

DATUEN ANALISIA ETA ENUNTZIATUEN LOGIKA *

Yosu YURRAMENDI

1. Sarrera

Ikertzaik askoren lanak datuen jasoketa batekin hasten dira sarritan. Datu pila baten aurrean ez dago ikertzailearentzako, dakigunez, inolako araubiderik horren gainean gorputz teoriko bat eraikitzeko. Dena dela, zeregin horretan saiatzea ez da lan antzua, datuak nolabait deskribatu, mamitu eta sailkatu behar bait dira.

Bada teknika multzo bat zeregin hoietarako lehenengo urratsak ematen laguntzeko balio duena: estatistika deskribatzailea deritzana. Hona hemen John Tukey estatistikari ospetsuak zer zioen aspalditxo "Annals of Mathematical Statistics, vol. 33, 1962" aldizkarian "The Future of Data Analysis" izenburu pean (2. orr.):

"For a long time I have thought I was a statistician, interested in inferences from the particular to the general. But as I have watched mathematical statistics evolve, I have had cause to wonder and to doubt".

Eta aurreraxeago:

"I have come to feel that my central interest is in data analysis, which I take to include, among other things: procedures for analyzing data, techniques for interpreting the results of such procedures, ways for planning the gathering of data to make its analysis easier, more precise or more accurate, and all the machinery and results of (mathematical) statistics which apply to analyzing data".

Beharbada, lerroaldi hauetan oso garbi geratu ez den zerbait, datu pila bat analisatzea zertan datzan jakitea da. Gehienetan, datuak zenbait objeturi behatutako ezaugarrietatik datoz; datuak analisatzea beraz, objetuen arteko erlazioak, ezaugarrien artekoak eta, objektu eta ezaugarrien artekoak azaltzean edo argitaratzean datza.

Datuen analisisan erabili behar diren prozedura eta tekniken nolokotasunari honela deritzo John Tukey-k aipatutako (6. orr.):

"Data analysis, and the parts of statistics which adhere

* El abstract inglés y la versión castellana de este artículo siguen en las páginas 275-291.

to it, must then take on the characteristics of a science rather those of mathematics, specifically:

(b-1) Data analysis must seek for scope and usefulness rather than security.

(b-2) Data analysis must be willing to err moderately often in order that inadequate evidence shall more often suggest the right answer.

(b-3) Data analysis must use mathematical argument and mathematical results as bases for judgment rather than as "bases for proof or stamps of validity".

Eta geroxeago:

"All sciences have much of art in their makeup. (...) As well as teaching facts and well-established structures, all sciences must teach their apprentices how to think about things in the manner of that particular science, and wath are its current beliefs and practices. Data analysis must do the same."

Hemen "structure" hitza axioma eta definizio multzo bat bezala ulertzen dugu. Arrazoinamendu -(logiko)- aren bidez bereizgarriak edo propietateak ondorioztatzen dira.

Idazki honen garaietatik gaurko egunetara kalkulugailuen garapen eta hedapena itzelak izan dira. Beren kalkuluetarako gaitasun eta oroi-menak, lehen bururatu ere ezin ziren eginkizunak ekarri ditu, eta honela, aztergarriak diren datu pilak gero eta haundiagoak dira.

Ikus dezagun orain, datuen analisiaren hastapenak zeintzu diren hamar bat urte beranduago, eskola frantzesako burua edo izan den Jean-Paul Benzecri-ren iritziz, bere inguruko lankide batzurekin argitaratutako "L'Analyse des Données. Tome I: La taxinomie. Tome 2: L'Analyse des correspondances. Edt. Dunod, 1973, 4^{ème} édition 1982" idazlanean (3-17 orr.):

"7^{er} Principe: Statistique n'est pas probabilité. Sous le nom de statistique mathématique, des auteurs (qui, je vous le dis en français, n'écrivent guère dans notre langue ...) ont édifié une pompeuse discipline, riche en hypothèses qui ne sont jamais satisfaites dans la pratique. Ce n'est pas de ces auteurs qu'il faut attendre la solution de nos problèmes

DATUEN ANALISIA ETA ENUNTZIATUEN LOGIKA

typologiques".

.....

"2^{ème} Principe: Le modèle doit suivre les données, non l'inverse. (...) Mais ce dont nous avons besoin c'est d'une méthode rigoureuse qui extraie des structures à partir des données.

.....

3^{ème} Principe: Il convient de traiter simultanément des informations concernant le plus grand nombre possible de dimensions.

.....

4^{ème} Principe: Pour l'analyse des faits complexes et notamment de faits sociaux, l'ordinateur est indispensable. Principe évidemment vrai ... mais qu'en eussent pensé nos pères les Gaulois il y a 15 ans?

.....

5^{ème} Principe: Utiliser un ordinateur implique d'abandonner toutes techniques conçues avant l'avènement du calcul automatique. Je dis technique, non sciences: (...)"

Ikusten denez, bi irakasleek badute zenbait puntu elkarren artean, beren aurkezteko estiloak zeharo desberdinak izan arren. Beren ustez, badirudi, datuen analisisa estatistika deskribatzaile klasikoaren itxura berria dela. Gure idazki honetan, hastapen hauei erantzuten dien eredu baten eraketa azalduko dugu, guztiz orokorra dena eta hain zuzen enuntziatuen logikarekin lotura duena. Bertan zenbait arazo planteatuko dugu.

2. Galdeketa eredu

Datu pila baten aurrean eredu bat eraikitzerakoan, ikuspegiren batetik edo bestetik, jasotako datuen berreraiketa baten ahalmena eman behar da. Gehienetan ordea, asmo hau nekez bete daiteke, datuen koda-ketak berak informazio guztiaren zati bat besterik ez bait du biltzen.

Hemen azalduko dugun ereduak galdeketa izena du, gizazientzia delakoetan maiz erabili ohi dena gogoratzen du eta.

Bedi Ω aztergai diren objektuen multzoa. Aurrerantzean finitua dela joko dugu:

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$$

Bedi X , Ω multzoaren ezaugarri baten agertzeko eren multzoa,

edo konparazioari jarraika, Ω -koek galdera bati eman liezizkioketen erantzunen multzoa. Hau ere finitutzat hartuko dugu:

$$X = \{x_1, x_2, \dots, x_p\}$$

Bedi X^* , Ω -koei egindako galdera bat:

$$X^*: \Omega \rightarrow X$$

hau da, Ω -ren aplikazio edo aldagai bat.

Bitez

$$\{X_n \mid n=1, \dots, r\} \quad \text{eta} \quad \{X_n^* \mid n=1, \dots, r\}$$

r galderari eman lekizkieken erantzunen multzo-sorta bat eta honi dago-kion aplikazio-sorta. Honela, ereduak suposatzen duena zera da: Ω -ko bakoitzak galdera bakoitzari erantzuten diola, eta gainera, era bakar batez.

Bedi

$$\{e_{nk}^* \mid n=1, \dots, r \text{ y } k=1, \dots, p_n\}$$

Ω -ren aplikazio-sorta bat $\{0,1\}$ multzoan, era honetan definituta:

edozein $\omega \in \Omega$, edozein $n = 1, \dots, r$ eta edozein $k = 1, \dots, p_n$ -entzako baldin

$$X_n^*(\omega) = X_{nk},$$

bada orduan

$$e_{nk}^*(\omega) = 1,$$

eta bestela

$$e_{nk}^*(\omega) = 0.$$

Ohar daitekenez, multzo hau erantzun guztien identifikatzaileena da, funtzio karakteristikoena alegia. Multzo honen bereizgarri edo propietate bat hau da:

edozein $\omega \in \Omega$, eta edozein $n = 1, \dots, r$ -entzako

$$\sum_{k=1}^{p_n} e_{nk}^*(\omega) = 1$$

DATUEN ANALISIA ETA ENUNTZIATUEN LOGIKA

E ikurraz izendatutako dugu multzo hau.

Noski,

$$\bigcup_{n=1}^r X_n$$

eta E multzoen artean korrespondentzia biuniboko bat dago.

Eredu hau eraikitzerakoan, asmoa, esan bezala, Ω -koen artean izan daitezken erlazio nabarmenak (zehaztu beharko den zentzu batetan) argitaratzea da, hala nola erantzunen artekoak eta, objektu eta erantzunen artekoak ere, guztiok, erantzunei dagozkien identifikatzaileen bitartez.

Bedi

$$E^*: \Omega \rightarrow \{0, 1\}^P, \quad P = \sum_{n=1}^r P_n,$$

izanik, era honetan definituta:

edozein $\omega \in \Omega$ -rentzako

$$E^*(\omega) = (e_{nk}^*(\omega) \mid n=1, \dots, r; k=1, \dots, P_n).$$

Eraikitako eredia murriztu egingo dugu galdera edo aldagai dikotomikoen kasuetara, idazki honen xederako nahikoa izango bait da. Edozein $n = 1, \dots, r$ -rentzako X_n multzoa bi elementuz osatuta dago:

$$X_n = \{x_n^+, x_n^-\},$$

eta bi hauei dagozkien identifikatzaileak e^{*+} eta e^{*-} dira. Arestian aipatutako bereizgarria beste honetan bihurtu da:

edozein $\omega \in \Omega$ eta edozein $n = 1, \dots, r$ -rentzako

$$e_n^{*-}(\omega) = 1 - e_n^{*+}(\omega).$$

Hortaz, datorrenerako E^* aplikazioaren pareko hau hartuko dugu gogotan:

$$E^{*+}: \Omega \rightarrow \{0, 1\}^{\Gamma},$$

edozein $\omega \in \Omega$ -rentzako

$$E^{*+}(\omega) = (e_n^{*+}(\omega) \mid n=1, \dots, r).$$

3. Enuntziatuen logika

X_n , $n = 1, \dots, r$ multzoak definitzean, esan bezala, r galderari eman lekizkieken erantzunak definitzen ditugu. Hauek bi baldin badira, erraz pentsa liteke kualitate beten jabe izatea adierazi dezakela batek, eta ez izatea besteak.

Enuntziatuen logikan (logikaren mailarik apalena bestalde) bi enuntziatuen X_k , X_l arteko loturak edo erlazioak, $(X_k \times X_l)$ multzoaren azpimultzoak dira. $(X_k \times X_l)$ bidekarketa cartesiarra 4 elementuz osatuta dago, eta dakigunez, 2^4 lotura edo erlazio desberdin definitu litezke.

Arruntenak, adibidez, hauek dira.

- a) Konjuntzioa: $X_k \& X_l = \{(x_k^+, x_l^+)\}$
- b) Disjuntzioa: $X_k \vee X_l = \{(x_k^+, x_l^+), (x_k^+, x_l^-), (x_k^-, x_l^+)\}$
 $= \{(x_k^-, x_l^-)\}^c$
- c) Baldintza: $X_k \rightarrow X_l = \{(x_k^+, x_l^+), (x_k^-, x_l^+), (x_k^-, x_l^-)\}$
 $= \{(x_k^+, x_l^-)\}^c$
- d) Baldintza bikoitza: $X_k \leftrightarrow X_l = X_k \rightarrow X_l \cap X_l \rightarrow X_k$
 $= \{(x_k^+, x_l^+), (x_k^-, x_l^-)\}$
- e) Bateriaezinezko disjuntzioa¹: $X_k \nleftrightarrow X_l = \{(x_k^+, x_l^-), (x_k^-, x_l^+)\}$

Bestalde, $\{0, 1\}^2$ -ren aplikazioak $\{0, 1\}$ multzoan 2^4 dira ere, eta beraz, korrespondentzia biuniboko bat egokitu daiteke erlazio eta aplikazioen artean parekaketa hau egin ondoren:

- (x_k^+, x_l^+) eta (1, 1);
- (x_k^+, x_l^-) eta (1, 0);
- (x_k^-, x_l^+) eta (0, 1);
- (x_k^-, x_l^-) eta (0, 0).

Adibideetan azaldutako erlazioak honela adierazi litezke deritzaten funtzio logikoren bitartez

$$\begin{aligned} \text{edozein } (e_1, e_2) &\in \{0, 1\}^2 \text{ -rentzako} \\ \text{a) } f_{\&}(e_1, e_2) &= e_1 e_2 \end{aligned}$$

DATUEN ANALISIA ETA ENUNTZIATUEN LOGIKA

- b) $f_{\vee}(e_1, e_2) = 1 - (1 - e_1)(1 - e_2)$
- c) $f_{\rightarrow}(e_1, e_2) = 1 - e_1(1 - e_2)$
- d) $f_{\leftrightarrow}(e_1, e_2) = 1 - (e_1 - e_2)^2$
- e) $f_{\nleftrightarrow}(e_1, e_2) = (e_1 - e_2)^2$

$\prod_{n=1}^r X_n$ biderkaketa cartesiarraren azpimultzoak r enuntziatuen arteko erlazioak adierazten dituzte, eta horietako batzuek badute zentzurik hizkuntza mailan. Aurreko adibideetako erlazioak, esaterako, multzo berri honetan daude murgilduta:

- a) $X_k \& X_l = \{(x_1, \dots, x_r) | x_n \in X_n \ (n=1, \dots, r) \ \& \ x_k = x_k^+ \ \& \ x_l = x_l^+\}$
- b) $X_k \vee X_l = \{(x_1, \dots, x_r) | x_n \in X_n \ (n=1, \dots, r) \ \& \ x_k = x_k^+ \ \vee \ x_l = x_l^+\}$
- c) $X_k \rightarrow X_l = \{(x_1, \dots, x_r) | x_n \in X_n \ (n=1, \dots, r) \ \& \ x_k = x_k^- \ \vee \ x_l = x_l^+\}$
- d) $X_k \leftrightarrow X_l = (X_k \rightarrow X_l) \cap (X_l \rightarrow X_k)$
- e) $X_k \nleftrightarrow X_l = (X_k \leftrightarrow X_l)^c$

Bi dimentsiotan egin dugunaren antzeko parekaketa baldin badago- kie $\{0,1\}^r$ -ko elementuei, adibideetako erlazioak funtzio logiko hauen bitartez adierazi litezke:

edozein $(e_n \mid n = 1, \dots, r) \in \{0,1\}^r$ -rentzako

- a) $f(e_1, \dots, e_r) = e_k e_l$
- b) $f(e_1, \dots, e_r) = 1 - (1 - e_k)(1 - e_l)$
- c) $f(e_1, \dots, e_r) = 1 - e_k(1 - e_l)$
- d) $f(e_1, \dots, e_r) = 1 - (e_k - e_l)^2$
- e) $f(e_1, \dots, e_r) = (e_k - e_l)^2$

Logikaren ikuspegitik arazo bat enuntziatuen mailan planteatzen de- nean, beti aurkitu daiteke f funtzio logiko bat arazo horren ispilu dena (definitutako korrespondentzia biunibokoaren zentzuan).

Soluzioa $\{0,1\}^r$ -ren $f^{-1}(1)$ azpimultzoa da. Dakigunez, $f^{-1}(1) = \emptyset$ baldin bada erlazioari kontresana deitzen diogu; aitzitik, $f^{-1}(1) = \{0,1\}^r$ bada, tautologia.

Datuen analisiaren ikuspegitik ordea, arazoa bestelakoa da. Erantzun- en identifikatzaileen bitartez, $\{0,1\}^r$ -ren azpimultzo batetara heltzen gara, eta arazoa zein funtzio logikoen soluzioa den jakitean datza, hala nola azpimultzo horren parte berezi batzuri zeintzu dagozkie (berezita- sun hori aurrerago geratuko da zehaztuta).

Ω -ren tamaina (kardinala) nahiko handia denean (demagun $N = 500$) eta ezaugarrien kopurua ere ($r = 30$), Ω -ren irudia E^{*+} aplikazioaren bitartez, aski handia izaten da bi edota hiru ezaugarrien arteko erlazioak ($(X_k \times X_l)$ eta antzekoak ez bestek) errazki azaltzeko (hura hauen soluzioa izateko alegia), nahiz 2^{30} -ren ondoan ($\{0,1\}^{30}$ -ren kardinala) oso txikia izan; datuen ulerkuntzak bestalde, erlazio ximpleak eskatzen ditu (esaterako, konjuntzioak, disjuntzioak, gai gutxietako baldintzak, eta abar).

Aurkako bi oharpen hauen artean soluziobide bat bilatu beharra dago. Estatistikaren hildotik abiatuz hurbilketa bat egin ohi da: erlazio ximpleak bilatzen dira Ω edota bere parte berezi batzuren gehiengoarentzako ("gehiengo" hitza zehaztu beharko den zentzu batetan).

Zer nolako loturak aztertu nahi izateak metodoaren muinean egon behar du. Adibidez, bi galderaren arteko baldintzazko erlazioak aztertu nahi bagenitu, ez litzateke okerra izango

$$\sum_{w \in \Omega} (1 - e_k^{*+}(w))(1 - e_l^{*+}(w))$$

bezalako zenbakiak kalkulatzeko galdera bikote bakoitzarentzako (Ω osoan beteko balitz erlazio hori, batura horrek N balioko luke), eta haundienak hartzea; esaterako, datuei Guttman-en eskala² egokitzea ohizko praktika izan da orain dela urte batzue arte.

Beste adibide bat, bi galderaren arteko konjuntziozko erlazioena da. Bi galderaren elkartasun maila, biei batera baiezkoa eman dutenen kopuruan oinarritu liteke:

$$n_{kl}^{++} = \sum_{w \in \Omega} e_k^{*+}(w)e_l^{*+}(w)$$

Sarritan ordea, ez da sendoa izaten ezezkoiei baliorik ez ematearen arrazoia: bietako galderaren bat alderantziz planteatuz gero, elkartasun maila desberdina izango zen. Halako kasu baten aurrean jokabiderik zuzenena izaten da galderaren bi erantzunak gogotan hartzea. Horrela bi galderari dagozkien lau gertakizunek, $X_k \times X_l$ -koek alegia, ezagubide berdina izaten dute. Bestalde, ez da informazio gehiagorik galtzen.

Dagikunez:

$$n_{kl}^{++} + n_{kl}^{+-} + n_{kl}^{-+} + n_{kl}^{--} = N$$

| | | |
|---------|---------------|---------------|
| | x_1^+ | x_1^- |
| x_k^+ | n_{kl}^{++} | n_{kl}^{+-} |
| x_k^- | n_{kl}^{-+} | n_{kl}^{--} |

Datuen analisiaren metodoek indartsuak izan behar dute erlazioak aditzera emateko ikertzaleari: ezaugarrien azterketa orokor batetara bideratu behar dira, eta agerian utzi behar dute beren arteko loturak, antzekotasunak eta desberdintasunak.

4. Adierazpide geometrikoa

Atal honetan, eraikitako eredia espazio geometriko batetan kokatuko dugu. Kokatze honen erakuntza guztiz orokorra izan liteke beste zenbait eredurentzako.

Bedi

$$J = \bigcup_{n=1}^r X_n$$

erantzunen multzoa.

Aipatutako elkartasun maila aplikazio hau da: $k: J \times J \rightarrow N$, edozein $(j, j') \in J \times J$ -rentzako

$$k(j, j') = \sum_{w \in \Omega} e_j^*(w) e_{j'}(w)$$

Definizio honetatik aise ondorioztatzen dira datozen propietateak:

- a. Aplikazio hau simetrikoa da: edozein $(j, j') \in J \times J$ -rentzako $k(j, j') = k(j', j)$
- b. Edozein $j \in J$ -rentzako $k(j, j) \in \Omega$ -n izandako bere maiztasuna da.
- c. Edozein $j \in J$ eta edozein X_n , $n = 1, \dots, r$, $k(x_n^+, j) + k(x_n^-, j) = k(j, j)$
- d. Edozein X_n , $n = 1, \dots, r$ -rentzako $k(x_n^+, x_n^-) = 0$
- e. Edozein X_n , $n = 1, \dots, r$ -rentzako $k(x_n^+, x_n^+) + k(x_n^-, x_n^-) = N$

Elkartasun mailak hobeto ulertze arren egokia izaten da bakoitza proporzionalki adieraztea:

edozein $j \in J$ eta edozein $j' \in J$ -entzako:

$$f_{j'}^j = k(j, j') / k(j, j); (f_{j'}^j \geq 0 \text{ beraz})$$

Horrela:

a. edozein X_n , $n = 1, \dots, r$ eta edozein $j \in J$ -entzako:

$$f_{X_n^+}^j + f_{X_n^-}^j = 1$$

b. edozein X_n , $n = 1, \dots, r$ -rentzako:

$$f_{X_n^+}^X = 1; f_{X_n^-}^X = 0; f_{X_n^+}^X = 0; f_{X_n^-}^X = 1.$$

Erantzun baten profila $|\mathbb{R}^{2r}$ -ko puntu hau da:

edozein $j \in J$ -rentzako

$$f_j^j = (f_{j'}^j | j \in J) \in |\mathbb{R}^{2r}$$

Erantzun baten masa zenbaki positibo hau da:

edozein $j \in J$ -rentzako

$$f_j = k(j, j) / N.$$

Bereizgarri hau dute bestalde: edozein X_n , $n = 1, \dots, r$ -rentzako

$$f_{X_n^+}^X + f_{X_n^-}^X = 1; \quad \sum_{j \in J} f_j = r \quad \text{beraz.}$$

Honela, erantzunen multzoari $|\mathbb{R}^{2r}$ -ko puntu-sistema bat egokitu zaio:

$$N(J) = \{(f_j^j, f_j) | j \in J\},$$

$f_j^j | j \in J$ -ren koordinatuak izanik, eta f_j bere masa edo sistema horretan $j \in J$ -ri dagokion garrantzia (bere maiztasunaren arauera).

Puntu-sistema honen grabitate-zentrua (ulerbide mekanikoa) edo batezbesteko profilaren (ulerbide estatistikoa) koordinatuak hauexek dira: edozein $j \in J$ -rentzako:

$$1/r \times \sum_{j' \in J} f_{j'}^j \times f_j^{j'} = 1/r \times \sum_{j' \in J} \frac{k(j', j')}{N} \times \frac{k(j', j)}{k(j', j')} = 1/r \times \frac{rk(j, j)}{N} = f_j$$

edo laburki esanda: $f_j = (f_j | j \in J) \in |\mathbb{R}^{2r}$ da $N(J)$ -ren batezbesteko profila.

Bestalde, galdera bakoitzaren bi erantzunen profilen batezbestekoa

hauxe bera da ere:

edozein X_n , $n = 1, \dots, r$ eta edozein $j \in J$ -entzako:

$$f_{x_n^+} \times f_j^{x_n^+} + f_{x_n^-} \times f_j^{x_n^-} = \frac{k(x_n^+, x_n^+)}{N} \times \frac{k(x_n^+, j)}{k(x_n^+, x_n^+)} + \frac{k(x_n^-, x_n^-)}{N} \times \frac{k(x_n^-, j)}{k(x_n^-, x_n^-)}$$

$$= \frac{k(j, j)}{N} = f_j$$

Sistema honen bereizgarri bat beraz, hauxe da: edozein galderaren bi erantzunen profilak eta batezbeteko profila, zuzen batean daude.

Sistema honen barrutik erantzunen arteko erlazioak argitaratu nahi ditugu.

Ohar daitekenez, bi galderaren artean baldintza bikoitzazko erlazioa edo bateraezinezko disjuntziozkoa baldin badago, orduan beren erantzunen profilak binaka hartuta nolabait, berdinak dira. Eta alderantziz, bi erantzunen profilak berdinak baldin badira, beren aurkakoak bai eta ere, eta beraz, dagozkien bi galderen arteko erlazioa baldintza bikoitzazkoa da (baiezkoak edota ezezkoak badira berdinak) ala bateraezinezko disjuntziozkoa (bestela).

Honela esan dezakeguna zera da: bi erantzunen hurbiltasuna espazio honetan, beren baliokidetasunaren legez dagoela burututa. Hurbiltasuna ordea, ez dago oraindik zeharo definituta: bi profilen arteko distantzia definitu behar da erantzunak espazio metriko batetan murgiltzeko.

Halako distantzia bati bereizgarri hau betetzea eskatu lekiok: bi galdera baldintza bikoitzazko erlazioan badaude (edo bateraezinezko disjuntziozkoan) orduan bi galderak batutzeak (espazioaren dimentsioak laburtzen) ez ditu aldatu behar gainontzeko erantzunen arteko distantziak.

Eskakizun honi euklidearra izatearena gehitzen baldin bazaio (ezaguna. zaigun errealitate bati lotzeko, eta gainera, erosotasun matematiko bat lortzeko), orduan holako distantzia bat definitzen da:

edozein $(j', j'') \in J \times J$ -rentzako

$$d^2(j', j'') = \sum_{j \in J} (f_j^{j'} - f_j^{j''})^2 / f_j$$

Ikusten denez, batugaiak erantzunen garrantziaz daude neurtuta. Distantzia honi χ^2 -ren distantzia deritzaio, bi galderaren kasuan estatistika

klasikoan χ^2 -ren testarena gogoratzen duelako.

Erantzunen adierazpide geometrikoa azalduz gero, datorrenean bere analisirako bi metodo matematiko orokor aurkeztuko ditugu: faktore-analisisa eta sailkaketa³.

5. Faktore-analisisa.

Puntu sistema baten faktore-analisiaren metodoa mende honetan burutu zen psikometriaren alorrean Thurstone-n eskutik (1930. urtearen gerroztik), adimen orokorrari buruzko Spearman-en teoriak formalizatu.

Ikuspuntu asko izan da alor honetan, baina datuen analisiarenatik honela aurkeztu liteke: $N(J)$ puntu-sistemarentzako ardatz nagusia bilatu \mathbb{R}^{2^J} espazioan, hau da, erantzunak edozein ardatzean proiektatu (definitutako distantziaren zentzuan), eta inertzia (ulerbide mekanikoa) edo bariantza (ulerbide estatistikoa) haundiena duena bilatu (karratu minimoen hastapena aplikatu); ondoren, ardatz honekiko perpendikularren artean arazo bera planteatu, eta honelaxe aritu espazioaren dimentsioa agortu arte.

Erantzunek ardatz hauetan hartzen dituzten balioei faktoreak deritze, eta nolabait esateko, beren koordenatu berriak dira \mathbb{R}^{2^J} espazioan. Koordenatu kartesiar hauek bereizgarri batzuk dituzte; esaterako, gehiengotan faktoreen arteko lehenengo batzuren artean ia informazio guztia (bariantzaren zentzuan) mamitzen da.

Ω -ko elementuak ere proiektatu daitezke espazio honetan (beraz, ardatz nagusietan), beren profil hauek definituz:

$$\text{edozein } w \in \Omega \text{ eta edozein } j \in J \text{-entzako } f_j^w = e_j^*(w).$$

Galdeketa-ereduan faktore-analisisa $k: J \times J \rightarrow \mathbb{N}$ aplikazioan oinarritu beharrean, $k': \Omega \times J \rightarrow \{0, 1\}$ aplikazioan oinarritu bagenu, edozein $(w, j) \in \Omega \times J$ -rentzako $k'(w, j) = e_j^*(w)$ izanik, J -ko elementuen profilak, masak eta elkarren arteko distantziak arestian definitutakoaren antzeko era batez definituta (Ω -ko elementuekiko), hala nola Ω -koenak (J -ko elementuekiko), badirudi erantzunen beste faktore-analisi bat lortuko gurekela. Izatez, bi analisiak antzekoak dira eta badute lotura matematikorik. Gainera, Ω -ko elementuen faktoreak berdinak dira bi analisisietan.

Frogatu egiten da ere, bigarren faktore-analisan $N(\Omega)$ eta $N(J)$ puntu-sistemen ardatz nagusiak berberak direla, eta beraz, batera adierazi daitezkeela erantzunak eta objektuak.

DATUEN ANALISIA ETA ENUNTZIATUEN LOGIKA

Praktikan emaitza hauek garrantzi haundikoak dira: alde batetik, ardatz bakoitzean zein objektu zein erantzunetik daude hurbil (χ^2 -ren distantziaren legez) begiratu liteke eta honela elkartasun maila erlatiboak neurtu; bestetik, nahikoa da emaitzak lortzeko lehenengo analisiari ekitea, honen kalkuluen dimentsioak askoz laburragoak izaten dira eta.

Erantzunen puntu-sistema ardatz nagusi batetan (edota bi ardatz nagusik osatzen duten planu batetan) proiektatzen denez gero, erantzunen arteko antzekotasunak eta, egokiro azaldu daitezke. Adibidez, galderen artean Guttman-en eskala baten egitura balego, lehenengo faktoreak ordenatuko lituzke erantzun guztiak; are gehiago, gainontzeko faktoreak lehenengoaren funtzio polinomikoak izango lirateke.

Faktoreak kalkulatzeko, r -ren balioa oso txikia den kasutan ezik, ordenagailuen erabilpena ezinbestekoa gertatzen da; eskuz egitea zeharo larza eta luzea izango zen.

6. Sailkaketa

Sailkaketa metodoen helburua multzo bateko elementuak sail homogeno gutxitan elkartzea da.

Oraingo honetan elementu bakoitza sail bakar batetan dagoen kasutara mugatuko gara, partiketen⁴ kasutara alegia.

Metodo hauek bitan banatu litezke: hierarkikoak eta ez hierarkikoak.

Metodo ez hierarkikoek partiketa bat osatzen dute zuzenean, sail kopurua zein den aldeaz aurretik erabakiz gero. Metodo hierarkikoek berriaz, partiketa segida bat osatzen dute, partiketak fintasunaren⁵ arauera ordenatuta daudelarik.

Esan daiteke, azken metodo hauek XVII. eta XVIII. mendeetako taxonomiariek burutu zituztela.

Metodo hauen adibide bat, ondoko algoritmo ximple bezain emankorrean oinarritzen da:

1.- Sail bakunen partiketa hartu, alegia, multzoaren elementu bakoitzak sail bat osatzen duenarena.

2.- Bi sail batu erizpide bati jarraituz, eta beraz, partiketa berri bat eraiki (noski, zaharra finagoa da berria baino).

3.- Partiketa berria multzo osoarena baldin bada, gelditu; bestela, 2. pausora itzuli.

Algoritmo hau eta antzekoak aplikatzerakoan, auzia erizpidean dago.

Kontutan hartuaz partiketa baten aurrean, alde batetik sailen barruko inertzia edo bariantza kalkulatu daitekela, eta bestetik sailen artekoa (sail bakoitzeko ordezkari bat definituz gero), erizpide bat (Ward (1963)--ena deritzana) hauxe izan liteke: sailen barruko (arteko) bariantza txikiena (haundiena) emango duten bi sailak batu.

Ω eta J multzoen gain sailkaketa bana eraiki bada, sail bakoitzaren ordezkarien arteko elkartasun mailak azalduaz, hasierakoen (partiketa bakenen) kondentsazio bat lortuko da, eta informazioa mamituta agertuko. Multzo baten eta bestearen sailak ("bereziak" izendatu ditugu idazki honen hasieran) elkarren artean konparatu eta interpretatu litezke. Honela Ω -ren parte edo sail bateko elementu guztiek era berean erantzun baldin badiete galdera batzuri (erantzunek J -ren sail "berezi" bat osatzen dutelarik), orduan esan liteke Ω -ren sail horretan konjuntziazko erlazio bat betezen dela.

Baldintzazko erlazioen arazoak leku berezi bat dute sailkaketaren metodoen barruan. Ω -ko elementuen partiketa bat aldeztetik ezartzen bada, hain zuzen galdera batek (edo gehiago) definitzen duena, bere sailen eta gainontzeko galderen sailen artean halako erlazioak bilatu daitezke. Estatistikan, arazo hau bereizketa⁶ izenaren bidez da ezaguna.

Azkenik, sailkaketa eta faktore-analisiaren arteko lotura dugu. $N(J)$ puntu-sistema aztertzeke distantzia berbera erabili baldin bada bi metodoetarako, biok elkarren osagarriak gertatzen dira, hau da, faktoreak (erlazio nabarmenenak) eta sailak elkarren artean interpretatzeko ezagutzen dira zenbait adierazle. Xehetasun hauetan, hala nola idazki honetan zehar ikutu baizik ez ditugun bestetan sartu nahi duen irakurleak, tresna aberastak aurkitu litzake ondoko bibliografian.

OHARRAK

¹ Bietako galderaren bat, eta bakar bat, alderantziz planteatuz gero, baldintza bikoitzazko erlazioa azaltzen da.

² Delako Guttman-en eskala bat r galderaren gain $(r-1)$ baldintzazko erlazioa ximple dira, orden egoki batetan horrelako egitura dutelarik:

$$(x_1 \rightarrow x_2) \cap (x_2 \rightarrow x_3) \cap \dots \cap (x_{r-1} \rightarrow x_r).$$

³ Factor analysis and classification.

DATUEN ANALISIA ETA ENUNTZIATUEN LOGIKA

⁴ Bitez A eta IP bi multzo. IP A -ren partiketa bat da, baldin eta soilik baldin:

a. edozein $P \in IP$ -rentzako, $P \subset A$ y $P \neq \emptyset$;

b. edozein $P \in IP$ eta edozein $P' \in IP$ -rentzako, $P \neq P'$ izanik,
 $P \cap P' = \emptyset$ bada;

eta

c. $\bigcup_{P \in IP} P = A$;

⁵ Bitez IP eta IP' A multzoaren bi partiketa. IP , IP' baino finagoa da baldin eta soilik baldin edozein $P \in IP$ -rentzako badago $P' \in IP$, $P \subset P'$ izanik.

⁶ discrimination.

BIBLIOGRAFIA

BENZECRI, J. P. *Histoire et Préhistoire de l'Analyse des Données*,
París: Dunod, 1982.

BENZECRI, J. P. et al. *L'Analyse des Données* (Tome I: *La Taxinomie*;
Tome II: *L'Analyse des Correspondances*), París: Dunod, 1^{ère} éd.
1973; 4^{ème} éd. 1982.

TUKEY, J. W. "The Future of Data Analysis", *Annals of Mathematical Statistics*, Vol. 33 (1962).

TUKEY, J. W. *Exploratory Data Analysis*, Addison Wesley, 1977.

Facultad de Filosofía y CCEE

Universidad del País Vasco, San Sebastián