denean. Egia esan, (7.14)  $n$  ekuazioko eta  $n^2$  ezezaguneko sistema bat da. Horregatik, ebakitzailere-hurbilpen arrakastatsu bat eraikitzea da posibilitate horien artean hautatzeko irizpide egokiena aukeratzeko.

Quasi-Newton metodoek  $\{\mathbf{A}_k\}$  segida bat eraikitzen dute, non  $\mathbf{A}_k$  matrizeak  $\mathbf{J}(\mathbf{x}_k)$  matrize jacobiarra ahalik eta hoberen hurbiltzen baitu. Alegia,  $k$ -garren iterazioan  $\mathbf{f}(\mathbf{x})$  ordezkatzeko dugu hurbilpen lineal honekin:

$$\mathbf{M}_k(\mathbf{x}) = \mathbf{f}(\mathbf{x}_k) + \mathbf{A}_k(\mathbf{x} - \mathbf{x}_k),$$

zeina eredu afina baita.

### 7.3.1. Broyden-en metodoa

Broydenek ideia erraz bat erabili zuen  $\mathbf{A}_{k+1}$ -en hurbilpen egokia lortzeko (ikus [6]): iterazio berrian,  $\mathbf{A}_k$  matrizeak jacobiarri buruz ematen digun informazio gehiena gordetzea. Alegia,  $\mathbf{A}_{k+1}$  aukeratzeko orduan, saiatuko gara minimizatzen eredu afinaren aldakuntza  $k$ -garren eredu afinarekiko, hots:

$$\begin{aligned} \mathbf{M}_{k+1}(\mathbf{x}) - \mathbf{M}_k(\mathbf{x}) &= [\mathbf{f}(\mathbf{x}_{k+1}) + \mathbf{A}_{k+1}(\mathbf{x} - \mathbf{x}_{k+1})] - [\mathbf{f}(\mathbf{x}_k) + \mathbf{A}_k(\mathbf{x} - \mathbf{x}_k)] \\ &= \mathbf{f}(\mathbf{x}_{k+1}) + \mathbf{A}_{k+1}[(\mathbf{x} - \mathbf{x}_k) - (\mathbf{x}_{k+1} - \mathbf{x}_k)] - [\mathbf{f}(\mathbf{x}_k) + \mathbf{A}_k(\mathbf{x} - \mathbf{x}_k)] \\ &= \mathbf{f}(\mathbf{x}_{k+1}) + \mathbf{A}_{k+1}(\mathbf{x} - \mathbf{x}_k) - \mathbf{A}_{k+1}(\mathbf{x}_{k+1} - \mathbf{x}_k) - \mathbf{f}(\mathbf{x}_k) - \mathbf{A}_k(\mathbf{x} - \mathbf{x}_k) \\ &= \mathbf{f}(\mathbf{x}_{k+1}) - \mathbf{f}(\mathbf{x}_k) - \mathbf{A}_{k+1}(\mathbf{x}_{k+1} - \mathbf{x}_k) + (\mathbf{A}_{k+1} - \mathbf{A}_k)(\mathbf{x} - \mathbf{x}_k) \\ &= \mathbf{y}_k - \mathbf{A}_{k+1}\mathbf{s}_k + (\mathbf{A}_{k+1} - \mathbf{A}_k)(\mathbf{x} - \mathbf{x}_k) \quad (\text{eta (7.14) kontuan hartuz}) \\ &= (\mathbf{A}_{k+1} - \mathbf{A}_k)(\mathbf{x} - \mathbf{x}_k). \end{aligned}$$

Hortaz, minimizatu behar duguna hau izango da:

$$(\mathbf{A}_{k+1} - \mathbf{A}_k)(\mathbf{x} - \mathbf{x}_k).$$

Orain,  $\mathbf{x} \in \mathbb{R}^n$  guztietarako, beti adieraz dezakegu honela  $\mathbf{x} - \mathbf{x}_k$  kendura:

$$\mathbf{x} - \mathbf{x}_k = \alpha \mathbf{s}_k + \mathbf{t},$$

non  $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k (\neq \mathbf{0})$  eta  $\mathbf{t}^T \mathbf{s}_k = 0$  (hots,  $\mathbf{t} \in \mathbb{R}^n$  da  $\mathbf{s}_k \in \mathbb{R}^n$  bektorearekiko  $n - 1$  dimentsioko azpiespazio ortogonalaren edozein bektore). Ondorioz,  $\mathbf{M}_{k+1}(\mathbf{x}) - \mathbf{M}_k(\mathbf{x})$  minimizatzeke, hau minimizatu behar dugu:

$$\alpha(\mathbf{A}_{k+1} - \mathbf{A}_k)\mathbf{s}_k + (\mathbf{A}_{k+1} - \mathbf{A}_k)\mathbf{t}. \quad (7.15)$$

Lehenengo gaiaz ezin dugu ezer egin, zeren (7.14), ebakitzaileren ekuazioaren arabera, hau gertatzen baita:

$$(\mathbf{A}_{k+1} - \mathbf{A}_k)\mathbf{s}_k = \mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k, \quad (7.16)$$

hau da, emaitza-bektorea konstantea da ( $k$ -garren iterazioan datu horiek guztiak finkatuta daude). Orain, (7.15) adierazpenaren bigarren gaiaz, hau eska dezakegu:

$$(\mathbf{A}_{k+1} - \mathbf{A}_k)\mathbf{t} = \mathbf{0} \quad \forall \mathbf{t} \in \mathbb{R}^n \text{ non } \mathbf{s}_k^T \mathbf{t} = 0.$$

Hori lor dezakegu  $\mathbf{A}_{k+1}$  egoki bat aukeratuz. Berdintza horrek hau betetzera behartzen du:

$$\mathbf{A}_{k+1} - \mathbf{A}_k = \mathbf{u} \mathbf{s}_k^T \quad \left( = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \mathbf{s}_k^T = \begin{bmatrix} u_1 \mathbf{s}_k^T \\ \vdots \\ u_n \mathbf{s}_k^T \end{bmatrix} \right) \quad (7.17)$$

non  $\mathbf{u} \in \mathbb{R}^n$ . Alegia,  $\mathbf{A}_{k+1} - \mathbf{A}_k$  bat heinakoa izan behar du. Orain, (7.14) ebakitzaileren ekuazioa betetzea (7.16) betetzearen baliokidea da, eta, orduan,  $\mathbf{u}$  bektoreak hau bete behar du:

$$(\mathbf{u} \mathbf{s}_k^T) \mathbf{s}_k = \mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k,$$

ondorioz, ebakitzaileren ekuazioa betetzeko hau izango dugu:

$$\mathbf{u} = \frac{\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k}{\mathbf{s}_k^T \mathbf{s}_k}.$$

Beraz, (7.17) ekuazioko  $\mathbf{A}_{k+1}$  bakanduz eta  $\mathbf{u}$  bektore hori ordezkaturaz, hau lortzen da:

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \mathbf{u} \mathbf{s}_k^T = \mathbf{A}_k + \frac{\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k}{\mathbf{s}_k^T \mathbf{s}_k} \mathbf{s}_k^T. \quad (7.18)$$

Ondorioz, (7.18) eguneratzeak ondoz ondoko bi eredu afinen diferentzia txikiagotzen du, eta ebakitzaileren ekuazioa betetzen du.

### Newtonen sistema hurbilduaren ebazpena

Bat heineko Broydenen eguneratzea erabiltzen badugu, (7.18),  $n$  funtzio-balioztatze bakarrik erabili behar dugu iterazio bakoitzean (diferentzia finituen kasuan,  $n^2 + n$  behar ditugu). Bestalde,  $\mathbf{J}(\mathbf{x}_k)\mathbf{p} = -\mathbf{f}(\mathbf{x}_k)$  Newtonen sistemaren ordez, sistema hau ebatzi behar dugu:

$$\mathbf{A}_k\mathbf{p} = -\mathbf{f}(\mathbf{x}_k),$$

eta sistema hori ebazteko,  $O(n^3)$  eragiketak erabili behar ditugu oraindik. Arazo hori gainditzeko, matrizeen propietate bat erabiliz, formula erraz baten bidez lortuko dugu  $\mathbf{A}_k$ -ren alderantzizkoa.

**7.3. lema.** (*Sherman-Morrison-Woodbury*) Izan bitez  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  bektoreak eta demagun  $\mathbf{A} \in \mathbb{R}^{n \times n}$  matrizea ez dela singularra. Orduan,  $\mathbf{A} + \mathbf{u}\mathbf{v}^T$  ez da singularra, baldin eta soilik baldin hau betetzen bada:

$$\sigma = 1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq 0.$$

Gainera, hau betetzen da:

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{1}{\sigma} \mathbf{A}^{-1} \mathbf{u}\mathbf{v}^T \mathbf{A}^{-1}.$$

Lema hori (7.18) adierazpenari zuzenean aplikatzeak hau ematen digu:

$$\begin{aligned} \mathbf{A}_{k+1}^{-1} &= (\mathbf{A}_k + \mathbf{u}\mathbf{s}_k^T)^{-1} = \mathbf{A}_k^{-1} - \frac{\mathbf{A}_k^{-1} \mathbf{u}\mathbf{s}_k^T \mathbf{A}_k^{-1}}{1 + \mathbf{s}_k^T \mathbf{A}_k^{-1} \mathbf{u}} \\ &= \mathbf{A}_k^{-1} - \frac{\mathbf{A}_k^{-1} \left( \frac{\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k}{\mathbf{s}_k^T \mathbf{s}_k} \right) \mathbf{s}_k^T \mathbf{A}_k^{-1}}{1 + \mathbf{s}_k^T \mathbf{A}_k^{-1} \left( \frac{\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k}{\mathbf{s}_k^T \mathbf{s}_k} \right)} \\ &= \mathbf{A}_k^{-1} - \frac{(\mathbf{A}_k^{-1} \mathbf{y}_k - \mathbf{s}_k) \mathbf{s}_k^T \mathbf{A}_k^{-1}}{\mathbf{s}_k^T \mathbf{s}_k + \mathbf{s}_k^T \mathbf{A}_k^{-1} \mathbf{y}_k - \mathbf{s}_k^T \mathbf{s}_k} \end{aligned}$$

eta, azkenik, *Broydenen alderantzizko formula* hau lortzen da:

$$\mathbf{A}_{k+1}^{-1} = \mathbf{A}_k^{-1} + \frac{(\mathbf{s}_k - \mathbf{A}_k^{-1} \mathbf{y}_k) \mathbf{s}_k^T \mathbf{A}_k^{-1}}{\mathbf{s}_k^T \mathbf{A}_k^{-1} \mathbf{y}_k}. \quad (7.19)$$

Adierazpen horrek  $O(n^2)$  eragiketa behar ditu eta,  $\mathbf{s}_k = -\mathbf{A}_k^{-1} \mathbf{f}(\mathbf{x}_k)$  kalkulatzeko,  $O(n^2)$  eragiketa egin behar ditugu. Beraz,  $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{A}_k^{-1} \mathbf{f}(\mathbf{x}_k)$  iterazio berria kalkulatzeko, guztira  $O(n^2)$  eragiketa egin behar ditugu,  $O(n^3)$  erabili beharrean.

Froga daiteke Broydenen metodoaren konbergentzia lokala superlineala dela; ikus [9].

**7.2. algoritmoa. Quasi-Newton metodoa, Broydenen eguneratzeaz.**

**0. urratsa.** SARRERA. Sartu:  $\mathbf{f}(\mathbf{x})$ ,  $\mathbf{x}_0$  (hasierako puntua),  $\mathbf{A}_0^{-1}$ ,  $ze > 0$  (zehaztasun erlatiboa),  $emax > 0$  (errorearen tolerantzia) eta  $imax$  (iterazioen kopuru maximoa). Kalkulatu  $\mathbf{f}(\mathbf{x}_0)$ . Jarri  $k = 0$ .

**1. urratsa.** Kalkulatu hau:

$$\mathbf{s}_k = -\mathbf{A}_k^{-1} \mathbf{f}(\mathbf{x}_k). \quad (7.20)$$

**2. urratsa.** Kalkulatu honako puntu hau:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k.$$

**3. urratsa.** Kalkulatu  $\mathbf{f}(\mathbf{x}_{k+1})$ .

**4. urratsa.** Baldin  $\|\mathbf{s}_k\|/\|\mathbf{x}_{k+1}\| < ze$  edo  $\|\mathbf{f}(\mathbf{x}_{k+1})\| < emax$  edo  $k \geq imax$  bada, gelditu egiten da. Bestela, hurrengo urratsera doa.

**5. urratsa.** Kalkula itzazu  $\mathbf{y}_k = \mathbf{f}(\mathbf{x}_{k+1}) - \mathbf{f}(\mathbf{x}_k)$  eta, gero,  $\mathbf{A}_{k+1}^{-1}$  honela:

$$\mathbf{A}_{k+1}^{-1} = \mathbf{A}_k^{-1} + \frac{(\mathbf{s}_k - \mathbf{A}_k^{-1} \mathbf{y}_k) \mathbf{s}_k^T \mathbf{A}_k^{-1}}{\mathbf{s}_k^T \mathbf{A}_k^{-1} \mathbf{y}_k}.$$

**6. urratsa.** Egin  $k = k + 1$  eta joan 1. urratsera.

**7. urratsa.** IRTEERA. Emaitza:  $\mathbf{x}_{k+1}$ .

Ohartarazi behar da  $\mathbf{A}_0$  hasierako matrizea identitate-matrizea edo hasierako jacobiarra edo bere diagonalak izan daitekeela; edo, algoritmo horretaz esperientzia badugu, matrize egoki bat.

**7.4. adibidea. Izan bedi**

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} x_1 + x_2 - 3 \\ x_1^2 + x_2^2 - 9 \end{bmatrix},$$

zeinak  $[0, 3]^T$  eta  $[3, 0]^T$  erroak baititu. Izan bedi  $\mathbf{x}_0 = [1, 5]^T$  eta aplikatu aurreko algoritmoa hau hartuz:

$$\mathbf{A}_0^{-1} = \mathbf{J}(\mathbf{x}_0)^{-1} = \begin{bmatrix} 1 & 1 \\ 2 & 10 \end{bmatrix}^{-1} = \begin{bmatrix} 1.2500 & -0.1250 \\ -0.2500 & 0.1250 \end{bmatrix}.$$

*Ebazpena.* Algoritmoaren 1. eta 2. urratsak egiten dira:

$$\mathbf{s}_0 = -\mathbf{A}_0^{-1} \mathbf{f}(\mathbf{x}_0) = - \begin{bmatrix} 1.2500 & -0.1250 \\ -0.2500 & 0.1250 \end{bmatrix} \begin{bmatrix} 3 \\ 17 \end{bmatrix} = \begin{bmatrix} -1.625 \\ -1.375 \end{bmatrix}.$$

3.ak zera ematen digu:

$$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{s}_0 = \begin{bmatrix} 1 \\ 5 \end{bmatrix} + \begin{bmatrix} -1.625 \\ -1.375 \end{bmatrix} = \begin{bmatrix} -0.625 \\ 3.625 \end{bmatrix}.$$

4.aren bidez, erraz ikus daiteke soluzioa oraindik ez dugula aurkitu hemen:

$$\mathbf{f}(\mathbf{x}_1) = \begin{bmatrix} 0 \\ 4.53125 \end{bmatrix}.$$

Beraz, algoritmoaren 5. urratsa burutzerakoan, emaitza hauek lortzen dira:

$$\mathbf{y}_0 = \mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_0) = \begin{bmatrix} 0 \\ 4.53125 \end{bmatrix} - \begin{bmatrix} 3 \\ 17 \end{bmatrix} = \begin{bmatrix} -3 \\ -12.46875 \end{bmatrix},$$

$$\mathbf{z} = -\mathbf{A}_0^{-1}\mathbf{y}_0 = -\begin{bmatrix} 1.2500 & -0.1250 \\ -0.2500 & 0.1250 \end{bmatrix} \begin{bmatrix} -3 \\ -12.46875 \end{bmatrix} = \begin{bmatrix} 2.19141 \\ 0.80859 \end{bmatrix},$$

$$p = -\mathbf{s}_0^T \mathbf{z} = -\begin{bmatrix} -1.625 \\ -1.375 \end{bmatrix} \begin{bmatrix} 2.19141 \\ 0.80859 \end{bmatrix} = 4.6729,$$

$$\begin{aligned} \mathbf{C} &= \frac{(\mathbf{s}_0 + \mathbf{z})\mathbf{s}_0^T \mathbf{A}_0^{-1}}{p} = \frac{\left( \begin{bmatrix} -1.625 \\ -1.375 \end{bmatrix} + \begin{bmatrix} 2.19141 \\ 0.80859 \end{bmatrix} \right) \begin{bmatrix} -1.625 & -1.375 \end{bmatrix} \begin{bmatrix} 1.2500 & -0.1250 \\ -0.2500 & 0.1250 \end{bmatrix}}{4.6729} \\ &= \begin{bmatrix} -0.2045 & 0.0038 \\ 0.2045 & -0.0038 \end{bmatrix}, \end{aligned}$$

$$\mathbf{A}_1^{-1} = \mathbf{A}_0^{-1} + \mathbf{C} = \begin{bmatrix} 1.2500 & -0.1250 \\ -0.2500 & 0.1250 \end{bmatrix} + \begin{bmatrix} -0.2045 & 0.0038 \\ 0.2045 & -0.0038 \end{bmatrix} = \begin{bmatrix} 1.0455 & -0.1212 \\ -0.0455 & 0.1212 \end{bmatrix}.$$

Berriro 2. urratsa aplikatuz, zera lortzen da:

$$\mathbf{s}_1 = -\mathbf{A}_1^{-1}\mathbf{f}(\mathbf{x}_1) = -\begin{bmatrix} 1.0455 & -0.1212 \\ -0.0455 & 0.1212 \end{bmatrix} \begin{bmatrix} 0 \\ 4.53125 \end{bmatrix} = \begin{bmatrix} 0.5492 \\ -0.5492 \end{bmatrix}.$$

Eta 3. urratsaren bitartez, hau:

$$\mathbf{x}_2 = \mathbf{x}_1 + \mathbf{s}_1 = \begin{bmatrix} -0.625 \\ 3.625 \end{bmatrix} + \begin{bmatrix} 0.5492 \\ -0.5492 \end{bmatrix} = \begin{bmatrix} -0.07575 \\ 3.07575 \end{bmatrix}.$$

Ikus dezakegunez,  $\mathbf{x}_2 \approx [0 \ 3]^T$ .  $\square$

Praktikan, inplementazio hau ondo dabil, baina desabantaila bat dauka: nekez detektatzen da  $\mathbf{A}_{k+1}$ -en baldintzapen txarra.  $\mathbf{A}_k$ -ren  $\mathbf{QR}$  faktORIZAZIOAK arazo hori ez duenez eta erabiltzen dituen eragiketen kopurua antzekoa denez ( $O(n^2)$ ), gaur egun, (7.19) formula ordezkatu egin du. Metodo horretan,  $\mathbf{A}_k = \mathbf{QR}$  faktORIZAZIOA dugunez, erraz lortzen da  $\mathbf{A}_{k+1} = \mathbf{A}_k + \mathbf{u}\mathbf{s}_k^T = \mathbf{QR} + \mathbf{u}\mathbf{s}_k^T = \overline{\mathbf{Q}} \overline{\mathbf{R}}$  faktORIZAZIOA. Horretarako, Givensen biraketak (transformazio ortogonalak) erabiltzen dira; ikus 6.3. atala.



## 7.4. Murrizketarik gabeko optimizazioa

Izan bedi  $F(\mathbf{x})$  aldagai anitzeko funtzio eskalar bat eta  $\mathbf{x} \in \mathbb{R}^n$ . Orduan, murrizketarik gabeko optimizazioan, honelako problema dugu:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}).$$

Orokorki, ez dakigu zenbat minimo lokal daukan eta nola aurkitu minimo globala era eraginkorren, hots, minimo guztietako  $F(\mathbf{x})$  funtzioaren balio txikiena ematen duena. Bestalde, bistan dago baliokidetasun hau:

$$\max_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \iff \min_{\mathbf{x} \in \mathbb{R}^n} -F(\mathbf{x}).$$

Aurrez ikusitako sistema linealen/ez-linealen teoria oso ondo datorkigu murrizketarik gabeko optimizazio-problemak ebazteko.

**7.3. Aldagai anitzetarako Taylorren teorema.** *Izan bedi  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$  eta demagun  $F(\mathbf{x}) \in \mathcal{C}^3$  funtzioa dugula. Orduan,  $\mathbf{d} = (d_1, d_2, \dots, d_n)^T$  urrats bektore baterako, Taylorren garapenak hau ematen digu:*

$$F(\mathbf{x} + \mathbf{d}) = F(\mathbf{x}) + \nabla F(\mathbf{x})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 F(\mathbf{x}) \mathbf{d} + O(\|\mathbf{d}\|^3). \quad (7.21)$$

Teorema horretan,  $\nabla F(\mathbf{x})$  gradiente bektorea da eta  $\nabla^2 F(\mathbf{x})$  matrize hessiarra, eta honela adieraz ditzakegu:

$$\nabla F(\mathbf{x}) = \begin{bmatrix} \frac{\partial F}{\partial x_1} \\ \frac{\partial F}{\partial x_2} \\ \dots \\ \frac{\partial F}{\partial x_n} \end{bmatrix}, \quad \nabla^2 F(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 F}{\partial x_1^2} & \frac{\partial^2 F}{\partial x_1 x_2} & \dots & \frac{\partial^2 F}{\partial x_1 x_n} \\ \frac{\partial^2 F}{\partial x_2 x_1} & \frac{\partial^2 F}{\partial x_2^2} & \dots & \frac{\partial^2 F}{\partial x_2 x_n} \\ \dots & \dots & \ddots & \dots \\ \frac{\partial^2 F}{\partial x_n x_1} & \frac{\partial^2 F}{\partial x_n x_2} & \dots & \frac{\partial^2 F}{\partial x_n^2} \end{bmatrix}.$$

Bestalde, teoremako baldintzetan  $F$ -ren  $\mathbf{d}$  norabidearekiko deribatua  $\mathbf{x}$  puntuan honela kalkula daiteke:

$$\nabla F(\mathbf{x})^T \mathbf{d} = \sum_{i=1}^n d_i \frac{\partial F}{\partial x_i}.$$

Gainera, zera dugu:

$$\frac{1}{2} \mathbf{d}^T \nabla^2 F(\mathbf{x}) \mathbf{d} = \frac{1}{2} \sum_{i=1, j=1}^n \frac{\partial^2 F}{\partial x_i \partial x_j} d_i d_j.$$

Hortaz,  $\mathbf{x}^*$  minimo lokala bada, hots, puntu horren ingurunean  $F$  funtzioak puntu horretan balio txikiena hartzen badu, orduan  $\mathbf{d}$  guztietarako hau dugu:

$$F(\mathbf{x}^* + \mathbf{d}) = F(\mathbf{x}^*) + \nabla F(\mathbf{x}^*)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 F(\mathbf{x}^*) \mathbf{d} + O(\|\mathbf{d}\|^3) \geq F(\mathbf{x}^*).$$

Orduan,  $\mathbf{x}^*$  minimoa izateko hau da baldintza beharrezkoa:

$$\nabla F(\mathbf{x}^*) = \mathbf{0}.$$

Berdintza hori betetzen duten puntuei *puntu kritikoak* deritze.

Gainera,  $\mathbf{x}^*$  puntu kritikoak minimoa izateko bete behar duen *baldintza nahikoa*  $\nabla^2 F(\mathbf{x}^*)$  matrize hessiarra definitu positiboa izatea da.

Ikus ditzagun baldintza horiek teorema honetan.

#### 7.4. teorema: murrizketarik gabeko minimizazioaren baldintzak.

*Demagun  $F \in \mathcal{C}^3$ . Orduan zera dugu:*

- $\mathbf{x}^*$  puntuan  $F$  funtzioak minimo lokal bat izateko, baldintza beharrezkoa da  $\mathbf{x}^*$  puntu kritikoa izatea eta  $\nabla^2 F(\mathbf{x}^*)$  matrize hessiarra erdidefinitu positiboa izatea.
- $\mathbf{x}^*$  puntuan  $F$  funtzioak minimo lokal bat izateko, baldintza nahikoa da  $\mathbf{x}^*$  puntu kritikoa izatea eta  $\nabla^2 F(\mathbf{x}^*)$  matrize hessiarra definitu positiboa izatea.

*Frogantza.* Izan bedi  $\|\mathbf{d}\|$  oso txikia,  $\|\mathbf{d}\|^2 \ll \|\mathbf{d}\|$  betetzeko. Orduan,  $\nabla F(\mathbf{x}^*) \neq \mathbf{0}$  bada, beti da posible  $\mathbf{d}$  norabide bat aurkitzea  $\nabla F(\mathbf{x}^*)^T \mathbf{d} < 0$  gertatzeko, eta, hortaz,  $F(\mathbf{x}^* + \mathbf{d}) < F(\mathbf{x}^*)$  eta  $\mathbf{x}^*$  ez da minimoa izango.

Bestalde,  $\mathbf{x}^*$  puntu kritikoa minimo lokal hertsia izateko  $\mathbf{d}$  norabide guztietarako, non  $0 < \|\mathbf{d}\| \ll 1$ , hau bete behar da:

$$F(\mathbf{x}^* + \mathbf{d}) = F(\mathbf{x}^*) + \frac{1}{2} \mathbf{d}^T \nabla^2 F(\mathbf{x}^*) \mathbf{d} + O(\|\mathbf{d}\|^3) > F(\mathbf{x}^*)$$

eta hori gertatuko da  $\nabla^2 F(\mathbf{x}^*)$  matrize hessiarra definitu positiboa bada.

Baldintza beharrezkorako,  $\mathbf{x}^*$  minimo lokala bada,  $F(\mathbf{x}^* + \mathbf{d}) \geq F(\mathbf{x}^*)$  gertatu behar da, eta, beraz, matrize hessiarrak ezin du  $\mathbf{d}^T \nabla^2 F(\mathbf{x}^*) \mathbf{d} < 0$  bete inolako  $\mathbf{d}$ -tarako, hots, matrize hessiarra erdidefinitu positiboa izan behar da.  $\square$

##### 7.4.1. Oinarrizko metodoak

Puntu kritikoa izateko baldintzak ekuazio ez-linealen sistema bat ematen du:

$$\mathbf{f}(\mathbf{x}) \equiv \nabla F(\mathbf{x}) = \mathbf{0}.$$

Hortaz, puntu kritiko bat aurkitzea da ekuazio ez-linealen sistema baten ebazpenaren kasu berezi bat, eta Newtonen metodoa eta bere antzekoak zuzen erabil daitezke. Argi dago

kasu honetan  $J(\mathbf{x}) = \nabla^2 F(\mathbf{x})$ , hots, matrize jacobiarra matrize hessian delako. Baina, orain, jacobiarra matrize simetrikoa da, eta minimoaren ondoan definitu positiboa izaten da; hori garrantzitsua da.

## Newtonen metodoa

### 7.3. algoritmoa. Murrizketarik gabeko minimizaziorako Newtonen metodoa.

**0. urratsa.** Jarri  $k = 0$ .

**1. urratsa.**  $\nabla^2 F(\mathbf{x}_k)\mathbf{d}_k = -\nabla F(\mathbf{x}_k)$  ebatziz  $\mathbf{d}_k$  lortzen da.

**2. urratsa.** Jarri  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$ .

**3. urratsa.** Baldin bukaerako baldintzak betetzen badira (esate baterako,  $\|\nabla F(\mathbf{x}_k)\| < \varepsilon$ ), emaitza  $\mathbf{x}_{k+1}$  da. Bestela,  $k = k + 1$  jarri eta segi 1. urratsean.

$\mathbf{x}$  puntu baterako esango dugu  $F$ -rekiko  $\mathbf{d}$  norabidea *beherakorra* dela,  $\nabla F(\mathbf{x})^T \mathbf{d} < 0$  bada. Orduan (7.21) Taylorren formula kontuan hartuz,  $\lambda > 0$  nahiko txiki baterako,  $F(\mathbf{x} + \lambda \mathbf{d}) < F(\mathbf{x})$  dugu.

Bestalde,  $\mathbf{x}^*$  puntu kritikoa zeladura-puntua izan daiteke eta, orduan, puntu horretan existituko dira  $\mathbf{d}$  urrats beherakorrak ( $F(\mathbf{x}^* + \mathbf{d}) < F(\mathbf{x}^*)$ ) eta  $\mathbf{d}$  urrats gorakorrak ( $F(\mathbf{x}^* + \mathbf{d}) > F(\mathbf{x}^*)$ ). Hori gertatuko da  $\mathbf{x}^*$  puntu kritikoa matrize hessian indefinitua denean.

## Metodo mota bat

Murrizketarik gabeko metodo gehienek honelako iterazioak erabiltzen dituzte:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \lambda_k \mathbf{d}_k, \quad \text{non} \quad \mathbf{d}_k = -\mathbf{A}_k^{-1} \nabla F(\mathbf{x}_k),$$

$\lambda_k > 0$  nahiko txikia hartzen da  $F(\mathbf{x}_k + \lambda_k \mathbf{d}_k) < F(\mathbf{x}_k)$  betetzeko.

$\mathbf{A}_k$  simetriko definitu positiboa bada, orduan nahitaez  $\mathbf{d}_k$  norabide beherakorra da, zeren  $\nabla F(\mathbf{x}_k)^T \mathbf{d} = -\nabla F(\mathbf{x}_k)^T \mathbf{A}_k^{-1} \nabla F(\mathbf{x}_k) < 0$ . Geroago,  $\mathbf{A}_k$  hautatzeko aukera desberdinak har ditzakegu kontuan.

Bilaketa metodo bat justifikatzeko, merezi du Newton metodoaren abantailak eta desabantailak aztertzea. Abantailak: konbergentzia lokal koadratikoa. Desabantailak:

- Metodoak matrize hessianaren existentzia eskatzen du.

- Okerrago, beharrezkoa da hessianra balioztatzea.
- Iterazio bakoitzean ekuazio linealen sistema bat ebatzi behar da.
- Minimoto kanpoan hessianra ez-definitu positiboa izan daiteke.
- Konbergentzia ezin da kontrolatu (konbergitzen da minimo batera?).

Horrelako arazoak gainditzeko, honako metodo hauek ikusiko ditugu.

### 7.4.2. Metodo globalak

Orain arte erabilitako metodoak lokalak izan dira, hots,  $\mathbf{x}_0$  hasierako puntuak  $\mathbf{x}^*$  errore nahiko hurbil egon behar du, metodo horiek konbergenteak izateko. Baina, hori gertatzen ez denean, zein metodo erabil dezakegu?

#### Gradiente metodoa

Dakigunez,  $\nabla F(\mathbf{x})^T \mathbf{d} = \|\nabla F(\mathbf{x})\|_2 \|\mathbf{d}\|_2 \cos(\theta)$  dugu ( $\theta$  da  $\nabla F(\mathbf{x})$  eta  $\mathbf{d}$  bektoreen arteko angelua), eta  $\|\mathbf{d}\|_2 = 1$  bada,  $\nabla F(\mathbf{x})^T \mathbf{d}$  biderkadurak balio txikiena hartuko du  $\theta = -\pi$  denean, hots,  $\mathbf{d} = -\nabla F(\mathbf{x}) / \|\nabla F(\mathbf{x})\|_2$  denean. Orokorki,  $\mathbf{d} = -\nabla F(\mathbf{x})$  norabidea aukeratzen denean, *gradientearen metodoa* (edo *beherapen azkarreneko metodoa*) dugu. Orduan,  $\mathbf{d} = -\nabla F(\mathbf{x})$  norabidean  $F$  funtzioaren balioa txikiagotzen da,  $\nabla F(\mathbf{x}) \neq \mathbf{0}$  izanik, eta honela aurkitzen da norabide horretan eman behar duen urratsaren luzera, alegia,  $\alpha$ -rako aukera zehatzena ( $\alpha > 0$  izanik):

$$\min_{\alpha} F(\mathbf{x} + \alpha \mathbf{d}), \quad (7.22)$$

baina, orokorki, kalkulu hori garestiegia da beste metodo ez-zehatz eta eraginkor batzuekin konparatzen badugu. Aurretiko metodo horrek  $F$ -ren minimo batera jotzen du konbergentzia linealarekin, eta, batzuetan oso astiro, lineala da.

Gradiente metodoan  $\mathbf{A}_k = \mathbf{1}$  da.

#### Funtzio koadratiko baterako konbergentzia-ratioa

Izan bedi  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{b}^T \mathbf{x}$  funtzio koadratikoa, non  $\mathbf{Q}$  simetrikoa eta definitu positiboa baita. Gradiente  $\nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{b}$  da, eta  $\mathbf{x}^*$  minimoa  $\mathbf{Q} \mathbf{x} = \mathbf{b}$  sistema linealaren soluzio bakarra.

Demagun  $\mathbf{x}^*$ -ren kalkuluan beherapen azkarreneko metodoaz  $k$ -garren iterazioan gaudela, hots,  $\mathbf{x}_k$  dugu. Izan bedi  $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$ , hortaz, hurrengo iterazioa kalkulatzeko, hau egin

behar dugu:

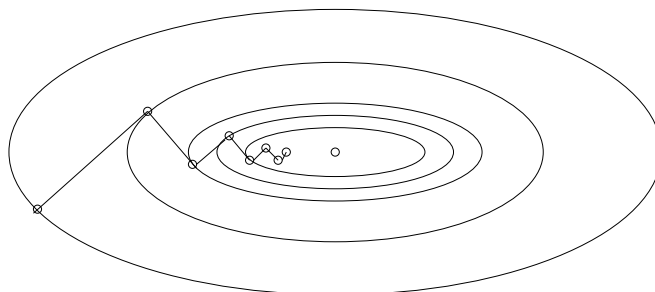
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \mathbf{g}_k$$

eta  $\alpha$  kalkulatu behar da. Horretarako, bilaketa lineal zehatza (7.22) erabiliko dugu, gure kasuan,  $f(\mathbf{x}_k - \alpha \mathbf{g}_k)$  funtzioa  $\alpha$ -rekiko minimizatu behar dugu. Deribatuz eta zerora berdinduz, hau lortzen dugu (frogapena ariketa gisa geratzen da):

$$\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}. \quad (7.23)$$

Ondorioz, funtzio koadratiko definitu positibo baterako beherapen azkarreneko iterazioa hau da:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} \mathbf{g}_k. \quad (7.24)$$



### 7.1. irudia. Gradiente metodoaren urratsak.

Irudian bi dimentsioko funtzio koadratiko baterako metodo honek sortutako iterazio-segida tipiko bat ikus daiteke.  $f$ -ren sestra-kurbak elipsoideak dira, eta haien ardatzak  $\mathbf{Q}$  matrizearen autobektore unitarioen norabideen gainean daude. Kontuan izan iterazioak soluziorantz sigi-sagan doazela.

Konbergentzia-ratioa kalkulatzeko, norma haztatu bat erabiliko dugu,  $\mathbf{Q}$ -norma hau:  $\|\mathbf{x}\|_{\mathbf{Q}}^2 = \mathbf{x}^T \mathbf{Q} \mathbf{x}$ . Orduan,  $\mathbf{Q} \mathbf{x}^* = \mathbf{b}$  erlazioa erabiltzen badugu, hau erraz froga dezakegu:

$$\frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{Q}}^2 = f(\mathbf{x}) - f(\mathbf{x}^*). \quad (7.25)$$

**7.5. teorema: konbergentzia-ratioa.**  $\mathbf{Q}$  matrizea simetrikoa eta definitu positiboa bada eta beherapen azkarreneko (7.24) iterazioa erabiltzen badugu, hau dugu:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{Q}}^2 = \left( 1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k)(\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k)} \right) \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2. \quad (7.26)$$

Frogantza. (7.25) erabiliz, zera dugu:

$$\|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{Q}}^2 = 2(f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})).$$

Bestalde, lehenengo  $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \mathbf{g}_k$  erabiliz, hau dugu:

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) &= \left(\frac{1}{2}\mathbf{x}_k^T \mathbf{Q} \mathbf{x}_k - \mathbf{b}^T \mathbf{x}_k\right) - \left(\frac{1}{2}\mathbf{x}_{k+1}^T \mathbf{Q} \mathbf{x}_{k+1} - \mathbf{b}^T \mathbf{x}_{k+1}\right) \\ &= \alpha \mathbf{g}_k^T \mathbf{Q} \mathbf{x}_k - \frac{1}{2}\alpha^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k - \alpha \mathbf{b}^T \mathbf{g}_k = \alpha \mathbf{g}_k^T (\mathbf{Q} \mathbf{x}_k - \mathbf{b}) - \frac{1}{2}\alpha^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k, \end{aligned}$$

eta gero,  $\mathbf{g}_k = \mathbf{Q} \mathbf{x}_k - \mathbf{b}$  adierazpena eta (7.23) bilaketa lineal zehatza kontuan hartuz, hau lortzen da:

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) &= \alpha \mathbf{g}_k^T \mathbf{g}_k - \frac{1}{2}\alpha^2 \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \\ &= \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} (\mathbf{g}_k^T \mathbf{g}_k) - \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k)^2} \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} - \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k)} = \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}. \end{aligned}$$

Ondorioz,

$$\|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{Q}}^2 = \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}.$$

Bestalde,  $\mathbf{g}_k = \mathbf{Q} \mathbf{x}_k - \mathbf{b} = \mathbf{Q} \mathbf{x}_k - \mathbf{Q} \mathbf{x}^* = \mathbf{Q}(\mathbf{x}_k - \mathbf{x}^*)$  denez,

$$\|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 = (\mathbf{x}_k - \mathbf{x}^*)^T \mathbf{Q} (\mathbf{x}_k - \mathbf{x}^*) = \mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k.$$

Azkenik, aurreko bi berdintzak erabiliz, zera lortzen dugu:

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{Q}}^2 &= \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} = \left(1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k) \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2}\right) \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 \\ &= \left(1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k) (\mathbf{g}_k^T \mathbf{Q}^{-1} \mathbf{g}_k)}\right) \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2. \quad \square \end{aligned}$$

Adierazpen horrek deskribatzen du  $f$ -ren beharrezko zehatza iterazio bakoitzean. Erlazionatuko dugu  $\mathbf{Q}$ -ren baldintzazko zenbakiarekin. Horretarako, Kantarovich-en desberdintza behar dugu. Frogapena [18] testuliburuan aurki daiteke.

**7.6. lema: Kantarovichen desberdintza.** *Izan bedi  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  matrize simetriko eta definitu positibo bat. Orduan,  $\forall \mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{y} \neq \mathbf{0}$ , zera betetzen da:*

$$\frac{(\mathbf{y}^T \mathbf{y})^2}{(\mathbf{y}^T \mathbf{Q} \mathbf{y})(\mathbf{y}^T \mathbf{Q}^{-1} \mathbf{y})} \geq \frac{4Mm}{(M+m)^2},$$

$M$  eta  $m$ , hurrenez hurren,  $\mathbf{Q}$  matrizearen autobalio handiena eta txikiena izanik.

Ondorioz, hori (7.26) adierazpenean aplikatuz, hau lortzen dugu:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{Q}}^2 \leq \left(1 - \frac{4Mm}{(M+m)^2}\right) \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2 = \left(\frac{M-m}{M+m}\right)^2 \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{Q}}^2,$$

eta (7.25) erabiliz, azkenik, zera aurkitu dugu:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{M-m}{M+m}\right)^2 (f(\mathbf{x}_k) - f(\mathbf{x}^*)). \quad (7.27)$$

$\mathbf{Q}$ -ren zenbakizko baldintza  $r = \kappa_2(\mathbf{Q}) = M/m$  dugunez, aurreko adierazpenaz hau lortzen da:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{r-1}{r+1}\right)^2 (f(\mathbf{x}_k) - f(\mathbf{x}^*)). \quad (7.28)$$

Emaitza hori zabal diezaiokegu edozein  $f \in \mathcal{C}^2$  helburu funtzio ez-lineali, baina kasu hone-tan  $\mathbf{Q} = \nabla^2 f(\mathbf{x}^*)$ . Aurreko desberdintzak erakusten du gradiente metodoaren konbergentzia-ratioa oso motela izan daitekeela, nahiz eta matrize hessiarraren baldintzazko zenbakia oso handia ez izan. Adibidez,  $\kappa_2(\mathbf{Q}) = 800$  eta  $f(\mathbf{x}^*) = 0$  badira, (7.28) desberdintzak iradoki-tzen du funtzioaren balioa 0.08 bider bakarrik gutxituko dela, bostehun iterazio egin ondoren gradientearen metodoa erabiltzen badugu.

### Newtonen metodo zehaztugabeak

Iterazio batean matrize hessianra definitu positiboa ez bada, Newtonen metodoak ez du lortuko norabide beherakor bat. Horrelako matrize hessianra duten problema handietarako  $\mathbf{d}_k$  norabide beherakorra aurkitzeko, Newtonen metodoaren gisako iterazioa erabiltzen dugu matrize hessianra aldatuz jarraian azaltzen den moduan.

Lehenengo, Gill and Murray-ren matrize hessiarrari Choleskyren faktORIZAZIO aldatua-ren algoritmoa aplikatuz (ikus [11])  $\mathbf{A}_k = \nabla^2 F(\mathbf{x}_k) + \mathbf{D} = \mathbf{R}^T \mathbf{R}$  lortuko dugu,  $\mathbf{D}$  matrize diagonal izanik eta bere gaiak ez-negatiboak.

Matrize hessianra definitu positiboa bada,  $\mathbf{D} = \mathbf{0}$  da, eta ez dugu aldatu behar.

Aldiz,  $\mathbf{D} \neq \mathbf{0}$  bada, Gerschgorin-en teorema kontuan hartuz (ikus 5.7. teorema)  $\mu$  positibo txikiena aukeratzen dugu  $\mathbf{A}_k = \nabla^2 F(\mathbf{x}_k) + \mu \mathbb{1}$  definitu positiboa izateko (hots, autobalio txikiena zero baino handiagoa izateko). Horretarako, nahikoa da hau hartzea:

$$\bar{\mu} = \min_{1 \leq i \leq n} \left\{ a_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \right\},$$

eta  $\mu = |\bar{\mu}| + \eta$  jarri, non  $\eta = \sqrt{\varepsilon_M}$ .

Azkenik, hau hartuko dugu:

$$\mathbf{A}_k = \nabla^2 F(\mathbf{x}_k) + \mu_k \mathbb{1}, \text{ non } \mu_k = \max\{\mu, \|\mathbf{D}\|_\infty\}.$$

Horrela kalkulaturako  $\mathbf{A}_k$  matrizea benetako matrize hessianretik simetriko definitu positibo hurbila da. Metodo hau Dennis & Schnabel-en [9] testuliburuan aurki daiteke.

### Quasi-Newton metodoak

Metodo hauek deribatuen kalkuluaren kostua ekiditeko erabiltzen dira. Ekuazio ez-linealetako sistemak ebazteko Broydenen metodoa erabili dugu. Baina, metodo hori ez datorkigu ondo  $\nabla^2 F(\mathbf{x})$  matrize hessiarra hurbiltzeko, zeren ez baitu gordetzen simetria (ikus 7.18):

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \mathbf{u}_k \mathbf{s}_k^T, \quad \text{non} \quad \mathbf{u}_k = \frac{\mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k}{\mathbf{s}_k^T \mathbf{s}_k}$$

eta orokorrean  $\mathbf{u}_k \neq \mathbf{s}_k$  gertatzen baita.

Simetria gorde nahi badugu eta ebakitzaileren ekuazioa ( $\mathbf{A}_{k+1} \mathbf{s}_k = \mathbf{y}_k$ ) betetzea, orduan hasierako  $\mathbf{A}_0$  simetrikoa izan behar da, eta eguneratzeko formula honelakoa:

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \sigma \mathbf{v}_k \mathbf{v}_k^T, \quad \text{non} \quad \sigma = \pm 1,$$

$\mathbf{v}_k$  bektore baterako. Hots, oraingo matrizea aldatzen dugu bat heineko beste matrize simetrikoko bat batuz.  $\mathbf{A}_{k+1}$  behartuz ebakitzaileren ekuazioa betetzera, adierazpen bakar hau lortuko dugu:

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \frac{\mathbf{q}_k \mathbf{q}_k^T}{\mathbf{s}_k^T \mathbf{q}_k}, \quad \text{non} \quad \mathbf{q}_k = \mathbf{y}_k - \mathbf{A}_k \mathbf{s}_k.$$

Frogapena ariketa bezala geratzen da. Formula horri *symmetric rank-one (SR1)* izenarekin ezagutzen dugu.

Bestalde, nahiz eta  $\mathbf{A}_k$  definitu positiboa izan  $\mathbf{A}_{k+1}$  indefinitua izan daiteke. Baina, ebakitzaileren metodoaren bidez ez da simetria bakarrik gorde dezakeguna,  $\mathbf{A}_k$  definitu positibo izatea ere finka dezakegu, eta, ondorioz,  $\mathbf{d}_k$  norabide beherakorra izatea. Hori lortzen da *Broyden-Fletcher-Goldfarb-Shanno (BFGS)* eguneratzea erabiltzen badugu. Eguneratze horretan, matrize hessiarren alderantzikoa ( $\nabla^2 F(\mathbf{x}_k)^{-1}$ ) hurbiltzen da  $\mathbf{H}_k$  matrizeen bidez, adierazpen honen bitartez eguneratuz:

$$\mathbf{H}_{k+1} = (\mathbb{1} - \rho_k \mathbf{s}_k \mathbf{y}_k^T) \mathbf{H}_k (\mathbb{1} - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T, \quad \text{non} \quad \rho_k = \frac{1}{\mathbf{y}_k^T \mathbf{s}_k}.$$

Orain, aurreko adierazpenean Sherman-Morrison-Woodbury formula aplikatuz,  $\nabla^2 F(\mathbf{x}_k)$  matrize hessiarren hurbilpena ( $\mathbf{A}_k$ ) honela eguneratu dezakegu:

$$\mathbf{A}_{k+1} = \mathbf{A}_k - \frac{\mathbf{A}_k \mathbf{s}_k \mathbf{s}_k^T \mathbf{A}_k}{\mathbf{s}_k^T \mathbf{A}_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}.$$

Azken bi formula horietan erraz ikusten denez, aldaketaren heina 2 da.

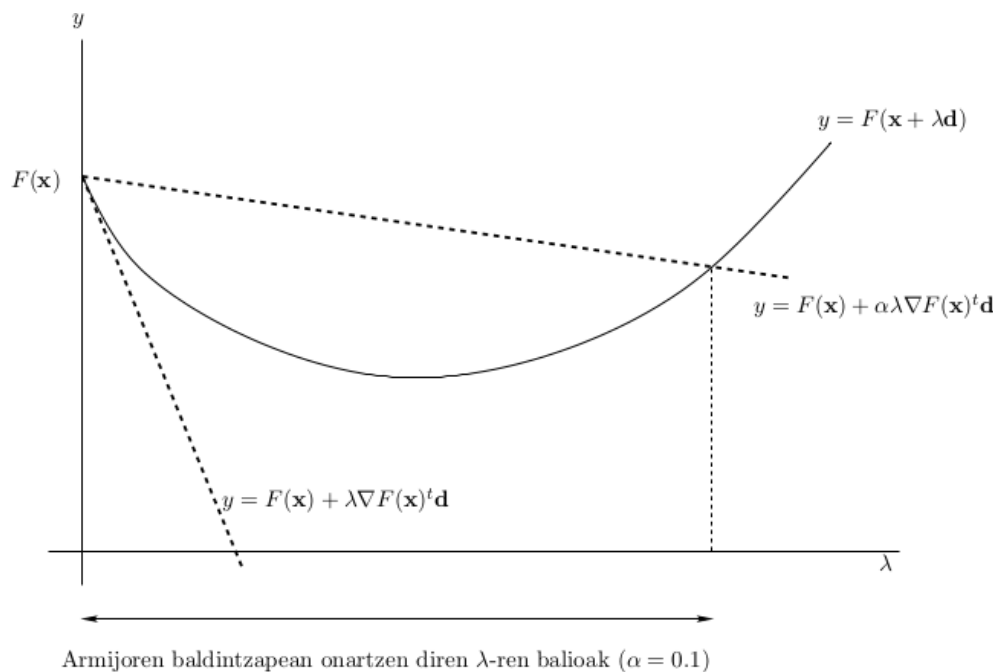
### Konbergentzia globala: Armijo-Goldstein-en baldintzak (edo Wolfe-ren baldintzak)

Armijoren baldintza erabiltzen da funtzioaren balioaren beherapen nahiko bat izateko. Alegia,  $\alpha \in (0, 1)$  hartuz, eta  $\lambda > 0$ -ri hau betetzea exijituz:

$$F(\mathbf{x} + \lambda \mathbf{d}) \leq F(\mathbf{x}) + \alpha (\nabla F(\mathbf{x})^T \mathbf{d}) \lambda. \quad (7.29)$$



Ikus 7.2. irudia.

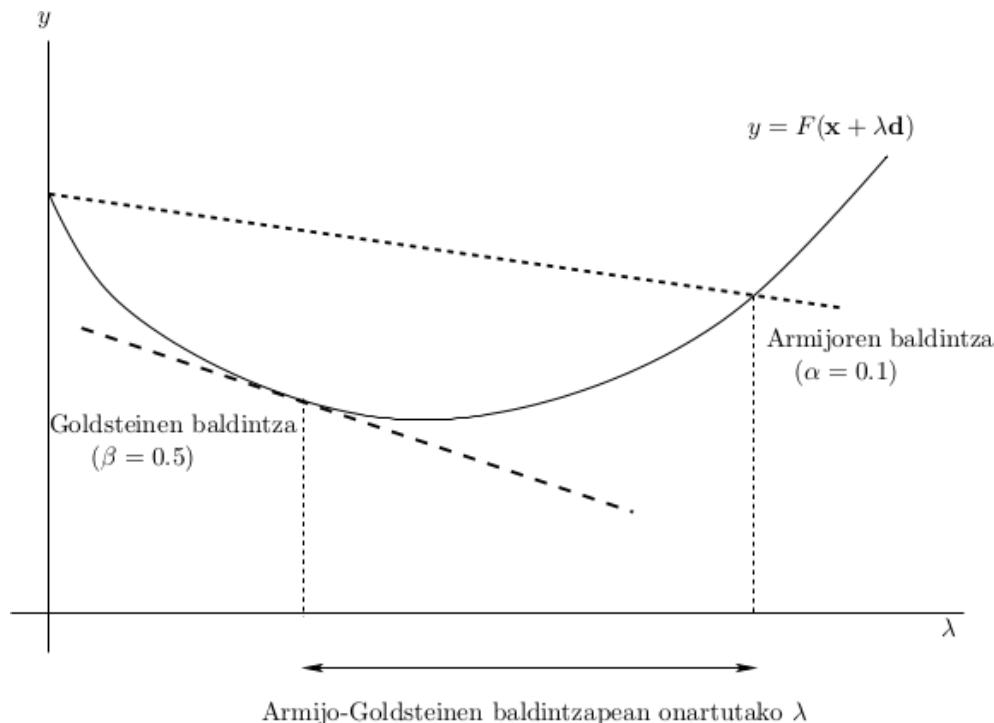


### 7.2. irudia. Armijoren baldintza.

Jakin badakigu  $\lambda$  nahiko txiki baterako  $F(\mathbf{x} + \lambda\mathbf{d}) < F(\mathbf{x})$  betetzen dela, baina ez txiki-  
 kiegia; hori da, hain zuzen ere, Goldsteinen baldintzarekin saihesten dena. Baina, guk nahi  
 dugu  $\mathbf{d}$  norabidean  $\mathbf{x} + \lambda\mathbf{d}$  puntuko  $F$ -ren beherapena aurreko  $\mathbf{x}$  puntuan dugunaren ratio  
 bat baino handiagoa izatea, hots:

$$\nabla F(\mathbf{x} + \lambda\mathbf{d})^T \mathbf{d} \geq \beta \nabla F(\mathbf{x})^T \mathbf{d}, \quad (7.30)$$

non  $\beta \in (\alpha, 1)$  (ikus 7.3. irudia). Hori da Goldsteinen baldintza. Atzeranzko estrategia  
 erabiltzen badugu, baldintza hori bide batez betetzen da.



### 7.3. irudia. Armijo-Goldsteinen baldintzak.

Honako teorema honek frogatzen du (ikus [29, 30])  $\mathbf{d}$  norabide beherakor bat izanez gero,  $\mathbf{x} + \lambda\mathbf{d}$  puntuak egonik, (7.29) eta (7.30) Armijo-Goldsteinen baldintzak (edo Wolferen baldintzak) betetzen dituztela.

**7.7. teorema.** *Izan bedi  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  etengabe diferentziagarria (hots,  $C^1$  klasekoa). Izan bitez  $\mathbf{x}_k, \mathbf{d}_k \in \mathbb{R}^n$ , non  $\nabla F(\mathbf{x}_k)^T \mathbf{d}_k < 0$  baita (hots,  $\mathbf{d}_k$  norabide beherakorra da), eta demagun  $\{F(\mathbf{x}_k + \lambda\mathbf{d}_k) \mid \lambda > 0\}$  behe bornatua dela. Orduan, baldin  $0 \leq \alpha \leq \beta < 1$ , badaude  $\lambda_u > \lambda_l > 0$  konstanteak, non  $\mathbf{x}_{k+1} = \mathbf{x}_k + \lambda_k \mathbf{d}_k$  puntuak (7.29) eta (7.30) baldintzak betetzen baititu,  $\lambda_k \in (\lambda_l, \lambda_u)$  bada.*

[9]-ko 6.3.2. teoreman aurki daiteke teorema honen frogantza eta Wolfe-n teoremaren enunziazioa (6.3.3. teorema) eta frogantza [29, 30]. Teorema horrek iterazio guztietan  $\mathbf{d}_k$  beherakorra bada eta bilaketa linealaren bitartez kalkulaturako  $\lambda_k$  Armijo-Goldsteinen baldintzak betetzen baditu, gradientea zerora konbergitzen da, edo funtzioa ez da behe-bornatua eta  $\infty$ -ra jotzen du. Hau da, teorema horrek baldintza horietan konbergentzia globala bermatzen du.  $\nabla F(\mathbf{x}_k)^T \mathbf{d}_k = 0$  ere izan daiteke, baina kasu hori saihestu dezakegu (ikus 123 orrialdea [9] liburuan).

Gainera, testuliburu berean Dennis and Moré-k erakusten dute  $k_0$  batetik aurrera  $\mathbf{d}_k$  norabide beherakorra Newtonen urratsetik nahiko hurbil badago (hots,  $\mathbf{d}_k \approx -\nabla^2 F(\mathbf{x}_k)^{-1} \nabla F(\mathbf{x}_k)$ ),

orduan  $k_0$ -tik aurrera  $\lambda_k = 1$  onargarria dela eta,  $\nabla F(\mathbf{x}^*) = \mathbf{0}$  bada, konbergentzia superlineala izango da. Hori gertatzen da quasi-Newton metodoetan.

**7.5. adibidea.** Izan bitez  $F(x_1, x_2) = x_1^4 + x_1^2 + x_2^2$ ,  $\mathbf{x} = (x_1, x_2)^T = (1, 1)^T$ ,  $\mathbf{d} = (-3, -1)^T$ , eta  $\alpha = 0.1$  (7.29) adierazpenean eta  $\beta = 0.5$  (7.30) -ean.

*Ebazpena.* Hau dugunez:

$$\nabla F(\mathbf{x})^T \mathbf{d} = (6, 2)(-3, -1)^T = -20 < 0,$$

$\mathbf{x}$  horretan  $\mathbf{d}$  norabide beherakorra da  $F(\mathbf{x})$ -rekiko. Orain, demagun  $\mathbf{x}(\lambda) = \mathbf{x} + \lambda \mathbf{d}$  dela. Baldin  $\lambda = 1$  bada,  $\bar{\mathbf{x}} = \mathbf{x}(1) = \mathbf{x} + 1 \cdot \mathbf{d} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1 \cdot \begin{bmatrix} -3 \\ -1 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \end{bmatrix}$  dugu, eta, orduan:

$$\nabla F(\bar{\mathbf{x}})^T \mathbf{d} = (-36, 0)(-3, -1)^T = 108 > -10 = \beta \nabla F(\mathbf{x})^T \mathbf{d}.$$

Beraz,  $\bar{\mathbf{x}}$  (7.30) betetzen du, baina:

$$F(\bar{\mathbf{x}}) = 20 > 1 = F(\mathbf{x}) + \alpha \lambda \nabla F(\mathbf{x})^T \mathbf{d}$$

eta, ondorioz, ez du (7.29) betetzen. Halaber,  $\lambda = 0.1$  bada,  $\bar{\mathbf{x}} = \mathbf{x}(0.1) = (0.7, 0.9)^T$  puntuak (7.29) betetzen du, baina (7.30) ez. Aldiz,  $\lambda = 0.5$  bada,  $\bar{\mathbf{x}} = \mathbf{x}(0.5) = (-0.5, 0.5)^T$  bi baldintzak betetzen ditu, (7.29) eta (7.30). Beraz, 7.3. irudian  $\lambda = 0.1$  eskualde onargarriaren ezker aldean dago,  $\lambda = 1$  eskualde horren eskuinaldean, eta  $\lambda = 0.5$  eskualde onargarrian bertan.  $\square$

## Atzeranzko bilaketa lineala

Orain zehaztuko dugu nola aukeratu behar dugun  $\lambda_k$ . Gaur egungo estrategia da hasieran  $\lambda_k = 1$  hartzea, eta  $\mathbf{x}_k + \mathbf{d}_k$  ez bada onargarria, «atzera egitea» (hots,  $\lambda_k$  txikiagotzea)  $\mathbf{x}_k + \lambda_k \mathbf{d}_k$  onargarri bat lortu arte. «Onargarriak» esan nahi du (7.29) eta (7.30) baldintzak betetzea. Bigarren baldintza betetzen da atzeranzko estrategia erabiltzeagatik, zeren horrek urrats txikiak saihesten baititu. Hori dela eta, ez da agertzen algoritmoan.

### 7.4. algoritmoa. Atzeranzko bilaketa lineala.

**0. urratsa.** Eman  $\alpha \in (0, 1/2)$ ,  $0 < l < u < 1$  eta  $\lambda_k = 1$ ;

**1. urratsa.**  $F(\mathbf{x}_k + \lambda_k \mathbf{d}_k) > F(\mathbf{x}_k) + \alpha \lambda_k \nabla F(\mathbf{x}_k)^T \mathbf{d}_k$  gertatzen den bitartean, egin  $\lambda_k := \rho \lambda_k$ , non  $\rho \in [l, u]$ ;

**2. urratsa.**  $\mathbf{x}_{k+1} = \mathbf{x}_k + \lambda_k \mathbf{d}_k$ .

Praktikan,  $\alpha$  nahiko txikia denez, funtzioaren balioari beheratze oso txikia eskatzen baitzaio, gure algoritmoan  $\alpha = 10^{-4}$ . Bertsekas-ek  $\rho = 1/2$  hartzea proposatzen du (ikus [3]). Beraz, horrela bada,  $\lambda_k$  aukeratzeko  $\{1, 2^{-1}, 2^{-2}, 2^{-3}, \dots\}$  segidari jarraituko diogu, eta helduko da  $2^{-m}$  ( $m \in \mathbb{N}$ ) balio batera, non hau betetzen baita (Armijoren baldintza):

$$F(\mathbf{x}_k + \lambda_k \mathbf{d}_k) \leq F(\mathbf{x}_k) + 10^{-4} \lambda_k \nabla F(\mathbf{x}_k)^T \mathbf{d}_k.$$

### 7.4.3. Minimo karratu ez-linealak

Honelako problemak dira:

$$\min F(\mathbf{x}) = \frac{1}{2} \sum_{j=1}^m r_j^2(\mathbf{x}),$$

non  $r_j : \mathbb{R}^n \rightarrow \mathbb{R}$  funtzio leuna baita ( $r_j \in C^1$ ) eta  $m \geq n$ . Gainera,  $r_j$  bakoitzari *hondar* deritzogu, eta  $\mathbf{r}(\mathbf{x}) = (r_1(\mathbf{x}), r_2(\mathbf{x}), \dots, r_m(\mathbf{x}))^T$  bektoreari *hondar-bektore*;  $F(\mathbf{x})$  honela idatz dezakegu:

$$F(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2.$$

Bestalde,  $F(\mathbf{x})$ -ren deribatuak honela idatz daitezke  $\mathbf{r}$  funtzioaren matrize jacobiarren bitartez:

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \frac{\partial r_1}{\partial x_1} & \frac{\partial r_1}{\partial x_2} & \cdots & \frac{\partial r_1}{\partial x_n} \\ \frac{\partial r_2}{\partial x_1} & \frac{\partial r_2}{\partial x_2} & \cdots & \frac{\partial r_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial r_m}{\partial x_1} & \frac{\partial r_m}{\partial x_2} & \cdots & \frac{\partial r_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla r_1(\mathbf{x})^T \\ \nabla r_2(\mathbf{x})^T \\ \vdots \\ \nabla r_m(\mathbf{x})^T \end{bmatrix}.$$

$$\nabla F(\mathbf{x}) = \sum_{j=1}^m r_j(\mathbf{x}) \nabla r_j(\mathbf{x}) = \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x}). \quad (7.31)$$

$$\begin{aligned} \nabla^2 F(\mathbf{x}) &= \sum_{j=1}^m \nabla r_j(\mathbf{x}) \nabla r_j(\mathbf{x})^T + \sum_{j=1}^m r_j(\mathbf{x}) \nabla^2 r_j(\mathbf{x}) \\ &= \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \sum_{j=1}^m r_j(\mathbf{x}) \nabla^2 r_j(\mathbf{x}). \end{aligned} \quad (7.32)$$

Berdintza horiek frogatzea ariketa gisa geratzen da.

### Gauss-Newtonen metodoa

$F$  funtzioa deribagarria bada, bere minimoan  $\nabla F(\mathbf{x}) = \mathbf{0}$  ekuazioa bete behar du. Gainera, baldintza hori  $\mathbf{x}^*$  puntuan betetzen bada eta  $\nabla^2 F(\mathbf{x}^*)$  definitu positiboa bada (hori gertatzen da  $\mathbf{J}(\mathbf{x}^*)$  hein betekoa bada), orduan  $\mathbf{x}^*$  puntuan  $F$  funtzioak minimo bat heltzen

du. Beraz,  $\nabla F(\mathbf{x}) = \mathbf{0}$  sistema ez-lineala ebazteko Newtonen metodoa erabiltzen badugu, iterazio bakoitzean honelako sistema bat ebatzi beharko dugu:

$$\nabla^2 F(\mathbf{x})\mathbf{d} = -\nabla F(\mathbf{x}). \quad (7.33)$$

Gauss-Newtonen metodoa Newtonen metodoaren aldaera bezala ikus daiteke. Kasu honetan,  $\mathbf{d}_k$  bilaketa-norabidea kalkulatzeko (7.33) Newtonen sisteman, (7.32) deribatuaren bigarren batugaia (hots,  $\sum_{j=1}^m r_j(\mathbf{x})\nabla^2 r_j(\mathbf{x})$ ) kendu eta gero sistema ebatziz lortzen da. Alegia,  $\mathbf{d}_k^{GN}$  kalkulatu da sistema hau ebatziz:

$$\mathbf{J}_k^T \mathbf{J}_k \mathbf{d} = -\mathbf{J}_k^T \mathbf{r}_k, \quad \text{non } \mathbf{J}_k = \mathbf{J}(\mathbf{x}_k) \text{ eta } \mathbf{r}_k = \mathbf{r}(\mathbf{x}_k). \quad (7.34)$$

Hauek dira metodo horren abantaila batzuk:

- $\nabla^2 F(\mathbf{x}_k) \approx \mathbf{J}_k^T \mathbf{J}_k$  hurbilpenak  $\nabla^2 r_j(\mathbf{x})$  matrize hessianaren kalkulua saihesten du.
- Askotan,  $\mathbf{J}^T \mathbf{J}$  gaia kendutako gaia baino askoz esanguratsuagoa da, bai  $r_j$  hondarrak txikiak direlako, bai ia linealak direlako; eta, orduan,  $\|\nabla^2 r_j\|$  txikia da. Hori dela eta, oso sarri, metodo honek Newtonen metodoaren ia portaera berdina du, eta konbergentzia lokal azkarra (ia koadratikoa). Izan ere,  $\mathbf{x}^*$  soluzioan hondar-bektorea nulua denean (hots,  $\mathbf{r}(\mathbf{x}^*) = \mathbf{0}$ ), orduan, baldin  $\mathbf{x}_0$  hasierako bektorea  $\mathbf{x}^*$ -tik nahiko hurbil badago, konbergentzia koadratikoa dela frogatu daiteke; ikus [9].
- $\mathbf{J}_k$  hein betekoa denean eta  $\nabla F(\mathbf{x}_k) \neq \mathbf{0}$ ,  $\mathbf{d}_k^{GN}$  norabidea norabide beherakorra da, zeren  $\mathbf{J}_k^T \mathbf{J}_k$  definitu positiboa baita. Orduan,  $\mathbf{x}_{k+1} = \mathbf{x}_k + \lambda \mathbf{d}_k^{GN}$  kalkulatzeko bilaketa lineala erabil dezakegu  $\lambda$  egoki bat aurkituz (hots, Wolferen baldintzak egiaztatuz), non  $F$ -ren beherapena nahikoa baita.
- Bistan denez,  $\mathbf{d}_k^{GN}$  aurkitzeko (7.34) ekuazio normalak ditugu, eta, dakigunez, balio-kideak dira sistema hori ebaztea eta honako minimo karratu lineal hauen problema ebaztea:

$$\min_{\mathbf{d}} \|\mathbf{J}_k \mathbf{d} + \mathbf{r}_k\|^2.$$

Beraz,  $\mathbf{d}_k^{GN}$  norabidea kalkulatu dezakegu ikusitako  $QR$  faktORIZAZIOA erabiliz, eta horrela ez dugu izango  $\mathbf{J}_k^T \mathbf{J}_k$  kalkulatzeko beharrik.

Metodo horrek desabantaila hauek ditu:

- Problema nahiko ez-linealak direnean edo hondar nahiko handiak dituenean, geldiro doa konbergentzia lineal lokala.
- Problema oso ez-linealak direnean edo hondar oso handiak dituenean, ez da lokalki konbergente.
- Ez dago ondo definituta  $\mathbf{J}_k$  jacobiarra zutabe hein betekoa ez denean.

### Levenberg-Marquardt-en metodoa

$\mathbf{J}_k$  hein betekoa ez denean  $\mathbf{J}_k^T \mathbf{J}_k$  singularra izango, da eta, orduan,  $\mathbf{d}_k$  beherakorra dela finkatzeko,  $\mathbf{J}_k^T \mathbf{J}_k$  matrizea aldatuko dugu  $\mathbf{D}$  matrize diagonal batez. Horretarako, Choleskyren faktORIZAZIO aldatua erabiliko dugu, eta  $\mathbf{J}_k^T \mathbf{J}_k + \mathbf{D}$  definitu positiboa izatea lortuko. Ondorioz, sistema bateragarri zehaztu honetaz lortutako  $\mathbf{d}_k$  norabidea beherakorra izango da:

$$(\mathbf{J}_k^T \mathbf{J}_k + \mathbf{D})\mathbf{d} = -\mathbf{J}_k^T \mathbf{r}_k.$$

Gero,  $\mathbf{x}_{k+1} = \mathbf{x}_k + \lambda \mathbf{d}_k$  kalkulatzeko, bilaketa lineala erabil dezakegu  $\lambda$  egoki bat aurkituz, konbergentzia globala finkatzeko.

Beste aukera bat da  $\mu$  egoki bat aurkitzea  $\mathbf{J}_k^T \mathbf{J}_k + \mu \mathbf{1}$  definitu positiboa izateko eta  $\mathbf{d}$  beherakorra; metodo horri Levenberg-Marquardten metodoa deritzogu. Orduan, norabide beherakorra honela kalkulaten da:

$$(\mathbf{J}_k^T \mathbf{J}_k + \mu \mathbf{1})\mathbf{d} = -\mathbf{J}_k^T \mathbf{r}_k \quad (7.35)$$

eta sistema hau ebaztea eta minimo karratu linealen problema hau ebaztea baliokideak dira:

$$\min_{\mathbf{d}} \frac{1}{2} \left\| \begin{bmatrix} \mathbf{J}_k \\ \sqrt{\mu} \mathbf{1} \end{bmatrix} \mathbf{d} + \begin{bmatrix} \mathbf{r}_k \\ \mathbf{0} \end{bmatrix} \right\|_2^2. \quad (7.36)$$

$\mu$  hori kalkulatzeko, *konfiantza eskualdeko metodo* bat erabil dezakegu, ikus [21].

## 7.5. Problemak

1. Ebatzi sistema hau Newtonen metodoaren 7.1. algoritmoa erabiliz,  $\mathbf{x}_0 = (x_0, y_0, z_0)^T = (0.5, 0.5, 0.5)^T$  hartuz:

$$\begin{aligned} f_1(x, y, z) &= x^2 + y^2 + z^2 - 1 = 0 \\ f_2(x, y, z) &= 2x^2 + y^2 - 4z = 0 \\ f_3(x, y, z) &= 3x^2 - 4y + z^2 = 0. \end{aligned}$$

Bigarren urratsa egiteko,  $LU$  metodoa erabili.

Baldin  $\mathbf{f}(x, y, z) = (f_1(x, y, z), f_2(x, y, z), f_3(x, y, z))^T$  bada, bukatu iteratzeko prozesua  $\|\mathbf{f}(x, y, z)\|_\infty < 6 \cdot 10^{-5}$  denean.

2. Ebatzi sistema hau Newtonen algoritmoaz,  $\mathbf{x}_0 = (x_0, y_0)^T = (1.5, 1.5)^T$  hartuz:

$$\begin{aligned} f_1(x, y) &= x^2 + y^2 - 2 = 0 \\ f_2(x, y) &= e^{x-1} + y^3 - 2 = 0. \end{aligned}$$

Baldin  $\mathbf{f}(x, y) = (f_1(x, y), f_2(x, y))^T$  bada, bukatu iteratzeko prozesua  $\|\mathbf{f}(x, y)\|_\infty < 10^{-4}$  denean.

3. Izan bedi sistema hau:

$$\begin{aligned} f_1(x, y) &= x^2 - y - 0.2 = 0 \\ f_2(x, y) &= y^2 - x - 0.3 = 0. \end{aligned}$$

Ebatzi sistema hori Newtonen metodoaz,  $\mathbf{x}_0 = (x_0, y_0)^T = (1.2, 1.2)^T$  hartuz; egin bi iterazio bakarrik. Orain  $\mathbf{x}_0 = (x_0, y_0)^T = (-0.2, -0.2)^T$  hartu, eta egin Newtonen metodoaren lehenengo bi iterazioak.

4. Ebatzi bigarren problema diferentzia finituzko Newtonen metodoaz.  
 5. Ebatzi bigarren problema Broydenen metodoaz.  
 6. Kalkulatu Broydenen metodoaren lehenengo hiru iterazioak sistema honetarako:

$$\begin{aligned} f_1(x, y) &= x_1 + x_2 - 3 = 0 \\ f_2(x, y) &= x_1^2 + x_2^2 - 9 = 0, \end{aligned}$$

$\mathbf{x}_0 = (x_0, y_0)^T = (2, 7)^T$  eta  $\mathbf{A}_0 = \mathbf{J}(\mathbf{x}_0)$  hartuz. Zein da  $\|\mathbf{x}_3 - \mathbf{x}_2\|_\infty$  zehaztasuna? Eta errorea (hots,  $\|(f_1(\mathbf{x}_3), f_2(\mathbf{x}_3))^T\|$ )?

7. Izan bedi sistema hau:

$$\begin{aligned} f_1(x, y) &= x^2 + y^2 - 2 = 0 \\ f_2(x, y) &= xy - 1 = 0. \end{aligned}$$

- (a) Egiaztatu  $(1, 1)$  eta  $(-1, -1)$  sistema horren soluzioak direla.

(b) Zelako arazoak ager daitezke Newtonen metodoa aplikatzen saiatzen bagara soluzio horiek kalkulatzeko?

8. Aurkitu sistema ez-lineal hauen soluzio bat, Newtonen metodoa erabiliz:

$$(a) \quad \begin{aligned} x_1^2 - 10x_1 + x_2^2 + 8 &= 0, \\ x_1x_2^2 + x_1 - 10x_2 + 8 &= 0; \end{aligned} \quad \text{hartu } \mathbf{x}_0 = (1.2, 1.2)^T.$$

$$(b) \quad \begin{aligned} 3x_1^2 - x_2^2 &= 0, \\ 3x_1x_2^2 - x_1^3 - 1 &= 0; \end{aligned} \quad \text{hartu } \mathbf{x}_0 = (0.5, 0.5)^T.$$

Ebatzi sistema horiek Newtonen metodoaz, eta bukatu  $\|\mathbf{f}(\mathbf{x})\|_\infty < 10^{-3}$  denean. Sistema bakoitzean, zenbat iterazio erabili ditugu? Zein da zehaztasun erlatiboa? Zein da konbergentziaren ordena?

9. Ebatzi aurreko problemako (a) eta (b) sistemak Broydenen metodoaz, eta bukatu  $\|\mathbf{f}(\mathbf{x})\|_\infty < 10^{-3}$  denean. Sistema bakoitzean, zenbat iterazio erabili ditugu? Zein da zehaztasun erlatiboa? Zein da konbergentziaren ordena?

10. (a) Froga ezazu  $f(\mathbf{x}) = x_1^2 - x_2^4 + 1$  funtzioak koordenatu-jatorrian zela-puntu bat duela. Hots, koordenatu-jatorria puntu kritikoa dela, baina ez dela minimo bat ezta maximo bat ere.

(b) Zer gertatzen da Newtonen metodoa aplikatzen badugu zela-puntu hori kalkulatzeko?

11. Izan bedi optimizazioko problema hau:

$$\min 5x^2 + 5y^2 - xy - 11x + 11y + 11.$$

(a) Aurkitu minimo izateko lehenengo ordenako baldintza beharrezkoak betetzen dituen puntu bat (hots, puntu kritikoa bat).

(b) Froga ezazu puntu hori minimo globala dela.

(c) Zein izango litzateke gradiente metodoaren konbergentzia-ratioa problema honi aplikatzen badiogu?

(d)  $\mathbf{x}_0 = (0, 0)^T$  puntutik abiatuz, gradiente metodoko zenbat iterazio erabili behar ditugu gehienez, funtzio errorea  $10^{-11}$  bider gutxitzeko?

12. Izan bedi helburu-funtzio hau:  $f(\mathbf{x}) = 3x_1^2 + 2x_1x_2 + x_2^2$ ,  $\mathbf{x}_0 = (1, 1)^T$ .

(a) Zein da  $f$ -rako beherapen azkarreneko norabidea  $\mathbf{x}_0$ -tik?

(b)  $(1, -1)^T$  norabide beherakorra da?

(c) Lehenengo ataleko norabide beherakorrerako, eta atzeranzko bilaketa lineala erabiliz, aurkitu  $\lambda$  onargarri bat  $\lambda_0 = 1$ ,  $\rho = 0.5$  eta  $\alpha = 0.1$  hartuz.

13. Izan bedi optimizazioko problema hau:  $\min f(x, y) = 3x^2 + y^4$ .



- (a) Gradiente-metodoaren iterazio bat aplikatu  $\mathbf{x}_0 = (1, -2)^T$  hartuz, eta,  $\lambda$  aurkitzeko, atzeranzko bilaketa lineala erabili  $\lambda = 1$ ,  $\rho = 0.5$  eta  $\alpha = 0.1$  hartuz.
- (b) Errepikatu (a)-ren kalkulua, baina orain  $\lambda = 1$ ,  $\rho = 0.1$  eta  $\alpha = 0.1$  hartuz. Konparatu kalkuluetan kostatu dena aurreko atalean kostatu denarekin.
- (c) Orain, egin Newtonen metodoaren iterazio bat (a) ataleko baldintza berdinekin  $\lambda$  onargarria kalkulatzeko. Zer gertatu da? Konparatu kalkuluetan kostatu dena aurreko ataletan kostatu denarekin.
14. Izan bedi helburu-funtzio hau:  $f(x, y) = x^2 + y^2 + xy - 3x$ .
- (a) Aurkitu  $f$ -ren minimo lokal bat.
- (b) Zergatik da (a)-ren soluzioa minimo global bat?
- (c) Gradientearen metodoa aplikatzen badugu (a) ataleko minimoa aurkitzeko, zein izango litzateke helburu-funtzioaren konbergentzia-ratioa?
15. Aurkitu gradiente-metodoaren konbergentzia-ratioa aplikatzen diogunean problema honi:
- $$\min f(x, y) = x^2 + 1.999xy + y^2.$$
16. Rosenbrock-en funtzioa  $\phi(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$  da, eta, erraz ikusten den bezala,  $(1, 1)^T$ -ean minimo bakarra heltzen du. Aplikatu gradientearen metodoa,  $\mathbf{x}_0 = (0, 0)^T$ -tik hasiz, iterazio bat bakarrik eginez, eta atzeranzko bilaketa erabiliz,  $\lambda = 1$ ,  $\rho = 0.5$  eta  $\alpha = 0.1$  hartuz. Zein da gutxi gorabehera konbergentzia-ratioa? Zenbat iterazio behar ditu, gutxi gorabehera, funtzio-errorea  $10^{-6}$  bider gutxitzeko?
17. Rosenbrock-en  $\phi(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$  funtzioaren minimo bakarra kalkulatzeko, erabili Newtonen metodoa eta BFGS metodoa bi iteraziotan zehar. Bi kasuetan, atzeranzko bilaketa erabiliz,  $\lambda = 1$ ,  $\rho = 0.5$  eta  $\alpha = 0.1$  hartuz, eta  $\mathbf{x}_0 = (0, 0)^T$ . Zein metodorekin hurbildu da hobeto soluzioa?
18. Doitu  $d = \phi(t, \mathbf{x}) = x_1 e^{-x_2 t}$  funtzioa  $t_i = -1, 0, 1, 2$  puntuetan  $d_i = 2.7, 1, 0.4, 0.1$  balioei. Horretarako, erabili Gauss-Newtonen bi iterazio  $\mathbf{x} = (1, 1)^t$  puntutik. Bilaketa linealerako, erabili atzeranzko bilaketa,  $\lambda = 1$ ,  $\rho = 0.5$  eta  $\alpha = 10^{-4}$  hartuz.
19. Froga ezazu 7.4.3. ataleko (7.35) eta (7.36) metodoak baliokideak direla  $\mathbf{d}_k$  kalkulatzeko.
20. Izan bedi  $\mathbf{r} : \mathbb{R} \rightarrow \mathbb{R}^3$ , non  $r_i = e^{t_i x} - y_i$ ,  $i = 1, 2, 3$ , eta  $f(x) = \frac{1}{2} \mathbf{r}(x)^T \mathbf{r}(x)$ , non datu hauek baititugu:

$$(t_1, y_1) = (1, 2) \quad (t_2, y_2) = (2, 4) \quad (t_3, y_3) = (3, \cdot).$$

Jo dezagun  $y_3$ -ren balioak eta  $x_0$  hasierako puntuak balio hauek har ditzaketela:

- (a)  $y_3 = 8$  eta  $x_0 = 1$ .
- (b)  $y_3 = -1$  eta  $x_0 = 0$ .

Ebatzi problema hau (a) eta (b) kasuetan, Newtonen metodoaz eta Gauss-Newtonen metodoaz, urratsaren luzera aldatu barik (hots,  $\lambda$  erabili gabe). Metodo bakoitzeko, bukatu iterazioak  $\|\nabla f(\mathbf{x}_k)\|_\infty < 10^{-5}$  denean. Zer gertatzen da kasu bakoitzean? Zergatik?

### MATLABez programatzeko problemak:

21. Ekuazio ez-linealen sistematarako, 7.1. Newtonen metodoaren algoritmoa dugu. Eraiki M-fitxategi bat goiko algoritmoa inplementatzeko. Fitxategi horren izena `snewton.m` izango da.

- (a) Gogora ezazue goiburu hau izan behar duela:
- ```
function [x,nf,err,ze,i]=snewton(x0,emax,zemax,imax),
```
- non `nf` funtzio-balioztatze kopurua baita.
- (b)  $\mathbf{f}(\mathbf{x})$  bektorea kalkulatzeko, funtzio bat definitu behar da M-fitxategi baten bidez, `fun.m` izenekoa (ez ahaztu, bektoreak zutabeak direla idazkera matrizialean).
- (c) Era antzeko batean kalkulatu da  $\mathbf{f}(\mathbf{x})$ -ren deribatua; hots,  $\mathbf{J}(\mathbf{x})$  matrizea kalkulatzeko, M-fitxategi bat eraikitzen da, `Jfun.m` izenekoa.
- (d) Lehenengo urratsean  $\mathbf{x}_k$  punturako  $\mathbf{f}$  funtzioaren eta  $\mathbf{J}$  deribatuaren balioak kalkulatu dira iterazio horretan.
- (e) Algoritmo horren 2. urratsean, (7.3) Newtonen sistema ebatzi behar dugu. Urrats horretan, erabili  $LU$  deskonposizioa eta aurreranzko eta atzeranzko ebazpenak.
- (f) Sistema ebatzi eta gero,  $\mathbf{x}_{k+1}$  berria izango dugu; hor (4. urratsean) erabiliko dugun norma infinitu-norma izango da.  $\|\mathbf{p}_k\|_\infty/\|\mathbf{x}_{k+1}\|_\infty < ze$  edo  $\|\mathbf{f}(\mathbf{x}_{k+1})\|_\infty < emax$  edo  $i = imax$  denean, gelditu egiten da; bestela, 1. urratsera joango da.
- (g) Izan bedi  $\mathbf{f}(\mathbf{x}) = \begin{bmatrix} x^2 + y^2 - 2 \\ e^{x-1} + y^3 - 2 \end{bmatrix}$ . Ebatzi  $\mathbf{f}(\mathbf{x}) = \mathbf{0}$  sistema Newtonen metodoaz, zuk eraikitako kodearen bitartez. Hartu  $\mathbf{x}_0 = (2, 2)^T$ ,  $ze = 10^{-6}$ ,  $emax = 10^{-6}$  eta  $imax = 20$ .

22. (a) Eraiki `J=dfjfun(x)` goiburuko M-fitxategi bat diferentzia finituzko metodoak erabiliz  $\mathbf{J}(\mathbf{x})$  matrize jakobiarra hurbiltzeko; ikus 7.2.1. atala, non  $h = \sqrt{\varepsilon_M}$  (MATLABen `eps`= $\varepsilon_M$  da).

(b) Aurreko kodea ondo ibiltzen dela egiaztatu eta gero, eraiki beste kode bat, non  $\mathbf{J}(\mathbf{x})$  matrize jakobiar hurbildu hori baita. Haren goiburua hau izango da:

```
function [x,nf,err,ze,i]=dfnewton(x0,emax,zemax,imax),
```

non `nf` funtzio-balioztatze kopurua baita.

Funtzio horrek Newtonen metodoaren funtzioarena egingo du, baina berari matrize jakobiar hurbildua emanaz, ez benetakoa. Berak jakobiarraren balio hurbildua lortuko du, `dfjfun.m` erabiliz.

- (c) Ebatzi aurreko ariketaren azken atala, parametro berdinekin, `dfnewton` funtzioa erabiliz. Gero, kodean, aldatu  $h$ -ren balioa  $h = 0.1, 0.01, 0.001, 0.0001$  desberdinetarako, eta, taula batez, konparatu bere eraginkortasuna: iterazio kopurua eta funtzio-balioztatze kopurua.
23. Implementatu 7.2. algoritmoa: Quasi-Newton metodoa, matrize hessianra Broydenen metodoaz eguneratuz. Haren goiburua hau izango da:
- ```
function [x,nf,err,ze,i]=qnewton(x0,emax,zemax,imax),  
non nf funtzio-balioztatze kopurua baita.
```
- Ebatzi 21. ariketako (g) atala, parametro berdinekin, `qnewton` funtzioa erabiliz.
24. Konparatu 21.(g) ataleko problemarako 21.-23. ariketetako hiru kodeen eraginkortasuna, taula bat sortuz eta zutabe bakoitzean metodo bakoitzerako funtzio-balioztatze kopurua eta iterazio kopurua idatziz. Zer ondorio atera dezakezu?

# Bibliografia

- [1] U.M. Ascher eta Ch. Greif, *A First Course in Numerical Methods*. SIAM, 2011.
- [2] K.E. Atkinson, *An Introduction to Numerical Analysis*. John Wiley and Sons, 1989.
- [3] D.P. Bertsekas, *Constrained Optimization and Lagrange Multiplier Methods*. Academic Press, New York, 1982.
- [4] D.P. Bertsekas, *Nonlinear Programming. (Second edition)* Athena Scientific, Belmont, USA, 2003.
- [5] R.P. Brent, *Algorithms for Minimization without Derivatives*. Prentice-Hall, Englewood Cliffs, NJ, 1973.
- [6] C.G. Broyden, *A Class of Methods for Solving Nonlinear Simultaneous Equations*. Mathematics of Computation (American Mathematical Society) 19 (92): 577–593.
- [7] S.C. Chapra, *Applied Numerical Methods with MATLAB for Engineers and Scientists*. McGraw-Hill, AEB, 2008.
- [8] T.J. Dekker, *Finding a zero by means of successive linear interpolation*. In B. Dejon; P. Henrici, *Constructive Aspects of the Fundamental Theorem of Algebra*, Wiley-Interscience, Londres, 1969.
- [9] J.E. Dennis eta R.B. Schnable, *Numerical methods for unconstrained optimization and nonlinear equations*. Prentice-Hall, 1983. SIAMek berrinprimatuta, 1996.
- [10] I.S. Duff, A.M. Erisman eta J.K. Reid, *Direct Methods for Sparse Matrices*. Oxford University Press, 1986.
- [11] P.E. Gill, W. Murray eta M.H. Wright, *Practical Optimization*. Academic Press, Londres, 1981.
- [12] P.E. Gill, W. Murray eta M.H. Wright, *Numerical Linear Algebra and Optimization, Volume 1*. Addison-Wesley, 1991.
- [13] J.A. George eta J.W. Liu, *Computer Solution of Large Sparse Positive Definite Systems*. Prentice-Hall, Englewood Cliffs, NJ, 1981.

- [14] A. Gilat eta V. Subramaniam, *Numerical Methods for Engineers and Scientists: An Introduction with Applications Using MATLAB*. John Wiley & Sons, AEB, 2008.
- [15] G.H. Golub eta C.F. Van Loan, *Matrix Computations*. The John Hopkins University Press, Baltimore, 1996.
- [16] D.J. Higham eta N.J. Higham, *MATLAB Guide*. SIAM, 2005.
- [17] N.J. Higham, *Accuracy and Stability of Numerical Algorithms*. SIAM, 2002.
- [18] D.G. Luenberger, *Linear and Nonlinear Programming*. Addison-Wesley, AEB, 1989.
- [19] J.F. Mathews eta K.D. Fink, *Métodos Numéricos con MATLAB*. Pearson, Prentice Hall, Madrid, 2005.
- [20] C.B. Moler, *Numerical Computing with MATLAB*. SIAM, Filadelfia, AEB, 2004.
- [21] J. Nocedal eta S.J. Wright, *Numerical Optimization*. Springer Series in Operations Research, New York, 1999.
- [22] J.M. Ortega eta W.C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York, 1970.
- [23] M. Overton, *Numerical Computing with IEEE Floating Point Arithmetic*. SIAM, Filadelfia, 2001.
- [24] J.M. Sanz-Serna, *Diez lecciones de Cálculo numérico*. Secretariado de Publicaciones e Intercambio Editorial, Universidad de Valladolid, 2010.
- [25] J. Stoer eta R. Bulirsch, *Introduction to Numerical Analysis*. Springer-Verlag, 2002.
- [26] G.W. Stewart, *Introduction to Matrix Computations*. Academic Press, New York, 1973.
- [27] G.W. Stewart, *Matrix Algorithms: Basic Decompositions*. SIAM, Filadelfia, 1998.
- [28] L.N. Trefethen eta D. Bau, *Numerical Linear Algebra*. SIAM, Filadelfia, 1997.
- [29] P. Wolfe, *Convergence conditions for ascent methods*. SIAM Review 11, 226–235.
- [30] P. Wolfe, *Convergence conditions for ascent methods. II: Some corrections*. SIAM Review 13, 185–188.

# Kontzeptuen Aurkibidea

- K*(*A*) azpiespazioa, 163
- L* matrizea, 118
- LDL<sup>T</sup>* faktORIZAZIOA, 135
- LU* faktORIZAZIOA, 122, 124
- P* matrizea, 118
- QR* faktORIZAZIOA, 155, 158, 160, 200, 213
- R<sup>T</sup>R* faktORIZAZIOA, 135
- SVD* deskonposizioa, 177
- U* matrizea, 118
- $\gamma$ -funtzio Lipschitz jarraitua, 190
- MATLAB, 5
  
- Alborapena, 57
- Alderantzizko interpolazio koadratikoa (IQI), 104
- Algoritmo
  - egonkorra, 68
  - ezegonkorra, 68
- Algoritmoa, 43
- Armijo-Goldstein-en baldintzak, 210
- Armijoren baldintza, 208
- Atzeranzko bilaketa lineala, 211
- Atzeranzko ebazpena, 122
- Atzeranzko ordezkatzeko prozesua, 119
- Aurreranzko ebazpena, 122, 195
- Aurreranzko ordezkatzeko prozesua, 120
- Autobalio errealak, 133
- Autobektoreen oinarri ortonormala, 133
  
- Bakartze-metodoak, 83
- Baldintzazko zenbaki orokortua, 179
- Baldintzazko zenbakia, 126
- Balio singulararrak, 177
- Balio singularretako deskonposizioa, 176
- Beharpen azkarreneko metodoa, 204
- Bektore singulararrak, 177
- Berretzailea, 55
  
- Biderkatzaileak, 116
- Bilaketa lineal zehatza, 205
- Biribiltzea, 56
- Biribiltzearen unitatea, 51, 57, 58
- Bisekzio-metodoa, 85, 88
- Bit, 55
- Broyden-en alderantzizko formula, 198
  
- Cholesky-ren faktORIZAZIOA, 133
- Choleskyren faktORIZAZIO ALDATUA, 207
- Cramer-en metodoa, 115
  
- Definitu positiboa, 133
- Deskonposizio ortogonal osoa, 175, 176
- Diagonal menperatzailea, 125, 134, 145
- Diferentzia finituak, 3
- Diferentzia finituzko Newtonen metodoa, 192
- Digitu esanguratsuak, 57
  
- Ebakitzailearen ekuazioa, 196
- Ebakitzailearen metodo aldatua, 101
- Ebakitzailearen metodoa, 100
- Ekuaizio ez-linealen sistema, 187
- Ekuaizio normalak, 164, 213
- Eragindako matrize-norma, 127
- Errore
  - absolutua, 44
  - erlatiboa, 44
- Errorearen hedapena, 67, 68
- Espekto-erradioa, 129
- Eulerren metodoa, 3
- Ezabapen-matrizeak, 122
- Ezabatze gaussiarra, 122
- Ezabatze-prozesua, 121
  
- Frobeniusen norma, 128
- Funtzioaren erroa, 83

- Gauss-Newtonen metodoa*, 213
- Gauss-Seidel-en iterazioa*, 141
- Gerschgorin-en teorema*, 134
- Givens-en biraketak*, 161
- Goldsteinen baldintza*, 209
- Gradientearen metodoa*, 204
  
- Hondar optimoa*, 175
- Hondar-bektorea*, 163, 165, 212
- Horner-en metodoa*, 62
- Householderren*
  - bektorea*, 156
  - islapenak*, 155, 166
  - matrizea*, 155, 157, 175
- Hurbilpen lineala*, 187
- Hurbiltze-ordena*, 64
  
- Identitate-matrizea*, 118
- Iterazio-funtzioa*, 96
  
- Jacobi-ren iterazioa*, 139
  
- Kantorovichen desberdintza*, 206
- Konbergentzia*
  - p ordenakoa gutxienez*, 95
  - koadratikoa*, 94, 96
  - lineala*, 94, 100
  - superlineala*, 94
- Konbergentzia superlineala*, 198
  
- Levenberg-Marquardt-en metodoa*, 214
  
- Makinaren errorea*, 57, 194
- Mantisa*, 55, 56
- MATLAB*, xi
- MATLABeko*
  - \ ikurra*, 11, 116
  - / ikurra*, 11, 116
  - ; ikurra*, 8
  - azpifuntzioak*, 20
  - Booleko balioak*, 25
  - egiturak*
    - for...end*, 25
    - if...else...end*, 26
    - if...elseif...else...end*, 26
    - if...end*, 26
  - switch*, 27
  - while*, 28
- era bektorizatua*, 25
- eragiketa*
  - .\*, ./, .\ eta .^ gaiEz gai*, 11
  - arimetikoak*, 5
  - matrizialak*, 10
- eragile logikoak*, 25
- erlazio eragileak*, 24
- funtzio*
  - funtzioak*, 33
  - anonimoak*, 32
  - argumentua*, 35
  - erabilgarriak*, 6
  - M fitzategiak*, 18
  - nagusia*, 20
  - parametroak aldatzea*, 35
- funtzioak*
  - chol(A)*, 137
  - cond(A,p)*, 133
  - fzero*, 105
  - inv(A)*, 133
  - nnz(A)*, 138
  - norm(X,p)*, 128
  - roots*, 105
  - svd(A)*, 133
  - disp*, 21
  - fprintf*, 22
  - fprintfen formatuak eta kontrolak*, 22
  - input*, 21
- grafikoak*, 13
- instrukzioak*
  - = esleipena*, 7
  - beep*, 30
  - break*, 29
  - disp*, 29
  - feval*, 18
  - fplot*, 14
  - hold off*, 15
  - hold on*, 15
  - linspace*, 9
  - logspace*, 9
  - lookfor*, 19
  - meshgrid*, 15

- pause, 30
- plot3, 14
- plot, 13
- print, 15
- subplot, 15
- tic, 30
- toc, 30
- M* fitzategia, 17
- MAT fitzategi bitarra, 23
- matrizeak, 7
- Matrize
  - behe-triangeluar unitatea, 122
  - behe-triangeluarra, 119
  - goi-triangeluarra, 119
  - jacobiarra, 187, 212
  - ortogonala, 119
- Matrize baten
  - baldintza, 130, 165
  - baldintzazko zenbakia, 130
- Matrize-
  - norma, 126
  - norma bateragarriak, 127
  - normen trinkotasuna, 126
- Metodo
  - egonkorren konbergentzia, 142
  - globalak, 204
  - grafikoak, 81
  - irekiak, 83
  - lokala, 100
- Minimo karratu ez-linealak, 212
- Minimo karratu linealak, 162, 213
- Moore-Penrose-ren sasiaderantzizko matrizea, 178
- Mullerren metodoa, 102
- Newtonen metodo aldatua, 195
- Newtonen metodo orokortua, 100
- Newtonen metodoa, 95, 188, 213
  - erro anizkoitza daudenean, 100
  - periodikotasuna, 100
- Norabide beherakorra, 210
- $O()$  notazioa, 64
- Oinarri-soluzioa, 173, 174
- Ondo baldintzatua, 126, 131
- Ordenagailuaren errorea, 56
- Ordenagailuaren zehaztasuna, 56
- Permutazio-matrizea, 118
- Pibota, 116, 120
- Pibotatze baztergarria, 125
- Pibotatze partziala, 121
- Pibotatzea, 117
- Problemaren baldintza, 126
- Puntu finko erakargarria, 93
- Puntu finkoaren
  - iterazioa, 90
  - konbergentziaren analisia, 90
- Puntu higikorra, 55
- Quasi-Newton metodoak
  - BFGS eguneratzea, 208
  - SR1 eguneratzea, 208
- QuasiBroyden-en metodoa, 196
- Regula falsi metodoa, 88
- Sasiaderantzizko matrizea, 164, 178
- Sistema baten baldintzazko zenbakia, 129
- Sistema lineal gaindeterminatua, 163
- Sistema linealak: ebazpen zuzenak eta iterati-  
boak, 115
- Soluzioa
  - analitikoa, 1, 2
  - zehatza, 1
  - zenbakizkoa, 2
- Trunkatze-errorea, 64
- Txarto baldintzatua, 126, 131
- Wolferen baldintzak, 208, 213
- Zarata, 70
- Zatiki bitarrak, 53
- Zehaztasun bakuna, 58
- Zehaztasun bikoitza, 57
- Zenbaki bitarrak normalizatzea, 55
- Zenbaki-sistema
  - bitarra, 51
  - hamartarra, 51
- Zenbakizko metodoak, 2



*Zeroin algoritmoa, 105*

*Zifra esanguratsuen ezeztapena, 61*

*Zutabe hein beteko matrize baten QR faktori-  
zazioa, 166*

*Zutabe hein beteko matrizea, 164*

*Zuzen ukitzailea, 98*

**UNIBERTSITATEKO ESKULIBURUAK**  
MANUALES UNIVERSITARIOS

INFORMAZIOA ETA ESKARIAK • INFORMACIÓN Y PEDIDOS

UPV/EHUko Argitalpen Zerbitzua • Servicio Editorial de la UPV/EHU  
argitaletxea@ehu.eus • editorial@ehu.eus  
1397 Posta Kutxatila - 48080 Bilbo • Apartado 1397 - 48080 Bilbao  
Tfn.: 94 601 2227 • [www.ehu.eus/argitalpenak](http://www.ehu.eus/argitalpenak)

eman ta zabal zazu



Universidad  
del País Vasco

Euskal Herriko  
Unibertsitatea