



STATISTIKARAKO SARRERA

Koerlazioa, erregresioa eta datu
anizkoitzen analisia

Plangintza berriei egokitua

Karmele Fernandez Agirre

ESTATISTIKARAKO SARRERA

Koerlazioa, erregresioa eta
datu anizkoitzen analisia

Plangintza berriei egokitua

Karmele Fernandez Agirre

© Karmele Fernandez Agirre

© Udako Euskal Unibertsitatea

I.S.B.N.: 84-86967-70-8

Lege-gordailua: BI-19-96

Inprimategia: BOAN S.A. Padre Larramendi 2. BILBO

Banatzaileak: UEU. General Concha 25, 6. BILBO
Zabaltzen: Igarabide, 88. DONOSTIA

AURKIBIDEA

HITZAURREA	XIII
I. EZAUGARRI ESTATISTIKO BAKUNAK, BANAKETAK, TAULAK, ADIERAZPIDE GRAFIKOAK	
<i>I.1. ESTATISTIKAREN HISTORIA LABURRA. ESTATISTIKA ETA PROBABILITATEA</i>	3
<i>I.2. OROKORTASUNAK</i>	5
<i>I.2.1. Populazioa eta lagina</i>	5
<i>I.2.2. Unitate estatistiko edo indibiduoak</i>	6
<i>I.3. EZAUGARRI ESTATISTIKOAK</i>	6
<i>I.3.1. Ezaugarri estatistiko bakunak</i>	6
<i>I.3.1.1. Ezaugarri kuantitatiboak</i>	7
<i>I.3.1.2. Ezaugarri kualitatiboak</i>	7
<i>I.4. MAIZTASUN-BANAKETA BAKUNAK</i>	8
<i>I.4.1. Maiztasun-banaketak (absolutuak, erlatiboak eta metatuak)</i>	8
<i>I.4.2. Taulak</i>	9
<i>I.5. MAIZTASUN-BANAKETA BAKUNEN ADIERAZPIDE GRAFIKOAK</i>	13
<i>I.5.1. Barra-diagramak</i>	13
<i>I.5.2. Histogramak</i>	15
<i>I.5.3. Maiztasun-poligonoak</i>	17
<i>I.5.4. Diagrama linealak</i>	18
<i>I.5.5. Beste adierazpide grafikoak</i>	18
II. EZAUGARRI BAKUNEN BALIO TIPIKOAK	
<i>II.1. MOMENTUAK</i>	23
<i>II.1.1. Momentu arruntak edo jatorriarekikoak</i>	23
<i>II.1.2. Momentu zentratuak edo batezbestekoarekikoak</i>	24

II.2. MAIZTASUN-BANAKETA BAKUNEN BALIO	
TIPIKO EDO ESTADISTIKOAK	26
II.2.1. Zentru-joeraren edo posizioaren balio tipikoak	26
II.2.1.1. Batezbesteko aritmetikoa	27
II.2.1.2. Batezbesteko geometrikoa	30
II.2.1.3. Batezbesteko harmonikoa	31
II.2.1.4. Moda	31
II.2.1.5. Mediana eta koantilak	34
II.2.2. Sakabanatzearen balio tipikoak	36
II.2.2.1. Bariantza	37
II.2.2.2. Batezbesteko desbidazioa	38
II.2.2.3. Ibiltartea	39
II.2.3. Itxuraren balio tipikoak	40
II.2.3.1. Asimetri koefizientea	40
II.2.3.2. Kurtosi koefizientea edo zapaltasun/zorroztasun-koefizientea	41
II.3. ALDAGAIEN TRANSFORMAZIO LINEALAK	42
II.3.1. Aldagai zentratua	43
II.3.2. Aldagai tipifikatua	44
II.4. KONTZENTRAZIO-NEURKETAK	44
II.4.1. LORENZ-en kurba	45
II.4.2. GINI-ren indizea	47

III. EZAUGARRI ESTADISTIKO BIKOITZAK, BANAKETAK, TAULAK, ADIERAZPIDE GRAFIKOAK

III.1. EZAUGARRI ESTADISTIKO BIKOITZAK	51
III.2. MAIZTASUN-BANAKETA BIKOITZAK	51
III.2.1. Maiztasun-banaketa bikoitzak eta bazter-maiztasunak	51
III.2.2. Maiztasun-banaketa bikoitzak: Taulak	53
III.3. MAIZTASUN-BANAKETA BIKOITZEN	
ADIERAZPIDE GRAFIKOAK	55
III.3.1. Sakabanatze-diagrama edo puntu-hodeia	55
III.4. MAIZTASUN-BANAKETA BALDINTZATUAK	57

III.5. TAULA BATEN DEPENDENTZIA EDO	
INDEPENDENTZIA	60
III.6. KONTINGENTZI TAULA BATEN ERRENKADA	
ETA ZUTABEEN BATEZBESTEKO SOSLAIK	62
IV. EZAUGARRI BIKOITZEN BALIO TIPIKOAK	
IV.1. MOMENTUAK	65
IV.1.1. Momentu arruntak edo jatorriarekikoak	65
IV.1.2. Momentu zentratuak edo batezbestekoarekikoak	66
IV.2. MAIZTASUN-BANAKETA BIKOITZEN BALIO	
TIPIKO EDO ESTADISTIKOAK	67
IV.2.1. Kobariantza: S_{xy}	67
IV.2.2. Koerlazio-koefizientea: r_{xy}	68
IV.3. ALDAGAIEN TRANSFORMAZIO LINEALAK	70
IV.4. KOERLAZIO GABEKO ALDAGAIEN BI PROPIETATE	73
IV.5. SAKABANATZE- ETA KOERLAZIO-MATRIZEAK	74
IV.6. OROKORPENA	75
V. KOERLAZIOA ETA ERREGRESIOA	
V.1. SARRERA	79
V.2. ERREGRESIOA R^2-n	79
V.2.1. Batezbestekoaren erregresioa	80
V.2.2. Karratu txikiaren erregresioa	80
V.3. KARRATU TXIKIEN ERREGRESIO LINEALA R^2-n	81
V.3.1. Karratu txikiaren erregresio zuzena	81
V.3.2. Karratu txikiaren erregresio linealaren propietateak	84
V.3.3. Hondar bariantza eta mugatze-koefizientea	85
V.3.4. Erregresio-koefizientearen eta koerlazio-koefizientearen zeinuaren azterketa. Doikuntzaren egokitasuna	86
V.4. ERREGRESIO LINEALA R^n-n	89
V.4.1. Sarrera	89
V.4.2. Erregresio hiperplanoa	90

V.4.3. Propietateak	91
V.4.4. Erregresio hiperplanoaren koefizienteak lortzeko metodoa	92
V.4.5. β erregresio partzial estandarizatuen koefizienteak	94
V.4.6. Hondar bariantza. Mugatze-koefizientea eta koerlazio-koefiziente anizkoitza	95
V.4.7. Edozein aldagai azalduaren erregresioaren orokorpena	97
V.4.8. Erregresio-koefiziente desberdinen zeinuaren azterketa. Doikuntzaren egokitasuna	97
V.5. KOERLAZIO PARTZIALA	98
V.5.1. Koerlazio partziala R^3 -n	99
V.5.2. Koerlazio partziala R^n -n	101
V.5.3. Koerlazio partzial eta erregresio-koefizienteen arteko erlazioa	103

V. A. Eranskina: ALDAGAI ESTADISTIKO "n"-KOITZEN MATRIZE-ESTADISTIKOAK

V.A.1. DATU-MATRIZEAK	107
V.A.2. BATEZBESTEKO-BEKTOREA: PROPIETATEAK.....	108
V.A.3. KOBARIANTZA MATRIZEA	109
V.A.3.1. Kobariantza matrizearen propietateak	110
V.A.4. KOERLAZIO-MATRIZEA	112
V.A.4.1. Koerlazio-matrizearen propietateak	114
V.A.5. DERIBAZIO BEKTORIALA	114

VI. ZENBAKI-INDIZEAK

VI.1. INDIZE SINPLEAK	117
VI.2. PONDERAZIO GABEKO INDIZE KONPLEXUAK	118
VI.2.1. Batezbesteko aritmetiko sinplearen metodoa	118
VI.2.2. Batezbesteko agregatu sinplearen metodoa	118

VI.3. INDIZE KONPLEXU PONDERATUAK	124
VI.3.1. Balioen, prezioen eta kopuruen indizeak	124
VI.3.1.1. LASPEYRES-en indizeak	127
VI.3.1.2. PAASCHE-ren indizeak	127
VI.3.1.3. FISHER-en indizeak	129
VI.3.1.4. Propietate eta erlazio batzuk: batera- garritasuna eta alderantzizkotasuna.....	131
VI.3.1.5. Kalkulua	132
VI.4. INDIZE KONPLEXUEN ERAIKETAN SORTZEN	
DEN ZENBAIT ERAGOZPEN	134
VI.4.1. Aldagaien hautapena	134
VI.4.2. Somatutako leku eta denboraren hautapena	134
VI.4.3. Talde eta azpitaldeen hautapena	134
VI.4.4. Oinarri-denboraren hautapena	134
VI.4.5. Formula eta ponderazioen hautapena	135
VI.4.6. Indizearen esangura eta zabaldura	135
VI.5. ERAGOZPEN BEREZI BATZUK	136
VI.5.1. Oinarri-aldaketa indize sinpleetan	136
VI.5.2. Berriztapen eta loturak indize konplexuetan	136
VI.6. ZENBAKI-INDIZEEN APLIKAZIOAK	138
VI.6.1. Kontsumo-prezioen indizeak	138
VI.6.2. Moneta-unitate arruntetan dauden magnitudeen deflazioa	139

VII. DESKRIBAPEN ESTADISTIKOAREN ADIERAZPIDE GEOMETRIKOAK

VII.1. OROKORTASUNAK	143
VII.2. X DATU-MATRIZEAREN BI ADIERAZPIDE GEOMETRIKO	143
VII.3. BI ALDAGAI ETA BI OHARPEN	144
VII.4. BI ALDAGAI ETA HIRU OHARPEN.....	151
VII.5. ERREGRESIO BAKUNA ALDAGAI ESTADISTIKOEN BALIO-MULTZOEN ESPAZIOAN	157

**VIII. DOIKUNTZA ORTOGONALA ETA
KOBARIANTZA MATRIZEAREN AUTODIREKZIOAK**

VIII.1. SARRERA	163
VIII.2. ANALISIA \mathbb{R}^2 ESPAZIOAN. DOIKUNTZA ORTOGONALAREN ZUZENA	163
VIII.3. DOIKUNTZA ORTOGONALAREN ZUZENAREN LORPENA	167
VIII.4. HODEI DUALAREN AURKEZPEN GRAFIKOA	170
VIII.5. TRANTSIZIO ERLAZIOAK	172
VIII.6. BATERAKO AURKEZPEN GRAFIKOA ETA INTERPRETAZIOA	173
VIII.7. DOIKUNTZA ORTOGONALAREN ZUZENA ALDAGAI NORMATUENTZAKO	174

IX. DATU ANIZKOITZEN ANALISIA

IX.1. SARRERA	181
IX.2. ANALISI OROKORRA	182
IX.2.1. \mathbb{R}^n espazioaren \mathbb{R}^q azpiespazio baten bidez egindako doikuntza	183
IX.2.2. \mathbb{R}^m espazioaren \mathbb{R}^q azpiespazio baten bidez egindako doikuntza	186
IX.2.3. \mathbb{R}^n eta \mathbb{R}^m espazioen arteko erlazioa	187
IX.2.4. X datu-taularen berreraketa	189
IX.3. OSAGAI NAGUSIZKO ANALISIA	191
XII.3.1. \mathbb{R}^n espazioan egindako analisisa	192
XII.3.2. \mathbb{R}^m espazioan egindako analisisa	193
XII.3.3. Baterako aurkezpen grafikoak	197
XII.3.4. Adibidea	199
IX.4. KORRESPONDENTZI ANALISI FAKTORIALA	207
IX.4.1. Hodeiak, masak eta distantziak	210
IX.4.2. \mathbb{R}^n espazioan egindako analisisa	213
IX.4.3. \mathbb{R}^m espazioan egindako analisisa	214
IX.4.4. \mathbb{R}^n eta \mathbb{R}^m espazioen arteko erlazioa	215

<i>IX.4.5. Maiztasun taularen berreraiketa</i>	217
<i>IX.4.6. Interpretaziorako laguntzak</i>	218
<i>IX.4.7. Korrespondentzia Anitzeko Anlisi Faktoriala</i>	220
<i>IX.5. ELEMENTU GEHIGARRIAK</i>	223
<i>IX.6. SAILKAPEN-METODOAK</i>	224
<i>IX.6.1. Goranzko Sailkapen Hierarkikoa</i>	224
<i>IX.6.1.1. Bariantzaren metodoa</i>	225
<i>IX.6.2. Adibidea</i>	227

HITZAURREA

Eskuarteon daukazun liburua, Ekonomia eta Enpresa-Zientzien Fakultateko "Estatistikarako Sarrera" asignaturaren testuliburu gisa, sortuz eta moldatuz joan garen liburu batzuren ondotik dator.

Duela 13 urte, 1983. urtean alegia, *Estatistikaren Hastapenak* izenaz lehen apunte-liburua argitaratu genuen, eta ondoren, *Estatistika Deskribatzailea. Koerlazioa, Erregresioa eta Datu Anizkoitzen Analisia* izenburuaz 1989. eta 1993. urteetan asignaturaren egitarauai egokituz eta gai batzuk berrituz bi argitalpen atera genituen.

Orain, bi urte t'erdi igaro ondoren *Estatistikarako Sarrera. Koerlazioa, Erregresioa eta Datu Anizkoitzen Analisia* kaleratzen dugu. Bi urte t'erdi epe laburra da, baina Sarrikoko ikastetxean aldaketa asko gertatu dira. Hauetatik garrantzitsuena plangintza berrien martxan jartzea izan da.

Liburu honetan, irakasgaiak funtsean aurrekoak badira ere lizentziatura berriei egokiturik daude. Ekonomia lizentziaturaren "Estatistikarako Sarrera" asignaturak 6 kreditu ditu, hau da, 60 ordu eta Administrazio eta Enpresa Zuzendaritza lizentziaturarenak 3 kreditu, hau da, 30 ordu.

Irakasgai guztiek, Ekonomia lizentziaturako asignaturaren egitaraua betetzen dute.

Administrazio eta Enpresa Zuzendaritza lizentziaturaren egitaraua lehen VI irakasgaiak beteko dute, gainera gai hauetan zenbait frogapen ez da emango. VII, VIII eta IX irakasgaiak "Estatistikarako Sarrera" asignatura gainditzeko, behar ez badira ere, lizentziatura honetan dagoen zenbait berezitasunetan, oso garrantzizkoak izan daitezke.

Bestalde, testuliburu hau Unibertsitateko beste ikastetxe batzuetan, euskaraz azaltzen diren Estatistika asigaturetan, erabilgarria izan daiteke; adibidez, Enpresa Unibertsitate Eskoletan, Gizarte-Graduatuen Eskolan nahiz beste fakultate eta eskoletan. Irakasle eta ikasle askori laguntzeko ilusioz lan egin dugu.

Azken 13 urte hauetan, asignaturaren egitarauetan egindako talde-lana izan da eta meritua ez da gurea bakarrik. Talde hauetan eduki ditugun lankide guztiak biziki eskertu nahi ditugu eta bereziki guztiok oroimenean daukagun J.M. Piris.

Bestalde, euskararen aldetik lagun askoren laguntza eskergaitza izan da, gehienak UEUren ingurukoak direlarik. Eskerrak guztiei eta edizio honetan hainbeste lagundu diguten UEUko Mari Karmen Menika, Nekane Intxaurtza eta Izaro Gurrutxaga.

Egileak
Sarriko, 1996.eko urtarrila

- I. EZAUGARRI ESTATISTIKO BAKUNAK,
BANAKETAK, TAULAK,
ADIERAZPIDE GRAFIKOAK**
- I.1. ESTATISTIKAREN HISTORIA LABURRA.
ESTATISTIKA ETA PROBABILITATEA*
- I.2. OROKORTASUNAK*
 - I.2.1. Populazioa eta lagina*
 - I.2.2. Unitate estatistiko edo indibiduoak*
- I.3. EZAUGARRI ESTATISTIKOAK*
 - I.3.1. Ezaugarri estatistiko bakunak*
 - I.3.1.1. Ezaugarri kuantitatiboak*
 - I.3.1.2. Ezaugarri kualitatiboak*
- I.4. MAIZTASUN-BANAKETA BAKUNAK*
 - I.4.1. Maiztasun-banaketak (absolutuak,
erlatiboak eta metatuak)*
 - I.4.2. Taulak*
- I.5. MAIZTASUN-BANAKETA BAKUNEN
ADIERAZPIDE GRAFIKOAK*
 - I.5.1. Barra-diagramak*
 - I.5.2. Histogramak*
 - I.5.3. Maiztasun-poligonoak*
 - I.5.4. Diagrama linealak*
 - I.5.5. Beste adierazpide grafikoak*

I.1. ESTADISTIKAREN HISTORIA LABURRA. ESTADISTIKA ETA PROBABILITATEA

Estatistika hitza gutxienez hiru zentzu desberdinetan erabil daiteke.

Esan ohi da:

- “estatistika batzuk” izen arrunt bezala erabiliz edo.
- “..... txosten estatistiko bat.....” adjetibotzat hartuz edo.
- “.... Estatistika ikastea....” izen propio bezala azkenik.

Hitz honen adierazpen orokor, zahar eta ez-espezializatua hauxe da:

MULTZO HANDIEN “Ikasketa” edo “erabilera”, MULTZO ZENTZUAN HARTURIK.

Honela, **hasieran** Estatistika artea eta teknika izan zen.

Mundu neolitikoa herriak eta salerosketak gero eta ugariagoak egin zirenez, Estatistikaren beharra sortu zen.

Gizartearen hasiera historikotik praktikatu da beraz. Israeldar, erromatar, txinatar eta inken aparteko inperioarenak, adibide eder batzuk besterik ez dira.

Baina “Cinquecento” italiarrean (XVI. mendean) izena eta bere buruaren kontzientzia hartzen hasi zela dirudi, kontabilitatearekin nahiko lotuta hasi ere.

XVIII. mendetik aurrera argiago eta garbiago agertzen da, beharbada. Estatu-zientziaren alboan hasi zen, gehienbat, Alemaniak emandako bultzada medio.

Eta honetan, oraindik Estatistika ez da zientzia, eta gertaera aleatorioetatik eta probabilitate-teoriatik urrun zegoen artea zen.

Baina geroztik “estatistikariei” arazoak agertzen hasi zitzaizkien:

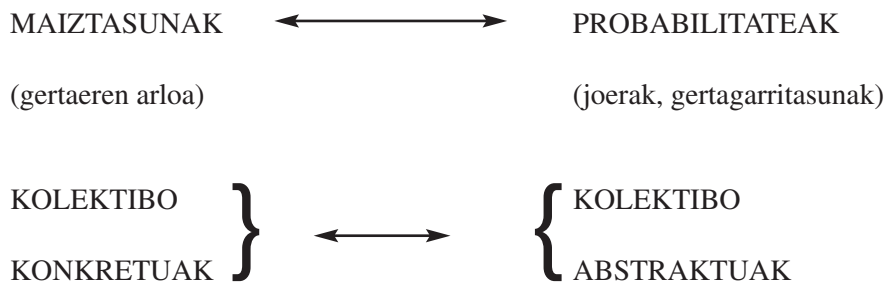
- nola lortu ondorio zientifikoak estatistiketatik?
- nola lortu informazio fidagarriak kolektibo aldakor edo zentsatzeko zail direnetatik?

Britaniar giroan XIX. mendearen amaieran eta XX.aren hasieran, oztopo hauetatik **irteten zen heinean**, Estatistika zientzia bilakatu zen.

Orain Estatistika probabilitate-teoriarekin esentziazko harremanetan dagoen zientzia da.

Honela, hurrengo galderetan ikusiko dugun maiztasun kontzeptutik probabilitate kontzeptura pasatu zen.

Eskematikoki ikusiz:



Fenomenoetatik bere kausetarainoko bidea egiten duen ezagupide honi, **INDUKZIOA** deritzo Estatistikaren alorrean eta, hementxe, **INFERENTZIA ESTATISTIKOA**.

Estatistika zientziaren oinarri trinkoan bermatutako estimazioak egiten dira gaur egun, beraien fidagarritasuna ere neurtzen delarik.

Hau dena **INFERENTZIA ESTATISTIKOARI** dagokio.

Apunte-liburu honetan ikusten diren gaiak Estatistika Deskribatzaileari dagozkio; hau da, oinarrizko Estatistika ikusiko dugu edo Estatistikaren hastapenak, eta halaber, Estatistika Deskribatzaile Anizkoitza edo Datu-Analisia IX. gaian.

Dakusagun, bada, gaian sartu aurretik, ESTADISTIKA DESKRIBATZAILEAREN definizio bat.

ESTADISTIKA DESKRIBATZAILEA, datu-multzo HANDIAK aurkez ditzakeen teknika bat da, multzo hau eskurakoi bihurtuz, beraren egitura gardenduz eta beraren barne-harremanak neurtuz.

1.2. OROKORTASUNAK

1.2.1. Populazioa eta lagina

Ikerketa estatistiko bakoitzean objektu nagusi diren pertsona, ondasun, saioaren emaitza (adibidez: datuaren aurpegia), probintzia baten herri, hauteskundearen alderdi politiko, edota fabrikazio-unitateek osatzen duten multzoari POPULAZIO edo UNIBERTSO deritzogu.

Populazioa ondo definitzeko edozein elementu partikular populaziokoa den ala ez jakin ahal izatea, beharrezkoa zaigu.

Populazioak **bukaezinak** (adb: boltsa bateko erauzketak) ala **bukakorrak** (eskualde bateko enpresa kooperatiboak) izan daitezke.

LAGINA: Lagina azpipopulazio errepresentagarri bat besterik ez da; populazio batez konklusio fidagarriak atera nahi baditugu, ahal dugun azpipopulazio errepresentagarriena aukeratu beharko dugu lagintzat.

LAGIN errepresentagarri hau erauzten duen prozesuari **zorizko laginketa** deritzo.

Lagina erauztean populazioaren elementu guztiei posibilitate berdina ematen badiegu, **zorizko lagin sinplea** deritzogu honela aukeratutako laginari.

Lagina erauzteko beste era eta jokabide asko daude, geruzatuz, mordoketaz... baina hau **LAGINKETA-TEORIA**ri dagokio.

Lagina osatzen duten elementuen kopuruari laginaren **tamainua** deritzo.

Sinbolikoki:

Ω = Populazioa

Ω' = Lagina $\Omega' = \{\omega_1, \omega_2, \dots, \omega_N\} \subset \Omega$

N = Laginaren tamainua.

1.2.2. Unitate estatistiko edo indibiduoak

Populazioa edo lagina osatzen duen elementu bakoitzari **unitate estatistiko** edo indibiduo deritzogu.

Indibiduo izena, estatistika deskribatzailearen jatorri demografikoari zor diogu.

1.3. EZAUGARRI ESTATISTIKOAK

Ikerketa estatistiko batean helburua ez dugu populazioa edo lagina exhaustiboki ikertzea, beraren ezaugarri batera edo batzuetara mugatuko gara baizik.

Ikertu nahi dugun ezaugarria bakarra bada, **ezaugarri estatistiko bakuna** deritzo.

Populazioaren bi edo ezaugarri gehiago batera ikertu nahi baditugu, **ezaugarri estatistiko anizkoitza** deritzo.

Ezaugarri estatistikoek (bakun ala anizkoitzek) populazio edo laginen multzoko indibiduo bakoitzean **balio indibidual** bat har dezakete.

Lagineko indibiduo-multzoari, ezaugarriaren **balio-multzo** bat dagokio.

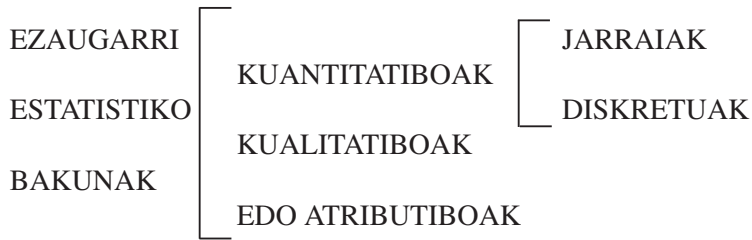
Sinbolikoki:

$\Omega = \{ \omega_1, \omega_2, \dots, \omega_N \} \longrightarrow \{ x_1, x_2, \dots, x_n \}$

$\Omega = \{ \omega_1, \omega_2, \dots, \omega_N \} \longrightarrow \{ (x_1 y_1), (x_2 y_2) \dots (x_n y_n) \}$

1.3.1. Ezaugarri estatistiko bakunak

Ezaugarri estatistiko bakunen sailkapena bere balio multzo desberdinen arabera:



1.3.1.1. Ezaugarri kuantitatiboak

Ezaugarriaren balio indibidualak populazioan edo laginean **zenbakiak** baldin badira, edota balio-multzoa zenbakizkoa bada, ezaugarri kuantitatiboa deritzo.

Ezaugarriari dagokion balio-multzoa, **multzo diskretua** denean, **ezaugarri diskretua** deritzo.¹

Balio biren artean edozein balio har dezakenean ezaugarri jarraia deritzo.

1.3.1.2. Ezaugarri kualitatiboak

Ezaugarriak har ditzakeen balio indibidualak kategoria edo atributuak direnean, edota **balio-multzoa** kategoria edo atributuek osatzen dutenean (diskretuak beraz), ezaugarri kualitatiboa deritzo.

Ondo klasifikatu ahal izateko, ezaugarri kualitatiboak hiru baldintza hauek bete behar ditu:

- **Ondo definitua**, hau da, kategoria edo atributu bakoitzak zer ulertarazten duen argiro azaldu behar du.

- **Esklusibotasuna**, hau da, indibiduo bat ezin liteke bi kategoriatan egon.

- **Exhaustibotasuna**, hau da, indibiduo guztiek nonbait klasifikatuak egon behar dute.

ADIBIDEAK: Eskualde bateko enpresa kooperatiboen lagin bat aztertzean:

- Langile-kopurua, ezaugarri kuantitatibo diskretua da;

- Produktzioa, ezaugarri kuantitatibo jarraia da;

1. Balio-multzoko elementuak, puntu isolatuak direnean, balio-multzo diskretu bat dugu.

- Lekutasuna, ezaugarri kualitatiboa da.

I.4. MAIZTASUN-BANAKETA BAKUNAK

I.4.1. Maiztasun-banaketak (absolutuak, erlatiboak eta metatuak)

Ezaugarri diskretuetan:

Suposa dezagun N indibiduo osatutako lagin bat, non ezaugarri bakun batek “m” balio edo kategoria desberdin hartzen dituen.

Hots: \longrightarrow
 $\{ \omega_1, \omega_2, \dots, \omega_N \} \longrightarrow \{ x_1, x_2, \dots, x_j, \dots, x_m \}$
 $\{ \omega_1, \omega_2, \dots, \omega_N \} \qquad \{ A, B, \dots, J, \dots, M \}$

x_j balioaren MAIZTASUN ABSOLUTUA, sinbolikoki n_j ikurraz adieraziko dugu, eta laginean x_j balioa hartzen duten indibiduen kopurua da.

Non:

$$\sum_{j=1}^m n_j = N$$

x_j balioaren MAIZTASUN ERLATIBOA, sinbolikoki f_j ikurraz adieraziko dugu, non:

$$f_j = \frac{n_j}{N} \qquad \sum_{j=1}^m n_j = N$$

hau da, laginean x_j balioa hartzen duten indibiduen kopurua erlatiboki kontutan hartuz lortzen den proportzioa da.

Maiztasun erlatiboak portzentaietan ere adierazten dira, batzuetan nahiko komenigarria baita horrela egitea.

Orduan:

$$P_j = \%f_j \cdot 100 \qquad \sum_{j=1}^m p_j = 100$$

Laginari dagokion **balio-multzoa ordenatu** ahal badugu (ezaugarria kuantitatiboa bada, beti; bestela ez beti), txikitik handira ordenatuko dugu.

Hots:

$$x_1 \leq x_2 \leq \dots \leq x_j \dots \leq x_m$$

Kasu honetan, x_j balioaren MAIZTASUN ABSOLUTU METATUA, sinbolikoki N_j , honela definituko dugu:

$$N_j = n_1 + n_2 + \dots + n_j \qquad N_m = \sum_{j=1}^m n_j = N$$

hau da, laginean x_j balioa eta txikiagoak hartzen dituzten indibiduen kopurua da.

Eta: x_j **balioaren MAIZTASUN ERLATIBO METATUA**, sinbolikoki F_j , honela definituko dugu:

$$F_j = \frac{n_1 + n_2 + \dots + n_j}{N} = \frac{N_j}{N} = f_1 + \dots + f_j$$

$$F_m = \sum_{j=1}^m f_j = 1$$

hau da, x_j balioa eta txikiagoak hartzen dituzten indibiduen kopurua N indibiduo guztiekiko kontsideratutako proportzioa (batekotan).

Edozein maiztasun erlatibo bezala, portzentaetan ere adierazten dira.

$$\text{Kasu honetan, } F_m = \sum_{j=1}^m p_j = 100$$

I.4.2. Taulak

Datu edo balioen ordenaketa eta elkarketari TABULAZIOA deritzo, eta, azkeneko datu-disposizioari, TAULA ESTADISTIKOA.

Lagin bati dagokion **maiztasun-banaketen taula** (ezaugarri bakun eta diskretua izanik) ondokoa da.

Lagina	Balio- -multzoa	MAIZTASUN-BANAKETA			
		M. absolutuak	M. absolutu metatuak	M. erlatiboak	M. erlatibo metatuak
ω_1	x_1 A	n_1	N_1	f_1	F_1
ω_2	x_2 B	n_2	N_2	f_2	F_2
..	
ω_j	x_j J	n_j	N_j	f_j	F_j
..	
ω_N	x_m M	n_m	N	f_m	1

Ezaugarri jarraietan

Nahiz eta **teorikoki** ezaugarria etengabea izan eta fenomeno askok ezaugarri mota honi erantzun, **praktikan** zehaztasun gehiago edo gutxiago duten tresnez neurtzen direnez gero, diskretu bukaezinak gertatzen zaizkigu.

Hots:

$$\{\omega_1, \omega_2, \dots, \omega_N\} \xrightarrow{x} \{x_1, x_2, \dots, x_N\} \subset (a, b) \subset \mathbb{R}$$

baina praktikan:

$$\{\omega_1, \omega_2, \dots, \omega_N\} \xrightarrow{x} \{x_1, x_2, \dots, x_N, \dots\} \subset (a, b) \subset \mathbb{R}$$

Lehen urratsean: ezaugarriari dagokion **balio-multzoa mailatan zatituko dugu**, mailen muturrak ondo definituz.

$$I = (a_0, a_1) \cup (a_1, a_2) \cup \dots \cup (a_{n-1}, a_n)$$

non:

$$a \leq a_0 \leq a_1 \leq a_2 \leq \dots \leq a_{n-1} \leq a_n \leq b$$

Mailen kopurua, tamainu desberdinekoak hartu ala ez..... azterketa bakoitzean erabaki beharreko gauzak dira, baina beti ere mailaketak hiru baldintza hauek bete beharko ditu:

- Ondo definituak izatea
- Esklusibotasuna
- Exhaustibotasuna

Bigarren urratsean: maila bakoitza balio batera mugatuko dugu, **klase-ordezkari** deritzoguna.

Normalki ordezkaria, maila tartearen **erdiko** puntua da, baina goi ala beheko muturrean ere izan liteke.

Honela bada, balio-multzoko maila bakoitza, balio batera mugatzen dugu (ordezkariaren baliora) eta laburpen honetaz informazio estatistikoa galdu arren, hasieran ezaugarri jarraia zena, ezaugarri diskretu bukakor bihurtzen dugu, eta dagokion maiztasun-banaketa kalkulatu ondoren, ondoan dagoen taula bezalako batean adieraz daiteke.

Lagina	Mailak edo klaseak	Klase-ordezkariak	M. absolutuak	M. absolutu metatuak	M. erlatiboak	M. erlatibo metatuak
ω_1	(a_0, a_1)	x_1	n_1	N_1	f_1	F_1
ω_2	(a_1, a_2)	x_2	n_2	N_2	f_2	F_2
..
	(a_{j-1}, a_j)	x_j	n_j	N_j	f_j	F_j
..
ω_N	(a_{m-1}, a_m)	x_m	n_m	N	f_m	1
$\sum_{j=1}^m n_j = N$					$\sum_{j=1}^n f_j = 1$	

ADIBIDEAK

a) X ezaugarri kualitatiboa da.

Eskualde bateko populazioaren banaketa produkzio-sektoreen arabera:

Sektoreak	Portzentaiak	Portzentaia metatuak
Nekazaritza	70	70
Industrigintza	25	95
Zerbitzuak	5	100

(Kasu honetan, sektore-kategoria bakoitzari dagozkien maiztasun erlatiboak portzentaiaetan bakarrik jarri ditugu).

b) X ezaugarria kuantitatiboa da.

b-1) Diskretu-taldekaturia

Enpresen banaketa, langile kopuruaren arabera herri bateko enpresetatik lortutako lagin batean.

Langile kopurua	Enpresa kopurua	
	m. absolutuak	m. abs. metatuak
0-tik 9-ra	3.876	3.876
10-tik 49-ra	2.028	5.904
50-tik 99-ra	976	6.880
100-tik 499-ra	689	7.569
500-tik 999-ra	320	7.889
1.000-tik 4.999-ra	304	8.193
5.000-tik 9.999-ra	295	8.488
10.000-tik aurrera	12	8.500
	N = 8.500	

(Kasu honetan, laginean langile kopuru ezaugarriari zegozkion balio desberdinak asko zirenez gero, taldekatu egin ditugu mailaka edo klaseka).

b-2) Ezaugarri jarraia

Klaseak (altuera m.)	Klase- -ordezkariak	Portzentaiak	Portzentaia metatuak
1'5 b. gutxiago		1	1
1'5-tik 1'55-era	1'525	7	8
1'55 1'60	1'575	12	20
1'60 1'65	1'625	18	38
1'65 1'70	1'675	30	68
1'70 1'75	1'725	20	88
1'75 1'80	1'775	10	98
1'80 b. gehiago		2	100

Soldaduzkara urte batean doazen gazteen banaketa, beren altueraren arabera:

(Kasu honetan, mailaketan, maila guztiek tamainu berdina dute, aurrenekoa eta azkenekoa kenduz).

1.5. MAIZTASUN-BANAKETA BAKUNEN ADIERAZPIDE GRAFIKOAK

Nahiz eta maiztasun-banaketen taulak informazio guztia hartu (balioak mailetan laburtuak izan ez badira behintzat), grafikoki adierazteak asko laguntzen du ikusarazten eta ulertzen.

Adierazpide grafiko asko erabili izan arren, batzuk bakarrik ikusiko ditugu.

1.5.1. Barra-diagramak

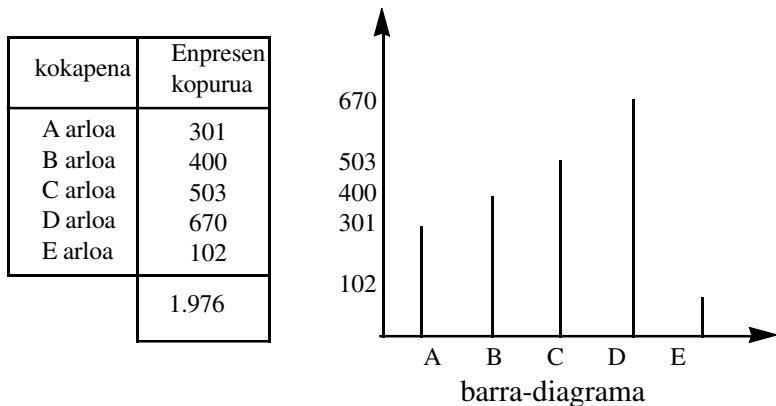
Oso baliagarriak zaizkigu ezaugarriak **diskretuak** ditugunean (kuantitatiboak zein kualitatiboak).

Abzisa-ardatzean balio edo kategoria desberdinak (ahal bada ordenaturik), ezarriko ditugu, eta ordenatu-ardatzean berriz laginari dagozkion maiztasunak (absolutuak, erlatiboak ala metatuak).

ADIBIDEAK

a) Ezaugarri kualitatibo diskretua.

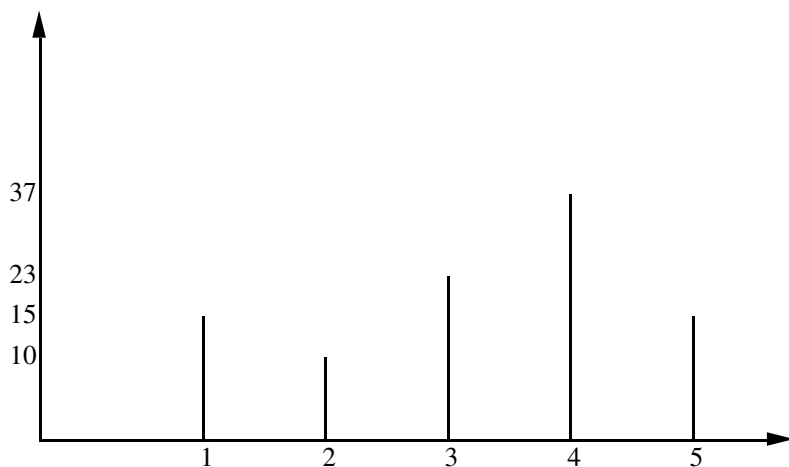
Enpresen banaketa bere kokapenaren arabera:



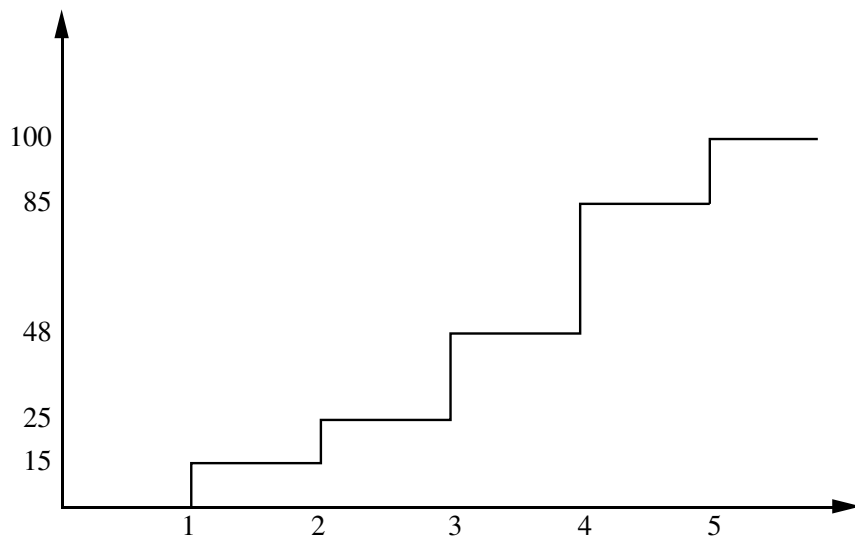
b) Ezaugarri kuantitatibo diskretua.

Loteen banaketa, beren baitan dituzten pieza akastunen kopuruaren arabera:

Txarren kopurua:	Maiztasun erlatiboa: (portzentaietan)	Maiztasun erlatiboa: Portzentaia metatuak
1	15	15
2	10	25
3	23	48
4	37	85
5	15	100



Maiztasun erlatiboen barra-diagrama

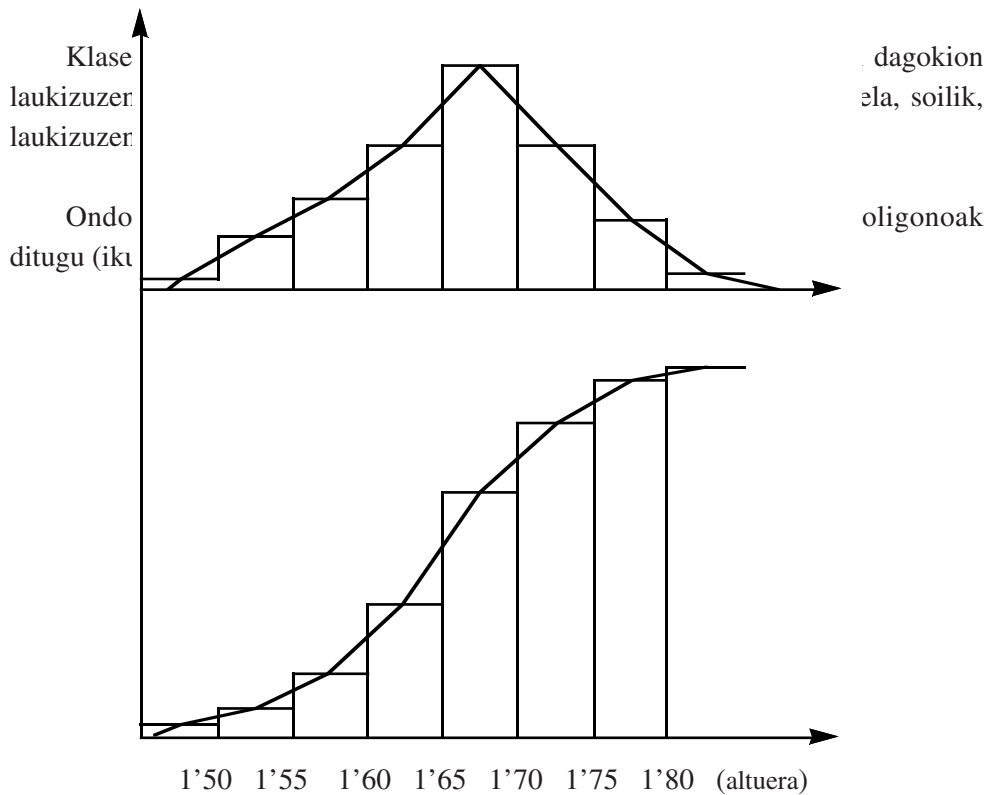


Maiztasun metatuen diagrama

1.5.2. Histogramak

Oso baliagarriak zaizkigu ezaugarriak jarraiak direnean.

Ezaugarriaren balio-multzoa mailetan zatitu ondoren eta maila edo klase bakoitza oinarritzat harturik, dagokion maiztasunaren (nahiz absolutua, erlatiboa edo metatua) **araberako azalera duen laukizuzen bat eraikitzen da.**



Histogramak eta maiztasun-poligonoak

Klaseen zabalerak, ordea, desberdinak baldin badira, laukizuzenen altuerak, maiztasun absolutuak edo erlatiboak zabaleraz zatituz, kalkulatu ditugu.

h_i , n_i eta c_i , i . klasearen altuera, maiztasuna eta zabalera, hurrenez hurren, izanik:

$$h_i = \frac{n_i}{c_i}, \text{ orduan } S_i \text{ azalera: } S_i = c_i \frac{n_i}{c_i} = n_i$$

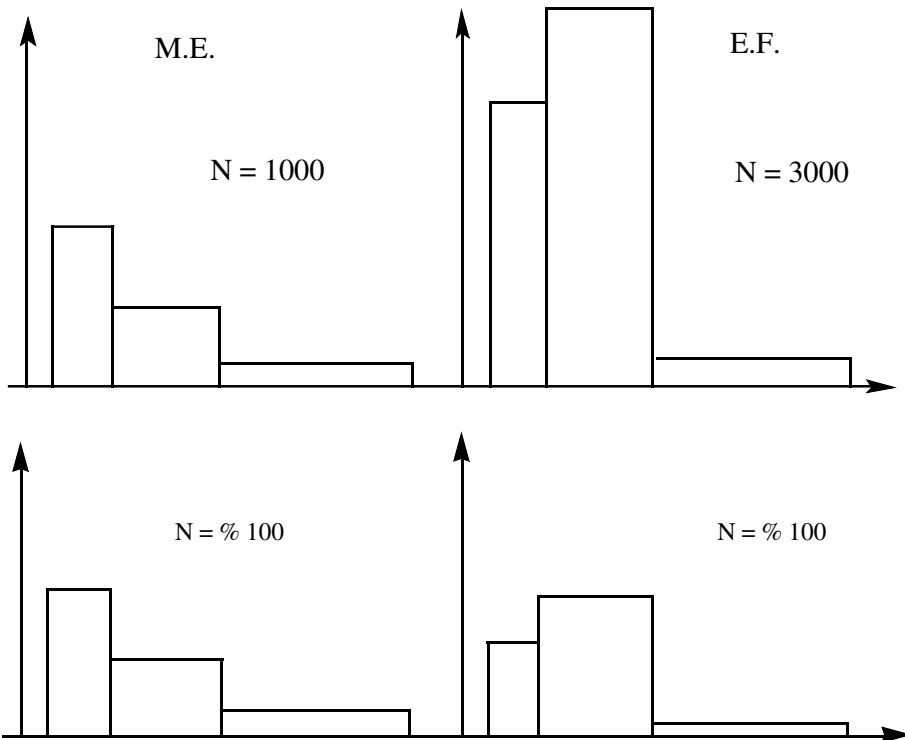
h_i altuera d_i bezala ere izendatzen da eta maiztasun-dentsitate bezala ezagutzen.

Adibidea:

Bi goi-mailako ikastetxetan, ikasleen adinak honela banatzen dira:

<u>Urteak (mailak)</u>	<u>M.E. Ikastetxea</u>	<u>E.F. Ikastetxea</u>
17 - 19	400	750
20 - 25	400	2.000
26 - 37	200	250

Adibidearen histogramak egin ditugu a) kasuan maiztasun absolutuak kontsideratuz eta b) kasuan, berriz, maiztasun erlatiboak portzentaietan.



Dakusagunez, maiztasun erlatiboak (ehunekotan edo ez) adieraztean totalaren ikusketa, N-rena hain zuzen, galtzen da, baina ordea bi banaketak konparagarriago bihurtzen dira.

1.5.3. Maiztasun-poligonoak

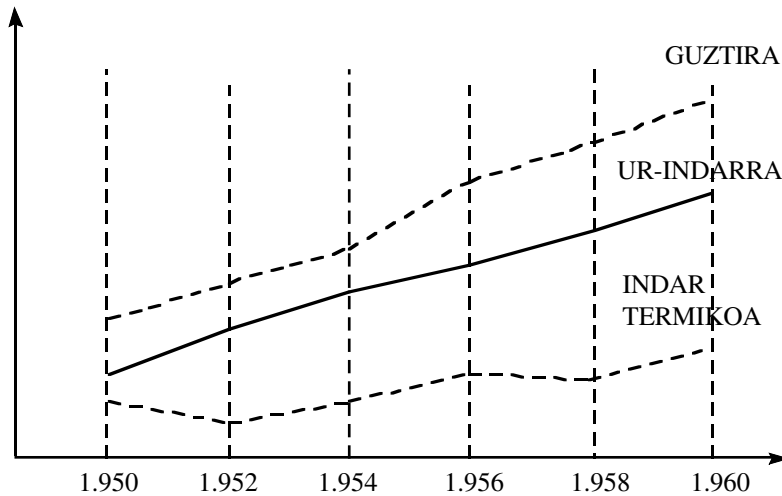
Maiztasun-poligonoa laukizuzenen goiko oinarriko erdiko puntuak marra batez lotuz eraikitzen da.

Honela egin dugu 16. orrialdeko histograman.

1.5.4. Diagrama linealak

Oso erabilgarriak izan ohi dira, balioen aldaketetan denborak izan lezakeen eragina aztertzeko.

Adibidez: Energi produkzioa 1.950etik, 1.960 urterarte:



1.5.5. Beste adierazpide grafikoak

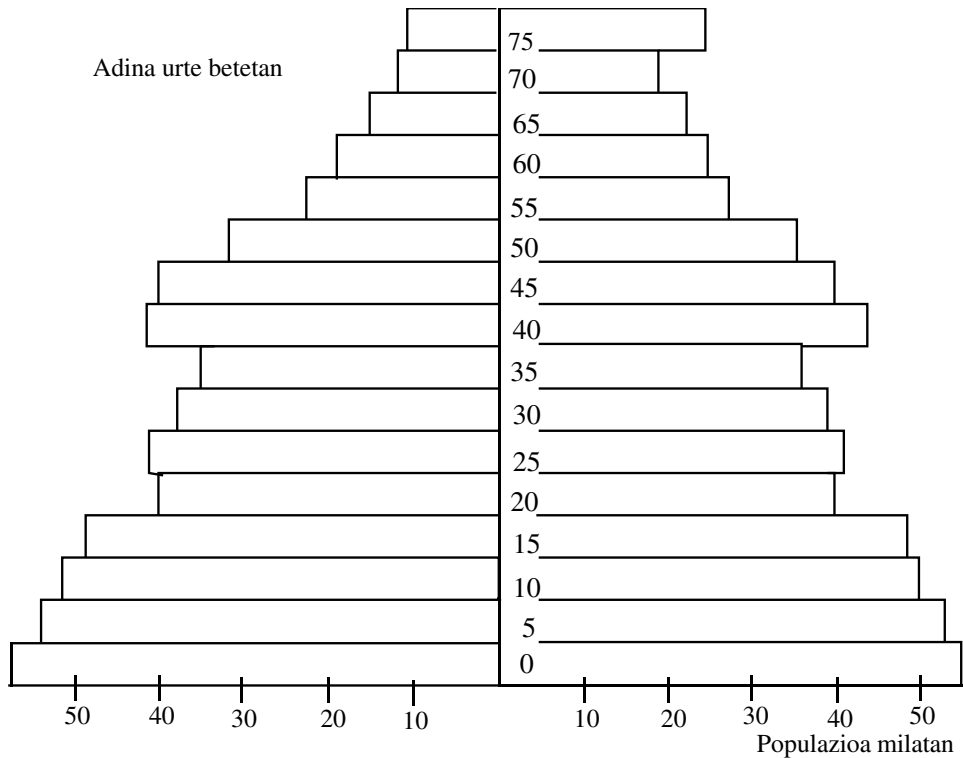
Beste batzuen artean, ikus ditzagun, azkenez, bi hauek:

- POPULAZIO-PIRAMIDEAK.

Histograma bikoitzak dira, eta honela eraikitzen dira: ordenatu-ardatzean adin-taldeak adierazten dira eta abzisa-ardatzean, ezkererantz eta eskuinerantz gizonezko eta emakumezkoen maiztasun erlatiboak populazio totalarekiko adierazten dira, dagozkien laukizuzenak eraikiz.

Maiztasun erlatiboak populazio totalarekiko kontsideratzean, gizonezkoen eta emakumezkoen laukizuzenak konpara daitezke adin-talde bakoitzerako.

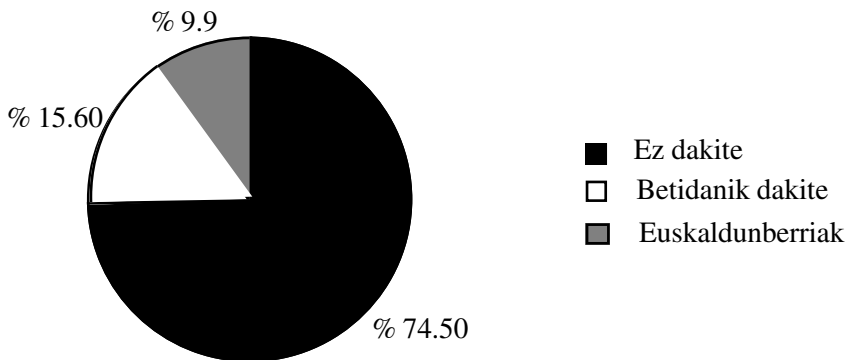
Halaber, maiztasun erlatiboak kontsideratzean, eskala berdinean eraikiz gero, tamainuaren aldetik oso desberdinak diren populazioak konpara daitezke.



Bizkaiko populazio-piramidea 1.975. urtean

- SEKTORE GRAFIKOAK

Grafiko hauek honela eraikitzen dira: zirkunferentziako 360 graduetatik, aldagaien balio ala kategoria bakoitzaren maiztasunari proportzionalki zati bat dagokio.



Euskararen ezaguera Sarrikoko Fakultatean 1979. urtean.

II. EZAUGARRI BAKUNEN BALIO TIPIKOAK

II.1. MOMENTUAK

II.1.1. Momentu arruntak edo jatorriarekikoak

*II.1.2. Momentu zentralak edo
batezbestekoarekikoak*

II.2. MAIZTASUN-BANAKETA BAKUNEN BALIO TIPIKO EDO ESTATISTIKOAK

II.2.1. Zentru-joeraren balio tipikoak

II.2.1.1. Batezbesteko aritmetikoa

II.2.1.2. Batezbesteko geometrikoa

II.2.1.3. Batezbesteko harmonikoa

II.2.1.4. Moda

II.2.1.5. Mediana

II.2.2. Sakabanatzearen balio tipikoak

II.2.2.1. Bariantza

II.2.2.2. Batezbesteko desbidazioa

II.2.2.3. Ibiltartea

II.2.3. Itxuraren balio tipikoak

II.2.3.1. Asimetri koefizientea

II.2.3.1. Kurtosi koefizientea edo

zapaltasun/zorroztasun-koefizientea

II.3. KONTZENTRAZIO-NEURKETAK

II.3.1. LORENZ-en kurba

II.3.2. GINI-ren indizea

II.1. MOMENTUAK

Datuek berekin duten informazioa, balio tipiko edo estatistikoen bidez sintetizatzea komeni da.

Noski, laburketa guztiei informazio-galtze bat dagokie; errakuntzarako motibo hau, bada, kontuan izan beharko dugu estatistiko bakoitza aztertzerakoan.

Baina, estatistiko edo balio tipikoak aurkeztu baino lehen, momentuak ikusi behar ditugu, balio tipikoak momentuak baitira edo momentuen bidez definitzen.

Momentuak bi motatakoak izan daitezke:

- arruntak edo jatorriarekikoak.
- zentralak edo batezbestekoarekikoak.

II.1.1. Momentu arruntak edo jatorriarekikoak

Maiztasun-banaketa bakun baten h-ordenako momentu arrunta edo jatorriarekikoa, sinbolikoki a_h ikurrak adieraziko da eta honela definituko non:

N = Laginaren edo populazioaren indibiduen kopurua.

x_i = Ezaugarriaren i . balioa.

m = Ezaugarriak hartzen dituen balio desberdinak; $m \leq N$ izango da.

$$a_h = \frac{\sum_{i=1}^N (x_i)^h}{N} = \frac{\sum_{i=1}^m (x_i)^h n_i}{N}$$

Momentu arrunten bi berezitasun ikusiko ditugu:

- 1) 0-ordenako momentuak 1 balioa hartzen du:

Hots:

$$a_0 = \frac{\sum_{i=1}^m (x_i)^0 n_i}{N} = 1$$

2) 1- ordenako edo lehenengo ordenako momentua batezbesteko aritmetikoa da:

Hots:

$$a_1 = \frac{\sum_{i=1}^m x_i n_i}{N} = \bar{x}$$

eta normalki \bar{x} ikurraz adierazten da.

Ikusten ari garen momentu arruntek nahiz ondoren ikusiko ditugunek, ez dute berekiko esangurarik, bide batez balio tipiko direnean izan ezik, adibidez, batezbesteko aritmetikoa, bariantza, e.a.

Halere, Estatistikan eragile matematiko bezala sartzen dira, teoriaren garapenean oso erabilgarriak baitira.

II.1.2. Momentu zentratuak edo batezbestekoarekikoak

Maiztasun-banaketa bakun baten h-ordenako momentu zentratua edo batezbestekoarekikoa, sinbolikoki m_h ikurraz adieraziko da eta honela definituko:

$$m_h = \frac{\sum_{i=1}^N (x_i - \bar{x})^h}{N} = \frac{\sum_{i=1}^m (x_i - \bar{x})^h n_i}{N}$$

Eragiketak eginez, ikus dezagun nola kalkulatzen diren momentu zentratuak momentu arrunten bidez:

$$\begin{aligned}
 m_h &= \frac{1}{N} \sum (x_i - \bar{x})^h n_i = \\
 &= \frac{1}{N} \sum \left(\binom{h}{0} x_i^h - \binom{h}{1} x_i^{h-1} \bar{x} + \binom{h}{2} x_i^{h-2} \bar{x}^2 - \dots + \right. \\
 &\quad \left. + (-1)^{h-1} \binom{h}{h-1} x_i \bar{x}^{h-1} + (-1)^{h-1} \binom{h}{h} \bar{x}^h \right) n_i = \\
 &= \binom{h}{0} a_h - \binom{h}{1} a_{h-1} \bar{x} + \binom{h}{2} a_{h-2} \bar{x}^2 - \dots + (-1)^{h-1} h \bar{x}^h + (-1)^h \bar{x}^h
 \end{aligned}$$

Hortik:

$$m_2 = a_2 - \bar{x}^2$$

$$m_3 = a_3 - 3a_2\bar{x} + 2\bar{x}^3$$

$$m_4 = a_4 - 4a_3\bar{x} + 6a_2\bar{x}^2 - 3\bar{x}^4$$

Momentu hauen artean bigarren ordenakoa (m_2), BARIANTZA delakoa da garrantzitsuena.

$$m_2 = \frac{\sum (x_i - \bar{x})^2 n_i}{N} = a_2 - \bar{x}^2 = \frac{\sum x_i^2 n_i}{N} - \left(\frac{\sum x_i n_i}{N} \right)^2$$

Eta azkenez, momentuak edozein punturekiko izan daitezke:

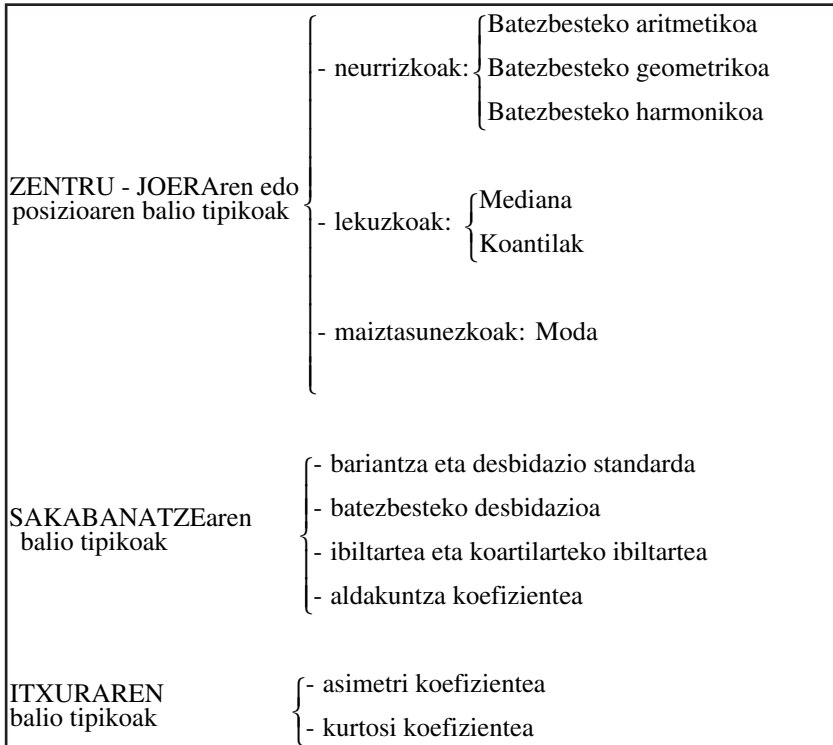
$$a_h, x_0 = \frac{\sum (x_i - x_0)^h n_i}{N} = \frac{\sum (x_i - x_0)^h n_i}{N}$$

$x_0 = 0$ baldin bada, momentu arruntak izango dira.

eta $x_0 = \bar{x}$ bada, zentratuak.

II.2. MAIZTASUN-BANAKETA BAKUNEN BALIO TIPIKO EDO ESTATISTIKOAK

Organigrama batean azalduko ditugu lehenik banaketa bakunetan erabili ahal diren estatistiko edo balio tipikoak, eta gero erabilienak banaka ikusten joango gara.



II.2.1. Zentru-joeraren edo Posizioaren balio tipikoak

Lagin batean ezaugarri bati dagozkion balioak ordenatzen baditugu, zentruan kokatzen diren balioak erabiltzen dira datu guztiak errepresentatzeko. Posizio-neurriak dira.

Balio batean datu-multzo baten ideia ematea argiro ikusten dugu, sarritan gehiegizko sinplifikazioa izanik.

Horregatik, zentruko estatistikoak beste estatistiko batzuekin (sakabanatzearenak..) osatzen dira, aztertzen ari garen populazio edo laginaren ideia erreala goa eduki nahi baldin badugu.

II.2.1.1. Batezbesteko aritmetikoa

Lagin batean $x_1, x_2, \dots, x_i, \dots, x_m$, zenbakizko balio-multzo bati dagozkion maiztasun absolutuak hurrenez hurren $n_1, n_2, \dots, n_i, \dots, n_m$ baldin badira, beraren batezbesteko aritmetikoa, \bar{x} , honako hau da:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^m x_i n_i \quad \text{non} \quad \sum_{i=1}^m n_i = N$$

Maiztasun erlatiboak, f_1, f_2, \dots, f_m , kontutan hartuz

$$\bar{x} = \sum_{i=1}^m x_i f_i$$

Ezaugarriak hartzen dituen balio guztiak desberdinak balira, maiztasun absolutuek 1 balioa hartuko lukete, hau da, $n_1 = n_2 = \dots = n_N = 1$ eta orduan:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

Maiztasun erlatiboak masak bezala kontsideratzen baditugu, batezbesteko aritmetikoak **grabitrate-zentrua** finkatuko du. Hau da, barra baten gainean eta x_1, x_2, \dots, x_m puntuetan maiztasunak bezalako pisuak kokatzen baditugu sistemak oreka lor dezan, esekitze-puntuak \bar{x} puntuan izan beharko du.

PROPIETATEAK:

1.- Edozein maiztasun-banaketatan, bere maiztasunaz biderkatuz, batezbestekoarekiko desbideratzeen batukaria zero egiten da.

Hots:

$$\sum_{i=1}^m (x_i - \bar{x}) n_i = 0$$

Frogapena:

$$\sum_{i=1}^m (x_i - \bar{x}) n_i = \sum_{i=1}^m x_i n_i - \bar{x} \sum_{i=1}^m n_i = N \bar{x} - N \bar{x} = 0$$

2.- Ezaugarri edo aldagai baten balio guztiak h parametroaz biderkatzen baditugu, batezbesteko aritmetikoa ere h-z biderkatuta edukiko dugu.

Hots:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^m x_i n_i \quad \text{izanez}$$

x_i balio bakoitza h x_i balioaz ordezkatzeko badugu:

$$\frac{1}{N} \sum_{i=1}^m h x_i n_i = h \frac{1}{N} \sum_{i=1}^m x_i n_i = h \bar{x}$$

3.- Aldagai bi edo gehiagoren batura den aldagai baten batezbesteko aritmetikoa, batezbesteko aritmetikoen batura da.

Hau da: $t = x + y + z$ baldin bada.

$\bar{t} = \bar{x} + \bar{y} + \bar{z}$ izango da.

Suposa dezagun oposizio-lehiaketa bat hiru azterketaz osatzen dela eta i. lehiakideak x_i, y_i, z_i puntuak ateratzen dituela, puntuazio totala $t_i = x_i + y_i + z_i$ izanez; N lehiakideren puntuazioen batezbesteko aritmetikoa ondokoa izango da:

$$\begin{aligned} \bar{t} &= \frac{\sum_{i=1}^N t_i}{N} = \frac{\sum_{i=1}^N (x_i + y_i + z_i)}{N} = \\ &= \frac{\sum_{i=1}^N x_i}{N} + \frac{\sum_{i=1}^N y_i}{N} + \frac{\sum_{i=1}^N z_i}{N} = \bar{x} + \bar{y} + \bar{z} \end{aligned}$$

BATEZBESTEKO ARITMETIKO PONDERATUA

Batzuetan batezbesteko aritmetikoa ateratzean, aldagaien garrantzi erlatiboa desberdina dela ohartzen gara.

Suposa dezagun hiru sagar mota saltzen direla 60, 80 eta 90 pezeta/kilo prezioetan hurrenez hurren.

Batezbesteko prezioa honela atera daiteke:

$$\bar{x} = \frac{60 + 80 + 90}{3} = 76'6, \text{ hau da, batezbesteko aritmetiko sinplea eginik.}$$

Baina mota bakoitzarentzat saltzen diren kiloak, hurrenez hurren, hauek baldin badira: 125, 72, 3

Batezbesteko ponderatua, zera izango da:

$$\bar{x}_{(\text{pon})} = \frac{60 \cdot 125 + 80 \cdot 72 + 90 \cdot 3}{125 + 72 + 3} = 67'6$$

Ikusten dugunez, batezbesteko aritmetiko ponderatua egitean, aldagaiaren balio bakoitza zenbaki batez biderkatzen dugu, balio-multzo barnean daukan garrantziaren arabera.

Zenbaki hauek pisuak edo ponderazioak dira eta sinbolikoki w_i deituko ditugu.

$$\text{Orduan } \bar{x}_{(\text{pon})} = \frac{\sum x_i w_i}{\sum w_i} \text{ edota } \bar{x}_{(\text{pon})} = \sum x_i w_i \quad \sum w_i = 1 \text{ denean}$$

Oharra: kontutan hartu behar da, bien formulak baliokideak badira ere ponderazioak ez direla maiztasunak. Hala ere, maiztasunak ponderazioak dira aldagaiaren balio desberdinek duten garrantzia maiztasunen arabera izango baita.

Beste adibide bat:

Ondoko taulan Hego Euskal Herriko lau herrialdeen kasuan hileroko gastuak/familia datuak dauzkagu.

Hileroko gastuak / familia, 1.968. urtean

	<u>Familiak</u>	<u>Gastuak (mila pezetatan)</u>
Araba	24	4,4
Gipuzkoa	53	5,3
Nafarroa	57	5,1
Bizkaia	84	5,1

$$\text{Batazbesteko aritmetikoa: } \bar{x} = \frac{19,9}{4} = 4,975$$

Baina, familien kopurua kontutan hartuz:

$$\bar{x}_{(\text{pon})} = \frac{4,4 \cdot 24 + 5,3 \cdot 53 + 5,1 \cdot 57 + 5,1 \cdot 84}{24 + 53 + 57 + 84} = 5,07$$

Noski, batezbesteko ponderatuen bidez errealitatearen azterketari gehiago hurbiltzen gatzaizkio, baina normalki datu guztiak ezagutzen ez ditugunez, batezbesteko aritmetiko sinplea egiten dugu.

II.2.1.2. Batezbesteko geometrikoa:

Ondoko berdintasunaren bidez adierazten dugu:

$$G = \sqrt[N]{x_1^{n_1} \cdot x_2^{n_2} \cdot \dots \cdot x_m^{n_m}} = \sqrt[N]{\prod_{i=1}^m x_i^{n_i}}$$

Arrazoi edo ehunekoen batezbesteko bat ateratzeko, batezbesteko aritmetikoa baino aproposagoa da eta ondoko propietatea betetzen du:

$$\text{Log } G = \frac{1}{N} \log \prod_{i=1}^m x_i^{n_i} = \frac{1}{N} \sum_{i=1}^m (\log x_i) n_i$$

Hots:

Batezbesteko geometrikoaren logaritmoa, aldagaiaren balioen logaritmoen batezbesteko aritmetikoa da.

II.2.1.3. Batezbesteko harmonikoa

Ondoko berdintasunaren bidez adierazten dugu

$$H = \frac{N}{\sum_{i=1}^m \frac{n_i}{x_i}} \quad \text{non} \quad N = \sum_{i=1}^m n_i$$

Ohar daiteke batezbesteko harmonikoa zera dela, x_i balioen alderantzizko balioen batezbesteko aritmetikoaren alderantzizko balioa.

Batezbesteko abiaduren batezbesteko bat ateratzeko, batezbesteko harmonikoa da egokiena.

Ikusitako batezbestekoak eta beste batzuk: batezbesteko koadratikoa, batezbesteko kubikoa, e.a. sar daitezke k ordenako batezbestekoaren kontzeptuan.

II.2.1.4. Moda

Maiztasun-banaketa batean gehien errepikatzen den aldagai edo ezaugarriaren balioari, hau da, maiztasunik handiena duenari, moda eta balio modala deritza.

$$\text{Hots:} \quad M_{(k)} = \sqrt[k]{\sum x_i^k f_i}$$

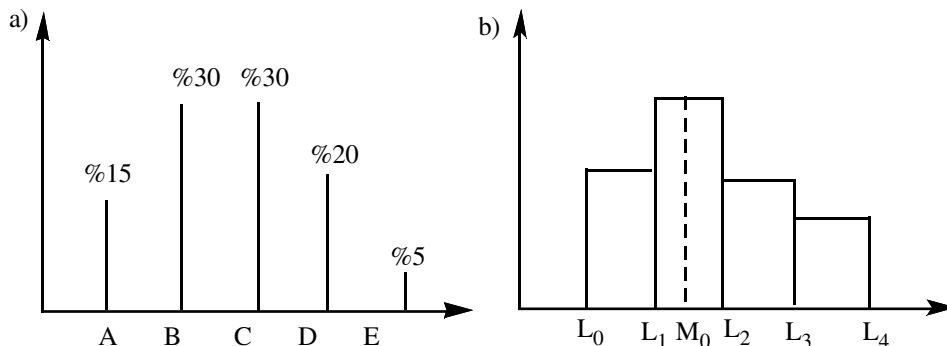
$$x_1, x_2, \dots, x_i, \dots, x_m \quad \text{edo} \quad A, B, \dots, I, \dots, M$$

balio edo kategorien multzo bati dagozkien maiztasun erlatiboak hurrenez hurren f_1, f_2, \dots, f_m baldin badira, maiztasun erlatibo gehien duen x_i balioari edo I kategoriari “moda” M_0 deritza.

Hots: $M_0 = x_i$ edo $M_0 = I$, n_i eta f_i maiztasun handienak izanik.

Banaketa, modagabea izan liteke $f_1 = f_2 = \dots = f_m$, edo moda bat baino gehiago ere eduki dezake, eta honen arabera, “bimodala”, “hirumodala”... orokorki moda-anitza izango da.

ADIBIDEAK



Kasu honetan banaketa bimodala dugu, B eta C arloak modak direlarik.

L₁, L₂ klasea da kasu honetan moda, klase honen ordezkaria, M₀, har daiteke modatzat.

Banaketa jarraia edo klasetan taldekatua bada, dakigunez, maiztasun handiena, M₀, duen klasean kokatuko da moda, eta klase-ordezkaria modatzat hartzen ez badugu, gutxi gora-beherako kalkuluak egin daitezke.

Klase edo tarte modala L_{i-1}, L_i izanez, balio modala estimatzeko bi irizpide desberdinetan oinarritutako formulak proposatuko ditugu.

1) Altueren alderaketaren irizpidea:

$$M_{(k)} = L_{i-1} + \frac{d_{i-1}}{d_{i-1} + d_{i+1}} \cdot c_i$$

non:

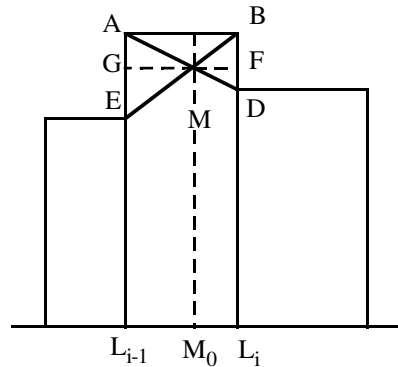
c_i = klaseen zabalera

$d_{i-1} = h_i - h_{i+1}$

$d_{i+1} = n_i - n_{i+1}$

h = i. klaseen altuera

Frogapena grafikoa



AME eta BMD triangeluak antzekoak dira, hau da, beraien altuerak oinarriekiko proportzionalak dira.

Hots:

$$\frac{MG}{MF} = \frac{AE}{BD}$$

Edo ondoko proportzio baliokidea:

$$\frac{MG}{MG + MF} = \frac{AE}{AE + BD} \quad \text{hau da,} \quad \frac{MG}{c_i} = \frac{d_{i-1}}{d_{i-1} + d_{i+1}}$$

eta, $M_0 = L_{i-1} + MG$ izatean, proposatu dugun formula frogatuta gelditzen da.

Klaseak zabalera berdinekoak balira, altueren alderaketa, zuzenki, maiztasunen alderaketa izango da.

Hau da:

$$d_{i-1} = n_i - n_{i-1}$$

$$d_{i+1} = n_i - n_{i+1}$$

2) Alderantzizko banaketa proportzionalaren irizpidea:

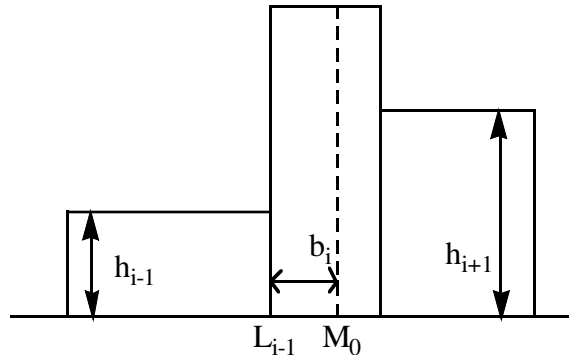
$$M_0 = L_{i-1} + \frac{h_{i+1}}{h_{i-1} + h_{i+1}} \cdot c_i$$

non: h_{i-1} = aurreko klasearen altuera
 h_{i+1} = hurrengo klasearen altuera

Ondoko klaseen altuerekiko alderantzizko banaketa proportzionala egitean, aurreko formula lortzen dugu.

Hots:

$$\frac{b_i}{h_{i+1}} = \frac{c_i - b_i}{h_{i-1}}$$



edo aurrekariak eta atzekariak batzean:

$$\frac{b_i}{h_{i+1}} = \frac{c_i}{h_{i+1} + h_{i-1}}$$

eta, $M_0 = L_{i-1} + b_i$ izatean, proposatu dugun formula frogatuta gelditzen da.

Klaseak zabalera berdinekoak balira, hau da, $c_i = c$ c zabalera biderkatzen eta zatitzen badugu eta $n_i = c \cdot h_i$ dela kontutan hartuz, ondoko hau izango da formula:

$$M_0 = L_{i-1} + \frac{h_{i+1} \cdot c}{(h_{i+1} + h_{i-1})c} = L_{i-1} + \frac{n_{i+1}}{h_{i+1} + h_{i-1}} c$$

II.2.1.5. Mediana eta koantilak

Mediana banaketaren balioak bi zati berdinetan zatitzen dituen lekuzko estatistikorik garrantzitsuena da.

Balio-multzoa ordenatuz gero, maiztasun absolutuen histograma azalera berdineko bi zatitan ebakitzen duen balioari mediana deritzo, eta hori zehazki betetzen duenik ez badago, aurrean dagoena hartzen dugu batzuetan medianatzat.

Ondoko baldintza betetzen duen balio bat da:

$$M_e = x_i / N_{i-1} < \frac{N}{2} \wedge N_i > \frac{N}{2}$$

Hala ere, banaketa diskretua bada eta indibiduen kopurua bakoitia, aurreko baldintza bete egingo da, baina ez beste edozein kasutan.

Banaketa diskretua eta N bikoiti baten aurrean bagaude.

$$M_e = \frac{x_i + x_{i+1}}{2} / N_i = \frac{N}{2}$$

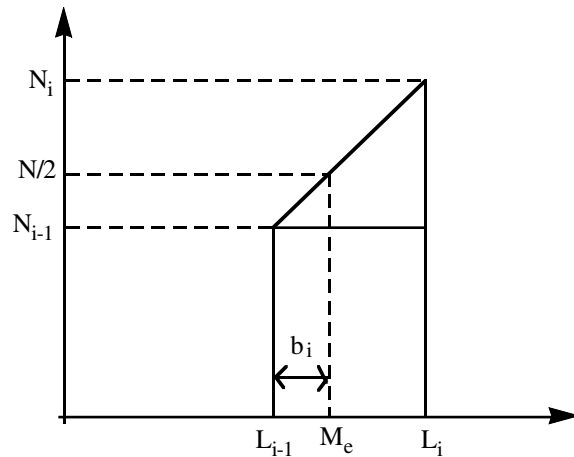
Hots: $M_e = \frac{\text{erdiko balioak}}{2}$

Banaketa taldekatuetan eta jarraietan: N/2 balioa duen tartea finkatuz gero eta tartean n_i balioak linealki gehitzen direla suposatuz kalkulatu dugu M_e.

Orduan: (L_{i-1}, L_i) M_e duen tartea izanez:

$$M_e = L_{i-1} + b_i$$

Irudian dauden bi triangeluak antzekoak izanez, zera idatz daiteke:



$$\frac{N_i - N_{i-1}}{c_i} = \frac{\frac{N}{2} - N_{i-1}}{b_i}$$

$$\text{eta: } b_i = \frac{N/2 - N_{i-1}}{N_i - N_{i-1}} c_i$$

Oharra: maiztasun erlatiboak erabiliz, era berdinean kalkula daiteke M_e ; kasu horretan,

$$\frac{N}{2} = \frac{1}{2}$$

KOANTILAK

Koantilak beste lekuzko neurriak dira. Mota desberdinetakoak dira: koartilak, dezilak eta zentilak.

Koartilak, dezilak eta zentilak medianaren antzera definitzen dira, baina bi parte hartu ordezkari (koartilak), hamar (dezilak) edo ehun (zentilak) parte hartuz.

Koartilak q_1, q_2, q_3 izango dira; (q_1, q_3) tartean banaketaren erdia daukagu, justu erdiko balioak, hau da, bazterreko balioak kenduz.

Dezilak $d_1 \dots d_9$ izango dira eta zentilak, $c_1 \dots c_{99}$.

Argiro ikusten da:

$$q_2 = d_5 = c_{50} = M_e \text{ dela.}$$

OHARRAK:

Erdiko baliorik garrantzitsuenak, batezbesteko aritmetikoa, moda eta mediana dira eta bereziki batezbesteko aritmetikoa, guztiak aldagaia edo ezaugarria neurtua dagoen unitateetan neurtuak direlarik.

Ondoko erlazio hau betetzen da:

$$M_0 \simeq 3 M_e - 2 \bar{x}$$

Eta gutxi gora-behera hiruretako edozein, besteen bidez kalkula daiteke.

II.2.2. Sakabanatzearen balio tipikoak

Datu-multzoaren ideia orokorra ematen ziguten zentzurako joeraren balio

tipikoek. Ideia hau osatzeko, zentzurako joeraren balioen inguruan datu-multzoa guztiz bilduta edo oso sakabanatuta dagoen adierazten diguten estatistiko batzuk ere behar ditugu.

Batezbesteko aritmetikoaren inguruan datuek duten sakabanatze-neurri bezala, estatistiko garrantzitsuenak hauek dira: Bariantza, Desbidazio Standarda eta Batezbesteko Desbidazioa.

II.2.2.1. Bariantza

Batezbestekoarekiko sakabanatze-neurririk garrantzitsuen da bariantza eta bera erabiliz lortzen ditugu desbidazio standarda eta aldakuntza koefizientea.

$\{x_1, x_2, \dots, x_i, \dots, x_m\}$ zenbakizko balio-multzo bati dagozkien maiztasun absolutuak hurrenez hurren $n_1, n_2, \dots, n_i, \dots, n_m$ baldin badira, bariantza, S_x^2 , deritzona hau da:

$$S_x^2 = \frac{1}{N} \sum_{i=1}^m (x_i - \bar{x})^2 n_i \quad \text{edota} \quad S_x^2 = \sum_{i=1}^m (x_i - \bar{x})^2 f_i$$

Bariantza neurri-unitate karratuetan daukagula erraz ikusten dugu formulatan.

Orduan, sakabanatze-neurri bezala bariantzaren erro karratua hartzea egokiagoa da, **Desbidazio Tipikoa edo Standarda**, hain zuzen.

$$S_x = \sqrt{\frac{1}{N} \sum_{i=1}^m (x_i - \bar{x})^2 n_i} = \sqrt{\sum_{i=1}^m (x_i - \bar{x})^2 f_i} = \sqrt{a_2 - \bar{x}^2}$$

eta ezaugarria neurtuta dagoen unitateetan, sakabanatze-neurketa daukagu.

Askotan, estatistika inferentzian estimatzaile bezala dituzten propietateengatik, S_x^{*2} Quasi-bariantza eta S_x^* Quasi-desbidazio standarda erabiltzen dira.

$$S_x^{*2} = \frac{\sum_{i=1}^m (x_i - \bar{x})^2 n_i}{N - 1} \quad \text{eta} \quad S_x^* = \sqrt{\frac{\sum_{i=1}^m (x_i - \bar{x})^2}{N - 1}}$$

Konputagailuko programa-paketeek S_x^{*2} eta S_x^* ematen dizkigute, Bariantza eta Desbidazio Standarda deitzen badituzte ere.

Desbidazio tipikoak batezbestekoarekiko suposatzen duen batekoa (edo ehunekoa) **aldakuntza koefizientea** deitzen da.

$$g_0 \text{ ikurraz adieraziko dugu sinbolikoki, non: } g_0 = \frac{S_x}{\bar{x}}$$

\bar{x} batezbestekoa zerora hurbiltzen bada, g_0 horrek ez du zentzu askorik (infiniturantz jotzen baitu kasu honetan).

Bariantzak garrantzizko propietate bat betetzen du, bigarren ordenako momentu zentralik txikiena baita, hain zuzen.

Hots:

$$\begin{aligned} a_{2, x_0} &= \frac{1}{N} \sum_{i=1}^m (x_i - \bar{x}_0)^2 n_i = \\ &= \frac{1}{N} \sum_{i=1}^m ((x_i - \bar{x}) + (\bar{x} - x_0))^2 n_i = \\ &= \frac{1}{N} \sum_{i=1}^m ((x_i - \bar{x})^2 + 2(x_i - \bar{x})(\bar{x} - x_0) + (\bar{x} - x_0)^2) n_i = \\ &= \frac{1}{N} \sum_{i=1}^m (x_i - \bar{x})^2 n_i + 2(\bar{x} - x_0) \sum_{i=1}^m (x_i - \bar{x}) \frac{n_i}{N} + \\ &+ \frac{1}{N} (\bar{x} - x_0)^2 \sum_{i=1}^m n_i = S_x^2 + (\bar{x} - x_0)^2 \end{aligned}$$

Eta argiro dagoenez, $x_0 = \bar{x}$ baldin bada, txikiena izango da, beste edozein kasutan bigarren batugaia positiboa baita.

II.2.2.2. Batezbesteko desbidazioa

Balio absolutuak harturik (konpentsaziorik izan ez dadin), batezbesteko batekiko balioek duten desbidazioen batezbesteko aritmetikoari, “batezbesteko desbidazioa” deritzo.

$$\text{Hau da: } D = \frac{1}{N} \sum_{i=1}^m |x_i - \text{batezbestekoa}| n_i$$

a) Aukeratutako batezbestekoa, batezbesteko aritmetikoa bada, batezbesteko desbidazioa hauxe izango da:

$$D_{\bar{x}} = \frac{1}{N} \sum_{i=1}^m |x_i - \bar{x}| n_i$$

Neurri honek indibiduen homogenotasuna adierazten digu, hau da, desbidazioak txikiak badira, batezbesteko desbidazioa txikia izango da eta, alderantziz, handiak badira, batezbesteko desbidazioa handia izango da.

b) Aukeratutako batezbestekoa mediana bada, batezbesteko desbidazioa hauxe izango da:

$$D_{M_e} = \frac{1}{N} \sum_{i=1}^m |x_i - M_e| n_i$$

Batezbestekoa mediana denean, batezbesteko desbidaziorik txikiena daukagu.

II.2.2.3. Ibiltartea

Ezaugarri zenbakizko batek hartzen dituen balio maximo eta minimoaren arteko diferentziari ibiltarte deritzogu.

Hots:

$$x_0 \leq x_1 \leq \dots \leq x_m \text{ izanik}$$

$$R = x_m - x_0$$

Koartil arteko ibiltarteak ere kontsidera daitezke, adibidez, $q_3 - q_1$ hirugarren eta lehenengo artekoena hain zuzen ere, edota, dezil arteko ibiltarteak e.a.

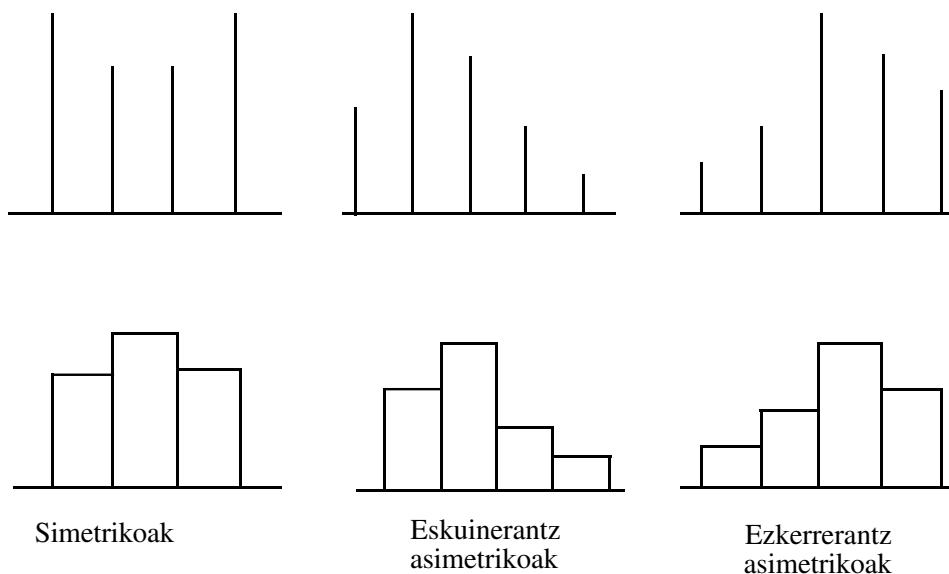
II.2.3. Itxuraren balio tipikoak

Aztertzen ari garen laginaren ezaugarria eta berari dagokion datu-multzoa, zentrurako joeraren eta sakabanatze-balioen bidez nahiko argi geratu bada ere, itxuraren balio tipikoak aztertzean gehiago osa daiteke; hauen artean asimetri koefizientea eta kurtosi koefizientea ikusiko ditugu.

II.2.3.1. Asimetri koefizientea

Maiztasun-banaketaren grafikoa simetrikoa denean, banaketa simetrikoa dela esango dugu, hau da, barra-diagrama edo histograma egitean simetrikoa bada, banaketaren elementuak batezbestekoarekiko bi aldeetan berdina kokatzen direla esan nahi du.

Aurrekoa betetzen ez bada, banaketa asimetrikoa izango da.



Banaketa bat moda batekoa eta simetrikoa baldin bada, batezbesteko aritmetikoa, moda eta mediana balio berean daude.

$$\text{Hots: } \bar{x} = M_0 = M_e$$

Asimetri koefizientea g_1 ikurraz adieraziko dugu sinbolikoki, non:

$$g_1 = \frac{m_3}{S_x^3}$$

$$m_3 = \sum_{i=1}^m (x_i - \bar{x})^3 \frac{n_i}{N} \text{ izanik, maiztasunen balioak } \bar{x} \text{ balioarekiko}$$

simetrikoki kokatuta badaude, $m_3 = 0$ izango da; baina maiztasunen balioak eskuinaldean **handiagoak** badira, hau da, $x_i > \bar{x}$ denean, orduan $m_3 > 0$ izango da, eta, alderantziz, ezker aldean handiagoak badira, hau da, $x_i < \bar{x}$ denean, orduan $m_3 < 0$.

Ikusten dugu, bada, m_3 asimetriaren balio adierazle bezala har dezakegula, baina neurri-unitateak ber hiru izanez zuzenkiago.

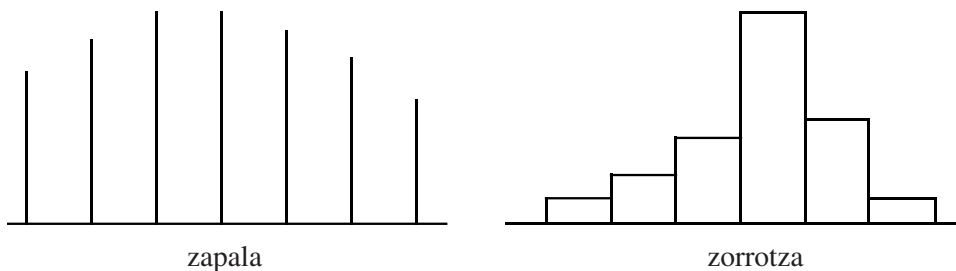
$$g_1 = \frac{m_3}{S^3}$$

erabiltzen da: $g_1 > 0$ bada, maiztasun-banaketa eskuinerantz asimetrikoa dela esaten da; $g_1 < 0$ bada, alderantziz, ezkererantz.

II.2.3.2. Kurtosi koefizientea edo zapaltasun/zorroztasun-koefizientea

Simetriarekin batera, itxuraren beste balio tipiko bat ikusiko dugu, kurtosi koefizientea hain zuzen.

Maiztasun-banaketaren grafikoa kontutan hartuz, konkretuki, maiztasun-poligonoa ezaugarri jarraia kasuan, edo barra-diagramaren goiko puntuak lotuko lituzkeen lerroa ezaugarri diskretuaren kasuan, kasu hauek dauzkagu:



Batezbesteko aritmetikoa eta bariantza berdina daukaten bi banaketa aztertuz, itxuraren aldetik bata oso zapala da eta oso zorrotza bestea.

Banaketaren maiztasunak batezbesteko aritmetikoaren inguruan handiak baldin badira, momentu zentratu bikoitiek, (adibidez, laugarren ordenakoa), balio handiagoak hartuko dituzte.

Orduan, m_4 zapaltasunaren adierazle bezala har daiteke, baina neurri-unitateak ber lau izanez, zuzenkiago g_2 ikurras adierazten den kurtosi koefizientea, zapaltasunaren adierazle bezala erabiliko dugu.

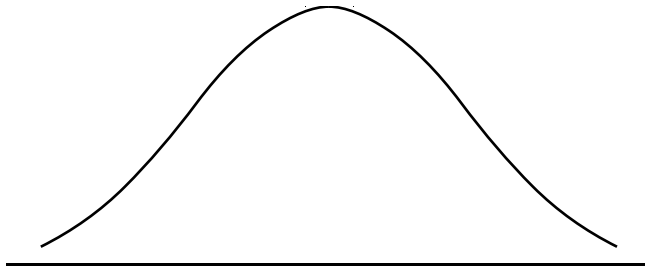
“Normal” deitutako banaketaren kurtosi koefizientea 3 da eta balio hau kurtosia aztertzean konparazio-puntu bezala hartuko dugu.

$$g_2 = \frac{m_4}{S_x^4} = \frac{\sum_{i=1}^m (x_i - \bar{x})^4 n_i}{N S_x^4}$$

g_2 , 3 baino handiagoa baldin bada, banaketa zorrotza izango da eta 3 baino txikiagoa baldin bada, zapala izango da.

g_2-3 , kurtosi koefizientetzat hartzen badugu, kurtosi positiboa edo negatiboa kontsidera daiteke.

Oharra: Aipatu dugun banaketa “Normala” 3. kurtsoan ikasiko da, inferentzia estatistikoan, oinarrizko probabilitate-eredua baita.



banaketa Normala : $g_2 = 3$

II.3. ALDAGAIEN TRANSFORMAZIO LINEALAK

X aldagaia U aldagaiaren transformatu lineala izan daiteke, hau da:

$$X = a U + b$$

$b = 0$ baldin bada, unitate-aldaketa daukagu soilik
 $a = 1$ baldin bada, jatorri-aldaketa edo traslazioa daukagu soilik, zeina berdin b baita.

Batzutan, jatorria eta unitatea aldatzea komenigarria da, honela, kalkuluak errazten baitira.

Ikus ditzagun, bi aldagaien estatistiko garrantzitsuenen arteko erlazioak:

$$\begin{aligned} \bar{x} &= \frac{1}{N} \sum_i x_i n_i = \frac{1}{N} \sum_i (a u_i + b) n_i = \\ &= a \cdot \frac{1}{N} \sum_i u_i n_i + b \frac{1}{n} \sum_i n_i = a \bar{u} + b \end{aligned}$$

Ikusten dugunez X eta U aldagaien batezbesteko aritmetikoak ez dira berdinak; pentsa dezagun estatistiko hau, ezaugarriak neurtzen ditugun unitateetan neurri bat dela.

Bariantza eta desbidazio tipikoen arteko erlazioak:

$$\begin{aligned} S_x^2 &= \frac{1}{N} \sum_i (x_i - \bar{x})^2 n_i = \frac{1}{N} \sum_i (a u_i + b - a \bar{u} - b)^2 n_i = \\ &= a^2 \cdot \frac{1}{N} \sum_i (u_i - \bar{u})^2 n_i = a^2 \cdot S_u^2 \end{aligned}$$

eta $S_x = |a| S_u$

Eta orokorki, momentu zentratuen arteko erlazioak:

$$\begin{aligned} m_h(x) &= \frac{1}{N} \sum_i (x_i - \bar{x})^h n_i = \frac{1}{N} \sum_i (a u_i + b - a \bar{u} - b)^h n_i = \\ &= a^h \frac{1}{N} \sum_i (u_i - \bar{u})^h n_i = a^h m_h(u) \end{aligned}$$

Ikusten dugunez, unitate-aldaketak eragina dauka bariantza eta edozein momentu zentratuan baina traslazioak ez.

II.3.1. Aldagai zentratua

Aldagai baten batezbesteko aritmetikoa zero baldin bada, **aldagai zentratua** deitzen da; aldiz, batezbesteko aritmetikoa zero eta desbidazio tipikoa bat balioak baldin badira, **aldagai tipifikatua** deritzogu.

Edozein X aldagairen kasuan bere batezbestekoa kenduz $Z = X - \bar{x}$, aldagai zentratua lortzen dugu.

$$\begin{aligned}\bar{z} &= \frac{1}{N} \sum_i z_j n_i = \frac{1}{N} \sum_i (x_i - \bar{x}) n_i = \\ &= \frac{1}{N} \sum_i x_j n_i - \frac{1}{N} \bar{x} \sum_i n_i = \bar{x} - \bar{x} = 0\end{aligned}$$

Kontutan hartu behar dugu zentratze-eragiketak ez duela unitate-aldaketarik suposatzen, jatorri-aldaketa edota aldagaiaren translazioa ardatzean suposatzen du bakarrik.

Aldagaiaren grabitate-zentrua, translazio honetan, jatorrira eraman dugu, $\tilde{z} = 0$ baita.

II.3.2. Aldagai tipifikatua edo standardizatua

Modu berean, edozein X aldagairen kasuan batezbestekoa kenduz eta desbidazio tipikoaz zatituz $T = (X - \bar{x}) / S_x$, aldagai tipifikatua lortzen dugu.

$$\begin{aligned}\bar{t} &= \frac{1}{S_x} \sum_i (x_i - \bar{x}) \frac{n_i}{N} = 0 \\ S_t^2 &= \frac{1}{N} \sum_i (t_i - \bar{t})^2 n_i = \frac{1}{N} \sum_i \frac{(\bar{x}_i - \bar{x})^2}{S_x^2} n_i = \\ &= \frac{1}{S_x^2} \frac{1}{N} \sum_i (\bar{x}_i - \bar{x})^2 n_i = \frac{1}{S_x^2} S_x^2 = 1\end{aligned}$$

Tipifikazio-eragiketak, zentratzeaz gainera unitatez aldatzea suposatzen du. Aldagai bakoitzaren balio zentratuak, bere desbidazio tipikoaz zatitzean, honen arabera neurturik ditugu, jatorrizko neurri-unitateak desagertu egiten direlarik.

II.4. KONTZENTRAZIO-NEURKETAK

Batezbesteko aritmetikoa egitean zenbakitzailean ateratzen zaigun kopuruak, abiapuntu estatistiko batetatik ikusita, batzutan ez du zentzu argirik eta interesgarririk. Adibidez, pertsona batzuren altuera-banaketa aztertzean, zenbakitzaile hori altueren batuketa izango da; baina, beste kasu batzutan, bereziki

ezaugarri sozioekonomikoetan, aztergarria izaten da, adibidez, alokairu-banaketa aztertzean, alokairu guztien kopurua edo “alokairu-masa” da.

Kontzentrazio-neurketek, zehazki, “masa” baten uniformetasuna neurtzea dute helburutzat.

Normalki errenta eta alokairu-banaketetan erabili ohi dira, baina beste edozein aldagairen banaketatan erabil daitezke.

Suposa dezagun alokairu-banaketa batetan langile guztiek alokairu berdina jasotzen dutela; orduan, banatzearen uniformetasuna osoa litzateke. Baina, ordez, alokairu-masa guztia langile batek jasoko balu, **uniformetasunik eza** ere osoa izango litzateke edo **kontzentrazioa maximoa** dela esango genuke kasu honetan.

Banaketak uniformetasun absoluturantz jotzen duenean, maiztasun-banaketaren batezbesteko aritmetikoa guztiz adierazgarria da, eta, alderantziz gertatuko da, kontzentrazioa maximoa denean.

II.4.1. LORENZ-en kurba

Grafikoki banaketen kontzentrazioa ikusarazten digun neurrian, interesgarria zaigu.

Lorenz-en kurba eraikitzeko, ondoko taulan ikusten ditugun zutabeak kalkulatu beharko ditugu.

x_i	n_i	$x_i n_i$	N_i	M_i	P_i	Q_i
x_1	n_1	$x_1 n_1$	N_1	M_1	P_1	Q_1
x_2	n_2	$x_2 n_2$	N_2	M_2	P_2	Q_2
.
.
x_i	n_i	$x_i n_i$	N_i	M_i	P_i	Q_i
.
.
x_k	n_k	$x_k n_k$	N	M	100	100

$$\sum_i n_i = N \qquad \sum_i x_i n_i = M$$

Azken zutabeak honela kalkulatzen dira:

$$N_i: (n_i) \text{ zutabea metatuz} \\ \text{hots: } N_i = n_1 + n_2 + \dots + n_i$$

$$M_i: (x_i n_i) \text{ zutabea metatuz} \\ \text{hots: } M_i = x_1 n_1 + x_2 n_2 + \dots + x_i n_i$$

$$P_i = \frac{N_i}{N} \cdot 100 \text{ maiztasunen (langileena esate baterako) portzentaia metatua}$$

$$Q_i = \frac{M_i}{M} \cdot 100 \text{ masaren portzentaia metatua.}$$

$\{(P_i, Q_i)\}_{i \in I}$ bikote guztiak sistema cartesiar batetan grafikoki adieraziz, eta (P_1, Q_1) , $(0,0)$ jatorriarekin eta beste bikote guztiak elkarren artean zuzenen bitartez lotuz LORENZ-en kurba lortzen dugu.

Adibidez: Suposa dezagun enpresa batetako langileen banaketa bere asteroko alokairuaren arabera (milaka pezetatan) ondoko taularen lehenengo bi zutabeetan emana dela:

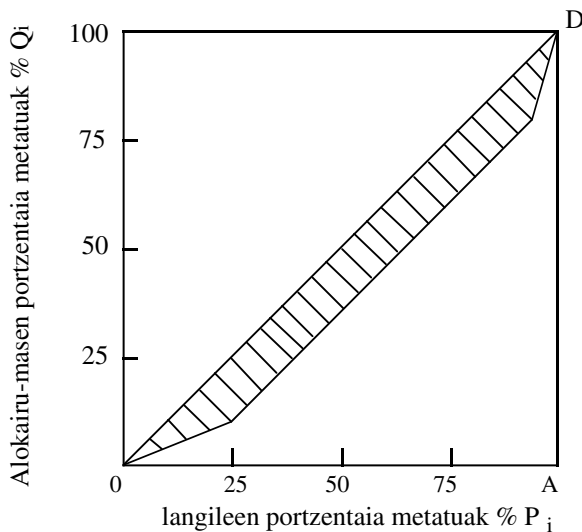
Alokairuak: x_i	Langile kopurua: n_i	Alokairu- -masa: $x_i n_i$	Metaketak N_i	M_i	Portzentaia metatuak: P_i	Q_i
8	10	80	10	80	20	10
15	20	300	30	380	60	47,5
20	15	300	45	680	90	85
24	5	120	50	800	100	100

$$N = 50 \quad M = \sum_{i=1}^k x_i n_i = 800$$

Adibideari dagokion LORENZ-en kurba ondoko hau izango da.

Eskalak, bi ardatzetan berdina direnez gero, LORENZ-en kurba, karratu baten barruan aurkitzen da beti. Jatorritik irteten den OD bere diagonal erabiliko dugu uniformetasunaren erreferentzia moduan (P_i, Q_i) multzoko puntu guztiak, OD-n aurkitzen badira, orduan uniformetasun osoa baitugu.

Hots: langileen arteko % 10-ak, alokairu-masaren % 10-a jasotzen du, % 30-ak, alokairu-masaren % 30-a, e.a.



Zenbat eta tarte handiagoa egon kurba eta diagonalaren artean, orduan eta **banaketa kontzentratuagoa** izango da, (edota zenbat eta kurba diagonaletik hurbilago egon, orduan eta **banaketa uniformeagoa**).

Abzisa ardatzean P_i kokatzen badugu eta ordenatu ardatzean Q_i , LORENZ-en kurba, karratuaren beheko triangeluan edukiko dugu (alokairuak txikienetatik handienetarantz ordenatu ditugunez gero, adibidez langileen % 25-ak ezin du alokairu-masaren % 25-a inoiz lortu).

II.4.2. GINI-ren indizea

LORENZ-en kurba banaketen kontzentrazioaren argitzaile bada ere, zenbakizko indize batetaz kontzentrazioa neurtzea, komenigarria zaigu; esate baterako, erabilgarria izango zaigu banaketak konparatzerakoan.

Helburu honi erantzuten dio **Gini-ren indizeak** edo **kontzentrazio-indizeak**.

Marraztutako azalera diagonal eta kurbaren artean dagoena bada, honela definituko dugu I_G , GINI-ren indizea:

$$I_G = \frac{\text{Marraztutako azalera}}{\triangle \text{ OAD azalera}}$$

$$\text{non } 0 \leq I_G \leq 1$$

Kontzentrazioa maximoa bada, orduan $\hat{\Delta}$ azalera = marraztutako azalera, eta $I_G = 1$.

Ordez, kontzentrazioa minimoa bada edota banaketa uniforme, orduan marraztutako azalera = 0 izango da eta $I_G = 0$.

Kalkulua: I_G GINI-ren indizea kalkulatzeko, azalera horiek zuzenki neurtuz kalkula litezke; baina beti era honetan erraza izaten ez denez gero, **gutxi gora-behera** ondoko formula erabiliz kalkulatzen da I_G :

$$I_G = \frac{\sum_{i=1}^{k-1} (P_i - Q_i)}{\sum_{i=1}^{k-1} P_i}$$

Kontura gaitzen, batukariak $k-1$ batugai bakarrik dituela, $p_k - q_k = 0$ baita.

Gutxi gora-beherako formula honetaz ere egiazta daiteke $0 \leq I_G \leq 1$ dela.

Adibidez: Lehenengo taularen datuekin jarraitzen badugu, aurreko formula erabiliz, I_G GINI-ren kontzentrazio-indizea ondoko hau izango da:

P_i	Q_i	$P_i - Q_i$
20	10	10
60	47'5	12'5
90	85	5
170		27'5

$$I_G = \frac{\sum_{i=1}^{k-1} (P_i - Q_i)}{\sum_{i=1}^{k-1} P_i} = \frac{27'5}{170} = 0'16$$

LORENZ-en kurba nahiz GINI-ren indizea aldagai sozioekonomikoetarako erabilgarriak dira. Adibide bezala hauek aipatuko ditugu: errenta pertsonalaren banaketa, nekazal-propietatearen banaketa, sektore batetako enpresa-produkzioaren banaketa eta abar.

**III. EZAUGARRI ESTADISTIKO BIKOITZAK,
BANAKETAK, TAULAK, ADIERAZPIDE
GRAFIKOAK**

III.1. EZAUGARRI ESTADISTIKO BIKOITZAK

III.2. MAIZTASUN-BANAKETA BIKOITZAK

*III.2.1. Maiztasun-banaketa bikoitzak eta
bazter-maiztasunak*

III.2.2. Maiztasun-banaketa bikoitzak: Taulak

III.3. MAIZTASUN-BANAKETA BIKOITZEN

ADIERAZPIDE GRAFIKOAK

III.3.1. Sakabanatze-diagrama edo puntu-hodeia

III.4. MAIZTASUN-BANAKETA BALDINTZATUAK

III.5. TAULA BATEN DEPENDENTZIA EDO

INDEPENDENTZIA

III.6. KONTINGENTZI TAULA BATEN ERRENKADA

ETA ZUTABEEN BATEZBESTEKO SOSLAIK

III.1. EZAUGARRI ESTATISTIKO BIKOITZAK

Orain arte ikusi ditugun fenomenoetan, indibiduo bakoitzari ezaugarri estatistiko bakoitzarekiko neurketa bat besterik ez zitzaion egiten.

Orain aipatuko ditugun ikerketetan, lagineko **indibiduo bakoitzari bi neurketa** egiten dizkiogu, bakoitza ezaugarri estatistiko bati dagokiona.

Ezaugarri estatistiko bikoitzak edo ezaugarri estatistiko bidimentsionalak ikusiko ditugu, bada, ondoren.

Adibidez: Urte batetan soldaduzkara doazen gazteei bi neurketa egiten zaizkie, edo beraien bi ezaugarri aztertzen dira batera: pisua eta altuera.

III.2. MAIZTASUN-BANAKETA BIKOITZAK

III.2.1. Maiztasun-banaketa bikoitzak eta bazter-maiztasunak

Suposa dezagun N indibiduoaz osatutako lagin bat, non X ezaugarri batek $x_1 x_2 \dots x_i \dots x_k$ balio desberdinak hartzen dituen.

Suposa dezagun, halaber, beste Y ezaugarri bakun batek lagin berdineko indibiduoetan $y_1 y_2 \dots y_j \dots y_l$ balio desberdinak hartzen dituela.

Orduan, (X,Y) ezaugarri bikoitzak, lagineko indibiduo bakoitzean balio-bikote bat hartzen du.

Hots:

$$\omega_p \rightarrow (x_i, y_j) \quad \begin{matrix} p = 1, \dots, N \\ i = 1, \dots, k \\ j = 1, \dots, l \end{matrix}$$

Eta (X,Y) ezaugarri bikoitzak, laginean hartuko lituzkeen **balio-bikote desberdinen** kopurua, $k \times l$ izango da. Hemendik aurrera horiek bakarrik hartuko ditugu kontuan.

Hots:

$$\{ \omega_1, \omega_2 \dots \omega_N \} \xrightarrow{(X,Y)} \{ (x_1 y_1) (x_1 y_2) \dots (x_i y_1) \dots \dots (x_i y_l) \dots (x_k y_1) \dots \dots (x_k y_l) \} = \{ (x_i y_j) \}$$

(x_i y_j) balio-bikotearen MAIZTASUN ABSOLUTUA

Sinbolikoki n_{ij} ikurraz adieraziko dugu, eta laginean (x_i y_j) balio-bikotea hartzen duten indibiduen kopurua da.

Non:

$$\sum_{i=1}^l \sum_{j=1}^k n_{ij} = n_{..} = N$$

(x_i y_j) balio-bikotearen MAIZTASUN ERLATIBOA

Sinbolikoki f_{ij} ikurraz adieraziko dugu, non:

$$f_{ij} = \frac{n_{ij}}{N} \quad \sum_{j=1}^k \sum_{i=1}^l f_{ij} = 1$$

eta portzentaetan:

$$p_{ij} = \%f_{ij} \cdot 100 \quad \sum_j \sum_i p_{ij} = 100$$

 x_i balioaren BAZTER-MAIZTASUN ABSOLUTUA

Sinbolikoki $n_{i.}$ ikurraz adieraziko dugu, non:

$$n_{i.} = \sum_{j=1}^k n_{ij} \quad \sum_j n_{i.} = n_{..} = N$$

 y_j balioaren BAZTER-MAIZTASUN ABSOLUTUA

Sinbolikoki $n_{.j}$ ikurraz adieraziko dugu, non:

$$n_{.j} = \sum_{i=1}^l n_{ij} \quad \sum_j n_{.j} = n_{..} = N$$

 x_i balioaren BAZTER-MAIZTASUN ERLATIBOA

Sinbolikoki $f_{i.}$ ikurraz adieraziko dugu, non:

$$f_{i.} = \frac{n_{i.}}{N} = \frac{\sum_j n_{ij}}{N}$$

y_j balioaren **BAZTER-MAIZTASUN ERLATIBOA**

Sinbolikoki $f_{.j}$ ikurrak adieraziko dugu, non:

$$f_{.j} = \frac{n_{.j}}{N} = \frac{\sum_i n_{ij}}{N}$$

eta

$$f_{..} = \sum_j f_{.j} = \sum_j f_{.j} = 1$$

III.2.2. Maiztasun-banaketa bikoitzak: Taulak

Lagin bati dagozkion maiztasun-banaketen taulak, ezaugarria bikoitza bada, ondoan azaltzen ditugunak dira.

X \ Y	y_1	y_2	y_j	y_l	
x_1	n_{11} f_{11}	n_{12} f_{12}	n_{1j} f_{1j}	n_{1l} f_{1l}	$n_{1.}$ $f_{1.}$
x_2	n_{21} f_{21}	n_{22} f_{22}	n_{2j} f_{2j}	n_{2l} f_{2l}	$n_{2.}$ $f_{2.}$
.				
x_i	n_{i1} f_{i1}	n_{i2} f_{i2}	n_{ij} f_{ij}	n_{il} f_{il}	$n_{i.}$ $f_{i.}$
.
x_k	n_{k1} f_{k1}	n_{k2} f_{k2}	n_{kj} f_{kj}	n_{kl} f_{kl}	$n_{k.}$ $f_{k.}$
Guztira	$n_{.1}$ $f_{.1}$	$n_{.2}$ $f_{.2}$			$n_{.j}$ $f_{.j}$			$n_{.l}$ $f_{.l}$	$n_{..} = N$ $f_{..} = 1$

BAZTER-BANAKETAK

X	M. Ab.	M. Er.
x_1	$n_{1.}$	$f_{1.}$
x_2	$n_{2.}$	$f_{2.}$
\vdots	\vdots	\vdots
x_i	$n_{i.}$	$f_{i.}$
\vdots	\vdots	\vdots
\vdots	\vdots	\vdots
x_k	$n_{k.}$	$f_{k.}$

Y	M. Ab.	M. Er.
y_1	$n_{.1}$	$f_{.1}$
y_2	$n_{.2}$	$f_{.2}$
\vdots	\vdots	\vdots
y_j	$n_{.j}$	$f_{.j}$
\vdots	\vdots	\vdots
\vdots	\vdots	\vdots
y_l	$n_{.l}$	$f_{.l}$

$$\sum_i n_{i.} = N$$

$$\sum_j n_{.j} = N$$

$$\sum_i f_{i.} = 1$$

$$\sum_j f_{.j} = 1$$

(X edo Y ezaugarriak, edo biak, jarraiak badira, ezaugarri bakunetan bezala mailakaturik diskretu bihurtuko ditugu).

Adibidea: 750 familia duen talde batean bi ezaugarri aztertu dira batera. Bata aktibo diren pertsonen kopurua, eta bestea hileroko soldata (milaka pezetatan).

Ondoko taulan dauzkagu maiztasun absolutuen banaketak.

Hileroko soldata milaka pezetatan

Pertsona aktiboak	(60,80)	(80,140)	(140,250)	Guztira
1	143	59	--	202
2	46	170	32	248
3	----	93	86	179
4	----	17	104	121
Guztira	189	339	222	750

Bazter-maiztasunaren esangura, bereziki azpimarratzea komeni zaigu.

Ikus dezagun nola maiztasun absolutuen kopuruetatik (ezaugarri bikoitzari dagozkionak), batuketa batzuekin, bazter-maiztasunen banaketak lortzen

ditugun: hau da, errenkadak batuz, ezaugarri batenak (pertsona aktiboak), eta zutabeak batuz, beste ezaugarriarenak (hileroko soldata).

Ohargarria da, bada, bazter-maiztasunen banaketak, maiztasun bakunen banaketa batzu besterik ez direla.

III.3. MAIZTASUN-BANAKETA BIKOITZEN ADIERAZPIDE GRAFIKOAK

Aldagai bikoitz baten maiztasun-banaketa ikustarazteko bi adierazpide grafiko ditugu, ESTEREOGRAMA eta SAKABANATZE-DIAGRAMA edo PUNTU-HODEIA hain zuzen. Hauetariko bigarrena ondoan azalduko dugu, bestea ez baita ia erabiltzen.

III.3.1. Sakabanatze-diagrama edo puntu-hodeia

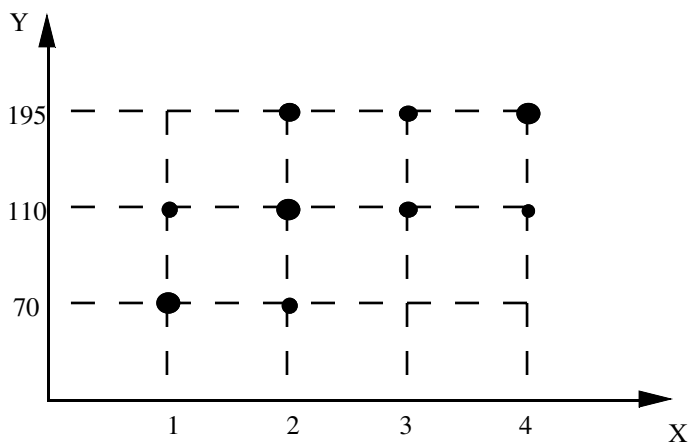
Izan bedi (X,Y) ezaugarri estatistiko bikoitz bat, eta izan bitez

$$\left\{ \begin{matrix} (x_i, y_j) \\ i = 1 \dots k \\ j = 1 \dots l \end{matrix} \right\} \quad (X,Y) \text{ ezaugarriak } N \text{ tamainuko lagin batetan hartzen dituen balio-bikote desberdinak.}$$

Orduan, ardatz cartesianarreko sistema batetan (x_i, y_j) balio-bikote bakoitza puntu batez adierazten badugu,

$$\left\{ \begin{matrix} (x_i, y_j) \\ i = 1, \dots, k \\ j = 1, \dots, l \end{matrix} \right\} \quad \text{multzo osoaren adierazpenari, sakabanatze-diagrama edo puntu-hodeia deritzo.}$$

Ondoan, 2.2. ataleko adibideari dagokion puntu-hodeia daukagu.



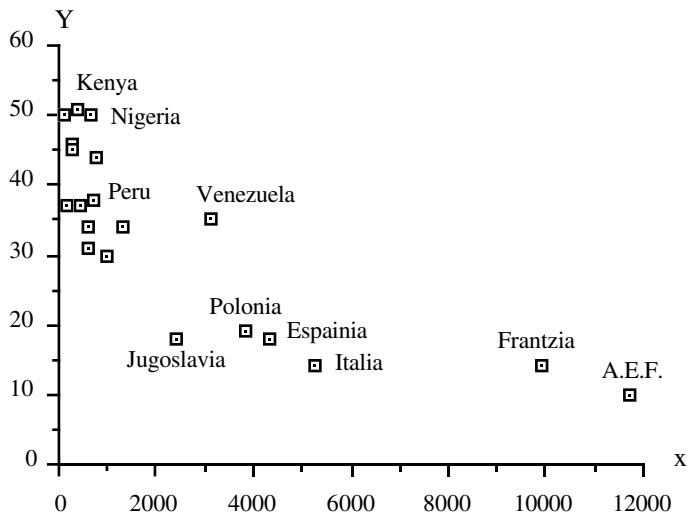
Ikusten dugunez, (x_i, y_j) balio-bikoteak errepikatzen direnean sakabanatze-diagramako puntuetariko batzu besteak baino lodiagoak dira, beraien azalerak eta balio-bikoteei dagozkien maiztasunak proportzionalak izanik.

Dena dela, taula bikoitzak edo datu taldekatuen taulak eta beraren dispersio-diagramak informazio berdintsua ematen digute.

Batez ere datu-multzo handiak ditugunean, daturik isolatuenak edo ezer gutxi errepikatzen diren datuenak dira taulak; eta horrelako kasuetan, puntu-hodeiak oso baliagarriak zaizkigu, multzoaren ibilbidea ikustarazteko.

X ezaugarria Nazio-Produktu Gordina eta Y jaiokortasuna izanik, ondoan 20 estaturi dagozkien datuak eta beraien puntu-hodeia ditugu.

<u>ESTATUA</u>	<u>IZENA</u>	<u>NPG</u>	<u>JAIOKORTASUNA</u>
1	Ethiopia	130	50
2	Birmania	160	37
3	Tanzania	270	46
4	Uganda	290	45
5	Kenya	380	51
6	Egipto	460	37
7	Thailandia	590	31
8	Filipinas	600	34
9	Nigeria	670	50
10	Peru	730	38
11	Maroko	740	44
12	Kolonbia	1010	30
13	Turkia	1330	34
14	Jugoslavia	2430	18
15	Venezuela	3130	35
16	Polonia	3830	19
17	Espainia	4340	18
18	Italia	5240	14
19	Frantzia	9940	14
20	A. Err. Federatua	11730	10



III.4. MAIZTASUN-BANAKETA BALDINTZATUAK

(X,Y) ezaugarri edo aldagaien banaketa bikoitzaren kasuan eta X aldagaiaren x_i balioarentzat, Y aldagaiaren maiztasun-banaketa kontsidera daiteke.

Y \ X	y_1	y_2	y_j	y_l	
.	
.	
.	
x_i	n_{i1}	n_{i2}	n_{ij}	n_{il}	n_i
.	
.	
.	

x_i balioa baldintzatzailea izanez, maiztasun-banaketa baldintzatuak ondokoak ditugu:

$Y/X = x_i$	$n_{j/i}$	$f_{j/i}$
y_i	$n_{1/i} = n_{i1}$	$\frac{n_{1/i}}{n_{i.}}$
y_2	$n_{2/i} = n_{i2}$	$\frac{n_{2/i}}{n_{i.}}$
\cdot	\cdot	\cdot
\cdot	\cdot	\cdot
\cdot	\cdot	\cdot
y_j	$n_{j/i} = n_{ij}$	$\frac{n_{j/i}}{n_{i.}}$
\cdot	\cdot	\cdot
\cdot	\cdot	\cdot
\cdot	\cdot	\cdot
y_l	$n_{l/i} = n_{il}$	$\frac{n_{l/i}}{n_{i.}}$

eta (1) x_i balioaz baldintzatutako Y aldagaiaren maiztasun absolutuen banaketa deitzen da.

Maiztasun-banaketa honen balioak, $n_{i.}$ bazter-maiztasunaz zatitzen baditugu (2), maiztasun-banaketa erlatiboa lortzen dugu, eta:

$$f_{i/j} = \frac{n_{j/i}}{n_{i.}} = \frac{n_{ij}}{n_{i.}}$$

X aldagaiak hartzen dituen balio guztiak kontsideratuz, beste hainbeste Y aldagaiaren banaketa baldintzatu edukiko ditugu.

Analogikoki, X aldagaiaren banaketa baldintzatuak edukiko ditugu, Y aldagaiak hartzen dituen balio guztiak baldintzatzailerak izanez.

Adibidea:

Ondoko taulan (X,Y) aldagaien banaketa daukagu. Y aldagaiak pieza baten luzera adierazten du.

X aldagaiak, ordea, akastuna (B) edo akasgabea (A) izatea.

X \ Y	71	72	73	n_i
B	3	1	5	9
A	2	4	6	12
n_j	5	5	11	21

Maiztasun-banaketa baldintzatuak (erlatiboak) hauek dira:

$f(y/x = B)$	$3/9$	$1/9$	$5/9$
$f(y/x = A)$	$2/12$	$4/12$	$6/12$

$f(x/y = 71)$	$f(x/y = 72)$	$f(x/y = 73)$
$3/5$	$1/5$	$5/11$
$2/5$	$4/5$	$6/11$

Horrelako banaketa baldintzatuak zera adierazten digute: aldagai baten balio bat finko kontsideratuz, beste aldagaiaren balioek dauzkaten proportzioak zeintzu diren.

Hau da: pieza horiek osatzen duten populazioan akastunen azpipopulazioa kontsideratuz, zer proportziotan daude luzera desberdinetakoak akastunen artean? Erantzuna $f(y/x = B)$ banaketak ematen digu.

Analogikoki, akasgabeen azpipopulazioa kontsideratzen badugu, $f(y/x = A)$ edukiko dugu.

Aldiz, Y aldagaiaren balio bat baldintzatzailea baldin bada, dagokion azpipopulazioa kontsideratuko dugu.

Adibidez, Y aldagaiak luzera milimetrotan adierazten badu eta 73 mm neurtzen dituztenen azpipopulazioa kontsideratuz, zer proportziotan dira akastunak ala akasgabeak pieza horiek? Erantzuna kasu honetan, $f(x/y = 73)$ banaketak ematen digu.

III.5. TAULA BATEN DEPENDENTZIA EDO INDEPENDENTZIA

X eta Y aldagaiak kontsideratuz, aldagai baten aldaketarekin bestearen banaketa aldatzen bada, biak elkarren artean dependenteak dira.

Dakusagun ondoko taula:

X \ Y	71	72	73	n_i	f_i
B	3	1	5	9	9/21
A	2	4	6	12	12/21
n_j	5	5	11	21	
f_j	5/21	5/21	11/21		1

X aldatzean, Y aldagaiaren banaketa aldatzen da, hau da:

y_j/x_i	$n(y/x = B)$	$n(y/x = A)$	$f(y/x = B)$	$f(y/x = A)$
71	3	2	1/3	1/6
72	1	4	1/9	1/3
73	5	6	5/9	1/2

Luzera horiek, bada, ez daude berdin banatuta akastunen artean eta akasgabeen artean, hots, nolabaiteko menpekotasun edo dependentzia ikusten da taula honetan.

Demagun, orain, honako taula hau:

X \ Y	71	72	73	$n_{i.}$	$f_{i.}$
B	3	6	9	18	18/30
A	2	4	6	12	12/30
$n_{.j}$	5	10	15	$30 = N = n_{..}$	
$f_{.j}$	5/30	10/30	15/30		

X aldatzean, Y aldagaiaren banaketa ez da aldatzen, hau da:

y_j/x_i	$n(y/x = B)$	$n(y/x = A)$	$f(y/x = B)$	$f(y/x = A)$
71	3	2	3/18	2/12
72	6	4	6/18	4/12
73	9	6	9/18	6/12

Ikusten dugunez, Y aldagaiaren banaketa baldintzatuak berdinak dira, luzerak proportzio berdinetan daude akastunen eta akasgabeen artean, eta Y aldagaiaren bazter-banaketa (5/30, 10/30, 15/30) ere berdina da.

Hau da, x_i balioaz baldintzatutako y_j balioaren maiztasun erlatibo baldintzatua eta y_j balioaren bazter-maiztasun erlatiboa berdinak dira:

$$\frac{n_{ij}}{n_{i.}} = \frac{n_{.j}}{N}$$

X aldagaiaren banaketa baldintzatuak eta bazter-banaketa kontsideratuz, gauza bera gertatzen da, hots, Y aldagaiaren balioak aldatzean, X aldagaiaren banaketa ez da aldatzen eta X eta Y independenteki banatzen dira.

Aurreko erlazioaren simetrikoa edukiko dugu, hau da:

$$\frac{n_{ij}}{n_{.j}} = \frac{n_{.i}}{N}$$

Bi erlazio horietatik erraz ateratzen den independentziaren baldintza beharrezkoa eta nahikoa ondokoa da:

$$f_{ij} = f_i \cdot f_j \quad \forall (i,j)$$

Hots, edozein (x_i, y_j) bikoterentzat, bazter-maiztasun erlatiboen biderkadura eta bikote horren maiztasun erlatiboa berdinak dira.

III.6. KONTINGENTZI TAULA BATEN ERRENKADA ETA ZUTABEEN BATEZBESTEKO SOSLAIK

X	Y	$Y_1 \dots y_j \dots y_l$	$n_{i.}$
x_1		$n_{11} \dots n_{1j} \dots n_{1l}$	$n_{1.}$
\vdots		$\vdots \quad \quad \quad \vdots$	\vdots
x_i		$n_{i1} \dots n_{ij} \dots n_{il}$	$n_{i.}$
\vdots		$\vdots \quad \quad \quad \vdots$	\vdots
x_k		$n_{k1} \dots n_{kj} \dots n_{kl}$	$n_{k.}$
$n_{.j}$		$n_{.1} \dots n_{.j} \dots n_{.l}$	N

Y eta X aldagaien maiztasun erlatiboen bazter-banaketak, Y aldagaiaren banaketa baldintzatuen eta X aldagaiaren banaketa baldintzatuen batezbesteko soslaiak dira. Y aldagaiaren banaketa baldintzatuen edo errenkada soslaien batezbestekoa hauxe da:

$$\sum_i \frac{n_{ij}}{n_{i.}} f_{i.} = f_{.j}$$

$\sum_i \frac{n_{ij}}{n_{i.}} f_{i.} = f_{.j}$ balioen batezbesteko ponderatua da, $f_{i.}$ maiztasuna x_i elementu baldintzatzaile bakoitzaren pisu erlatiboa izanik.

Halaber,
$$\sum_i \frac{n_{ij}}{n_{.j}} f_{.j} = f_{.j}$$

Hots: $f_{i.}, \frac{n_{ij}}{n_{.j}}$ balioen batezbesteko ponderatua da $f_{.j}$ maiztasuna y_j elementu baldintzatzaile bakoitzaren pisu erlatiboa izanik.

IV. EZAUGARRI BIKOITZEN BALIO TIPIKOAK

IV.1.MOMENTUAK

IV.1.1. Momentu arruntak edo jatorriarekikoak

IV.1.2. Momentu zentralak edo batezbestekoarekikoak

IV.2.MAIZTASUN-BANAKETA BIKOITZEN BALIO TIPIKO EDO ESTATISTIKOAK

IV.2.1. Kobariantza: S_{xy}

IV.2.2. Koerlazio-koefizientea: r_{xy}

IV.3.ALDAGAIEN TRANSFORMAZIO LINEALAK

IV.4.KOERLAZIO GABEKO ALDAGAIEN BI PROPIETATE

IV.5.SAKABANATZE- ETA KOERLAZIO-MATRIZEAK

IV.6.OROKORPENA

IV.1. MOMENTUAK

II. ikasgaian, kasu bakunerako gertatu zen bezala, balio tipikoak edo estatistikoak aurkeztu baino lehen, momentuak ikusi behar ditugu eta arrazoi berdinetatik, hots, balio tipikoak momentuak direlako edo momentuen bidez definitzen direlako.

IV.1.1. Momentu arruntak edo jatorriarekikoak

Maiztasun-banaketa bikoitz baten h . ordenako momentu arrunta edo jatorriarekikoa, ondoko formularen bidez definitzen da:

$$a_{h_1 h_2} = \frac{1}{N} \sum_i \sum_j x_i^{h_1} y_j^{h_2} n_{ij} \quad h_1 + h_2 = h$$

1. ordenako edo lehenengo ordenako momentu arruntak hauexek dira:

$$\begin{aligned} a_{10} &= \frac{1}{N} \sum_i \sum_j x_i n_{ij} = \frac{1}{N} \sum_i x_i \sum_j n_{ij} = \\ &= \frac{1}{N} \sum_i x_i n_{i.} = \bar{x} \end{aligned}$$

$$\begin{aligned} a_{01} &= \frac{1}{N} \sum_i \sum_j y_j n_{ij} = \frac{1}{N} \sum_j y_j \sum_i n_{ij} = \\ &= \frac{1}{N} \sum_j y_j n_{.j} = \bar{y} \end{aligned}$$

2. ordenako edo bigarren ordenako momentu arruntak hauexek dira:

$$\begin{aligned} a_{20} &= \frac{1}{N} \sum_i \sum_j x_i^2 n_{ij} = \frac{1}{N} \sum_i x_i^2 \sum_j n_{ij} = \\ &= \frac{1}{N} \sum_i x_i^2 n_{i.} = a_2(x) \end{aligned}$$

Hau da, X ezaugarriaren bazter-banaketa kontsideratuz, banaketa bakun honen bigarren ordenako momentu arrunta, aurrekoa da.

$$\begin{aligned} a_{02} &= \frac{1}{N} \sum_i \sum_j y_i^2 n_{ij} = \frac{1}{N} \sum_j y_j^2 \sum_i n_{ij} = \\ &= \frac{1}{N} \sum_j y_j^2 n_{.j} = a_2(y) \end{aligned}$$

Y ezaugarriaren bazter-banaketaren bigarren ordenako momentu arrunta, aurrekoa izango da.

$$a_{11} = \frac{1}{N} \sum_i \sum_j x_i y_j n_{ij}$$

Hirugarren ordenako momentuak hauexek izango lirateke:

$$a_{30}, a_{03}, a_{21}, a_{12}$$

IV.1.2. Momentu zentralak edo batezbestekoarekikoak

Maiztasun-banaketa bikoitz baten h-ordenako momentu zentrala edo batezbestekoarekikoa, ondoko formularen bidez definitzen da:

$$m_{h_1, h_2} = \frac{1}{N} \sum_i \sum_j (x_i - \bar{x})^{h_1} (y_j - \bar{y})^{h_2} n_{ij}$$

Bereziki, bigarren ordenakoak oso garrantzitsuak dira:

$$m_{20}, m_{02}, m_{11}$$

$$\begin{aligned} m_{20} &= \frac{1}{N} \sum_i \sum_j (x_i - \bar{x})^2 n_{ij} = \frac{1}{N} \sum_j (x_i - \bar{x})^2 \sum_i n_{ij} = \\ &= \frac{1}{N} \sum_i (x_i - \bar{x})^2 n_{i.} = S_x^2 \end{aligned}$$

Hori da, bestetik, X ezaugarriaren bariantza (ikus 37. orrialdea)

$$\begin{aligned} m_{20} &= \frac{1}{N} \sum_i \sum_j (y_i - \bar{y})^2 n_{ij} = \frac{1}{N} \sum_j (y_j - \bar{y})^2 \sum_i n_{ij} = \\ &= \frac{1}{N} \sum_j (y_j - \bar{y})^2 n_{.j} = S_y^2 \end{aligned}$$

Y ezaugarriaren bariantza da.

$$m_{11} = \frac{1}{N} \sum_i \sum_j (x_i - \bar{x})(y_i - \bar{y})n_{ij}$$

m_{11} momentua S_{xy} KOBARIANTZA deitzen dena da, balio tipiko garrantzitsuetariko bat, hain zuzen.

IV.2. MAIZTASUN-BANAKETA BIKOITZEN BALIO TIPIKO EDO ESTADISTIKOAK

(X,Y) ezaugarri estatistiko bikoitz bat bada, eta bere maiztasun-banaketa lagin batetan kalkulaturik baldin badugu, X eta Y-ren arteko erlazioa neurtzen duten estatistiko batzuk kalkulatu genitzake.

IV.2.1. Kobariantza: S_{xy}

X eta Y ezaugarriak, laginean duten erlazio-neurri bat da (non X eta Y ezaugarriak indibiduen gain neurtzen direneko unitateek eragina duten).

Esan dugunez, m_{11} bigarren ordenako momentu hau kobariantza da; ikus dezagun nola adierazten den momentu arrunten bidez:

$$\begin{aligned} S_{xy} = m_{11} &= \frac{1}{N} \sum_i \sum_j (x_i - \bar{x})(y_i - \bar{y})n_{ij} = \\ &= \frac{1}{N} \sum_i \sum_j (x_i y_j - x_i \bar{y} - y_j \bar{x} + \bar{x} \bar{y})n_{ij} = \\ &= \frac{1}{N} \sum_i \sum_j x_i y_j n_{ij} - \frac{1}{N} \bar{y} \sum_i \sum_j x_i n_{ij} - \end{aligned}$$

$$\begin{aligned}
& -\frac{1}{N}\bar{x} \sum_i \sum_j y_j n_{ij} + \frac{1}{N}\bar{x}\bar{y} \sum_i \sum_j n_{ij} = \\
& = a_{11} - \bar{y}\bar{x} - \bar{x}\bar{y} + \bar{x}\bar{y} = a_{11} - \bar{x}\bar{y}
\end{aligned}$$

IV.2.2. Koerlazio-koefizientea: r_{xy}

Estatistiko honek ere X eta Y aldagaiek laginean duten erlazioa neurtzen digu, baina kasu honetan berdin da X eta Y aldagaien balioak zein unitatetan neurtzen diren, neurketan erabiltzen diren unitateek ez baitute eraginik koerlazio-koefizientean.

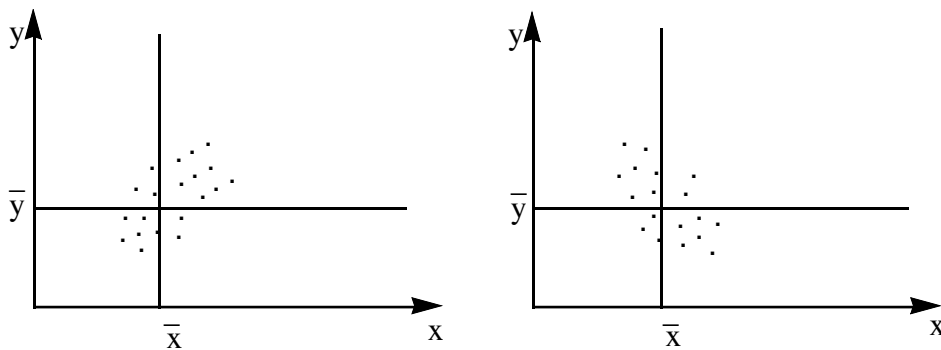
$$\begin{aligned}
r_{xy} &= \frac{S_{xy}}{S_x S_y} = \frac{\sum_i \sum_j (x_i - \bar{x})(y_j - \bar{y})n_{ij} / N}{\left(\frac{1}{N} \sum_i (x_i - \bar{x})^2 n_i\right)^{1/2} \left(\frac{1}{N} \sum_j (y_j - \bar{y})^2 n_{.j}\right)^{1/2}} = \\
&= \frac{a_{11} - \bar{x}\bar{y}}{(a_{20} - \bar{x}^2)^{1/2} (a_{02} - \bar{y}^2)^{1/2}}
\end{aligned}$$

Dakigunez S_x eta S_y desbidazio tipikoak, beti positiboak dira; orduan, S_{xy} kobariantza positiboa (negatiboa) baldin bada, r_{xy} koerlazio-koefizientea positiboa (negatiboa) izango da.

Intuitiboki honetaz ohar daiteke: puntu gehientzako $(x_i - \bar{x})$ eta $(y_j - \bar{y})$ kendurak biak positiboak edo negatiboak baldin badira, (hau da, (\bar{x}, \bar{y}) puntua jatorritzat hartzen badugu, sakabanatze-diagramako puntuak 1. eta 3. koadranteetan egongo dira) $S_{xy} > 0$ izango da; aldiz, puntu gehientzako zeinu kontrakoak baldin badira (puntuak 2. eta 4. koadranteetan egongo dira), $S_{xy} < 0$.

Hots: kobariantzaren zeinuak puntu-banaketaren ideia bat ematen digu.

Grafikoki:



Kobariantzak nahiz koerlazio-koefizienteak bi ezaugarrien artean dagoen erlazioa adierazten digute; baina ikusi dugu zein kasutan izango den $S_{xy} > 0$, hots, $(x_i - \bar{x})$ eta $(y_j - \bar{y})$ kendurak biak positiboak direnean (hau da, bi ezaugarriak gorakorrek dira) edo biak negatiboak direnean (hau da, bi ezaugarriak behakorrek). $S_{xy} < 0$ den kasuan, kenduren zeinuak kontrakoak direnez, ezaugarri bat gorakorra da eta behakorra bestea. Azkenez $S_{xy} = 0$ baldin bada, bi ezaugarriek ez dute erlazio linealik elkarren artean, hau da, koerlazio gabeak dira.

Ohar daiteke taula baten independentziak koerlazio eza halabehartzen duela, baina ez alderantziz.

Independentzia $(f_{ij} = f_i \cdot f_j) \not\Rightarrow$ koerlazio eza

Hots: Independentziaren kasuan

$$S_{xy} = \sum_i \sum_j (x_i - \bar{x})(y_j - \bar{y})f_{ij} = \sum_i (x_i - \bar{x})f_i$$

$$\sum (y_j - \bar{y})f_{.j} = m_1(x) \cdot m_1(y) = 0$$

r_{xy} koefizienteak ez du bakarrik aldagaien arteko erlazioa nolakoa den esaten (r_{xy} -ren zeinua S_{xy} -rena baita) baizik eta zer mailatan erlazonatuta dauden ere adierazten du.

Koerlazio-koefizienteak har ditzakeen balioak, $[-1, +1]$ tartekoak dira.

Zerora hurbiltzen den balio bat hartzen bada, bi ezaugarrien arteko erlazioa ahula da.

r_{xy} koefizientea +1 baliora hurbiltzen bada, erlazioa sakona eta zuzena da.

Aldiz, -1 baliora hurbiltzen bada, erlazioa sakona da, baina alderantzizkoa.

IV.3. ALDAGAIEN TRANSFORMAZIO LINEALAK

Ikus dezagun, bada, (X,Y) aldagaien ordean, haiekin erlazionatuta dauden (U,V) aldagaiak erabiltzen ditugunean, kobariantza aldatu egiten dela eta koerlazio-koefizientea ez.

Horrexegatik eta baita erregresio-teorian betetzen duen funtzioagatik ere, koerlazio-koefizientea, balio tipikorik garrantzitsuena da bi edo ezaugarri gehiago aztertzen ditugun ikerketetan.

Hots: $x_i = au_i + b$ eta $y_j = ev_j + f$ egiten badugu batezbesteko aritmetikoen arteko erlazioak, honako hauek dira:

$$\bar{x} = a\bar{u} + b \quad \text{eta} \quad \bar{y} = e\bar{v} + f$$

Ikus dezagun, koerlazio-koefizientea kalkulatzeko erabiltzen dugun desbidazio tipikoetan eta kobariantzan, zer eragina duen aldaketa honek.

$$S_x = |a|S_u \quad \text{eta} \quad S_y = |e|S_v$$

$$\begin{aligned} S_{xy} &= \frac{1}{N} \sum_i \sum_j (x_i - \bar{x})(y_j - \bar{y}) = \\ &= \frac{1}{N} \sum_i \sum_j (au_i + b - a\bar{u} - b)(ev_j + f - e\bar{v} - f)n_{ij} = \\ &= a \cdot e \cdot \frac{1}{N} \sum_i \sum_j (u_i - \bar{u})(v_j - \bar{v})n_{ij} = a \cdot e \cdot S_{uv} \end{aligned}$$

Jatorria aldatzeak eraginik ez duela ikusten da baina, bai baduela unitate-aldaketa horrek.

Aldiz, koerlazio-koefizientean ez du eraginik ez jatorri aldaketak, ezta unitate-aldaketak ere.

Hots:

$$r_{xy} = \frac{S_{xy}}{S_x S_y} = \frac{a \cdot e \cdot S_{uv}}{|a| S_u |e| S_v} = \frac{S_{uv}}{S_u S_v} = |r_{uv}|$$

Ikusten dugunez, bi aldagaien arteko koerlazio-koefizientea inbariantea da edozein transformazio linealen aurrean.

a edo e negatiboa balitz erlazioaren zentzua aldatuko litzateke.

KOERLAZIO-KOEFIZIENTEA, aldagai tipifikatuen arteko kobariantza da.

$$\text{Hots: } Z = \frac{(X - \bar{x})}{S_x} \quad \text{eta} \quad Z' = \frac{(Y - \bar{y})}{S_y} \quad \text{baldin badira:}$$

$$\begin{aligned} S_{zz'} &= \frac{1}{N} \sum_i \sum_j \left(\frac{x_i - \bar{x}}{S_x} \right) \left(\frac{y_j - \bar{y}}{S_y} \right) n_{ij} = \\ &= \frac{1}{S_x S_y} S_{xy} = r_{xy} \end{aligned}$$

Adibidea: Lagin batetan X eta Y ezaugarriak aztertu dira. Maiztasun-banaketaren taula ondokoa bada, kalkula ditzagun: batezbestekoak, bariantzak, kobariantza eta koerlazio-koefizientea, kalkulurik laburrena eginez.

X \ Y	1'735	1'740	1'745	n_i
239	7	9	9	25
243	8	7	5	20
247	11	9	10	30
251	4	7	14	25
n_j	30	32	38	100

Aldaketa hauek egiten ditugu: $u_i = \frac{x_i - 243}{4}$, $v_j = \frac{y_j - 1'740}{0'005}$

hots: $x_i = 4u_i + 243$ eta $y_j = 0'005 v_j + 1'740$

honela beste taula honetara iristen gara:

u \ v	-1	0	1	$n_{i.}$
-1	7	9	9	25
0	8	7	5	20
1	11	9	10	30
2	4	7	14	25
$n_{.j}$	30	32	38	100

Eta kalkuluak eginez:

$$\bar{u} = \frac{\sum_i u_i n_{i.}}{N} = 0'55$$

$$\bar{x} = 4 \bar{u} + 243 = 245'2$$

$$\bar{v} = \frac{\sum_j v_j n_{.j}}{N} = 0'08$$

$$\bar{y} = 0'005 \bar{v} + 1'740 = 1'7404$$

$$S_u^2 = 1'2475$$

$$S_x^2 = 4^2 S_u^2 = 19'96$$

$$S_v^2 = 0'6736$$

$$S_y^2 = 0'005^2 \cdot S_v^2 = 1'684 \cdot 10^{-5}$$

$$S_{uv} = 0'126$$

$$S_{xy} = 4 \cdot 0'005 \cdot S_{uv} = 0'00252$$

$$r_{xy} = \frac{S_{xy}}{S_x S_y} = \frac{S_{uv}}{S_u S_v} = 0'1374$$

Ikusten dugunez, X eta Y ezaugarrien arteko erlazioa zuzena da eta oso ahula; ez dute ia zerikusirik batak eta besteak.

IV.4. KOERLAZIO GABEKO ALDAGAIEN BI PROPIETATE

2.2. atalean esaten genuen bezala, $S_{xy} = 0$ denean ez dago inongo erlazio linealik X,Y aldagaien artean, horrexegatik $r_{xy} = 0$, hau da, koerlazio-koefiziente lineala zero izatean koerlazio gabeak direla esango dugu.

Dakusagun orain bada, koerlazio gabeko aldagaien bi propietateak.

1. X eta Y aldagaiak koerlazio gabeak direnean, beraien baturaren bariantza bariantzen batura da.

$$S_{xy} = 0 \quad S_{(x+y)}^2 = S_x^2 + S_y^2$$

Frogapena:

$$\begin{aligned} S_{(x+y)}^2 &= \frac{1}{N} \sum_i (x_i + y_i - \bar{x} - \bar{y})^2 = \frac{1}{N} \sum_i [(x_i - \bar{x}) + (y_i - \bar{y})]^2 = \\ &= \frac{1}{N} \sum_i [(x_i - \bar{x})^2 + (y_i - \bar{y})^2 + 2(x_i - \bar{x})(y_i - \bar{y})] = \\ &= \frac{1}{N} \sum_i (x_i - \bar{x})^2 + \frac{1}{N} \sum_i (y_i - \bar{y})^2 + 2 \frac{1}{N} \sum_i (x_i - \bar{x})(y_i - \bar{y}) = \\ &= S_x^2 + S_y^2 + 2S_{xy} \end{aligned}$$

Hots:

$$S_{xy} = 0 \quad \Rightarrow \quad S_{(x+y)}^2 = S_x^2 + S_y^2$$

2. X aldagaia Y,Z.... aldagaiekin koerlazio gabea denean, beraien edozein konbinazio linealekin koerlazio gabea da.

$$\left. \begin{array}{l} S_{x,y} = 0 \\ S_{x,z} = 0 \end{array} \right\} \Rightarrow S_{x,ay+bz+c} = 0$$

Frogapena:

$$\begin{aligned} S_{x,ay+bz+c} &= \frac{1}{N} \sum_i [(x_i - \bar{x})(ay_i + bz_i + c - a\bar{y} - b\bar{z} - c)] = \\ &= \frac{1}{N} \sum_i (x_i - \bar{x}) [a(y_i - \bar{y}) + b(z_i - \bar{z})] = a \frac{1}{N} \sum_i (x_i - \bar{x})(y_i - \bar{y}) + \\ &+ b \frac{1}{N} \sum_i (x_i - \bar{x})(z_i - \bar{z}) = a S_{x,y} + b S_{x,z} = 0 \end{aligned}$$

IV.5. SAKABANATZE- ETA KOERLAZIO-MATRIZEAK

Momentu zentralak honela ere adierazten dira:

$$l_{11} = m_{20} = S_x^2, \quad l_{22} = m_{02} = S_y^2, \quad l_{12} = l_{21} = m_{11} = S_{xy}$$

L eta **R** letraz laburki adierazten ditugun ondoko matrizeak, sakabanatze-edo kobariantza matrizea eta koerlazio-matrizeak dira.

$$\mathbf{L} = \begin{bmatrix} l_{11} & l_{12} \\ l_{21} & l_{22} \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} 1 & r_{12} \\ r_{21} & 1 \end{bmatrix}$$

Ikusten denez, simetrikoak dira biak.

Lehenago adibidekoak hauek dira:

$$\mathbf{L} = \begin{bmatrix} 19'96 & 0'0025 \\ 0'0025 & 1'684 \cdot 10^{-5} \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} 1 & 0'1374 \\ 0'1374 & 1 \end{bmatrix}$$

IV.6. OROKORPENA

Suposa dezagun n ezaugarri edo aldagai ditugula. Binaka hartuz, kobariantzak eta koerlazio-koefizienteak atera daitezke; orduan, \mathbf{L} eta \mathbf{R} matrizeak hauexek izango lirateke:

$$\mathbf{L} = \begin{bmatrix} l_{11} & l_{12} & \cdots & l_{1n} \\ l_{21} & l_{22} & \cdots & l_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} 1 & r_{12} & \cdots & r_{1n} \\ r_{21} & 1 & \cdots & r_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ r_{n1} & r_{n2} & \cdots & 1 \end{bmatrix}$$

non $l_{11}, l_{22}, \dots, l_{nn}$, n aldagaien bariantzak diren hurrenez hurren; $l_{1n} = l_{n1}$, n eta 1 aldagaien arteko kobariantza da.

$r_{1n} = r_{n1}$, n eta 1 aldagaien arteko koerlazio-koefizientea da eta diagonaleko 1 balioa aldagai bakoitzak berekiko duen koerlazioa da.

V. KOERLAZIOA ETA ERREGRESIOA

V.1. SARRERA

V.2. ERREGRESIOA R^2 -n

V.2.1. Batezbestekoaren erregresioa

V.2.2. Karratu txikiaren erregresioa

V.3. KARRATU TXIKIENEN ERREGRESIO LINEALA R^2 -n

V.3.1. Karratu txikiaren erregresio zuzena

V.3.2. Karratu txikiaren erregresio linealaren propietateak

V.3.3. Hondar bariantza eta mugatze-koefizientea

V.3.4. Erregresio-koefizientearen eta koerlazio-koefizientearen zeinuaren azterketa.

Doikuntzaren egokitasuna

V.4. ERREGRESIO LINEALA R^n -n

V.4.1. Sarrera

V.4.2. Erregresio hiperplanoa

V.4.3. Propietateak

V.4.4. Erregresio hiperplanoaren koefizienteak lortzeko metodoa

V.4.5. β erregresio partzial estandarizatuen koefizienteak

V.4.6. Hondar bariantza. Mugatze-koefizientea eta koerlazio-koefiziente anizkoitza

V.4.7. Edozein aldagai azalduaren erregresioaren orokorpena

V.4.8. Erregresio-koefiziente desberdinen zeinuaren azterketa. Doikuntzaren egokitasuna

V.5. KOERLAZIO PARTZIALA

V.5.1. Koerlazio partziala R^3 -n

V.5.2. Koerlazio partziala R^n -n

V.5.2. Koerlazio partziala eta erregresio-koefizienteen arteko erlazioa

V.1. SARRERA

Bi edo aldagai gehiago aztertzen dituen analisi-estatistiko baten aurrean, horien elkar aldaketa oso garrantzitsua da, beraien arteko berdintasun edo desberdintasunak ez baitira kasualitatearen ondorioak izango, baizik eta gehienetan, beraien artean erlazio funtzionalaren ondorioak izango dira.

Elkar-aldaketa zehazteko ditugun teknikak, bi motatakoak dira:

1. Erregresio tekniken bidez, $Y = f(X_1, \dots, X_n, a, b, \dots)$ funtzio matematikoa lortuko dugu, zeinak, X_1, X_2, \dots, X_n aldagaien balioak ordezkatzuz gero, Y -ren, aldagai dependentearen, balio teorikoak emango baitizkigu.
2. Koerlazio-tekniken bidez: aldagaien arteko kobariantza azalduko dizkigun koefiziente batzu lortuko ditugu.

V.2. ERREGRESIOA R^2 -n

Maiztasun banaketa bikoitzean, aldagaiaren balio bakoitzari bestearen zenbait balio dagozkio. Hala ere, aldagai horietako bat independentetzat hartzen badugu eta erlazio matematikoa estimatzen badugu, berorren balio bakoitzari, aldagai dependentearen balio bakarra dagokio.

Funtzio matematiko hori optimizazio erizpide bati jarraituz lortuko dugu, dispersio-diagrama edo puntu-hodeia kontutan harturik.

Erregresio-metodo asko daude:

1. Batezbestekoaren erregresioa.
2. Karratu txikiaren erregresioa.
Erregresio lineala.
Erregresio parabolikoa.

Egin dugun erregresioaz, aldagai baten bidez bestea auresateko baliatuko gara, hauxe baita, azken finean, aldagaien erlazioen azterketen helburua.

Edozein metodori jarraituz, bi erregresio desberdin kontsidera daitezke: Y -rena X -ekiko eta X -ena Y -rekiko.

V.2.1. Batezbestekoaren erregresioa

Y-ren X-ekiko erregresioa kontsideratuz, batezbestekoaren erregresioak aldagaiaren x_i $i = 1, \dots, k$ balio bakoitzari, $(Y/X = x_i)$ banaketa baldintzatuaren \bar{y}/x_i batezbesteko balioa elkartzen dio. Kasu honetan, x_i balio desberdinentzat, non $i = 1, \dots, k$ den, Y aldagaiaren batezbesteko balio baldintzatuak izango dira:

$$\bar{y} / x_1 = \sum_{j=1}^l y_j / x_1 f(y_j / x_1)$$

$$\bar{y} / x_2 = \sum_{j=1}^l y_j / x_2 f(y_j / x_2)$$

$$\vdots = \quad \quad \quad \vdots$$

$$\bar{y} / x_k = \sum_{j=1}^l y_j / x_k f(y_j / x_k)$$

$(x_i, \bar{y}/x_i)$, $i = 1, \dots, k$ puntuak asko eta hurbilak balira, kurba bat osatuko lukete. Kurba hori, batezbestekoaren erregresioarena edo erregresio enpirikoarena deritzo, hain zuzen.

V.2.2. Karratu txikiaren erregresioa

(X eta Y) bi aldagai estatistikoen puntu-hodeiaren arabera, beraien artean $Y = f(X, a, b)$ moduko erlazio funtzional bat dagoela pentsa dezakegu, horrela, X aldagai independentea edo azaltzailea izango litzateke eta Y, berriz, aldagai dependentea edo azaldua. Puntu-hodeiaren itxuraren arabera, $f(X, a, b)$ funtzio bat edo beste hartuko dugu.

Karratu txikiaren zentzuan, karratu txikiaren erregresioaren arabera funtzio hoberena doituko dugu. Doikuntza hau, a eta b parametroak lortzean datza, non hauek puntu-hodeiari hobekien doitzen zaion kurbarenak diren. Hau da, hurrengo balioa minimo egiten duen kurba aukeratu nahi dugu.

$$\sum_{i=1}^m [y_j - f(x_i, a, b)]^2$$

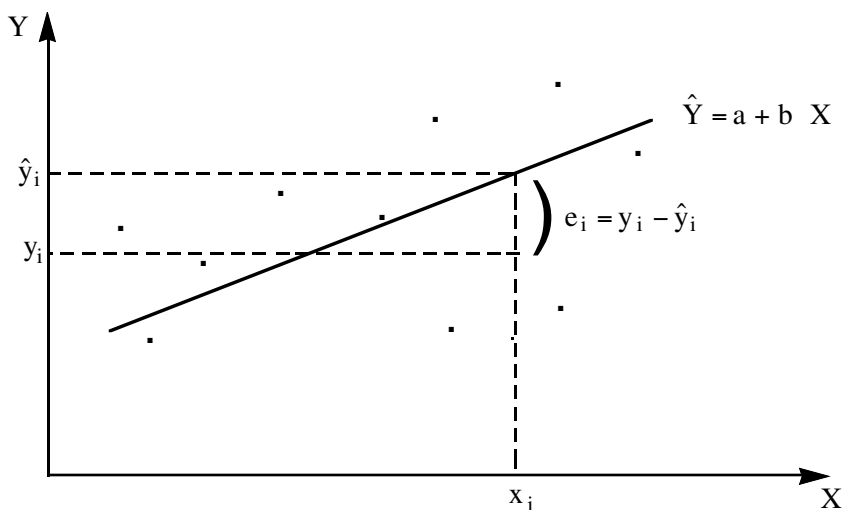
Kasu sinpleena funtzio lineala aukeratzen dugun kasua da, hau da, $f(X, a, b) = a + bX$.

V.3. KARRATU TXIKIENEN ERREGRESIO LINEALA R^2 -n

V.3.1. Karratu txikien erregresio zuzena

(x_i, y_i) , $i = 1, \dots, m$ puntu-hodeia emanik, eta $\hat{Y} = a + bX$ doikuntzaren funtzioa aukeraturik, x_i X-en balio bakoitzari, Y-ren bi balio dagozkie: puntu-hodeiko y_i balio erreala eta $X = x_i$ balioa aurreko zuzenean ordezkatuz lortzen dugun \hat{y}_i balio teorikoa edo estimatua.

1 definizioa. y_i balio erreala eta \hat{y}_i balio estimatuaren arteko diferentziari, errorea edo hondarra deritzo, eta e_i izendatuko dugu; hau da, $e_i = y_i - \hat{y}_i$, $i = 1, \dots, m$.



Karratu txikien eredua a eta b mugatzean datza, non erroreen karratuen batura minimo egiten duten.

Beraz, hurrengo problema askatu behar dugu:

$$\min_{a,b} \Phi(a,b) = \min_{a,b} \sum_{i=1}^m e_i^2 = \min_{a,b} \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \min_{a,b} \sum_{i=1}^m (y_i - a - bx_i)^2 \quad (5.1)$$

Φ funtzioak minimoa izateko beharrezko baldintza, a -rekiko eta b -rekiko deribatuak zero egitea da, hau da:

$$\frac{\partial \Phi}{\partial a} = 0 \Rightarrow ma + \left(\sum x_i \right) b = \sum y_i \quad (5.2)$$

$$\frac{\partial \Phi}{\partial b} = 0 \Rightarrow \left(\sum x_i \right) a + \left(\sum x_i^2 \right) b = \sum x_i y_i$$

non (5.2)-ko batukariak, eta baita aurrerantzean azalduko direnak, 1-tik m-rainokoak diren. (5.2) sistemari **ekuazio normalen sistema** deritzogu.

(5.2) sistema askatuz, ondorengo balioak lortzen ditugu:

$$a = \bar{y} - b\bar{x} \quad (5.3)$$

$$b = \frac{S_{xy}}{S_x^2}$$

non a, gai independentea den eta b, Y-ren X-ekiko erregresio zuzenaren X-en koefizientea, erregresio-koefizientea bezala ere ezagutzen dena.

(5.3)-n lortutako balioentzat Φ funtzioak minimoa izateko baldintza nahikoa betetzen du.

Ikus ezazu, (5.3)-n lortutako lehen ekuaziotik \bar{y} askatuz, erregresio zuzena grabitate zentrutik pasatzen dela.

(5.3)-n lortutako a-ren eta b-ren balioak, $Y = a + bX$ zuzenaren ekuazioan ordezkatzuz, normalean erabiltzen den ondoko adierazpena lortzen dugu:

$$\hat{Y} - \bar{y} = \frac{S_{xy}}{S_x^2} (X - \bar{x})$$

Ondoren daukagun \mathbf{X} matrizea

$$\begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{pmatrix}$$

eta $\mathbf{b} = (a, b)^T$; $\mathbf{y} = (y_1, \dots, y_m)^T$ bektorea kontutan izanik, orduan zera daukagu:

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} m & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \quad \text{eta} \quad \mathbf{X}^T \mathbf{y} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

eta horrela, (5.2) ekuazio normalen sistema matrizialki adieraz daiteke:

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y} \quad (5.4)$$

Baldin $\det(\mathbf{X}^T \mathbf{X}) \neq 0$, aurreko sistemaren ebazpidea existitzen da eta izango da:

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (5.5)$$

(5.5.) ebazpidea deribazio bektorialaren bitartez erraz lortzen da, ondoren ikusiko dugu bezala.

Izan bedi $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$ erroreen bektorea, non $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_m)^T$ estimatutako balioen bektorea den. Horrela, (5.1) problema era honetan adierazten da:

$$\min \mathbf{e}^T \mathbf{e} \quad (5.6)$$

Estimatutako balioen bektorea $\hat{\mathbf{y}} = \mathbf{X} \mathbf{b}$ denez, (5.6) hurrengo adierazpenaren baliokidea izango da.

$$\min_{\mathbf{b}} (\mathbf{y} - \mathbf{X} \mathbf{b})^T (\mathbf{y} - \mathbf{X} \mathbf{b}) = \min_{\mathbf{b}} (\mathbf{y}^T \mathbf{y} - 2 \mathbf{b}^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}) \quad (5.7)$$

eta \mathbf{b} -rekiko bektorialki deribatuz, (ikus A.7 eranskina), zera lortzen da:

$$\frac{\partial \mathbf{e}^T \mathbf{e}}{\partial \mathbf{b}} = -2 \mathbf{X}^T \mathbf{y} + 2 \mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{0} \quad (5.8)$$

nondik

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y} \quad (5.9)$$

lortzen den eta (5.4)-ren berdina den.

V.3.2. *Karratu txikiaren erregresio linealaren propietateak*

Aurreko atalean egindako garapenetatik hurrengo propietateak ondoriozta daitezke:

i) Erroreen batezbestekoa zero da: $\bar{e} = 0$

Frogapena: (5.4)-tik edo (5.8)-tik zera daukagu:

$$\mathbf{X}^T(\mathbf{y} - \mathbf{Xb}) = \mathbf{0} \quad (5.10)$$

non bi ekuazioko sistema den. Lehenengoak $\sum e_i = 0$ dela diosku, eta beraz, $\bar{e} = 0$.

ii) Estimaturako balioen batezbestekoa eta balio errealena berdinak dira:
 $\bar{\hat{y}} = \bar{y}$

Frogapena: i)-tik eta erroreen definiziotik berehala ondorioztatzen da.

iii) X aldagai azaltzailea eta e errore aldagaia elkarren artean koerlaziogabeak dira: $S_{xe} = 0$.

Frogapena: (5.10)-ko bigarren ekuaziotik, $\sum x_i e_i = 0$ lortzen dugu. Erroreek batezbestekoz zero dutenez, $S_{xe} = \frac{1}{m} \sum x_i e_i = 0$ da.

iv) \hat{Y} estimaturako aldagaia eta e errore aldagaia koerlaziogabeak dira:
 $S_{\hat{y}e} = 0$

Frogapena: $S_{\hat{y}e} = S_{(a+bx)e} = bS_{xe} = 0$

v) Bariantzaren deskonposaketa batukorra: $S_y^2 = S_{\hat{y}}^2 + S_e^2$

Frogapena: $Y = \hat{Y} + e$ dela kontutan izanik, eta (\hat{Y}, e) koerlaziogabeak direnez, baturaren bariantza bariantzen batura dela dakigu.

vi) $S_e^2 = (1 - r^2)S_y^2$

Frogapena (5.3)-ko b-ren adierazpenetatik eta $r = r_{xy} = \frac{S_{xy}}{S_x S_y}$ dela kontutan izanik, zera idatz dezakegu:

$$\begin{aligned}
 S_c^2 &= S_{(y-\hat{y})}^2 = S_{(y-a-bx)}^2 = S_{(y-bx)}^2 = S_y^2 + b^2 S_x^2 - 2b S_{xy} = \\
 &= S_y^2 + \frac{S_{xy}^2}{S_x^2 S_x^2} S_x^2 - 2 \frac{S_{xy}^2}{S_x^2} = S_y^2 - r^2 S_y^2 = (1 - r^2) S_y^2
 \end{aligned}
 \tag{5.11}$$

$$\text{vii) } S_{\hat{y}}^2 = r^2 S_y^2$$

Frogapena: **v)** eta **vi)** propietateetatik berehala ondorioztatzen da.

V.3.3. Hondar bariantza eta mugatze-koefizientea

Orain oso erraza zaigu koerlazio-koefizientearen karratua 0 eta 1 tartean dagoela frogatzea. **vi)** propietateetik eta bariantzak ez negatiboak direla kontutan izanik, $(1 - r^2)$ ez negatiboa dela ondorioztatu dezakegu, eta beraz, lehen azaldutakoa. Gainera, koerlazio-koefizientea -1 eta 1 tartean dagoela frogatzen da.

2 definizioa. Koerlazio-koefizienteren karratuari *mugatze-koefizientea* deritzo.

vii) propietatean ikusi dugunez, mugatze-koefizientea azaldutako bariantza totalarekiko portzentaia bezala ulertarazi daiteke.

Ikus ezazu, X-en Y-rekiko erregresioa egin nahi badugu, garapena lehenengokoaren berdina dela, baina X eta Y aldagaiaren betebeharrak elkar aldatuz, hortaz, orain $\hat{X} = a + bY$.

Bereziki, erregresio honen b koefiziente berria $b = \frac{S_{xy}}{S_y^2}$ izango litzateke eta normalean erabiltzen den zuzenaren adierazpena:

$$(\hat{X} - \bar{x}) = \frac{S_{xy}}{S_y^2} (Y - \bar{y})$$

Bi erregresioak elkarrekin lantzen baditugu, b_{xy} X-en Y-rekiko erregresioaren $(\hat{X}(Y))$ koefizientea izango da eta b_{yx} , berriz, Y-ren X-ekikoa $(\hat{Y}(X))$.

Erregresio linealaren koefizientea erregresio zuzenaren malda da. Orduan, b_{yx} -ek X-en aldaketekin Y-ren gehikuntza-tasa neurtzen digu. Hau da, b_{yx} -ek X aldagaia unitate batetan handitzen denean Y aldagaiaren aldakuntza adierazten digu.

Berehala froga dezakegu, mugatze-koefizientea bi erregresio-koefizienteen biderkadura dela. Hau da:

$$r^2 = b_{xy}b_{yx} \quad (5.12)$$

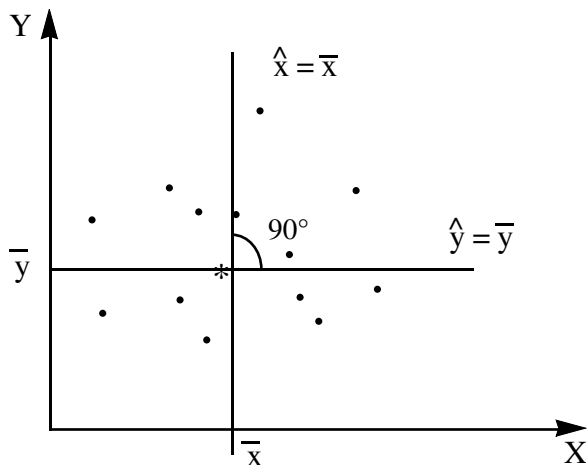
V.3.4. Erregresio-koefizientearen eta koerlazio-koefizientearen zeinuaren azterketa. Doikuntzaren egokitasuna

$$r_{xy} = r_{yx} = \frac{S_{xy}}{S_x S_y}, \quad b_{xy} = \frac{S_{xy}}{S_y^2}, \quad b_{yx} = \frac{S_{xy}}{S_x^2}$$

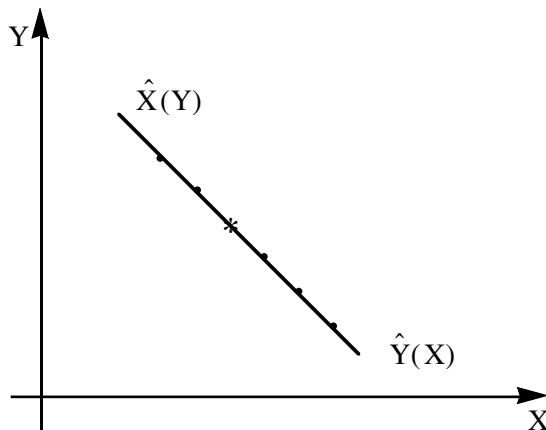
izanik eta bariantzak eta desbidazio tipikoak positiboak direnez, koefiziente hauek aldagaien arteko kobariantzaren zeinu berdina izango dute. Hau da, koerlazioa (eta beraz kobariantza) positiboa bada, bi zuzenen maldak positiboak lirerateke. Analogoki izango da kobariantza negatiboa denean.

Koerlaziogabe direnen kasuan ($S_{xy} = 0$), lortutako zuzenak $\hat{Y} = \bar{y}$ eta $\hat{X} = \bar{x}$ izango dira, eta beraz, grafikoki ikus daitezkeen bezala, ardatzen paraleloak izango dira eta beraien artean perpendikularrak. Koerlazio-koefizientea kontutan izanik, ondoren grafikoki aztertu ditzakegu:

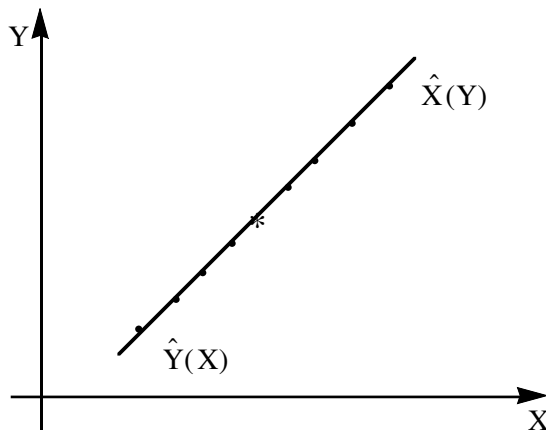
1. $r = 0$ ($S_{xy} = 0$, $S_y^2 = 0$). Kasu honetan, zuzenen arteko angelua 90° da. X aldagaiak ez du linealki Y azaltzen eta beraz, azaldutako bariantza zero da.



2. $r = -1$ ($S_{xy} = -S_x S_y$ eta $S_e^2 = 0$). Bi zuzenek bat egiten dute, (osatutako angelua 0° da) eta malda negatiboa dute.



3. $r = 1$ ($S_{xy} = S_x S_y$ eta $S_e^2 = 0$). Bi zuzenek bat egiten dute, (osatutako angelua 0° da) eta malda positiboa dute.



Ikus dezakegunez, 2. eta 3. kasuetan koerlazioa osoa da eta hortaz, hodeiko puntu guztiak zuzenaren gainean daude, $y_i = \hat{y}_i$, $i = 1, \dots, m$; kasu honetan doikuntza perfektua dela esango dugu.

Beraz, doikuntzaren egokitasuna r^2 -ren balioarekin zehazta dezakegu. Berau 1-etik hurbil badago, doikuntza ona izango da, baina zeretik hurbil badago doikuntza txarra dela esango dugu.

Adibidea:

X aldagaia lantegi batean egindako lan-orduak izanik eta Y aldagaia produzitu diren unitateak, kalkula dezagun koerlazio-koefizientea eta bi erregresio-zuzenak.

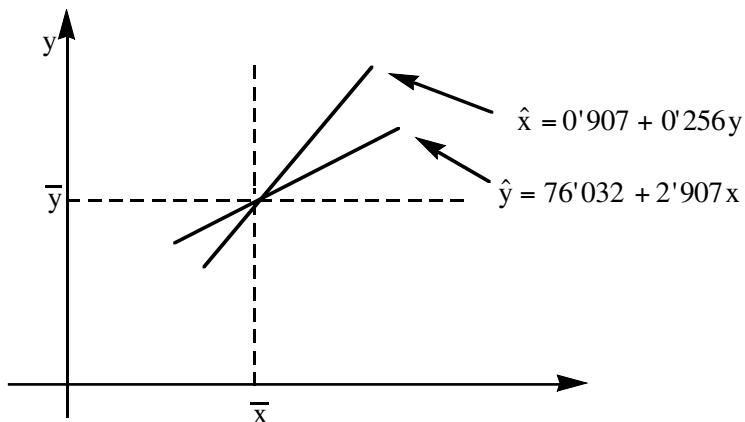
X	80	79	83	84	78	60	82	85	85	79	84
Y	300	302	315	330	300	250	300	340	315	330	310

Emaitzak hauek izan dira:

$$r = 0'863$$

$$\bar{x} = 79'9$$

$$\bar{y} = 308'36$$



Aurreko adibidean bi erregresio zuzenak egin baditugu ere, aldagaien arteko menpekotasuna alderantzizkoa da, hau da:

Bi galdera desberdin erantzuteko balio luketen zuzenak dira:

- 1) Zenbat lan-ordu beharko lirateke produzitutako unitateak kopuru jakina izateko?
- 2) Zein izango litzateke produzitutako unitateen kopurua egindako lan-orduak kopuru jakina bada?

V.4. ERREGRESIO LINEALA R^n -n

V.4.1. Sarrera

Atal honetan karratu txikiaren erregresio linealaren analisia n aldagaietara zabalduko dugu (X_1, \dots, X_n) .

Kasu honetan, m indibiduoentzako n dimentsiotako banaketa baten oharpenak ditugu, non x_{ij} -k i indibiduoaren ohartutako balioa j aldagaiarako adierazten duen.

Aldagai baten (azaldua) beste $n-1$ aldagaiekiko (azaltzaileak) karratu txikiaren erregresio lineala egin nahi dugu. Orokortasunik galdu gabe, X_1 aldagaia azaldutzat hartuko dugu eta

$\hat{X}_1(X_2, \dots, X_n)$ erregresio hiperplanoa doitu dugu:

$$\hat{X}_1 = a_1 + b_{12}X_2 + \dots + b_{1n}X_n \quad (5.13)$$

Ondorengo kalkuluak errazteko, R^2 -n erabilitako notazioa zabalduko dugu. Horrela, \mathbf{X} matrizea osatuko dugu:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{12} & \cdots & x_{1n} \\ 1 & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{m2} & \cdots & x_{mn} \end{pmatrix}$$

\mathbf{X} , beraz, aldagai independente bakoitzaren m obserbazioz osaturik dago, gehi 1-ez osatutako zutabe bat.

$$\mathbf{b} = (a_1, b_{12}, \dots, b_{1n})^T, \quad \mathbf{e} = \mathbf{x}_1 - \hat{\mathbf{x}}_1, \quad (5.14)$$

non $\mathbf{x}_1 = (x_{11}, \dots, x_{m1})^T$ X_1 aldagaiaren balio errealeko bektorea den eta $\hat{\mathbf{x}}_1 = (\hat{x}_{11}, \dots, \hat{x}_{m1})^T$, X_1 aldagaiaren estimatutako bektorea; beraz, $\hat{\mathbf{x}}_1 = \mathbf{X}\mathbf{b}$.

Hurrengo bektorea, $\bar{\mathbf{x}}$, batezbestekoen bektore bezala definituko dugu:

$$\bar{\mathbf{x}} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_n \end{pmatrix} \quad (5.15)$$

eta \mathbf{L} , (X_1, \dots, X_n) aldagaien arteko bariantza eta kobariantza matrizea izango da:

$$\mathbf{L} = \begin{pmatrix} S_1^2 & S_{12} & \cdots & S_{1n} \\ S_{12} & S_2^2 & \cdots & S_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ S_{1n} & S_{2n} & \cdots & S_n^2 \end{pmatrix} \quad (5.16)$$

non $S_{kj} = S_{x_k, x_j} = \text{Cov}(x_k, x_j)$ $k \neq j$, eta $S_j^2 = S_{x_j}^2 = \text{Var}(x_j)$, non $j = 1, \dots, n$.

Azkenik, \mathbf{R} koerlazio-matrizea izango da:

$$\mathbf{R} = \begin{pmatrix} 1 & r_{12} & \cdots & r_{1n} \\ r_{12} & 1 & \cdots & r_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n} & r_{2n} & \cdots & 1 \end{pmatrix} \quad (5.17)$$

non $r_{kj} = r_{x_k, x_j} = \text{Corr}(x_k, x_j) = \frac{S_{kj}}{S_k S_j}$ $k \neq j$.

\mathbf{L} eta \mathbf{R} matrize simetriko eta erdidefinitu positiboak dira, aurrerantzean erabiliko ditugun propietateekin (ikus V.A eranskina).

V.4.2. Erregresio hiperplanoa

(5.13) motako erregresio hiperplanoa doitzeko, karratu txikiaren ereduaren arabera, R^2 -n egindako minimizazio problema baliokidea planteatzen dugu matrizialki:

$$\min_{\mathbf{b}} \mathbf{e}^T \mathbf{e} = \min_{\mathbf{b}} (\mathbf{x}_1 - \mathbf{X}\mathbf{b})^T (\mathbf{x}_1 - \mathbf{X}\mathbf{b}) \quad (5.18)$$

non \mathbf{e} , \mathbf{x}_1 , \mathbf{b} eta \mathbf{X} aurretik definitutakoak diren.

(5.18)-n \mathbf{b} -rekiko bektorialki deribatuz eta zerori berdinduz, hurrengo ekuazio normalen sistema lortzen dugu:

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{x}_1 \quad (5.19)$$

eta $\det(\mathbf{X}^T \mathbf{X}) \neq 0$ bada, (5.18) askatzen duen \mathbf{b} parametroen bektorea izango da:

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{x}_1 \quad (5.20)$$

V.4.3. Propietateak

Orokorrean, n aldagaien kasuan, hurrengo erregresio propietateak ditugu:

i) Erroreak aldagai zentratuak dira: $\bar{e} = 0$.

Frogapena: Bi aldagaien kasuan bezala, n aldagaiekin (5.19)-tik zera daukagu:

$$\mathbf{X}^T (\mathbf{x}_1 - \mathbf{X} \mathbf{b}) = \mathbf{X}^T \mathbf{e} = \mathbf{0} \quad (5.21)$$

n ekuazio dituen sistema honen lehen ekuazioak $\sum e_i = 0$ dela diosku eta hemendik propietatea lortzen dugu.

ii) Balio estimatuen batezbestekoa eta balio errealeak berdinak dira:

$$\bar{\hat{x}}_1 = \bar{x}_1$$

Frogapena: **i)** -tik eta erroreen definiziotik ondoriozta daiteke.

iii) Erroreak eta edozein aldagai azaltzaile elkarren artean koerlaziogabeak dira: $S_{x_j e} = 0, \quad \forall j=2, \dots, n$.

Frogapena: (5.21)-ko beste $n-1$ ekuazioetatik ondorioztatzen da.

iv) Erroreak estimatutako balioekin koerlaziogabeak dira: $S_{\hat{x}_1 e} = 0$

Frogapena: $\hat{X}_1, X_2, \dots, X_n$ aldagaien konbinazio lineala dela jakinik, garbi ikus dezakegu.

v) Bariantzaren deskonposaketa batukorra: $S_{x_1}^2 = S_{\hat{x}_1}^2 + S_e^2$

Frogapena: $X_1 = \hat{X}_1 + e$ dela kontutan izanik eta (\hat{X}_1, e) aldagaiak koerlaziogabeak direnez, propietatea lortzen da.

V.4.4. Erregresio hiperplanoaren koefizienteak lortzeko metodoa

ii) propietateak zera daukagu:

$$a_1 = \bar{x}_1 - b_{12}\bar{x}_2 - \dots - b_{1n}\bar{x}_n \quad (5.22)$$

eta (5.14) hiperplanoa horrela adieraz dezakegu:

$$\hat{X}_1 - \bar{x}_1 = b_{12}(X_2 - \bar{x}_2) + \dots + b_{1n}(X_n - \bar{x}_n) \quad (5.23)$$

R^2 -n bezalaxe, lehendabizi hiperplanoa batezbestekoen bektoretik pasatzen dela ikusten da, hau da grabitate zentrotik. Bigarrenik, aldagaiak zentratuta daudela kontsideratzen badugu, erregresio-koefizienteak ez dira aldatzen. Orduan, orokortasunik galdu gabe, azken kasu honetarako koefizienteak kalkulatuko ditugu.

Aldagai zentratuentzat minimizazio prozesuaren ekuazio normalen sistema garatzen badugu, eta kasu honetan, honako matrize hauek

$$\mathbf{b} = (b_{12}, \dots, b_{1n})^T, \quad \mathbf{e} = \mathbf{x}_1 - \hat{\mathbf{x}}_1, \quad (5.24)$$

$$\mathbf{X} = \begin{pmatrix} x_{12} & \cdots & x_{1n} \\ x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \vdots \\ x_{m2} & \cdots & x_{mn} \end{pmatrix}$$

ditugarik, zera lortzen da:

$$\begin{aligned} b_{12} \sum x_{i2}^2 &+ b_{13} \sum x_{i2}x_{i3} + \cdots + b_{1n} \sum x_{i2}x_{in} = \sum x_{i1}x_{i2} \\ b_{12} \sum x_{i2}x_{i3} &+ b_{13} \sum x_{i3}^2 + \cdots + b_{1n} \sum x_{i3}x_{in} = \sum x_{i1}x_{i3} \\ \vdots & \quad \quad \quad \vdots \quad \quad \quad \vdots \\ b_{12} \sum x_{i2}x_{in} &+ b_{13} \sum x_{i3}x_{in} + \cdots + b_{1n} \sum x_{in}^2 = \sum x_{i1}x_{in} \end{aligned} \quad (5.25)$$

x_{ij} bakoitzak i indibiduoaren balioa j aldagai zentratutarako azaltzen duenez, (5.25) sistemako ekuazio bakoitza m -gatik zatitzen badugu, bariantza eta kobariantza terminotan adieraz dezakegu:

$$b_{1j} = -\frac{S_1 \mathbf{R}_{1j}}{S_j \mathbf{R}_{11}} \quad j = 2, \dots, n \quad (5.30)$$

Azkenik, koefizienteen balio guztiak ordezkatzuz lortzen dugun erregresio hiperplanoaren ekuazio orokorra ondokoa da:

$$\hat{X}_1 - \bar{x}_1 = -\frac{\mathbf{L}_{12}}{\mathbf{L}_{11}}(X_2 - \bar{x}_2) - \dots - \frac{\mathbf{L}_{1n}}{\mathbf{L}_{11}}(X_n - \bar{x}_n). \quad (5.31)$$

edota

$$\hat{X}_1 - \bar{x}_1 = -\frac{S_1 \mathbf{R}_{12}}{S_2 \mathbf{R}_{11}}(X_2 - \bar{x}_2) - \dots - \frac{S_1 \mathbf{R}_{1n}}{S_n \mathbf{R}_{11}}(X_n - \bar{x}_n). \quad (5.32)$$

V.4.5. β erregresio partzial estandarizatuen koefizienteak

Aurretik aldagaiak tipifikatu izan bagenitu, erregresio hiperplanoa horrela adieraziko litzateke:

$$\hat{T}_1 = \beta_{12} T_2 + \dots + \beta_{1n} T_n \quad (5.33)$$

non T_j , X_j hasierako aldagaiari dagokion aldagai tipifikatua den, eta β_{1j} notazioa, horrelako aldagaien erregresio-koefizienteentzako erabiliko dugun. Konturatu behar da, ez dagoela gai independenterik, aldagai tipifikatuak zentratuak baitaude. Gainera, aldagai hauek bariantzaz 1 dute, eta beraz, (5.30) erabiliz, erregresio-koefizienteak kasu honetan izango dira:

$$\beta_{1j} = -\frac{\mathbf{R}_{1j}}{\mathbf{R}_{11}} = \frac{b_{1j}}{S_1 / S_j} \quad j = 2, \dots, n. \quad (5.34)$$

3 definizioa. (5.34)-n lortutako koefizienteak *erregresio partzial standardizatuen koefizienteak* dira.

Oharra: R^2 -ren erregresioan, erregresio partzial standardizatuen koefizientea koerlazio-koefizientea $\beta_{12} = \beta_{21} = r_{12}$ da.

Koefiziente hauen garrantzia, adibide praktikoetan izango ditugun erantzunen interpretaziotan ikusiko da.

V.4.6. Hondar bariantza. Mugatze-koefizientea eta koerlazio-koefiziente anizkoitza

Erroreen batezbestekoa zero denez, bere bariantza matrizialki horrela adieraz dezakegu:

$$S_e^2 = \frac{1}{m} \mathbf{e}^T \mathbf{e} = \frac{1}{m} [\mathbf{x}_1 - \mathbf{Xb}]^T [\mathbf{x}_1 - \mathbf{Xb}]. \tag{5.35}$$

Aurreko biderkaketa garatuz eta \mathbf{b} , (5.19)-ko balioaz ordezkaturaz, (5.35) hurrengoaren baliokidea da:

$$S_e^2 = \frac{1}{m} [\mathbf{x}_1^T \mathbf{x}_1 - \mathbf{x}_1^T \mathbf{Xb}] \tag{5.36}$$

Ondorengo eragiketak errazteko, aldagaiak zentratuak daudela suposa dezakegu. Kasu orokorrari zabaltzea erraza izango litzateke, zeren badakigu aldagaien traslazioek hondar bariantzan ez dutela eragiten.

Horrela, datozen balio hauek (5.36)-ean ordezkutzen baditugu:

$$\begin{aligned} \mathbf{x}_1^T \mathbf{x}_1 &= mS_1^2 \\ \mathbf{x}_1^T \mathbf{X} &= (mS_{12} \dots mS_{1n}) \quad \text{eta} \\ \mathbf{b} &= \left(-\frac{L_{12}}{L_{11}} - \frac{L_{13}}{L_{11}} \dots - \frac{L_{1n}}{L_{11}} \right)^T \end{aligned} \tag{5.37}$$

$L_{11} \neq 0$ bada, hondar bariantzaren adierazpena lortzen dugu:

$$S_e^2 = S_1^2 \frac{L_{11}}{L_{11}} + S_{12} \frac{L_{12}}{L_{11}} + \dots + S_{1n} \frac{L_{1n}}{L_{11}} = \frac{|\mathbf{L}|}{L_{11}} = \frac{|\mathbf{R}|}{\mathbf{R}_{11}} S_1^2 \tag{5.38}$$

non $|\mathbf{L}|$ eta $|\mathbf{R}|$, \mathbf{L} eta \mathbf{R} matrizeen determinanteak diren.

1 oharra. \mathbf{L} eta \mathbf{R} matrize erdidefinitu positiboak direnez, $|\mathbf{L}|$, L_{11} eta $|\mathbf{R}|$, \mathbf{R}_{11} ezin dira negatiboak izan, eta hau hondar bariantzaren definizioarekin ados dago.

Mugatze-koefiziente anizkoitza

4 definizioa. X_1 -en X_2, \dots, X_n -rekiko erregresioaren $r_{1 \cdot 2, \dots, n}^2$ mugatze-koefiziente anizkoitza, azaldutako bariantzaren aldagai azalduaren bariantza totalarekiko proportzioa da.

Analitikoki,

$$r_{1 \cdot 2, \dots, n}^2 = \frac{S_{\hat{x}_1}^2}{S_{x_1}^2} = 1 - \frac{S_e^2}{S_{x_1}^2}, \quad (5.39)$$

Definizio honetatik, R^2 -n lortutako **vi** propietatearen orokorpena R^n -rako lortzen da.

$$\text{vi) } S_e^2 = (1 - r_{1 \cdot 2, \dots, n}^2) S_1^2,$$

eta beraz, $0 \leq r_{1 \cdot 2, \dots, n}^2 \leq 1$ ere betetzen da.

(5.38)-tik, (5.39) horrela adieraz daiteke:

$$r_{1 \cdot 2, \dots, n}^2 = 1 - \frac{|\mathbf{L}|}{\mathbf{L}_{11} S_1^2} = 1 - \frac{|\mathbf{R}|}{\mathbf{R}_{11}} \quad (5.40)$$

\mathbf{R} -ren determinantea bere lehenengo lerroko adjuntuetaz garatzen badugu, $|\mathbf{R}| = \mathbf{R}_{11} + r_{12} \mathbf{R}_{12} + \dots + r_{1n} \mathbf{R}_{1n}$, mugatze-koefizientea standardizatuak erregresio-koefizienteen eta koerlazio koefizienteen arteko biderkaduren batura izango da.

$$r_{1 \cdot 2, \dots, n}^2 = \beta_{12} r_{12} + \beta_{13} r_{13} + \dots + \beta_{1n} r_{1n} \quad (5.41)$$

5 definizioa. Mugatze-koefiziente anizkoitzaren erro karratu positiboari koerlazio-koefiziente anizkoitza deritzogu eta $r_{1 \cdot 2, \dots, n}$ bezala adierazten dugu.

Beraz,

$$r_{1 \cdot 2, \dots, n} = \sqrt{1 - \frac{S_e^2}{S_1^2}} = \sqrt{1 - \frac{|\mathbf{L}|}{\mathbf{L}_{11} S_1^2}} \quad (5.42)$$

Koerlazio-koefiziente anizkoitza, bi aldagaien arteko koerlazio simple edo totalaren arabera azaldu daiteke eta X_1 aldagai azalduaren eta beste aldagai azaltzaile guztien (X_2, \dots, X_n) arteko erlazioa azaltzen digu.

Bi aldagai genituenean, koerlazioak zuen zeinuaz hitzegiteak zentzua zuen, aldagaien kobariantza zuzena edo alderantzizkoa azaltzen baitzuen. Dimentsio haundiagoak ditugunean berriz, azaltzaile batek eragin positiboa izan dezake eta beste batzuk ordea, negatiboa, eta beraz, koerlazio-koefizientearen zeinuaren justifikaziorik ez dago.

Baina \hat{X}_1 aldagai azaltzaileen konbinazio lineala denez, $r_{1 \cdot 2 \dots n}$ X_1 eta \hat{X}_1 -ren arteko koerlazio totalatzat har dezakegu, eta horrela, zeinu positiboa baldin badu, aldagai baten eta bere estimatuaren arteko koerlazioa positiboa dela adieraziko du.

V.4.7. Edozein aldagai azalduaren erregresioaren orokorpena

Aldagai azaldua X_j , $j=1, \dots, n$, edozein aldagai denean, formulak aurrekoen berdinak dira, baina aldagai independentearen eta dependentearen betebeharrak era egokian aldatu behar dira. Horrela, orokorrean,

$$\hat{X}_j - \bar{x}_j = \sum_{k \neq j} b_{jk} (X_k - \bar{x}_k), \quad \text{non} \quad b_{jk} = -\frac{L_{jk}}{L_{jj}} \quad (5.43)$$

$$\hat{T}_j = \sum_{k \neq j} \beta_{jk} T_k, \quad \text{non} \quad \beta_{jk} = -\frac{R_{jk}}{R_{jj}} \quad \text{eta}, \quad (5.44)$$

$$r_{1 \cdot 2 \dots (j-1)(j+1) \dots n}^2 = 1 - \frac{|L|}{L_{jj} S_j^2} \quad (5.45)$$

Garbi dago, erregresioaren propietateak berridazten direla, aldagai azaltzaile desberdinei notazioa egokituz.

V.4.8. Erregresio-koefiziente desberdinen zeinuaren azterketa. Doikuntzaren egokitasuna

X_1 -ek azalduetako aldagaia izaten jarraituko du. Lehendabizi, (5.30) eta (5.34) parekatuz, b_{1j} eta β_{1j} -ren zeinuak, R_{1j} -ren zeinuaren arabera daudela ikus dezakegu, non biek zeinu berbera izango duten.

(5-23)-ko hiperplanoaren ekuazioan dakusagunez, x_j unitate batean gehitzen denean, beste aldagai azaltzaileak konstante mantenduz, X_1 aldagaiaren hazkuntza edo murriztapena b_{1j} baliokoa izango da. Partikularki, b_{1j} koefizienteen zeinuak bi aldagai hauen kobariantza, zuzena edo alderantzizkoa, erakusten du, besteak konstante mantenduz.

Aldagai tipifikatuekin lan eginez gero, β_{1j} koefizientea jatorrizko aldagaien neurri unitateekiko independentea izango da, baina bere zeinuaren esanahia b_{1j} -renaren berdina, honek aurretik azal dutako zeinuen berdintasuna justifikatzen duelarik. Bestalde, β koefiziente standardizatuen neurri unitateekiko independenteak izatean, beraien artean gonbaragarriak dira eta azaltzaileek azalduaren gain duten eragin erlatiboaren adierazle bezala erabiltzen dira, koefiziente standardizatuaren modulu haundiagoa duena eraginkorragoa delarik.

Bi aldagaiekin bezala, hiperplanoaren datuekiko doikuntzaren egokitasunaren analisia egin nahi dugu, eta mugatze-koefizientea adierazle egokia da. Lehen bezala, $r_{1 \cdot 2, \dots, n}^2$ -ren balioa 1 -etik hurbil badago, doikuntza oso ona izango da ($S_e^2 \approx 0$) eta 0 -tik hurbil badago, berriz, doikuntza txarra izango dugu ($S_e^2 \approx S_1^2$), 0 eta 1 har ditzakeen balioen muturrak izanik.

V.5. KOERLAZIO PARTZIALA

Har dezagun Jamaikako ron-aren prezioa eta apaizen alokairua bi aldagai estatistikoen elkar eboluzioa. Edozein arrazoi izanik, bi aldagaien arteko koerlazio-koefizientea kalkulatu nahi dugu, eta seguruenik balio positiboa eta altua izango da. Honek bien artean dependentzia dagoela esan nahi al du?

Inongo ekonomilariari ez litzaioke bururatuko Jamaikako ron-aren prezioa jeisteko, apaizen alokairua jeitsi behar dugunik!

Honek, aldagaiekin zerikusi zuzenik ez duten kausek eragiten dituzten koerlazioak eta koerlazio kausalak, hau da, aldagaien arteko dependentzia garbia azaltzen duten erlazioak-bereizi behar direla esan nahi du. Zentzu honetan, emandako adibidearen koerlazioak lehen motako koerlazioa azaltzen digu, hau da, beste hirugarren aldagai baten (edo gehiagoren) eraginagaitik eman dela, inflazioa adibidez. Aldagai honek, bere aldetik, beste bi aldagaiekin koerlazio kausal bat izan dezake. Hau da, inflazioa handitzen bada, prezioak eta alokairuak batera handitzea izan daiteke.

Beraz, atal honetan beste aldagaien eragina kontrolatzen duen koerlazio kontzeptu baten definizio bat eman nahi dugu.

Har dezagun beste adibide berri bat erreferentziatuz.

Amerikako zenbait herriren kriminaltasun-tasaren kausak aztertu nahi ditu ikertzaile batek. Horretarako, hurrengo lau aldagaien datuak hartu ditu:

- X_1 : koloreko populazioaren portzentaia.
- X_2 : 3000 dolar baino gutxiagoko sarrera duten populazioaren portzentaia.
- X_3 : herriaren tamainua.
- X_4 : kriminaltasun-tasa.

Datu horiekin ondorengo koerlazio-matrizea lortu da:

$$\mathbf{R} = \begin{pmatrix} 1 & 0.51 & 0.41 & 0.36 \\ & 1 & 0.29 & 0.60 \\ & & 1 & 0.49 \\ & & & 1 \end{pmatrix} \quad (5.46)$$

Aurreko matrizean (X_4) kriminaltasun-tasa eta (X_1) koloreko populazioaren portzentaiaren arteko koerlazioa, nahiz eta oso handia ez izan, positiboa dela ikusten da. (X_2) pobreak eta (X_3) herriaren masifikazioak X_4 -rekin ere koerlazio positiboa eta haundiagoa dute, eta era berean, aldagai horiek X_1 -ekin ere koerlazio positiboa dute.

Ikertzaileak kriminaltasun eta koloreko populazioaren arteko koerlazio totala benetakoa den edota X_2 eta X_3 bezalako beste aldagaien eraginagaitik den galdetu behar du. Horretarako, azken aldagai hauen eragina baztertu behar du eta ondoren, koerlazioen aldaketak aztertu. Aldagai batek beste bien arteko erlazioan duen eragina baztertzeko, baztertu nahi duguna aldaketarik ez lukeen ingurune batean kokatuz lortuko genuke, hau da, aldagai hori konstante egongo balitz. Ingurune hau sortzea ikerketa gehienetan posible ez denez koerlazio partziala zer den aztertuko dugu.

V.5.1. Koerlazio partziala R^{3-n}

6 definizioa. X_1 eta X_2 -ren arteko koerlazio partziala X_3 -ren eragina kenduta, $r_{12 \cdot 3}$ deritzogu. Koerlazioa, X_1 -en X_3 aldagaiarekiko eta X_2 -ren X_3 aldagaiarekiko erregresioetan sortzen diren erroreen arteko koerlazio totala da.

Hau da,

$$r_{12\cdot3} = \text{Corr}(e_{1\cdot3}, e_{2\cdot3}) = \frac{\text{Cov}(e_{1\cdot3}, e_{2\cdot3})}{\sqrt{\text{Var}(e_{1\cdot3})}\sqrt{\text{Var}(e_{2\cdot3})}} \quad (5.47)$$

non

$$e_{1\cdot3} = X_1 - \hat{X}_1(X_3) = (X_1 - \bar{x}_1) - \frac{s_{13}}{s_3^2}(X_3 - \bar{x}_3) \quad (5.48)$$

$$e_{2\cdot3} = X_2 - \hat{X}_2(X_3) = (X_2 - \bar{x}_2) - \frac{s_{23}}{s_3^2}(X_3 - \bar{x}_3)$$

Beraz, horrela definituriko koerlazio-koefiziente partzialak X_1 eta X_2 -ren arteko koerlazio linealaren maila neurtzen du, baina bietan X_3 -ren eragin lineala kendu ondoren.

(5.47) definizioa, r_{12} , r_{13} , r_{23} koerlazio totalen arabera azal dezakegu. Horretarako, aldagaien transformazio lineal zuzena egiten badugu, koerlazio-koefizientea ez dela aldatzen gogoratu behar dugu. Bereziki, aldagaiak tipifikaturik daudela pentsa dezakegu. Kasu honetan, (5.48)-ko adierazpenak hurrengoak izango dira:

$$\begin{aligned} e_{1\cdot3} &= T_1 - r_{13}T_3 \\ e_{2\cdot3} &= T_2 - r_{23}T_3 \end{aligned} \quad (5.49)$$

Horrela,

$$\begin{aligned} \text{Cov}(e_{1\cdot3}, e_{2\cdot3}) &= \text{Cov}(t_1 - r_{13}t_3, t_2 - r_{23}t_3) = \\ &= r_{12} - r_{13}r_{23} - r_{23}r_{13} + r_{13}r_{23} = r_{12} - r_{13}r_{23} \end{aligned} \quad (5.50)$$

$$\text{Var}(e_{1\cdot3}) = \text{Var}(t_1 - r_{13}t_3) = 1 - r_{13}^2 \quad (5.51)$$

$$\text{Var}(e_{2\cdot3}) = \text{Var}(t_2 - r_{23}t_3) = 1 - r_{23}^2 \quad (5.52)$$

Azkenik, (5.50), (5.51) eta (5.52) ordezkatuz (5.47)-an, hauxe lortzen dugu:

$$r_{12\cdot3} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{1 - r_{13}^2}\sqrt{1 - r_{23}^2}} \quad (5.53)$$

(5.53) erabiliz, eta gure adibideko aldagaiei aplikatuz, X_1 eta X_4 aldagaien arteko koerlazioa azter dezakegu, beste aldagaiek bakoitzean duten eragina baztertuz.

$$r_{14 \cdot 2} = \frac{r_{14} - r_{12}r_{24}}{\sqrt{1 - r_{12}^2} \sqrt{1 - r_{24}^2}} = \frac{0.36 - 0.60 \times 0.51}{\sqrt{1 - (0.60)^2} \sqrt{1 - (0.51)^2}} = 0.078$$

$$r_{14 \cdot 3} = \frac{r_{14} - r_{13}r_{34}}{\sqrt{1 - r_{13}^2} \sqrt{1 - r_{34}^2}} = \frac{0.36 - 0.41 \times 0.49}{\sqrt{1 - (0.41)^2} \sqrt{1 - (0.49)^2}} = 0.2$$

Erantzun hauekin, ikertzaileak $r_{14} = 0.36$ eta $r_{14 \cdot 2} = 0.078$ -ren arteko diferentzia aztertzen du, eta koloreko populazioaren eta kriminaltasun-tasaren arteko koerlazio totalaren zati handiena sarrera txikiko populazioaren portzentaiak bi aldagaietan duen eraginagaitik dela ondoriozta dezake.

Analisi honekin jarraituz, ikertzaileak bi aldagaien eragina batera atean kentzea pentsatzen du, pobrezia eta masifikazioa.

V.5.2. Koerlazio partziala R^n -n

7 definizioa. X_1 eta X_2 aldagaien arteko koerlazio partzialari, X_3, \dots, X_n aldagaien eragina baztertuz, $r_{12 \cdot 3 \dots n}$ deritzogu, eta $r_{12 \cdot 3 \dots n}$ koerlazioak X_1 -en eta X_2 -ren X_3, \dots, X_n -rekiko erregresioetan sortzen diren hondarren arteko koerlazioa da.

Hau da:

$$r_{12 \cdot 3 \dots n} = \text{Corr}(e_{1 \cdot 3 \dots n}, e_{2 \cdot 3 \dots n}) \quad (5.54)$$

X_1 eta X_2 aldagaien arteko koerlazio ohartuan, X_3 eta X_4 aldagaien eragina baztertzeke erabiltzen edo jarraitzen den prozesua, hau da: lehen horietatik batena baztertzea, eta ondoren lehenengo erregresioan lortutako erroreen bestearen eragina baztertzea baliokidea da. Hau errepikapen eredu da.

Logikoa denez, prozesu honen emaitza honakoa da:

$$r_{12 \cdot 34} = \frac{r_{12 \cdot 3} - r_{14 \cdot 3}r_{24 \cdot 3}}{\sqrt{1 - r_{14 \cdot 3}^2} \sqrt{1 - r_{24 \cdot 3}^2}} \quad (5.55)$$

Noski, emaitza berdina izango litzateke, lehendabizi X_4 -ren eragina

baztertuko bagenu; hau da:

$$r_{12 \cdot 34} = \frac{r_{12 \cdot 4} - r_{13 \cdot 4} r_{23 \cdot 4}}{\sqrt{1 - r_{13 \cdot 4}^2} \sqrt{1 - r_{23 \cdot 4}^2}} \quad (5.56)$$

Errepikapen eredu hau, erraza izan arren, aldagaien kopurua handitzen denean lankorra bihurtu daiteke.

Horregatik, \mathbf{L} edo \mathbf{R} -ren adjuntuen bitartez, baliokidea den matrizeen adierazpen bat erabiltzen da.

Orokorrean, $r_{12 \cdot 3 \dots n}$, X_1 eta X_2 aldagaien arteko koerlazioa bada, beste $n-2$ aldagaien eragina kenduta, zera daukagu:

$$r_{12 \cdot 3 \dots n} = -\frac{\mathbf{L}_{12}}{\mathbf{L}_{11}^{1/2} \mathbf{L}_{22}^{1/2}} = -\frac{\mathbf{R}_{12}}{\mathbf{R}_{11}^{1/2} \mathbf{R}_{22}^{1/2}} \quad (5.57)$$

Adibidearekin jarraituz,

$$r_{14 \cdot 23} = -\frac{\mathbf{R}_{14}}{\mathbf{R}_{11}^{1/2} \mathbf{R}_{44}^{1/2}} = -0.061$$

non \mathbf{R} matrizea, X_1 , X_2 , X_3 eta X_4 aldagaientzat (5.46)-ean azaldutako koerlazio-matrizea den.

Horrela, koloreko populazioaren portzentaia eta kriminaltasun-tasaren arteko koerlazio simple positiboa, aldagaien arteko zuzeneko elkar dependentziagaitik ez dela ikusten dugu, baizik eta masifikazio edota pobrezia bi aldagaietan duten eragin zuzenagaitik izan daiteke.

Erantzun hauek direla eta, zera pentsa dezakegu: beharbada pobrezia eta masifikazio maila berdina duten populazioekin beste ikerketa paralelo bat egingo bagenu, ez genuke aurkituko koloreko populazioa eta kriminaltasun-tasaren arteko koerlazio esanguratsurik. "Beharbada" esaten dugu, lortutako emaitzak ikertutako kolektibora mugatzen baitira, edonolako estrapolaziorik egin gabe.

2 oharra. Notazioak azaltzen ez duen arren, kontutan hartu behar da \mathbf{L} edo \mathbf{R} matrizeek, ikertzen den kasuari egokiak zaizkion dimentsio eta aldagaiak izan behar dituztela.

Horrela, formula berbera edozein aldagai kopurarekin erabil dezakegu.

Adibidez, X_1 -en X_4 -ren arteko koerlazio partziala, X_2 -ren eragina kenduta kalkulatu nahi badugu, X_1 , X_2 eta X_4 aldagaientzat definitutako \mathbf{R} matrizea kontsidera dezagun. Horrela:

$$\mathbf{R} = \mathbf{R}(X_1, X_2, X_4) = \begin{pmatrix} 1 & 0.51 & 0.36 \\ & 1 & 0.60 \\ & & 1 \end{pmatrix}$$

Eta aldagaien ordenaketa berria kontutan izanik, kalkulatu nahi dugun koerlazioa izango da:

$$r_{13 \cdot 2} = -\frac{\mathbf{R}_{13}}{\mathbf{R}_{11}^{1/2} \mathbf{R}_{33}^{1/2}} = 0.078$$

koerlazio-matrize berrian laugarren aldagaia hirugarrena izango baita.

Ikus daitekeenez, formula matrizialak, errepikapen formularekin batera etorri behar du; frogatzeko, adjuntoak garatu besterik ez da egin behar hiru aldagaien kasuan erraz frogatu ahal izateko. Kasu orokorra frogatu gabe geratzen da, kurtsoko asmoetatik kanpo baitago.

V.5.3. Koerlazio partzial eta erregresio-koefizienteen arteko erlazioa

Definizioan, $r_{12 \cdot 3 \dots n}$ -ren zeinuak X_1 eta X_2 -ren arteko kobariantza (zuzena edo alderantzizkoa) adierazten du, beste $n-2$ aldagaien aldaketek eraginik ez dutenean; hau da, hauek konstanteak mantentzen direnean. Analitikoki gainera, zeinu hori b_{12} eta β_{12} koefizienteen X_1 -en beste guztiekiko erregresio-koefizienteen zeinuen berdina izango du, eta beraz, interpretazio berbera, zeinu hori \mathbf{L}_{12} adjuntuak mugatzen baitu.

Azkenik, (5.43), (5.44) eta (5.57) adierazpenetatik hurrengo erlazioak ondorioztatzen dira:

$$b_{12}b_{21} = \beta_{12}\beta_{21} = r_{12 \cdot 3 \dots n}^2 \quad (5.58)$$

non b_{21} eta β_{21} bigarren aldagaiaren besteekiko erregresioaren koefizienteak diren.

***V.A. Eranskina. ALDAGAI ESTATISTIKO n -KOITZEN
MATRIZE-ESTATISTIKOAK***

V.A.1. DATU-MATRIZEAK

*V.A.2. BATEZBESTEKO-BEKTOREA:
PROPIETATEAK*

V.A.3. KOBARIANTZA MATRIZEA

V.A.3.1 Kobariantza matrizearen propietateak

V.A.4. KOERLAZIO-MATRIZEA

V.A.4.1. Koerlazio-matrizearen propietateak

V.A.5. DERIBAZIO BEKTORIALA

V.A.1. DATU-MATRIZEAK

Azkeneko ikasgaietan bi aldagaien arteko erlazioak aztertzeke tresna batzuk eman ditugu, baina hemendik aurrera gure ikasketako objektuak, hiru, lau,..... n aldagai izango dira eta ikasketak matrizialki egitea komenigarria da.

Has gaitezen, bada, matrizeak definitzen: X_1, X_2, \dots, X_n n aldagai kontsideratuz $X = [X_1, X_2 \dots X_n]$ beraiek osatzen duten zutabe-bektorea da:

Hots:

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$$

Aldagai horiek m indibiduen gain neurteta konkretu batzuk edo balio ohartu batzuk hartzen dute, hauexek X datu-matrizea osatzen dutelarik.

Hau da:

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m11} & x_{m2} & \dots & x_{mn} \end{bmatrix} = [x_{ij}]$$

Datu-matrizeak, X alegia, eta aldagai bikoitz baten maiztasun-banaketaren matrizeak ez dute zerikusirik, biak matrizeak direla baizik; lehenengo kasuan, x_{ij} , elementua i. indibiduoari dagokion j. aldagaiarekiko balioa da eta bigarreanean, ordea, $n_{ij}(x_i, y_j)$ balio bikotearen maiztasun absolutua.

j-garren aldagaiaren batezbesteko aritmetikoa eta desbidazio tipikoa

$$\bar{x}_j = \frac{\sum_{i=1}^m x_{ij}}{m} \quad \text{eta} \quad S_j = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_{ij} - \bar{x}_j)^2} \quad \text{izanik,}$$

Zentratze eta berreskalatze eragiketak

m x n ordenako \mathbf{X} datu-matrizea maiz aldatzen da, zentratze, berreskalatze, tipifikatze edo bere zutabeen (aldagaien) normetan tipifikatze eragiketen bitartez.

Datu-matrize **zentratua**, \mathbf{X} izendatuko duguna, horrela definitzen:

$$\mathbf{X} = \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1n} - \bar{x}_n \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \dots & x_{2n} - \bar{x}_n \\ \vdots & \vdots & \vdots & \vdots \\ x_{m1} - \bar{x}_1 & x_{m2} - \bar{x}_2 & \dots & x_{mn} - \bar{x}_n \end{bmatrix} = \left[x_{ij} - \bar{x}_j \right]_{\substack{j=1,\dots,n \\ i=1,\dots,m}}$$

Datu-matrizea $\frac{1}{\sqrt{m}}$ -gatik biderkatuz, \mathbf{X} izendatuko duguna, **zentratua eta berreskalatua** izango da:

$$\mathbf{X} = \begin{bmatrix} \frac{x_{11} - \bar{x}_1}{\sqrt{m}} & \frac{x_{12} - \bar{x}_2}{\sqrt{m}} & \dots & \frac{x_{1n} - \bar{x}_n}{\sqrt{m}} \\ \frac{x_{21} - \bar{x}_1}{\sqrt{m}} & \frac{x_{22} - \bar{x}_2}{\sqrt{m}} & \dots & \frac{x_{2n} - \bar{x}_n}{\sqrt{m}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{x_{m1} - \bar{x}_1}{\sqrt{m}} & \frac{x_{m2} - \bar{x}_2}{\sqrt{m}} & \dots & \frac{x_{mn} - \bar{x}_n}{\sqrt{m}} \end{bmatrix} = \left[\frac{x_{ij} - \bar{x}_j}{\sqrt{m}} \right]_{\substack{j=1,\dots,n \\ i=1,\dots,m}}$$

Zentratze eragiketak ez du datuen sakabanatze neurri simple edo konposatuetan eragiten, lekualdatze bat besterik ez baita. Horregatik, bi matrizeentzat notazio berbera erabiliko dugu. Hala ere, zentratutako eta \sqrt{m} -z berreskalatutako datu-matrizea, batzuetan erabiltzen da eta berreskalatzeak eragina izango du sakabanatzearen estatistikoetan.

Oharra: \mathbf{X} deitura mantentzea notazioan gehiegikeria da, baina ez du eragozpenik, kasu bakoitzean bere gai orokorra zehaztuko baita.

V.A.2. BATEZBESTEKO-BEKTOREA: PROPIETATEAK

\bar{x} -ikurraz adieraziko dugu eta aldagaien batezbestekoen zutabe-bektorea da.

$$\bar{\mathbf{x}} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_n \end{bmatrix} = [\bar{x}_j]$$

Aldagai bakun baten batezbestekoaren linealtasun-propietateak, batezbesteko bektorerako zabaltzen dira berehala.

1. Propietatea

$$\left. \begin{aligned} \mathbf{Y} &= \mathbf{A}\mathbf{X} + \mathbf{b} \\ \mathbf{Y} &= \mathbf{X} \mathbf{A}^T + \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \mathbf{b}^T \end{aligned} \right\} \bar{\mathbf{y}} = \mathbf{A}\bar{\mathbf{x}} + \mathbf{b}$$

non: \mathbf{A} matrizea $\mathbb{R}^n \rightarrow \mathbb{R}^n$ aplikazio lineal baten matrizea eta \mathbf{b} lekualdatze- edo traslazio-bektorea diren.

Honen ondorioz zera esan dezakegu:

Aldagaien jatorri-aldatzeak eta unitate- edo eskala-aldatzeak batezbestekoetan daukaten eragina, zera da: beraiek ere linealki aldatzea. Hau da, batezbestekoetan transformazio lineala daukagu.

$$\mathbf{Z} = \mathbf{X} + \mathbf{Y} \Rightarrow \bar{z} = \bar{x} + \bar{y}$$

Aldagaien bektore bat beste bi aldagaien bektoreen batura denean baturaren batezbesteko-bektorea bi batugaien batezbesteko-bektoreen batura da.

V.A.3. KOBARIANTZA MATRIZEA

$(X_1, X_2 \dots X_n)$ aldagai n -koitz baten bigarren ordenako momentu zentratuak bi eratakoak izan daitezke.

Aldagai bakun baten bazter-momentuak, adibidez:

$$m_{002..0} = \frac{1}{m} \sum_{i=1}^m (x_{i3} - \bar{x}_3)^2 = S_{x_3}^2$$

Aldagai bikoitz baten bazter-momentuak, adibidez:

$$m_{101.-0} = \frac{1}{m} \sum_{i=1}^m (x_{i1} - \bar{x}_1)(x_{i3} - \bar{x}_3) = S_{x_1 x_3}$$

Hots bariantzak edo kobariantzak dira.

Aldagai baten bariantza aldagaiak berekiko duen kobariantza bezala uler dezakegu, eta definizioz $S_{ij} = S_{ji}$ da.

Hau guztia kontutan hartuz aldagai n -koitzaren 2. ordenako momentu zentratu guztiek matrize simetriko bat osatzen dute, kobariantza matrizea hain zuzen.

\mathbf{L} ikurraz sinbolikoki adierazten dugu eta matrizialki honela lortzen da:

$$\mathbf{L} = \begin{bmatrix} S_1^2 & S_{12} & \dots & S_{1n} \\ S_{12} & S_2^2 & \dots & S_{2n} \\ \dots & \dots & \dots & \dots \\ S_{1n} & S_{2n} & \dots & S_n^2 \end{bmatrix} = \begin{bmatrix} l_{11} & l_{12} & \dots & l_{1n} \\ l_{12} & l_{22} & \dots & l_{2n} \\ \dots & \dots & \dots & \dots \\ l_{1n} & l_{2n} & \dots & l_{nn} \end{bmatrix}$$

$\mathbf{X} = [x_{ij} - \bar{x}_j]_{i,j}$, datu-matrize zentratua izanik, X_1, \dots, X_n aldagaien hurrengo bariantza eta kobariantza matrizea lortzen da:

$$\mathbf{L} = \frac{1}{m} \mathbf{X}^T \mathbf{X}$$

$\mathbf{X} = \left[\frac{x_{ij} - \bar{x}_j}{\sqrt{m}} \right]_{i,j}$, zentratu eta berreskalatutako datu-matrizea izanik, \mathbf{L} horrela lortzen da:

$$\mathbf{L} = \mathbf{X}^T \mathbf{X}$$

V.A.3.1. Kobariantza matrizearen propietateak

1. propietatea

$$\left. \begin{aligned} Y &= A X + b \\ Y &= X A^T + \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} b^T \end{aligned} \right\} \Rightarrow L(Y) = A L(X) A^T$$

Ikusten dugunez eskala aldatzeak badauka eragina kobariantza matrizean, baina ez jatorri aldatzeak.

Oharra: gogora dezagun $Y = aX + b$ izanik, bariantzen arteko erlazioa zera zela:

$$S_y^2 = a^2 S_x^2$$

2. propietatea:

Kobariantza matrizeak definitu positiboak edo erdidefinitu positiboak dira.

Matrize-algebran honek suposatzen du azpimatrize diagonal guztiak definitu positiboak edo erdidefinitu positiboak direla.

Datu-matrize tipifikatua

Tipifikatzeak edo norman tipifikatzeak aurretiko zentratzeaz gain, aldagaietan eskala aldatzea suposatzen du. Lehenengo kasuan, transformazioa egin ondoren, aldagaiek 1 baliodun desbidazio tipikoa dute, eta bigarrenean, m dimentsioko bektoreek osatzen dituzte, norma edo moduloa 1 dutelarik.

Horrela hauxe daukagu:

$$T = \begin{bmatrix} \frac{x_{11} - \bar{x}_1}{s_1} & \frac{x_{12} - \bar{x}_2}{s_2} & \dots & \frac{x_{1n} - \bar{x}_n}{s_n} \\ \frac{x_{21} - \bar{x}_1}{s_1} & \frac{x_{22} - \bar{x}_2}{s_2} & \dots & \frac{x_{2n} - \bar{x}_n}{s_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{x_{m1} - \bar{x}_1}{s_1} & \frac{x_{m2} - \bar{x}_2}{s_2} & \dots & \frac{x_{mn} - \bar{x}_n}{s_n} \end{bmatrix} = \left[\frac{x_{ij} - \bar{x}_j}{s_j} \right]_{\substack{j=1,\dots,n \\ i=1,\dots,m}} = [t_{ij}]_{i,j}$$

non $S_{ij}^2 = S_{tj} = 1$ den, edota:

$$\mathbf{N} = \begin{bmatrix} \frac{x_{11} - \bar{x}_1}{s_1 \sqrt{m}} & \frac{x_{12} - \bar{x}_2}{s_2 \sqrt{m}} & \dots & \frac{x_{1n} - \bar{x}_n}{s_n \sqrt{m}} \\ \frac{x_{21} - \bar{x}_1}{s_1 \sqrt{m}} & \frac{x_{22} - \bar{x}_2}{s_2 \sqrt{m}} & \dots & \frac{x_{2n} - \bar{x}_n}{s_n \sqrt{m}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{x_{m1} - \bar{x}_1}{s_1 \sqrt{m}} & \frac{x_{m2} - \bar{x}_2}{s_2 \sqrt{m}} & \dots & \frac{x_{mn} - \bar{x}_n}{s_n \sqrt{m}} \end{bmatrix} = \left[\frac{x_{ij} - \bar{x}_j}{S_j \sqrt{m}} \right]_{i=1, \dots, m}^{j=1, \dots, n} = [n_{ij}]_{i,j}$$

non $\|\mathbf{n}_j\| = 1$.

V.A.4. KOERLAZIO-MATRIZEA

Aldagai tipifikatuen arteko kobariantza matrizea da.
Hau da:

$$\mathbf{L}(T) = \frac{1}{m} \mathbf{T}^T \mathbf{T} = \frac{1}{m} \left[\frac{x_{ij} - \bar{x}_j}{S_j} \right]^T \left[\frac{x_{ij} - \bar{x}_j}{S_j} \right] = \mathbf{R}$$

$$\mathbf{R} = \begin{bmatrix} 1 & r_{12} & \dots & r_{1n} \\ r_{12} & 1 & \dots & r_{2n} \\ \dots & \dots & \dots & \dots \\ r_{1n} & r_{2n} & \dots & 1 \end{bmatrix}$$

$$\text{halaber: } \mathbf{R} = \mathbf{N}^T \mathbf{N} \quad \mathbf{N} = \left[\frac{x_{ij} - \bar{x}_j}{S_j \sqrt{m}} \right]$$

normatutako datu-matrizea izanik.

\mathbf{X} datu-matrizearen eta \mathbf{T} datu-matrize tipifikatuen arteko erlazioa hauxe da:

$$\mathbf{X} = \mathbf{T} \mathbf{D}_S^T + \begin{bmatrix} 1 \\ 1 \\ \cdot \\ \cdot \\ 1 \end{bmatrix} \bar{\mathbf{x}}^T$$

Halaber, X aldagaien bektorea eta T aldagai tipifikatuen artekoa:

$$X = D_S T + \bar{x}$$

Espresio honek adierazten duena zera da: transformazio lineal diagonal bat dugula (\mathbb{R}^n espazio bektorialean definitua), transformazio hau, D_S matrize diagonalak definitu baitu.

$$D_S = \begin{bmatrix} S_1 & 0 & \dots & 0 \\ 0 & S_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & S_n \end{bmatrix}$$

Dakusagun $L(X)$ eta $L(T)$ kobariantza matrizeen arteko erlazioa:

$$L(X) = D_S L(T) D_S^t$$

Hots:

$$L = D_S R D_S^t$$

Kobariantza eta koerlazio-matrizeen arteko erlazioa.

Estentsiboki:

$$L = \begin{bmatrix} S_1 & 0 & \dots & 0 \\ 0 & S_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & S_n \end{bmatrix} \begin{bmatrix} 1 & r_{12} & \dots & r_{1n} \\ r_{12} & 1 & \dots & r_{2n} \\ \dots & \dots & \dots & \dots \\ r_{1n} & r_{2n} & \dots & 1 \end{bmatrix} \begin{bmatrix} S_1 & 0 & \dots & 0 \\ 0 & S_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & S_n \end{bmatrix}^T$$

Dakusagun azkenik matrize hauen determinanteen eta adjuntuen arteko erlazioak.

L eta R matrizeen arteko elementuetan $S_{ij} = S_i S_j r_{ij}$ erlazioa betetzen dela ikusirik, ondoko erlazioak, erraz ateratzen ditugu:

$$\begin{array}{ll}
 |\mathbf{L}| = S_1^2 S_2^2 \dots S_n^2 |\mathbf{R}| & |\mathbf{L}| \text{ eta } |\mathbf{R}|, \mathbf{L} \text{ eta } \mathbf{R} \\
 & \text{matrizeen determinanteak izanik} \\
 \mathbf{L}_{11} = S_2^2 S_3^2 \dots S_n^2 \mathbf{R}_{11} & \mathbf{L}_{11} = l_{11} \text{ elementuaren adjuntua} \\
 & \mathbf{R}_{11} = r_{11} \text{ elementuaren adjuntua} \\
 \mathbf{L}_{12} = S_1 S_2 S_3^2 \dots S_n^2 \mathbf{R}_{12} & \mathbf{L}_{12} = l_{12} \text{ elementuaren adjuntua} \\
 & \mathbf{R}_{12} = r_{12} \text{ elementuaren adjuntua}
 \end{array}$$

V.A.4.1. Koerlazio-matrizearen propietateak

Gogoratu behar dugu (II.3.2.) koerlazioak, aldagai tipifikatuen arteko kobariantzak direla.

Honela, koerlazio-matrizea kobariantza matrize mota bat izatean, kobariantza matrizeen propietateak ditu.

Bestalde, koerlazio-matrizea, bi aldagaien arteko koerlazioa bezala (II.3) inbariantea da edozein transformazio linealen aurrean.

V.A.5. DERIBAZIO BEKTORIALA

Ondoren, eskema moduan, deribazio bektorialeko hurrengo formulak gogoratuko ditugu, R-ko beste hainbat funtzioekin erlazionatuz eta horrela bere ikasketa erraztuz.

$$\text{Karratua: } \frac{d(x^2)}{dx} = 2x \qquad \frac{\partial(\mathbf{x}^T \mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{x}$$

$$\text{Karratua: } \frac{\partial(xy)}{\partial x} = \frac{\partial(yx)}{\partial x} = y \qquad \frac{\partial(\mathbf{y}^T \mathbf{x})}{\partial \mathbf{x}} = \frac{\partial(\mathbf{x}^T \mathbf{y})}{\partial \mathbf{x}} = \mathbf{y}$$

$$\text{Biderkadura eskalar batez: } \frac{d(ax)}{dx} = a \qquad \frac{\partial(\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}^T$$

$$\text{Forma koadratikoa: } \frac{d(ax^2)}{dx} = 2ax \qquad \frac{\partial(\mathbf{x}^T \mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{A}\mathbf{x}, \text{ (asimetrikoa)}$$

$$\text{Eskalarra: } \frac{d(kf(x))}{dx} = k \frac{df(x)}{dx} \qquad \frac{\partial(kf(\mathbf{x}))}{\partial \mathbf{x}} = k \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$$

Deribazio matrizialaren formulen ikerketa sakonagoa egiteko, ikus adibidez, Martin Pliegoren *Introducción a la Estadística Económica y Empresarial*, (1994).

VI. ZENBAKI-INDIZEAK

VI.1. INDIZE SINPLEAK

VI.2. PONDERAZIO GABEKO INDIZE KONPLEXUAK

VI.2.1. Batezbesteko aritmetiko sinplearen metodoa

VI.2.2. Batezbesteko agregatu sinplearen metodoa

VI.3. INDIZE KONPLEXU PONDERATUAK

VI.3.1. Balioen, prezioen eta kopuruen indizeak

VI.3.1.1. LASPEYRES-en indizeak

VI.3.1.2. PAASCHE-en indizeak

VI.3.1.3. FISHER-en indizeak

VI.3.1.4. Propietate eta erlazio batzuk: bateragarritasuna eta alderantzizkotasuna

VI.3.1.5. Kalkulua

VI.4. INDIZE KONPLEXUEN ERAIKETAN SORTZEN DEN ZENBAIT ERAGOZPEN

VI.4.1. Aldagaien hautapena

VI.4.2. Somatutako leku eta denboraren hautapena

VI.4.3. Talde eta azpitaldeen hautapena

VI.4.4. Oinarri-denboraren hautapena

VI.4.5. Formula eta ponderazioen hautapena

VI.4.6. Indizearen esangura eta zabaldura

VI.5. ERAGOZPEN BEREZI BATZUK

VI.5.1. Oinarri-aldaketa indize sinpleetan

VI.5.2. Berriztapen eta loturak indize konplexuetan

VI.6. ZENBAKI-INDIZEEN APLIKAPENAK

VI.6.1. Kontsumo-prezioen indizeak

VI.6.2. Moneta-unitate arruntetan dauden magnitudeen deflazioa

VI.1. INDIZE SINPLEAK

Aldagai baten gora-beherakadak aztertzeko, **indize sinpleak** erabiltzen dira, eta aldagai batzuen gora-beherakadak “batera” aztertzeko **indize konplexuak**.

Hemen denboran zeharreko gora-beherakadak aztertuko dira. Erreferentzia gisa hartzen dugun denborari, **oinarria** deritzogu.

Erreferentziaren balioarekiko aldagaiaren balio bakoitzaren portzentaiak baino ez dira indize sinpleak.

Portzentaia hauen bidez, neurtzeko unitatea desagertzen da eta aldagai baten gora-beherakadak autonomoki azter daitezke. Honela indizeak aldagai-segiden gora-beherakadak errazten ditu (originalki unitate desberdinetakoak).

Izan bedi $x_0, x_1 \dots x_t$ aldagaiaren denborazko balio-segida bat; hots, 0, 1, t denboretan neurtutako balioak.

$t = 0$ oinarritzat harturik eta x_0 balio erreferentziala, indize sinpleak ondoko taulan agertzen diren bezala kalkulatu dira:

Denbora: t	Balio-segida: x_t	Indize Sinpleak
0	x_0	$I_0 = 100$
1	x_1	$I_1 = \frac{x_1}{x_0} \cdot 100$
2	x_2	$I_2 = \frac{x_2}{x_0} \cdot 100$
.	.	.
.	.	.
.	.	.
t	x_t	$I_t = \frac{x_t}{x_0} \cdot 100$

x_0 balio erreferentzial gisa “urte normal” bati dagokiona hartzea, oso garrantzitsua da.

Adibidez: Instintiboki uzta-bilketa bat ona edo txarra dela esatean, gureztat uzta-bilketa “normalak” direnekin konparaturik egiten dugu.

VI.2. PONDERAZIO GABEKO INDIZE KONPLEXUAK

n balioen aldakuntzak batera aztertu nahi ditugunean, indize konplexuak eraiki behar ditugu. Horiek n balio-segidak, balio-segida bakar batera laburtzen dituzte, non balioak, oinarritzat hartzen den urte batera erreferiturik dauden eta “konplexu” guztiaren mugimendua ikusarazten diguten.

Goazen, bada, ponderazio gabeko indize konplexuak kalkulatzeko bi metodo ikustera.

VI.2.1. Batezbesteko aritmetiko sinplearen metodoa

Indize konplexu hau, denbora bakoitzari dagozkion **indize sinpleen** batezbesteko aritmetikoa da.

Hots: t denboran x_1, x_2, \dots, x_n aldagaiei dagozkien $x_{1t}, x_{2t}, \dots, x_{it}, \dots, x_{nt}$ n balio-segidak baldin badira, non $t = 0, 1, 2, \dots, m$ den $I_{1t}, I_{2t}, \dots, I_{nt}$ indize sinpleak kalkulatu ondoren, indize konplexua formula honetaz kalkulaten den:

$$S_t = \frac{\sum_{i=1}^n I_{it}}{n} = \frac{\sum_{i=1}^n \frac{x_{it}}{x_{i0}}}{n} 100$$

Momentu batean, indize sinpleen elkarrekiko garrantzia ez hartzeagatik, ponderazio gabeko indizea da, orduan, **zenbat eta aldagai gehiagok parte hartu, orduan eta hobe gertatzen da metodo hau.**

Estatistika institutuan, urtez urteko aldagai kopuru handi batez Estatu espainolerako metodo honetaz edo metodo ponderatuen bidez kalkulaten diren indizeak, ez dira hain desberdinak.

VI.2.2. Batezbesteko agregatu sinplearen metodoa

Metodo honen bidez, une bakoitzerako aldagai guztien balioen batura egiten da eta lortzen den segidan indize sinpleak bakarrik kalkulaten dira.

Formula honetaz kalkulatzen dira:

$$B_t = \frac{\sum_{i=1}^n X_{it}}{\sum_{i=1}^n X_{i0}} = 100$$

Metodo hau, **ezin erabil daiteke aldagaiak unitate desberdinetan neurtuak direnean**, ez baitu zentzu askorik unitate desberdinetan neurtuak izan diren balioak batzea.

Ponderazio gabeko indize hau aurkez daiteke indize sinpleen indize ponderatu bezala ere, ponderazioak oinarri denborako balioak direlarik. (Ikus batezbesteko aritmetiko ponderatua 27. orrialdean)

$$B_t = \frac{\sum_{i=1}^n I_{it} X_{i0}}{\sum_{i=1}^n X_{i0}}$$

Ondoko taulan ponderazio gabeko indize konplexu hauen formulak dauzkagu hedaturik.

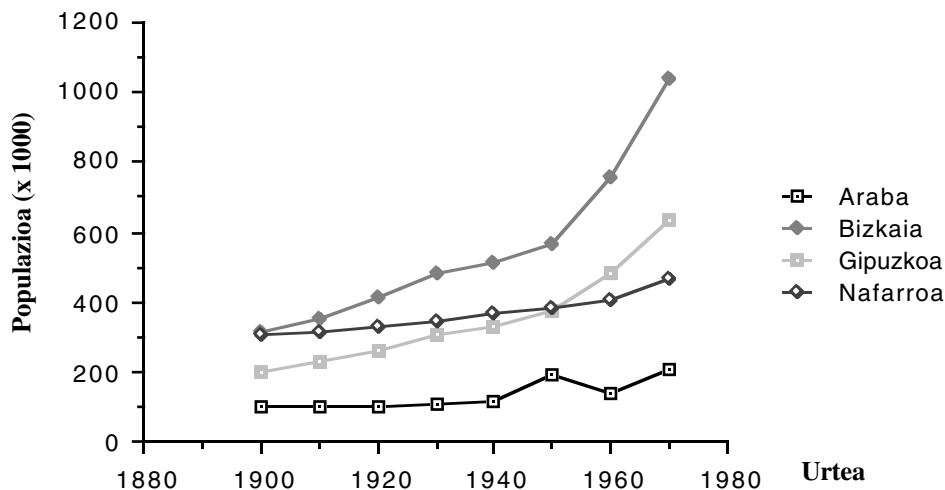
$\begin{matrix} x \\ t \end{matrix}$	$x_1, x_2, \dots, x_i, \dots, x_n$	I_i (Indize sinpleak)	S_t	$\sum_i x_{it} B_t$
0	$x_{10} x_{20} \dots x_{i0} \dots x_{n0}$	$I_{10} I_{20} \dots I_{n0}$	$S_0 = \frac{\sum_i I_{i0}}{n}$	$\sum_i x_{i0} B_0 = \frac{\sum_i x_{i0}}{\sum_i x_{i0}} 100$
1	$x_{11} x_{21} \dots x_{i1} \dots x_{n1}$	$I_{11} I_{21} \dots I_{n1}$	$S_1 = \frac{\sum_i I_{i1}}{n}$	$\sum_i x_{i1} B_1 = \frac{\sum_i x_{i1}}{\sum_i x_{i1}} 100$
⋮	⋮	⋮	⋮	⋮
t	$x_{1t} x_{2t} \dots x_{it} \dots x_{nt}$	$I_{1t} I_{2t} \dots I_{nt}$	$S_t = \frac{\sum_i I_{it}}{n}$	$\sum_i x_{it} B_t = \frac{\sum_i x_{it}}{\sum_i x_{it}} 100$
⋮	⋮	⋮	⋮	⋮
m	$x_{1m} x_{2m} \dots x_{im} \dots x_{nm}$	$I_{1m} I_{2m} \dots I_{nm}$	S_m	B_m

Adibidea: Hego Euskal Herriko populazioaren datu hauen bidez, populazio igoerak ikusiko ditugu herrialde bakoitzean.

Datuak:

Urtea	Araba	Bizkaia	Gipuzkoa	Nafarroa
t	x_1	x_2	x_3	x_4
1.900	96.385	311.361	195.850	307.669
1.910	97.181	349.932	226.684	312.235
1.920	98.668	409.550	258.557	329.875
1.930	104.176	485.205	302.329	345.883
(1) 1.940	112.876	511.135	331.753	369.618
1.950	188.012	569.188	374.040	382.932
1.960	138.934	754.383	478.337	402.042
1.970	204.323	1.043.311	631.003	464.867
1.980	258.200	1.179.105	692.586	507.367
1.990	272.101	1.153.622	675.529	527.318

1. Taula



1. Grafikoa

1. GRAFIKOAN IKUSTEN DIREN GAUZA NABARMENENAK

Grafiko honek, herrialde bakoitzeko **populazioaren hazkunde absolutua** adierazten digu.

Ondorio hauek atera genitzake:

1.900. urtean Nafarroako eta Bizkaiko populazioak berdintsuak ziren, Nafarroak XIX. mendean zuen populazioa askoz handiagoa izan arren.

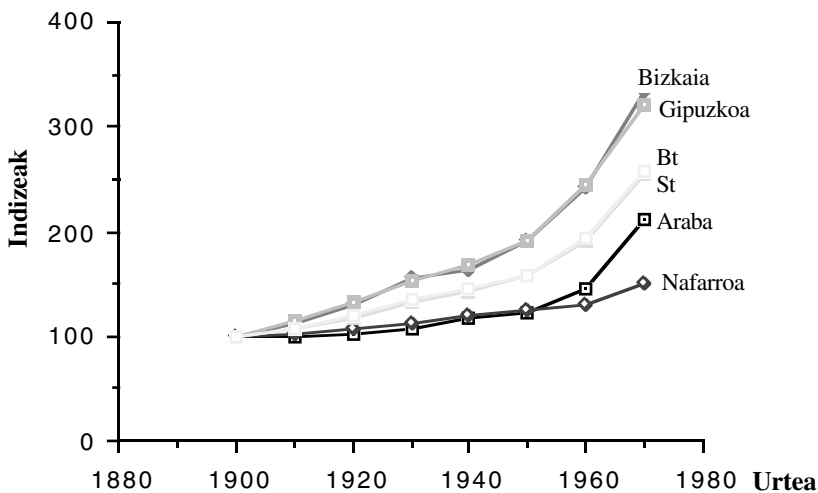
Hamarkada bakoitzean, ibilbidearen malda, gehitzen den biztanle-kopuruaren arabera da, (non gehikuntza hamarkadaren hasierako populazio absolutuarekiko baita).

Ikusten denez, Nafarroa eta Arabaren ibilbideak antzekoak dira alde batetik, eta bestetik, Bizkaia eta Gipuzkoarenak ere oso antzekoak dira.

Indizeak:

Urteak	I_{1t}	I_{2t}	I_{3t}	I_{4t}	S_t	$\sum_i x_{it}$	B_t
1.900	100'00	100'00	100'00	100'00	100'00	911'265	100'00
1.910	100'83	112'38	115'74	101'48	107'61	986'032	108'20
1.920	102'37	131'53	132'02	107'22	118'28	1.086'670	119'24
(2) 1.930	108'08	155'83	154'37	112'42	132'67	1.237'593	135'81
1.940	117'11	164'16	169'39	120'13	142'70	1.325'364	145'44
1.950	122'44	191'48	190'98	124'46	157'34	1.444'172	158'48
1.960	144'15	242'28	244'23	130'67	190'33	1.773'696	194'64
1.970	211'98	335'08	322'18	151'09	255'08	2.343'504	257'17

2. Taula



2. Grafikoa

2. GRAFIKOAN IKUSTEN DIREN GAUZA NABARMENENAK

Alde batetik grafiko honetan, herrialde bakoitzari dagokion populazio-indize sinpleen ibilbidea daukagu.

Indize sinple horien ibilbideek, **populazioaren gehikuntza erlatiboa, 1.900 urteko populazioarekiko**, adierazten digute, herrialde guztietako populazioa 1.900 urtean, 100 dela suposatuz.

1.900 urtetik zenbat eta gehiago urruntzen garen, orduan eta ibilbide horien maldek esangura gutxiago izango dute.

Beste batetik, Hego Euskal Herriko populazio-indize konplexu bien (S_t , B_t) ibilbideak ditugu (erdikoak).

Azkenik, hamarkada bakoitzeko urtez urteko hazkunde-tasak bilatu eta 3. grafikoan adierazten dira grafikoki.

Hots:

$$i_{t,t-1} = \frac{X_{i,t}}{X_{i,t-1}} = (1 + \alpha)$$

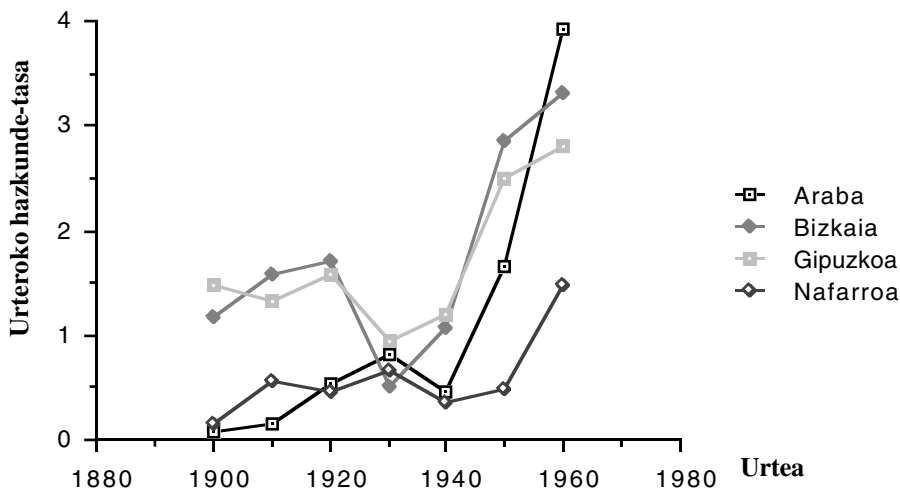
Hemen $i_{t,t-1}$ delakoa momentu baten aurrekoarekiko indize sinplea da. Hau da, $t-1$ uneko unitate bat, $1 + \alpha$ bihurtu da t unean, orduan hamarkada batean 0 uneko unitate bat, $(1 + \alpha)^{10}$ bihurtu da 10 urtetan, eta esenplu honetan datuak hamarkadaz ditugunez, **urteroko hazkunde-tasa** ondoko formulen bidez kalkulatzeko dugu:

$$\alpha + 1 = \sqrt[10]{\frac{X_{i,t}}{X_{i,t-10}}} \quad \text{eta} \quad \alpha = \sqrt[10]{\frac{X_{i,t}}{X_{i,t-10}}} - 1$$

Ondoko taulan, urteroko hazkunde-tasak adierazten dira:

$\alpha \times 100$	1.900/ 1.910	1.910/ 1.920	1.920/ 1.930	1.930/ 1.940	1.940/ 1.950	1.950/ 1.960	1.960/ 1.970
ARABA	0'08	0'15	0'54	0'81	0'45	1'65	3'93
BIZKAIA	1'17	1'59	1'71	0'52	1'08	2'86	3'30
GIPUZKOA	1'47	1'32	1'58	0'93	1'21	2'49	2'80
NAFARROA	0'15	0'55	0'47	0'67	0'35	0'49	1'47

3.Taula.



3. Grafikoa

3. GRAFIKOAN IKUSTEN DIREN GAUZA NABARMENENAK

Grafiko honetan, herrialde bakoitzari dagokion hazkunde-tasen ibilbidea daukagu.

Hazkunde-tasen ibilbideek **populazioaren gehikuntza erlatiboa adierazten** digute, kasu honetan gehikuntza **hamarkada bakoitzarekikoa** delarik.

Globalki 7 hamarkada horietan, Bizkaia eta Gipuzkoari dagozkien hazkunde-tasen ibilbideek dute azelerazio gehien.

Geldiera handia nabaritzen da lau herrialdeetan, 1.936 urte inguruan (gerra izan liteke arazo honen zergatia).

Lehenengo hamarkadan nahiko geldikorrak dira hazkunde-tasak bai Araban, bai Nafarroan; gero bietan hazkunde handia nabaritzen da, ordea, 1.930-1.940 hamarkadan, justu Bizkaia eta Gipuzkoan geldiera handia denean (migrazioa eta gerra izan litezke arazo honen zergatia?).

Azken hamarkadetan izugarritzko azelerazioa nabaritzen da lau herrialdeetan; kezkarriak Bizkaia eta Gipuzkoaren kasuak, azkenean Arabak, Bizkaia eta Gipuzkoaren saturazioa jasotzen duelarik.

Ariketa: Azken bi hamarkadako populazioaren datuak (1. taula) kontutan hartuz, lor itzazu dagozkien indizeak, hazkunde-tasak, eta adieraz itzazu 3 grafikoetan.

VI.3. INDIZE KONPLEXU PONDERATUAK

Orain arte, ponderazio gabeko indizeak aztertu ditugu. Baina, adibidez, arroza garia baino bi aldiz gehiago kontsumitzen bada, prezioen indize baten kalkuluan garrantzi handiagoa eman behar zaie arrozaren prezioaren aldaketei gariarenei baino.

Hau ponderaketaren bidez lortzen da.

Prezioen indize bat kalkulatzeko, ondasun bakoitza, kontsumitzen den kopuruaz ponderatzen bada, pautak batzuk jarraituz egingo dugu, eta artikulua edo ondasun baten zatitza aukeratuak izan diren pautak besterentzat ere errespetatu behar dira ahal bada.

Adibidez, hileroko janari prezioen indize batean, esnearen prezioa bere kantitateaz ponderatzen bada eta kopuru hau famili bati dagokiona baldin bada, (senitarte kopuru mugatu eta errenta-tarte baten barnekoa), orduan, ogiari dagokion ponderazioak ere irizpide berdinetan oinarritua izan behar du.

VI.3.1. Balioen prezioen eta kopuruaren indizeak

“p” delako ondasun baten prezioa bada eta “k” beraren kopurua (saldua, produktua...) orduan, “b” delako ondasunaren balioa dela esango dugu eta honela kalkulatu dugu:

$$b = p \cdot k$$

0, 1, ..., t, ..., m denboretan, ondasun baten prezioak

$p_0, p_1, \dots, p_t, \dots, p_m$ eta dagozkien kopuruak

$k_0, k_1, \dots, k_t, \dots, k_m$ baldin badira, orduan “denborazko balio-segida” hau izango da:

$$\begin{aligned}
 & \frac{b_t}{b_0} \\
 b_0 &= p_0 k_0 \\
 b_1 &= p_1 k_1 \\
 & \cdot \\
 & \cdot \\
 b_t &= p_t k_t \\
 & \cdot \\
 & \cdot \\
 b_m &= p_m k_m
 \end{aligned}$$

Baina denbora iragatean $k = kte$ bada, b_t denborazko balio-segidak, prezioaren gora-beherakadak adieraziko ditu bakarrik. Berdin, denbora iragatean $p = kte$ bada, b_t balio-segidak kopuruaren gora-beherakadak adieraziko ditu bakarrik.

Orduan, ondasun bati dagozkion hiru “denborazko balio-segida” desberdin idatz ditzakegu.

b_t	$b_t (k)$	$b_t (p)$
$p_0 k_0$	pk_0	$p_0 k$
$p_1 k_1$	pk_1	$p_1 k$
\cdot		
\cdot		
$p_t k_t$	pk_t	$p_t k$
\cdot		
\cdot		
$p_m k_m$	pk_m	$p_m k$

“O” ondasun bati dagozkion (balio, prezio eta kopuruaren) segida horiek balio errealetan adieraziak direnez gero, beste “ O_i ” ondasun batzuei dagozkienak ere batu ditzakegula, argi ikusten da.

Suposa dezagun $O_1 O_2 \dots\dots\dots O_i \dots\dots O_n$ n ondasun desberdin ditugula, (non O_i ondasun bakoitzaren balio, prezio eta kopuruaren denborazko balio-segidak ezagunak zaizkigun), n ondasun horiek, batera osatzen duten “konplexuaren indizeak” edo “indize konplexuak” kalkulatzekoan, balio-segida desberdinen batuketak, laburki, honela adieraz daitezke:

$$\begin{array}{ccc}
 \underline{\sum_{i=1}^n b_{it}} & \underline{\sum_{i=1}^n b_{it}(k)} & \underline{\sum_{i=1}^n b_{it}(p)} \\
 \underline{\sum_i p_{i0}k_{i0}} & \underline{\sum_i p_i k_{i0}} & \underline{\sum_i p_{i0}k_i} \\
 \underline{\sum_i p_{i1}k_{i1}} & \underline{\sum_i p_i k_{i1}} & \underline{\sum_i p_{i1}k_i} \\
 \vdots & \vdots & \vdots \\
 \underline{\sum_i p_{it}k_{it}} & \underline{\sum_i p_i k_{it}} & \underline{\sum_i p_{it}k_i} \\
 \vdots & \vdots & \vdots \\
 \underline{\sum_i p_{im}k_{im}} & \underline{\sum_i p_i k_{im}} & \underline{\sum_i p_{im}k_i}
 \end{array}$$

(Bi azkenetako segidetan, $p_i = k_{i0}$, \forall_i eta $k_i = k_{i0}$, \forall_i dira).

Beraz, **indize konplexuak** lortzeko konplexuari dagozkion (balio, prezio eta kopurua) hiru **denborazko balio segida** 1. horien **indize sinpleak kalkulatu** behar ditugu.

Hots oinarri-denbora $t = 0$ bada:

$$\begin{array}{ccc}
 \underline{I_t^b} & \underline{I_t^k} & \underline{I_t^p} \\
 \frac{\sum_i p_{i0}k_{i0}}{\sum_i p_{i0}k_{i0}} 100 & \frac{\sum_i (p_i)k_{i0}}{\sum_i (p_i)k_{i0}} 100 & \frac{\sum_i p_{i0}(k_i)}{\sum_i p_{i0}(k_i)} 100 \\
 \\
 \frac{\sum_i p_{i1}k_{i1}}{\sum_i p_{i0}k_{i0}} 100 & \frac{\sum_i (p_i)k_{i1}}{\sum_i (p_i)k_{i0}} 100 & \frac{\sum_i p_{i1}(k_i)}{\sum_i p_{i0}k_{i0}} 100 \\
 \vdots & \vdots & \vdots
 \end{array}$$

1. Lehenengoa prezio eta kopuruaren funtzioa da, bigarrena, kopuruaren funtzioa, eta hirugarrena prezioarena.

$$\begin{array}{ccc}
 \frac{\sum_i p_{it} k_{it}}{\sum_i p_{i0} k_{i0}} 100 & \frac{\sum_i (p_i) k_{it}}{\sum_i (p_i) k_{i0}} 100 & \frac{\sum_i p_{it} (k_i)}{\sum_i p_{i0} (k_i)} 100 \\
 \vdots & \vdots & \vdots \\
 \frac{\sum_i p_{im} k_{im}}{\sum_i p_{i0} k_{i0}} 100 & \frac{\sum_i (p_i) k_{im}}{\sum_i (p_i) k_{i0}} 100 & \frac{\sum_i p_{im} (k_i)}{\sum_i p_{i0} (k_i)} 100
 \end{array}$$

I_t^b I_t^k I_t^p aurreko indizeei, balio, prezio eta kopuruen **indize konplexu agregatuak** deritze, erabiltzen den agregazioa ponderatua bada ere.

Baina orain problema bat sortzen zaigu:

I_t^k “kopuruen indize konplexua” lortzeko, O_i ondasun bakoitzaren zein prezio hartuko dugu konstantetzat?

Edota:

I_t^p “prezioen indize konplexua” lortzeko, O_i ondasun bakoitzaren zein kopuru hartuko dugu konstantetzat?

Problema honetarako nahi beste ebazpide edo soluzio aurki daitezke, baina praktikan bi inposatu dira: LASPEYRES-ena eta PAASCHE-rena.

I_t^k eta I_t^p lortzeko, horiek urratutako bide desberdinak ikus ditzagun ondoren.

VI.3.1.1. LASPEYRES-en indizeak

I_t^k eta I_t^p prezio eta kopuruen indize konplexuak kalkulatzeko, O_i ondasun bakoitzaren kasuan konstantetzat hartzen ditugun kopuru eta prezioak, oinarri-denboran sortutakoak dira.

Hots, $k_i = k_{i0}$ eta $p_i = p_{i0}$

Honela lortutako indize konplexuei, **LASPEYRES-en prezio eta kopuruen indize konplexuak** deritzegu eta sinbolikoki L_t^p eta I_t^k adieraziko ditugu.

Hots:

$$\begin{array}{cc}
 \frac{L_t^p}{\sum_i p_{i0}(k_{i0})} & \frac{L_t^k}{\sum_i (p_{i0})k_{i0}} \\
 \frac{\sum_i p_{i0}(k_{i0})}{\sum_i p_{i0}(k_{i0})} 100 & \frac{\sum_i (p_{i0})k_{i0}}{\sum_i (p_{i0})k_{i0}} 100 \\
 \frac{\sum_i p_{i1}(k_{i0})}{\sum_i p_{i0}(k_{i0})} 100 & \frac{\sum_i (p_{i0})k_{i1}}{\sum_i (p_{i0})k_{i0}} 100 \\
 \vdots & \vdots \\
 \frac{\sum_i p_{it}(k_{i0})}{\sum_i p_{i0}(k_{i0})} 100 & \frac{\sum_i (p_{i0})k_{it}}{\sum_i (p_{i0})k_{i0}} 100 \\
 \vdots & \vdots \\
 \frac{\sum_i p_{im}(k_{i0})}{\sum_i p_{i0}(k_{i0})} 100 & \frac{\sum_i (p_{i0})k_{im}}{\sum_i (p_{i0})k_{i0}} 100
 \end{array} \quad (1)$$

Oinarrizko eragiketa batez, erraz ikus daiteke (1) formulak, ondoren datozen (2) formula horien baliokideak direla:

$$\begin{array}{cc}
 \frac{L_t^p}{100} & \frac{L_t^k}{100} \\
 \frac{\sum_i \frac{p_{i1}}{p_{i0}}(p_{i0} k_{i0})}{\sum_i (p_{i0} k_{i0})} 100 & \frac{\sum_i \frac{k_{i1}}{k_{i0}}(p_{i0} k_{i0})}{\sum_i (p_{i0} k_{i0})} 100 \\
 \vdots & \vdots
 \end{array}$$

$$\begin{array}{cc}
 \frac{\sum_i \frac{p_{it}}{p_{i0}} (p_{i0} k_{i0})}{\sum_i (p_{i0} k_{i0})} 100 & \frac{\sum_i \frac{k_{it}}{k_{i0}} (p_{i0} k_{i0})}{\sum_i (p_{i0} k_{i0})} 100 & (2) \\
 \vdots & \vdots & \\
 \frac{\sum_i \frac{p_{im}}{p_{i0}} (p_{i0} k_{i0})}{\sum_i (p_{i0} k_{i0})} 100 & \frac{\sum_i \frac{k_{im}}{k_{i0}} (p_{i0} k_{i0})}{\sum_i (p_{i0} k_{i0})} 100 &
 \end{array}$$

Ikusten denez, oinarri-denborako prezio edo kopuruaz biderkatzen eta zatitzen da zenbakitzailearen batugai bakoitzean.

Era honetara, batezbesteko aritmetiko ponderatu bezala, (2) formulak oso erabilgarriak zaizkigu: O_i ondasun bakoitzaren indize sinpleak erabiliz LASPEYRES-en prezio edo kopuruen indize konplexuak kalkulatzeko ditugu.

VI.3.1.2. PAASCHE-ren indizeak

I_t^p eta I_t^k prezio eta kopuruen indize konplexuak kalkulatzeko, O_i ondasun bakoitzaren kasuan indizea kalkulatzeko denborako, kopurua eta prezioa hartzen ditugu konstantetzat.

Hots: $k_i = k_{it}$ eta $p_i = p_{it}$

Honela lortutako indize konplexuei, **PAASCHE-ren prezio eta kopuruen indize konplexuak** deritzegu eta sinbolikoki P_t^p eta P_t^k adieraziko ditugu.

Hots:

$$\begin{array}{cc}
 \frac{P_t^p}{\quad} & \frac{P_t^k}{\quad} \\
 \\
 \frac{\sum_i p_{i1}(k_{i1})}{\sum_i p_{i0}(k_{i1})} 100 & \frac{\sum_i (p_{i1}) k_{i1}}{\sum_i (p_{i1}) k_{i0}} 100 \\
 \vdots & \vdots
 \end{array}$$

$$\begin{array}{ccc}
 \frac{\sum_i p_{it}(k_{it})}{\sum_i p_{i0}(k_{it})} 100 & \frac{\sum_i (p_{it}) k_{it}}{\sum_i (p_{it}) k_{i0}} 100 & \text{(1)} \\
 \vdots & \vdots & \\
 \frac{\sum_i p_{im}(k_{im})}{\sum_i p_{i0}(k_{im})} 100 & \frac{\sum_i (p_{im}) k_{im}}{\sum_i (p_{im}) k_{i0}} 100 &
 \end{array}$$

Lehen egin genuen bezala, zenbakitzailearen batugai bakoitzean, oinarri-denborako prezio edo kopuruaz biderkatzen eta zatitzen bada, (1) formuletatik datozen (2) formuletara iritsiko gara.

$$\begin{array}{ccc}
 \frac{P_t^p}{100} & \frac{P_t^k}{100} & \\
 \\
 \frac{\sum_i \frac{p_{i1}}{p_{i0}} (p_{i0} k_{i1})}{\sum_i (p_{i0} k_{i1})} 100 & \frac{\sum_i \frac{k_{i1}}{k_{i0}} (p_{i1} k_{i0})}{\sum_i (p_{i1} k_{i0})} 100 & \\
 \vdots & \vdots & \\
 \frac{\sum_i \frac{p_{it}}{p_{i0}} (p_{i0} k_{it})}{\sum_i (p_{i0} k_{it})} 100 & \frac{\sum_i \frac{k_{it}}{k_{i0}} (p_{it} k_{i0})}{\sum_i (p_{it} k_{i0})} 100 & \text{(2)} \\
 \vdots & \vdots & \\
 \frac{\sum_i \frac{p_{im}}{p_{i0}} (p_{i0} k_{im})}{\sum_i (p_{i0} k_{im})} 100 & \frac{\sum_i \frac{k_{im}}{k_{i0}} (p_{im} k_{i0})}{\sum_i (p_{im} k_{i0})} 100 &
 \end{array}$$

Ikusten denez, (2) formula hauek ere, indize sinpleen bidez konplexuak kalkulatzera eramaten gaituzte.

LASPEYRES eta PAASCHE-ren indizeak (1) formulen bidez batezbesteko gehikortu ponderatu bezala adierazten dira, eta (2) formulen bidez batezbesteko aritmetiko ponderatu bezala.

VI.3.1.3. FISHER-en indizeak

FISHER-en indize konplexu edo “indize ideal” deritzona, LASPEYRES eta PAASCHE-ren indizeen batezbesteko geometrikoa da.

Hots:

$$F_t^p = \sqrt{L_t^p P_t^p} \text{ FISHER – en prezioen indize konplexua.}$$

$$F_t^k = \sqrt{L_t^k P_t^k} \text{ FISHER – en kopuruen indize konplexua.}$$

VI.3.1.4. Propietate eta erlazio batzuk: bateragarritasuna eta alderantzizkotasuna.

Bateragarritasuna

Indizeen propietate garrantzitsu bat **bateragarri izatea da**, alegia, $B = P \cdot K$ berdintasuna betetzea B , P eta K balio, prezio eta kopuruen indizeak izanik.

LASPEYRES eta PAASCHE-ren indizeek ez dute betetzen erlazio hau, baina bai FISHER-enak.

Hots:

$$L_t^p \cdot L_t^k \neq 100 B_t$$

$$F_t^p \cdot F_t^k = 100 B_t$$

$$P_t^p \cdot P_t^k \neq 100 B_t$$

LASPEYRES eta PAASCHE-ren konbinaketaz ere betetzen da.

$$\text{Hots: } L_t^p \cdot P_t^k = 100B_t \quad P_t^p \cdot L_t^k = 100 B_t$$

Sinboloen ordezkari formula osoak ipintzen baditugu, aurreko erlazioak erraz frogatuko ditugu.

Alderantzizkotasuna

t denborako o denborarekiko indizea eta o denborako t denborarekiko indizea elkarren artean balio alderantzizkoak izateari deritzogu alderantzizkotasuna.

Indize sinpleek ondoko erlazio hau betetzen dute:

$$\frac{I_{t,0}}{100} \neq \frac{1}{I_{t,0} / 100}$$

LASPEYRES eta PAASCHE-ren prezio edota kopuru-indizeek ez dute betetzen, baina bai FISHER-enek.

$$\frac{L_{t,0}^p}{100} \neq \frac{1}{L_{t,0}^p / 100}$$

$$\frac{P_{t,0}^p}{100} \neq \frac{1}{P_{t,0}^p / 100}$$

$$\frac{F_{t,0}^p}{100} = \frac{1}{F_{t,0}^p / 100}$$

LASPEYRES eta PAASCHE-ren prezio edota kopuru-indizeen konbinaketaz ere betetzen da.

$$\frac{L_{t,0}^p}{100} = \frac{1}{P_{t,0}^p / 100} \quad \frac{P_{t,0}^k}{100} \neq \frac{1}{L_{t,0}^k / 100}$$

VI.3.1.5. Kalkulua

Ondoko taulan labore batzuen prezio eta kopuruak ditugu, 1.960 eta 1.964 urteetakoak hain zuzen.

	Arroza		Artoa		Garia	
	p	k	p	k	p	k
1.960	24	100	18	40	20	50
1.964	32	80	12	40	40	70

LASPEYRES-en prezioen indizeaz:

$$\sum_i p_{it} k_{i0} = 32(100) + 12(40) + 40(50) = 5.680$$

$$\sum_i p_{i0} k_{i0} = 24(100) + 18(40) + 20(50) = 4.120$$

$$L_{64}^p = \frac{5.680}{4.120} 100 = 138$$

Alegia, oinarri-denborako kopuruaren balioa, prezioen bidez %38 gehitu dela, baina gehienetan, “prezioak 1.960-1.964 tartean, %38 inguruan goratu direla” esango da.

Bestalde, PAASCHE-ren prezio-indizeaz:

$$\sum_i p_{it} k_{it} = 32(80) + 12(40) + 40(70) = 5.840$$

$$\sum_i p_{i0} k_{it} = 24(80) + 18(40) + 20(70) = 4.040$$

$$P_{64}^p = \frac{5.840}{4.040} 100 = 144$$

Indize honek LASPEYRES-enak bezalako esanahia du, ordea, eraiketa desberdina duenez, oraingo emaitza %44 da.

VI.4. INDIZE KONPLEXUEN ERAIKETAN SORTZEN DEN ZENBAIT ERAGOZPEN

VI.4.1. Aldagaien hautapena

Dakigun bezala, aldagaien multzo bat batera aztertzea, indize konplexuari dagokio.

Beraz, aipatu dugun multzoaren hautapena, lehendabiziko arazoa izango da. Adibidez, aldagaiak nekazal prezioak baldin badira, beren kopurua nahiko handia izango da, areago, kalitate desberdinak kontuan hartzen baditugu; beraz, gehienetan, ondasun garrantzitsuenetatik azpipopulazio bat hautatuko da.

Honek, produktuen eta beren kalitateen izendatze zehatz batetara eramaten gaitu, honela, denbora iragatean, datu-segidak erizpide berdinekoak izango baitira.

VI.4.2. Somatutako leku eta denboraren hautapena

Somatutako zenbakizko datuak hartzean arduraz produktuen kalitateak definitzea eta somaketen leku eta denborak jakinaraztea komeni da. Prezioen kasuan, lekuzat produkzio, salketa edo kontsumo-lekua har liteke; denboratzat une bat (aste, egun edo hilabete bat) har liteke. Ohituraz, maiztasun, adiera edo garrantzi gehieneko erizpideari jarraituz aukeratzen dira leku eta denborak; beraz diren kasu guztien arteko hautapen aleatorio bat ez da posible.

VI.4.3. Talde eta azpitaldeen hautapena

Indize konplexua kalkulatzean, honetaz gainera aldagai, talde eta azpitalde batzuren “konplexua” kalkulatzeko badugu, argibide gehiago izango dugu.

Adibidez, nekazal prezioen arloan bi talde handi kontsidera genitzake, sekain lurrena eta lur ureztatuena. Lehenengo taldean azpitaldeak har litezke, adibidez, laboreak, lekariak..., eta bigarrean, adibidez, barazkiak, fruituak... Beste sailkapen batetan, barne-kontsumorako eta esportaziorako produktuak izan zitezkeen. Talde eta azpitalde hauen hautapena, egin nahi den ikerketaren menpe egongo da.

VI.4.4. Oinarri-denboraren hautapena

Oinarriari dagozkion datuak, beti ehun bezala hartzen dira, orduan, esan genuen bezala, garrantzitsuena “urte normal” bat hartzea da.

Indizearen barnean dauden aldagaiak, aldaketa nabarmenik ez badute (adibidez, industri produktuetan gertatzen den bezala), urte normal batetakoak oinarritzat har daitezke; bestalde, aldaketa handiak baldin badituzte (adibidez nekazal produktuetan bezala), hiruzpalau urtetako erdineurri bat hartuko da oinarritzat.

Beste arazo garrantzitsu bat zera da, oinarri-denbora oraingo denboratik ez dela oso urrun egon behar, baina, dena dela, beste galdera batetan (oinarri aldaketari buruzkoan) honetaz gehiago hitz egingo da.

VI.4.5. Formula eta ponderazioen hautapena

Formula eta ponderazioek elkarbide bat dute. Ponderazioen informazioa lortu ezin bada, normalki, indize sinplearen batezbesteko aritmetikoa egingo da, baina, lortu ahal bada, gehienetan LASPEYRES, PAASCHE eta FISHER-en indizeak erabiliko dira.

LASPEYRES-en indizean, ponderazioa oinarri-denborakoa denez, informazio gutxiena behar du; PAASCHE-renak gehiago behar du, ponderazioa aldakorra baita; eta azkenik, FISHER-ena, dakigunez, beste bien batezbesteko geometrikoa da.

VI.4.6. Indizearen esangura eta zabaldua

Ondasunen kopuruaren menpe zabaldua dago, adibidez, %75 edo %90 bat izango da. Batzutan, hautatu gabe gelditu diren ondasunetariko informazioa baldin badugu, zuzenketak sar daitezke indizeetan.

Indizearen esangura, batezbestekoen ikuspegitik begiratuko dugu; indize konplexuarekin batera sakabanatzearen neurri bat edukitzen dute eta hau handiago da oinarri-denboratik urruntzen garenean.

LASPEYRES eta PAASCHE-ren indizeak batera kalkulatzen baditugu, emaitzen desberdintasuna gehitu egingo da oinarri-denboratik urruntzen garenean eta hau sakabanatzearen aipamen bezala har dezakegu.

Desberdintasun hau nahiko handia egiten denean, oinarrien berriztapena egitea komeniko da.

VI.5. ERAGOZPEN BEREZI BATZU

VI.5.1. Oinarri-aldaketa indize sinpleetan

Adibide batez ikus dezagun.

Urtea	I_t (oinarria = 1.955)	I_t (oinarria = 1.965)
1.955	100	36'4
1.956	116'7	42'4
1.957	150	54'5
1.958	150	54'5
1.959	158'3	57'6
1.960	125	45'5
1.961	100	36'4
1.962	133'3	48'5
1.963	166'7	60'5
1.964	200	72'7
1.965	275	100

Salmenten indize bat baldin bada eta oinarria 1.965. urtean berritzen bada, momentu horretan denborazko segida hautsi egiten da eta denborazko segida berrian kalkuluak atzerantz egiten ditugu erregela proportzional baten bidez (I_t -ren zenbaki bakoitza bider 100 / 275 biderkatzean).

Alegia, 1.955 urteko salmentak %100 baldin baziren, orain 1.965 urtekoak %100 dira.

Bi segida hauek adierazten digute: %33'3 gehiago saldu zela 1.962. urtean 1.955.ean baino eta %51'5 gutxiago saldu zela 1.962. urtean 1.965.ean baino.

VI.5.2. Berriztapen eta loturak indize konplexuetan

Produkzio, aktibitate eta gastuen aldaketaren ondorio bezala, berriztapena inposatzen da. Momentu batetan ondasunik erabilgarrienak zirenak, bigarren postu batetan gelditzen dira denbora iragatean edo desagertu egiten dira. Bestaldetik ondasun berriak agertzen dira, gero arrunt bihurtzen direlarik.

Horregatik ondasun multzoa ezin dezakegu utz aldatu gabe. Denbora iragatean, aldagaiak, oinarriak eta ponderazioak aukeratu beharko ditugu berriro.

Indize konplexuen berriztapenak, segidaren hausturaren arazoa sortzen du.

Ikus dezagun nola egiten den lotura kasu errealean batetan:

Estatistika-Institutuak (Estatu espainolekoak), kontsumo-prezioen indizea berriztatu egiten zuen eta bi segida daude, bata 1.936. urtean oinarritua eta bestea 1.958.ean.

Kontsumo-prezioen indizea:

<u>Urtea</u>	<u>Oinarria 1.936 = 100</u>	<u>Oinarria 1.958 = 100</u>
1.956	643'1	—
1.957	712'4	—
1.958	807'7	100
1.959	866'7	—
1.960	876'9	—
1.961	—	111'3
1.962	—	117'3

Lotura-eragiketarak segida bakar bat lortzen duenez gero, 1.936. urtean edo 1.958.ean oinarrituz egin liteke. Bietan indize sinpleetan ikusi genuen erregela proportzionalen bidez egingo da.

Bi segidak, bi indize desberdinak ditugun urtean (1.958 urtea kasu honetan) berdinak egiten dira. Oinarri bakartzat 1.936. urtea hartzen badugu, orduan 1.958.eko 100, bider 8'077 biderkatzean 807'7 bihurtzen da, eta honela egingo da segida osoaz; aldiz, oinarri bakartzat 1.958 urtea hartzen badugu, orduan, 1.936.ean oinarriturik dagoen segida, 8'077 balioaz zatituko da.

Ondoko taulan, loturik dauden bi segida ikusiko ditugu.

<u>Urtea</u>	<u>Oinarria 1.936 = 100</u>	<u>Oinarria 1.958 = 100</u>
1.956	643'1	79'5
1.957	712'4	82'2
1.958	807'7	100
1.959	866'7	107'3
1.960	876'9	108'6
1.961	899'0	111'6
1.962	947'2	111'4

Argiro, indize sinpleen alderantziz, bi segida hauek ez dira konparagarriak, aldagai eta ponderazio desberdinez eginak baitaude, baina, ez dugu beste biderik, indize konplexuari jarraitasuna eman nahi badiogu.

Ez dugu ahaztu behar, indizeak fenomeno baten eboluzioa ikusteko erabiltzen diren adierazle batzu besterik ez direla.

VI.6. ZENBAKI-INDIZEEN APLIKAPENAK

Indizeen formulak era orokor batean ikusi ondoren, bi multzotan sailka ditzakegu, prezio eta kopuruenak, baina, formula hauek, indize batzutan zehaztu egiten dira, adibidez, industri produkzioen indizeak, alokairuen indizeak, kontsumo-prezioen indizeak, balore mugigarrien indizeak, inportazioen indizeak, esportazioen indizeak eta abar, gehienak ekonomi aktibitatearen eboluzioa ikusteko baliagarriak direlarik.

Ikasgai honetan, kontsumo-prezioen indizea ikusiko dugu, garrantzitsuena baita.¹

VI.6.1. Kontsumo-prezioen indizeak

Hipotesi bezala familia baten bizimaila konstantetzat hartzen badugu, denbora-une desberdinetan familia horrek kontsumoan ordaintzen duenaren gain konparaketak egitea da indize honen helburua. Alegia, aurrekontu (familia baten bizimaila deritzona) baten aldaketak neurtu nahi ditu.

Egikera erabiliena “erosketa-otarrea” deitzen dena da; hau inkesta baten bidez, Famili Aurrekontuen Inkesta, alegia, aurrekontuan sartzen diren ondasunen kopuruek, zerbitzuek eta industri produktuek osatzen dute.

Hartzen den familia, “normala” da, hots, bere aurrekontua batezbestekoari hurbiltzen zaio.

Oinarri-denborako prezioaz eta oraingo prezioaz “otarre”-aren balioa kalkulatzeko, izango da arazo bakarra, eta bi hauen zatidura kontsumo-prezioen indizea da.

Izan bitez:

k_{10} k_{20} ... k_{n0}	“otarre”-ko ondasunen kopuruak (bizimaila presuposatzen dutenak).
p_{10} p_{20} ... p_{n0}	haien prezioak oinarri-denboran.
eta	
p_{11} p_{21} ... p_{n1}	haien prezioak oraingo unean.

1. Besteak UEU-k 1979. urtean argitaratu zuen ZENBAKI-INDIZEAK liburuxkan ikus daitezke.

Orduan **kontsumo-prezioen indizea**, ondoko formulaz edukiko dugu.

$$I_{1,0} = \frac{\sum_{i=1}^n p_{i1}k_{i0}}{\sum_{i=1}^n p_{i0}k_{i0}} \cdot 100$$

Ikusten denez, indize hau LASPEYRES-en indize bat da.

Lehen esan dugunez, denbora iragatean indizearen oinarria aldatuko ez balitz, eta honekin batera indizearen estruktura gaurkotuko ez balitz, indizea balio gabe geratuko litzateke. Zer esanik ez indizea kontsumo-prezioena izanik, azken urteetan jasaten ditugun aldaketa teknologiko eta ekonomikoen egoeraz baldintzatuta dagoela.

Espainiako Estatistika Institutuak KPI kalkulatzeko oinarri hauek finkatu ditu: 1936, 1958, 1976, 1983 eta 1985.eko uztaila.

Azken urte hauetako egoera, hauex da: Famili Aurrekontuen Inkesta berria 1990/91 urteetan burutu zen, oinarri aldaketa 1992 urtea oinarritzak harturik egin zelarik. Kontsumo Prezioen Indize berria 1993. urteko urtarriletik aurrera ateratzen da.

Dakusagunez lehen aldaketa 22 urte iragan ondoren izan zen, eta bigarrena 18 urte iragan ondoren; azkeneko urteetan gero eta hurbilago dago aldaketa, kontsumoaren estrukturan aldaketak beharturik.

VI.6.2. Moneta-unitate arruntetan dauden magnitudeen deflazioa

Zenbaki-indizeen aplikazio garrantzitsuenetako bat deflazioa da.

Prezio edo balioen moneta-segida bat deflaktatzea, segida hori moneta efektuaz askatzea da.

$\sum_i p_{it}k_{it}$ balio-segida, moneta-balio arruntetan dago eta balio errealetan edo moneta-balio konstanteetan ipini nahi badugu, deflaktatze eragiketa eginez $\sum_i p_{i0}k_{it}$ balio-segidan bihurtuko da.

Bidea, PAASCHE-ren prezio-indizeaz zatitzea izango litzateke.

Honela

$$\frac{\sum_i p_{it} k_{it}}{P_t^P} = \frac{\sum_i p_{it} k_{it}}{\sum_i p_{it} k_{it}} = \sum_i p_{i0} k_{it}$$

P_t^P ez daukagunean, beste “ad hoc” indizea hartuko da, baina segida desberdin batzu parekatu nahi baditugu, indize berbera erabiliko dugu.

Adibidea: Ondoko taulan 90-92 urte-tarterako ondasun baten salmentak pezeta arruntetan ditugu. Zein izango litzateke salmenta horien balioa pezeta konstanteetan, urte horietan Laspeyres-en prezio-indizeak: 100, 131.42, 162.43, hain zuzen, deflaktatzeko erabiltzen baditugu?

Urtea	I_t (oinarria = 1.955)	I_t (oinarria = 1.965)	
Urtea	Salmentak pta arrut.	Deflazio- -indizea	Salmentak 1990.eko pta. konstanteetan
1990	1500000	100	1500000
1991	1800000	131.42	1396654
1992	2100000	162.43	1300873

**VII. DESKRIBAPEN ESTADISTIKOAREN
ADIERAZPIDE GEOMETRIKOAK**

VII.1. OROKORTASUNAK

*VII.2. X DATU-MATRIZEAREN BI ADIERAZPIDE
GEOMETRIKO*

VII.3. BI ALDAGAI ETA BI OHARPEN

VII.4. BI ALDAGAI ETA HIRU OHARPEN

*VII.5. ERREGRESIO BAKUNA ALDAGAI ESTADIS-
TIKOEN BALIO-MULTZOEN ESPAZIOAN*

VII.1. OROKORTASUNAK

Ikasgai honetan batezbestekoa, bariantza eta kobariantza geometrikoki ikustea hartzen dugu helburutzat.

Matrize-estatistikoak azaltzen genuen ikasgaietan, estatistika anizkoitzean ikasketak matrizialki egitea beharrezko bihurtzen dela esaten genuen.

Errealitatea estatistikoki¹ ikasi nahi badugu, bada, formulazio matematiko bektoriala edota matriziala nahitaezkoak eta bide batez paregabeak direla esan daiteke, errealitatea bektoriala baita.

$\mathbf{X}_{(m \times n)}$ datu-matrizearen bektoreak, \mathbb{R}^n , \mathbb{R}^m espazioetako bektoreak alegia, aurkeztu hasiko gara; ondoren, \mathbb{R}^2 , \mathbb{R}^3 espazioetan zentratuz, aipatutako estatistikoak geometrikoki ikusiko ditugu.

Zoritxarrez n edo m 3 baino handiago direnean iruditzen ez gertatzen ezer gutxi ikusiko dugu, baina, \mathbb{R}^2 , \mathbb{R}^3 espazioak ondo aztertuta baditugu espazio anizkoitzean gertatzen dena errazago ulertuko dugu.

VII.2. X DATU-MATRIZEAREN BI ADIERAZPIDE GEOMETRIKO

Lehen aurkeztutako $\mathbf{X}_{(m \times n)}$ datu-matrizean bi bektore mota kontsidera daitezke:

1) \mathbb{R}^n espazioan, indibiduoek dituzten n koordenatuek n aldagaiekin arabera \mathbb{R}^n espazioko bektoreak osatzen dituzte.

2) Bestalde, \mathbb{R}^m espazioan, aldagaiek dituzten m koordenatuek m indibiduen gain \mathbb{R}^m espazioko bektoreak osatzen dituzte.

1. Estatistika Deskribatzailean kokaturik, eta konkretuki, Datu-Analisiaren atarian, hau da, probabilitate-indukzioan oinarritzen den Inferentzia Estatistikotik at.

Hots:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1j} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2j} & \cdots & x_{2n} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{i1} & x_{i2} & \cdots & x_{ij} & \cdots & x_{in} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{m1} & x_{m2} & \cdots & x_{mj} & \cdots & x_{mn} \end{bmatrix}$$

Errenkadak, \mathbb{R}^n espazioko bektoreak dira, i. indibiduorako edo oharpenerako $(x_{i1} \ x_{i2} \ \dots \ x_{ij} \ x_{in})$ koordinatuak direlarik. Errenkadetan, bada, m **indibiduo-puntu** edo oharpen-puntu ditugu, orain arte edozein hodeitan genituen puntuak alegia.

Zutabeak, aldiz, \mathbb{R}^m espazioko bektoreak dira j. aldagairako $(x_{1j} \ x_{2j} \ \dots \ x_{ij} \ \dots \ x_{mj})$ koordinatuak direlarik. Zutabetan n **aldagai-puntu** ditugu hemendik aurrera kontutan hartuko dugun puntu-hodei berri bat osatzen dutelarik.

VII.3. BI ALDAGAI ETA BI OHARPEN

Bi aldagai bi indibiduen gain neurtu ondoren honako emaitza hauek ditugu.

$$\begin{matrix} & X_1 & X_2 \\ \omega_1 & \begin{bmatrix} x_1(1) & x_2(1) \end{bmatrix} \\ \omega_2 & \begin{bmatrix} x_1(2) & x_2(2) \end{bmatrix} \end{matrix} = (\text{konkretuki}) = \begin{bmatrix} 1 & 5 \\ 3 & 2 \end{bmatrix} \text{ datu-matrizea}$$

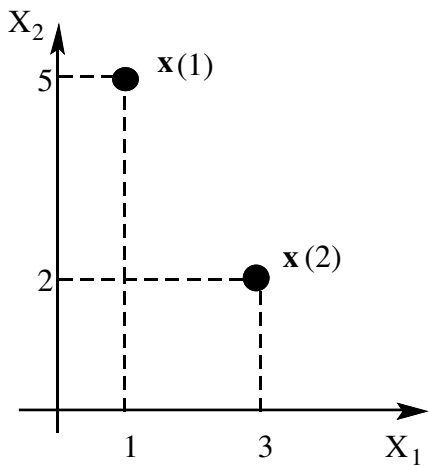
adibide bezala erabiliko dugu.

Kasu honetan $\mathbf{x}(1) = (1,5)^T$, $\mathbf{x}(2) = (3,2)^T$, \mathbb{R}^2 -ko bi bektoreak indibiduo-puntuak edo oharpen-puntuak dira. $\mathbf{x}_1 = (1,3)$, $\mathbf{x}_2 = (5,2)$, \mathbb{R}^2 -ko bi bektoreak, aldiz, aldagai-puntuak.

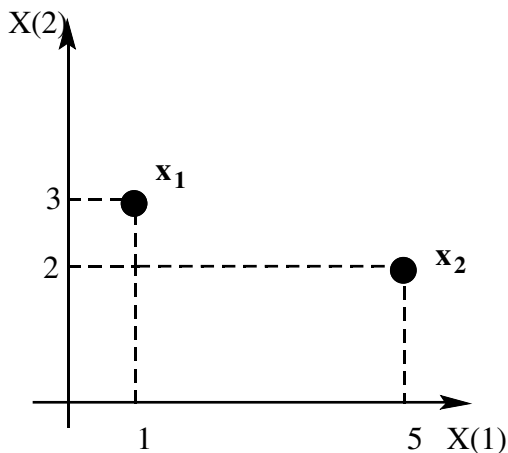
Konturatu behar dugu datu-matrizeko bektoreak ez direla desberdintzen errenkada-bektoreak edo zutabe-bektoreak izateagatik bakarrik, daukaten desberdintasuna benetan garrantzitsuagoa da: lehenengoak balio-multzo ez-homogenoak dira indibiduo baten gain aldagai desberdinen neurriak izanik; bigarrenak, ordea, balio-multzo homogenoak dira aldagai bakoitzaren neurriak izatean.

Datu-matrizearen (i,j) elementua adierazteko bi ikur ezberdin aukeratu ditugu. Hauexek dira: $x_{ij} = x_j(i)$

Biek X_j aldagaiak ω_i indibiduoaren gain hartzen duen balioa adierazten digute. Dakusagun, bada, adierazpide grafikoak:



Indibiduoaren Hodeia



Aldagaiaren Hodeia

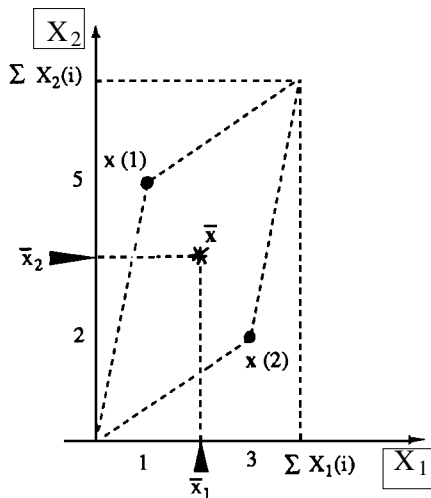
X_1, X_2 , aldagaien balioen batezbestekoek betebeharrak desberdinak betetzen dituzte bi hodei horietan

$$\bar{x}_1 = \frac{1}{m} \sum_i x_1(i) = 2$$

$$\bar{x}_2 = \frac{1}{m} \sum_i x_2(i) = 3.5$$

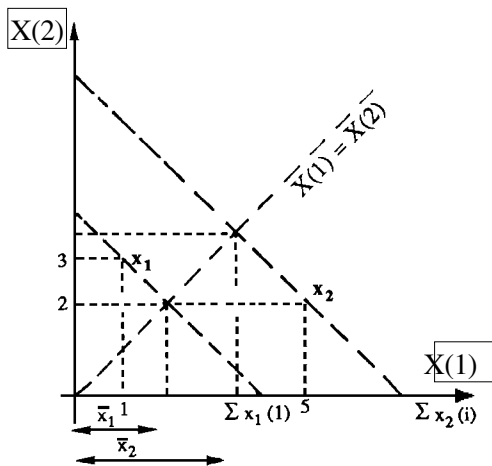
INDIBIDUOEN HODEIAN

Indibiduo-puntuaren grabitate-zentrua finkatzen dute.



ALDAGAIEN HODEIAN

Aldagai-puntuen osagaien batezbestekoa



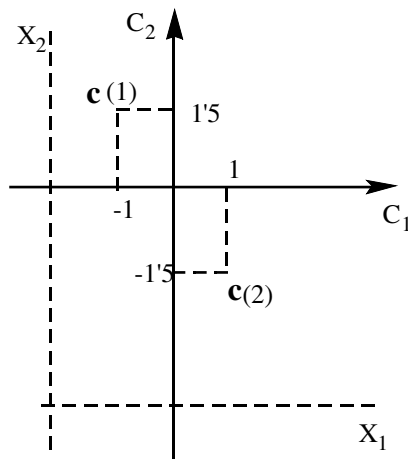
Aldagai-bektorearen batezbestekoa, $X(1) = X(2)$ erdikarian aldagai-puntuak daukan proiektatutako puntuaren osagai bakoitza bezala ikusten da.

Aldagaien balio bakoitzaren **zentratze-eragiketa**:

$C_i(i) = X_j(i) - \bar{x}_j$ zentratze-eragiketa eginez, **C** datu-matrize zentratua lortzen dugu.

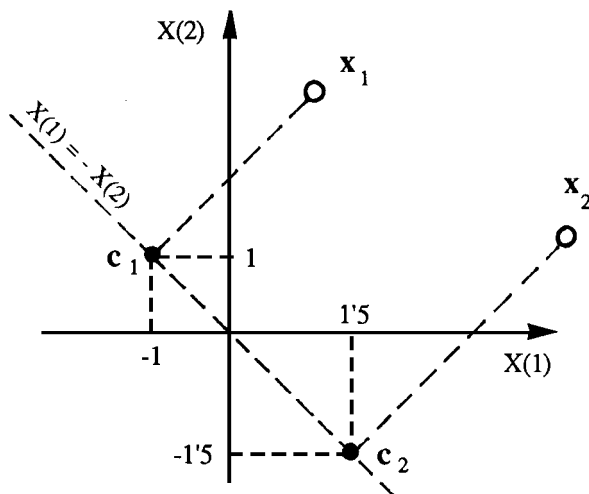
$$C = \begin{bmatrix} c_1(1) & c_2(1) \\ c_1(2) & c_2(2) \end{bmatrix} = \begin{bmatrix} -1 & 1,5 \\ 1 & -1,5 \end{bmatrix}$$

INDIBIDUOEN HODEIAN



Zentratzeak ardatzen aldaketa paralelo bat egitea suposatzen du hodei honetan.

ALDAGAIEN HODEIAN



Hodei honetan, ordea, batezbestekoz eta baturaz zero den $X(1) = -X(2)$ erdikarian aldagai-puntuak proiektatzea da zentratze-eragiketak suposatzen duena.

Aldagai batek bi oharpen dituenean, zentratzea $X(1) = -X(2)$ erdikarian proiektatzea da, bi balio zentratuak, berdinak balio absolutuz eta aurkakoak baitira. Ikus grafikoan.

Horrexegatik $X(2 \times 2)$ datu-matrizeak ezaugarri oso bereziak ditu.

Dakusagun:

1) Desbidazio standarda aldagai zentratuaren koordenatuen balio absolutua da.

$$S_j = |c_j(1)| = |c_j(2)|$$

$$\text{Kasu honetan: } S_1 = |-1| = |1| = 1$$

$$S_2 = |1'5| = |-1'5| = 1'5$$

Dakigunez:

$$S_1^2 = \frac{1}{m} \sum_i c_1^2(i)$$

$$S_2^2 = \frac{1}{m} \sum_i c_2^2(i)$$

$$S_1 = \frac{\|c_1\|}{\sqrt{m}}$$

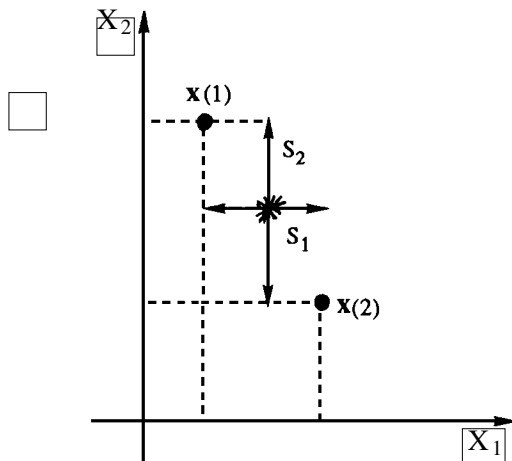
$$S_2 = \frac{\|c_2\|}{\sqrt{m}}$$

$m = 2$, hots bi indibiduo edo bi oharpen izanik

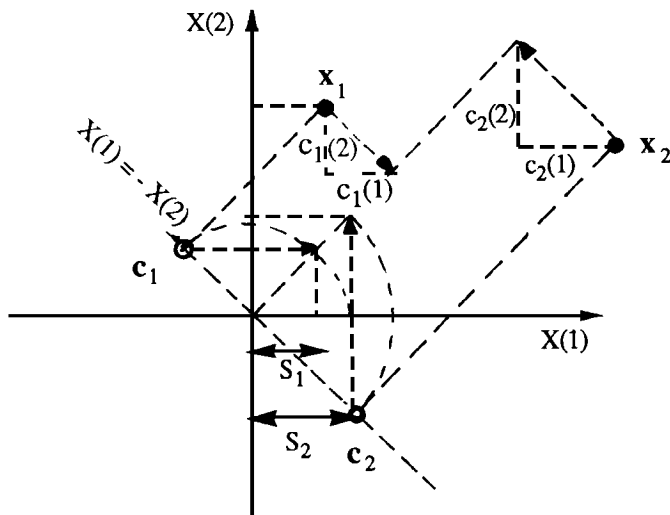
$$\mathbf{L} = \begin{bmatrix} S_1^2 & S_{12} \\ S_{12} & S_2^2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} -1 & 1 \\ 1'5 & -1'5 \end{bmatrix} \begin{bmatrix} -1 & 1'5 \\ 1 & -1'5 \end{bmatrix} = \begin{bmatrix} 1^2 & -1 \cdot 1'5 \\ -1 \cdot 1'5 & 1'5^2 \end{bmatrix}$$

Kobariantza matrizea lortzean, matrize hauen bigarren berezitasuna agertzen zaigu, baina horretaz geroxeago hitz egingo dugu.

Dakusagun, orain bada, desbidazio standarda nola agertzen den indibiduen hodeian.



Kasu berezi honetan ($X(2 \times 2)$) desbidazio standarda edo tipikoa “desbidazioa” balio absolutuan da.



$X(1) = X(2)$ erdikarian aldagai-puntua proiektatzean, proiektatutako puntuaren koordinatu bakoitza batezbestekoa zen bezalaxe, orain erdikarian aldagai-puntu zentratua proiektatzean, proiektatutako puntuaren koordinatu bakoitza desbidazio standarda da.

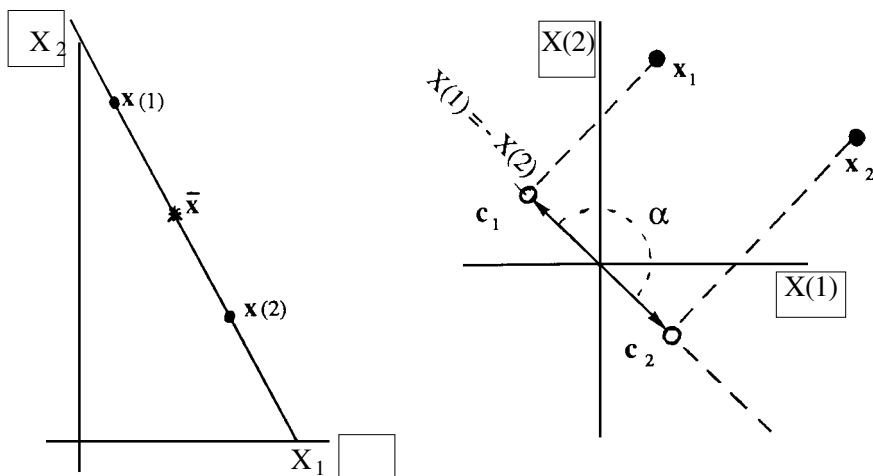
C_j bektoreak eta bere proiekzioak norma bera daukate:

$$\|c_j\| = \sqrt{m} S_j$$

Aldagai zentratuaren norma sakabanatzearen neurri on bat da, askotan erabilia izanik.

Azkenik azpimarratu behar dugu aldagai zentratuaren norma, jatorrizko aldagaiak erdikarira daukan distantzia euklidearra dela.

2) Bi oharpen ditugunean aldagaiak guztiz korrelatuak dira beti, hots koerlazio-koefizientea 1 edo -1 (gure adibidean bezala) izango da. Ikus **L** aurreko orrialdean.



$$r_{12} = \frac{S_{12}}{S_1 S_2} = \frac{\langle c_1, c_2 \rangle / m}{\frac{\|c_1\|}{\sqrt{m}} \frac{\|c_2\|}{\sqrt{m}}} = \frac{\langle c_1, c_2 \rangle}{\|c_1\| \|c_2\|} = \cos \alpha$$

Dakusagenez, aldagai zentratuen azpiespazio bektoriala $X(1) = -X(2)$ zuzena da, horrexegatik edozein aldagai bikotek osatzen duen angelua zero edo ehun eta laurogei gradukoa izango da.

Horrek, dakigunez, erregresio perfektua halabehartzen du.

Indibiduo-puntuek ezin dutela zuzen baten gainean ez izan ikusten dugu ezkerreko hodeian, lerrotatuak daude.

Aldagai zentratu bakoitza bestearen eskalarra dela ikusten dugu, ordea, eskuinaldekoan.

Ondoren, hau baino errepresentagarriago den kasu bat aztertuko dugu, bi aldagai eta hiru oharpenen kasua, alegia.

Ariketa: \mathbf{X} datu-matrizea $\begin{bmatrix} 1 & 4 \\ 5 & 6 \end{bmatrix}$ matrizea izanik. Ikus ezazu, grafikoki, \mathbf{c}_1 , \mathbf{c}_2 bektoreek zero graduko angelua osatzen dutela edota $r_{12} = 1$ dela.

VII.4. BI ALDAGAI ETA HIRU OHARPEN

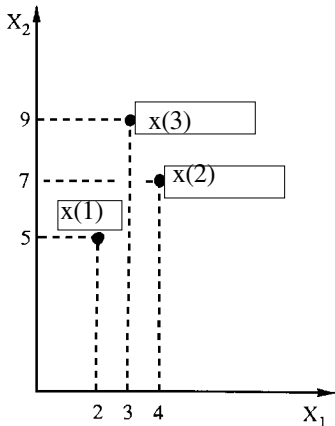
Ondoan, bi aldagai eta hiru indibiduoren kasuko matrizea daukagu.

$$\mathbf{X} = \begin{bmatrix} x_1(1) & x_2(1) \\ x_1(2) & x_2(2) \\ x_1(3) & x_2(3) \end{bmatrix}$$

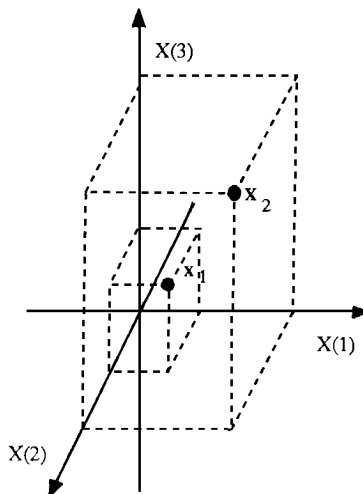
Eta adibide bezala kontsideratzen duguna.

$$\mathbf{X} = \begin{bmatrix} 2 & 5 \\ 4 & 7 \\ 3 & 9 \end{bmatrix}$$

Indibiduo-puntuen hodeia



Aldagai-puntuen hodeia

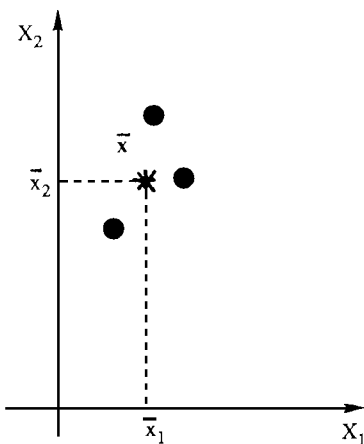


$$\bar{x}_1 = \sum_i x_1(i) / m = 3$$

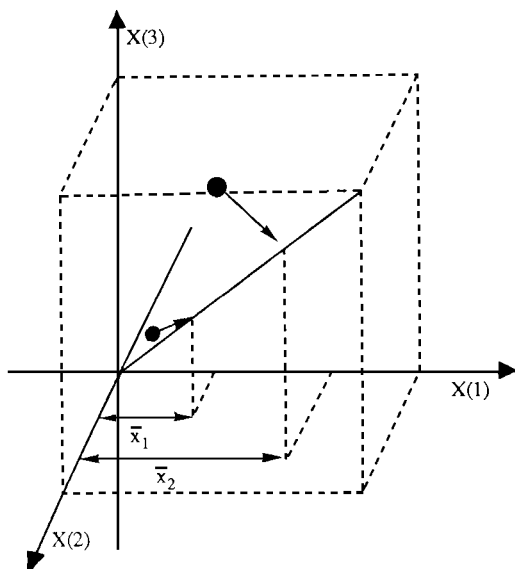
$$\bar{x}_2 = \sum_i x_2(i) / m = 7$$

Hots $\bar{\mathbf{x}} = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$

Dakusagun, bada, adierazpide grafikoak:



Batezbesteko bektorea indibiduen barizentrua edo grabitate-zentrua da.



Aldagai bektorearen batezbesteko osagaia $X(1) = X(2) = X(3)$ erdikarian proiektatutako puntuaren osagai bakoitza da.

Zentratze-eragiketa:

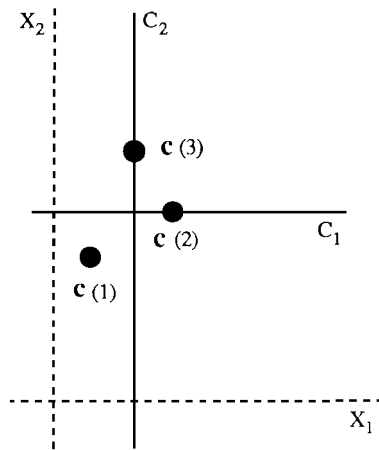
$$\bar{x}_1 = \sum_i x_1(i) / m = 3$$

$$\bar{x}_2 = \sum_i x_2(i) / m = 7$$

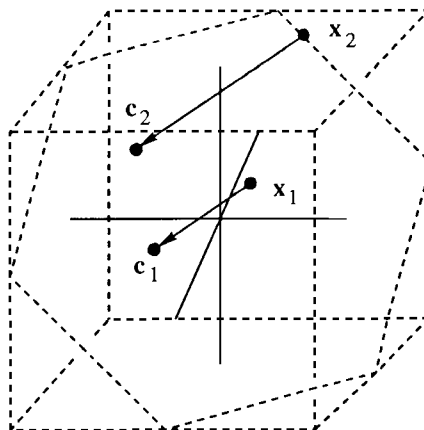
Hots $\bar{\mathbf{x}} = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$

datu-matrize zentratua lortzen dugu.

Grafikoki:



Berriro ere ardatzen aldaketa paralelo bat daukagu.

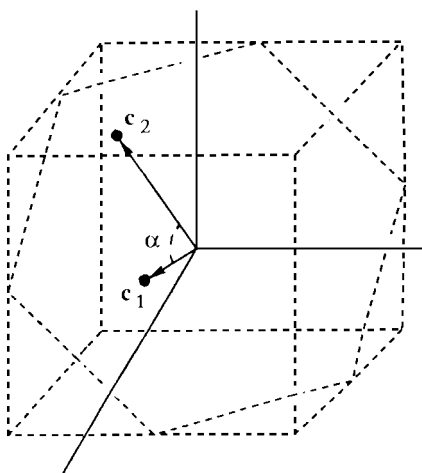


Lehen $X(1) = -X(2)$ erdikarian proiektatzen ziren aldagai-puntuak, hots, baturaz eta batezbestekoz zero diren bektoreen azpiespazioan.

Orain, bada, berdin gertatzen da baina kasu honetan baturaz eta batezbestekoz zero diren bektoreen azpiespazioa plano bat da, $X(1) + X(2) + X(3) = 0$ plano a alegia.

Irudian, kubo baten barruan dagoen hexagonoaren bidez ikusarazi nahi da azpiespazio hori, hau da, hexagonoak $X(1) + X(2) + X(3) = 0$ plano a ikusarazteko besterik balio ez duelarik.

Aldagai zentratuen grafika handituz aldagaien arteko koerlazioa geometrikoki ikusiko dugu.



Ikusi genuenez:

$$\|c_1\| = \sqrt{m} \quad S_1$$

$$\|c_2\| = \sqrt{m} \quad S_2$$

$$r_{12} = \frac{\langle c_1, c_2 \rangle}{\|c_1\| \|c_2\|} = \cos \alpha \quad \alpha = \arccos r_{12}$$

Aldagaien osagaiak ez dira berdinak balio absolutuan ez eta aldagaiek ez dute zergatik guztiz korrelatuak izan behar ere.

Darabilgun adibidean \mathbf{L} eta \mathbf{R} matrizeak hauek dira:

$$\mathbf{L} = \begin{bmatrix} 2/3 & 2/3 \\ 2/3 & 8/3 \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix} \quad \begin{array}{l} \text{hots } \alpha = 60^\circ \\ \alpha = \arccos 1/2 \end{array}$$

Ariketa: \mathbf{c}_2 bektorearen ordezkio berean dagoen $\mathbf{c}_2 = (-1, 0, 1)$ bektorea kontsideratuz, ikus ezazu $\mathbf{c}_1, \mathbf{c}_2$ bektoreen arteko angelua 60° koa dela.

Oharra: $\mathbf{c}_1, \mathbf{c}_2, 0$ puntuek alde berdineko triangelu bat osatzen dute, alde bakoitzak 2 balio duelarik.

DATU-MATRIZE TIPIFIKATUA ETA ZENTRU-NORMALIZATUA

Dakigunez aldagaien balioak tipifikatuak edo standardizatuak dituen datu-matrizeari tipifikatua deritzogu.

Adibidean:

$$\mathbf{C} = \begin{bmatrix} -1 & -2 \\ 1 & 0 \\ 0 & 2 \end{bmatrix}, \quad S_1 = \sqrt{\frac{2}{3}}, \quad S_2 = 2\sqrt{\frac{2}{3}} \quad \text{izanik}$$

$$\mathbf{T} = \begin{bmatrix} \frac{x_j(i) - \bar{x}_j}{S_j} \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{3}}{\sqrt{2}} & \frac{\sqrt{3}}{\sqrt{2}} \\ \frac{\sqrt{3}}{\sqrt{2}} & 0 \\ 0 & \frac{\sqrt{3}}{\sqrt{2}} \end{bmatrix}$$

Era berean, aldagai zentratuen balioak dagozkien normaz zatituz, datu-matrize zentru-normalizatua lortuko dugu.

Hau da:

$$\|\mathbf{c}_1\| = S_1 \sqrt{3} = \sqrt{2}$$

$$\|\mathbf{c}_2\| = S_2 \sqrt{3} = 2\sqrt{2}$$

$$\mathbf{N} = \left[\frac{x_j(i) - \bar{x}_j}{S_j \sqrt{m}} \right] = \begin{bmatrix} \frac{-1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & 0 \\ 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$$

Aldagai-puntu tipifikatuen norma \sqrt{m} da.

Hots:

$$\|t_j\| = \left\| \frac{x_j(i) - \bar{x}_j}{S_j} \right\| = \left\| \frac{\mathbf{c}_j}{S_j} \right\| = \frac{\|\mathbf{c}_j\|}{S_j} = \frac{\|\mathbf{c}_j\|}{\|\mathbf{c}_j\|/\sqrt{m}} = \sqrt{m}$$

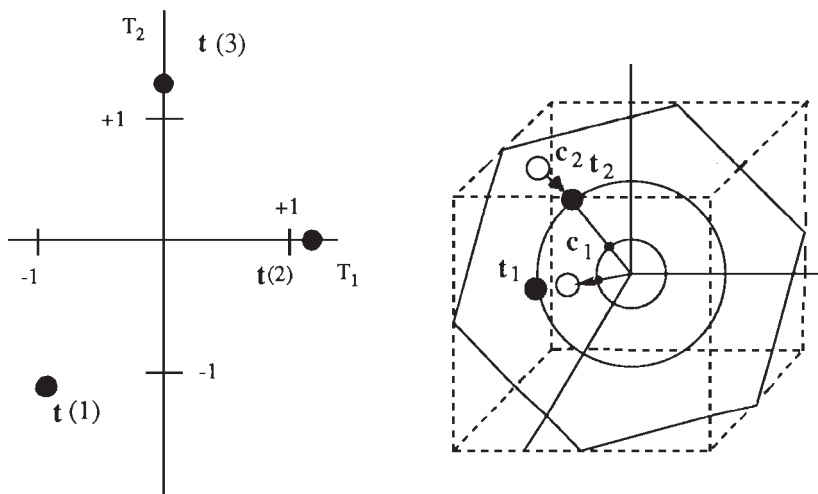
Noski, \mathbf{N} matrizeko zutabe-bektoreek 1 daukate normaz, bektoreak normalizatuak baitaude.

Edozein kasutan, aldagaiak zentratzeak translazio bat bakarrik suposatzen badu ere, tipifikatzeak edota zentratu ondoren normaltzeak eskala aldatzea suposatzen du.

Horrexegatik, \mathbb{R}^3 espazioan aldagai-puntu tipifikatuak erradioz $\sqrt{3}$ duen zirkulu batetan eta $m \geq 4$ denean erradioz \sqrt{m} duen esfera edo hiperesfera batetan kokatzen dira hurrenez hurren.

Aldagai-puntu zentru-normatuak, berriz, erradioz 1 neurtzen duen zirkulu, esfera, edo hiperesferan kokatzen dira.

Bi kasu horietan, balioak erlatibizatzen ditugu desbidazio standardaren arabera eta aldagai guztiak eskala berean ditugunez beraien direkzioak interesatzen zaizkigu.



Indibiduen balioak tipifikaturik ditugu ezkerreko hodeian.

Eskumaldeko hodeian erradioz $\sqrt{3}$ duen zirkunferentzian aldagai-bektore tipifikatuak ikusten ditugu, baita, aldagai-bektore normatuak erradioz 1 duen zirkunferentzian ere.

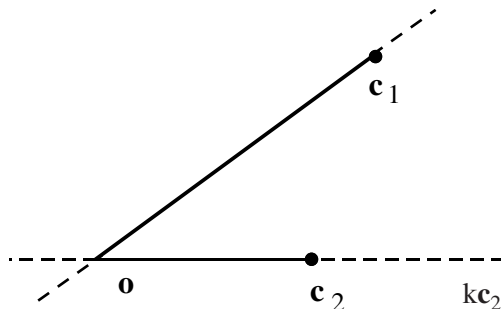
VII.5. ERREGRESIO BAKUNA ALDAGAI ESTATISTIKOEN BALIO-MULTZOEN ESPAZIOAN

X.2. atalean ikusi dugun bezala $\mathbf{X}(m \times n)$ datu-matrizearen zutabeak \mathbb{R}^m espazioko aldagai-bektoreak dira, m indibiduen gain aldagai bakoitzaren balioak hartu baititugu.

Aldagai-bektore zentratuak, ordea, beraien osagaien edo balioen batura zero denez, \mathbb{R}^{m-1} espazioko bektoreak dira. Konkretuki, X.3. eta X.4. ataletan, \mathbb{R}^1 eta \mathbb{R}^2 espazioetan kokaturik ikusi ditugun bi indibiduen eta hiru indibiduen kasuetarako, hurrenez hurren.

Bi aldagaien arteko erregresio bakuna (aldagai independente edo erregresore bat) beti, plano batean adieraz daiteke, alegia, \mathbb{R}^{m-1} dimentsioko bi aldagai-bektore zentratuek sortarazten duten espazio bidimentsionalean.

Bi aldagai eta m oharpen suposatuz, ikus ditzagun, bada, $X(m \times 2)$ datu-matrizeko bi aldagai zentratuak beraiek osatzen duten planoan adieraziak.



X_1 aldagaiaren X_2 aldagaiarekiko erregresioa suposatuz eta aldagai-bektore zentratuen bidez adieraziz:

$$\hat{c}_1 = b_{12}c_2$$

Hau da, aldagaiaren balio estimatuen bektorea c_2 bektorearen eskalarra da.

$b_{12} c_2$, bada, $O c_2$ segmentua kokatuta dagoen zuzenean kokatuko da, zuzen hau c_2 bektorearen eskalar guztien leku geometrikoa edo kc_2 direkzioa baita.

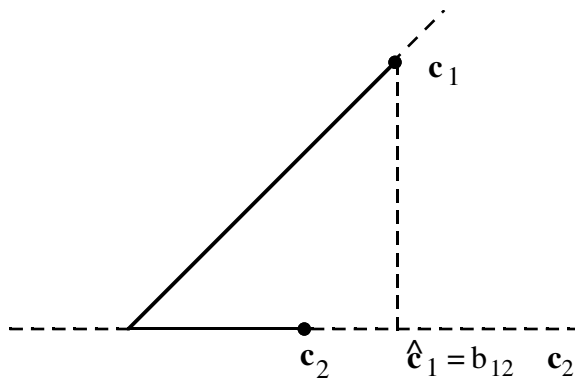
\mathbb{R}^m espazioan, c_1 eta $b_{12} c_2$ bektoreen arteko distantzia, norma euklidearraz kalkula daiteke.

$$d = + \sqrt{\sum_i [c_1(i) - b_{12} c_2(i)]^2} = \sqrt{\sum_i [c_1(i) - \hat{c}_1(i)]^2}$$

Eta $c_1(i) - \hat{c}_1(i)$ i . indibiduoaren hondarra edo errorea izanik, distantzia, errore karratuen baturaren erro karratua da, hau da, $\|e\|$ norma. Balio hau txikiena, eta bide batez, S_e^2 hondar-bariantza txikiena izatea da erregresioaren helburua.

Erregresioa, bada, geometrikoki honela planteatu daiteke:

$\hat{c}_1 = b_{12} c_2$ eta c_1 balio-multzoen arteko distantzia txikiena izatea da bilatzen dena. Balio hau, dakigunez, c_1 eta c_1 -ek, kc_2 direkzioan, daukan proiektzio ortogonalaren arteko distantziari dagokio.



Honela, geometrikoki, b_{12} erregresio-koefizientea lor dezakegu:

$$\text{proj. } \mathbf{c}_1 = b_{12} \|\mathbf{c}_2\|$$

eta dakigunez,

$$\text{proj. } \mathbf{c}_1 = \left\langle \mathbf{c}_1 \mid \frac{\mathbf{c}_2}{\|\mathbf{c}_2\|} \right\rangle = \frac{\langle \mathbf{c}_1 \mid \mathbf{c}_2 \rangle}{\|\mathbf{c}_2\|}$$

orduan,

$$b_{12} = \frac{\langle \mathbf{c}_1 \mid \mathbf{c}_2 \rangle}{\|\mathbf{c}_2\|^2} = \frac{\langle \mathbf{c}_1 \mid \mathbf{c}_2 \rangle / m}{\|\mathbf{c}_2\|^2 / m} = \frac{S_{12}}{S_2^2}$$

Azkenik, balio estimatuen bektorea

$$\hat{\mathbf{c}}_1 = \frac{S_{12}}{S_2^2} \mathbf{c}_2$$

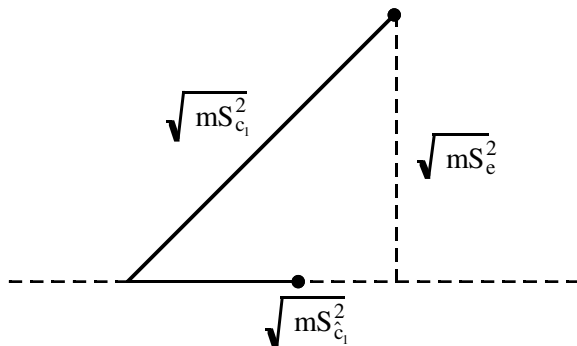
\mathbf{c}_1 bektorearen bektore-posizioaren osagaiak, $\hat{\mathbf{c}}_1$ bektorearekiko \mathbf{c}_1 -ren \mathbf{c}_2 -rekiko erregresioaren hondarrak dira.

Dakusagunez, \mathbf{e} bektorearen direkzioa \mathbf{c}_2 eta $\hat{\mathbf{c}}_1$ bektoreen direkzioarekiko ortogonal da. Honela geometrikoki erregresioaren bi propietate hauek ikusten ditugu:

$-r_{e\mathbf{c}_2} = r_{e\mathbf{x}_2} = 0$, hau da, hondarrak eta aldagai independentea koerlazio gabeak dira.

$-r_{e\hat{\mathbf{c}}_1} = r_{e\hat{\mathbf{x}}_1} = 0$, hau da, hondarrak eta balio estimatuak koerlazio gabeak dira.

Azkenik, hain garrantzitsua den bariantzaren deskonposaketa batukorra, geometrikoki ikus dezakegu:



Pitagoras-en teorema aplikatuz,

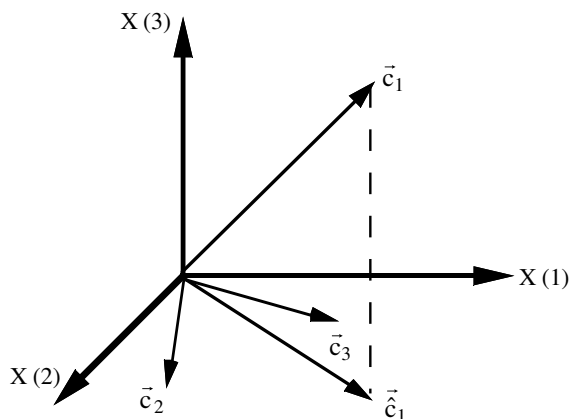
$$S_{c_1}^2 = S_{\hat{c}_1}^2 + S_e^2$$

edota

$$S_{x_1}^2 = S_{\hat{x}_1}^2 + S_e^2$$

Karratu txikiaren erregresioa, aldagaien kopurua edozein izanik ere, grafikoki berdintsu interpreta daiteke. Aldagai dependentearen balio zentratuen bektorea erregresore guztiek osatzen duten azpiespazioan proiektatuz egingo da.

Hurrengo grafikoan, X_1 aldagaiaren erregresioa X_2 , X_3 aldagaiekiko aurkezten dugu.



**VIII. DOIKUNTZA ORTOGONALA ETA
KOBARIANTZA MATRIZEAREN
AUTODIREKZIOAK**

VIII.1. SARRERA

VIII.2. ANALISIA \mathbb{R}^2 ESPAZIOAN. DOIKUNTZA
ORTOGONALAREN ZUZENA

VIII.3. DOIKUNTZA ORTOGONALAREN
ZUZENAREN LORPENA

VIII.4. HODEI DUALAREN AURKEZPEN
GRAFIKOA

VIII.5. TRANSIZIO ERLAZIOAK

VIII.6. BATERAKO AURKEZPEN GRAFIKOA
ETA INTERPRETAZIOA

VIII.7. DOIKUNTZA ORTOGONALAREN ZUZENA
ALDAGAI NORMATUENTZAKO

VIII.1. SARRERA

Gai honen helburua, hurrengo gaien aurkeztuko dugun Osagai Nagusizko Analisiaren garapena ulergarriagoa egitea da.

Osagai Nagusizko Analisiaren helburua, berriz, dimentsioa murriztuz, (m x n) datu-matrizearen ikasketa burutzea da. Honek, aurkezpen grafiko batzuk, indibiduen eta aldagaien interpretazioa eta berauen arteko erlazioak azaltzera eramango gaitu.

Demagun, adibidez, 100 automobil motaren ezaugarriei (motorra, leungailua, prezioa, garantia, konponketa zerbitzua...) dagozkien 20 aldagaien datuak jaso ditugula. Balio hauen (100 x 20) datu-matrizeak informazio aberatsa dauka baina automobilen arteko antzekotasunak eta ezberdintasunak ez dira erraz ikusiko.

Indibiduoak (automobilak) aldagai ezberdinetan hartzen dituzten balioen arabera parekatzea interesatzen zaigu, baina, aldagaien artean erlazioak egongo dira eta erlazio horietaz baliatuko gara daukaten dimentsioa murrizteko.

VIII.2. ANALISIA \mathbb{R}^2 ESPAZIOAN. DOIKUNTZA ORTOGONALAREN ZUZENA

Kasu honetan, bi aldagaien planotik bi ardatz nagusien planora iritsiko gara, lehen ardatzean hodeiaren informazio maximoa, posible den neurrian, proiektatzen delarik.

Lehen ardatz honi Doikuntza Ortogonalaren Zuzena (D.O.Z.) deituko zaio.

Jatorrizko datutatik abiatuz (x_{ij}), non $i = 1, \dots, m$ eta $j = 1, 2$, datu zentratu eta birreskalatuen matrizea osatuko dugu, horrela hurrengo eragiketa matematikoak erraztuko ditugu.

Honek (V.A.1.) ez ditu emaitzak aldatzen, egingo dugun analisia konparatiboa baita. Interesatzen zaiguna, aldagai ezberdinetan hartzen dituzten posizio erlatiboak aurkitzea da eta egindako transformazioak ez dauka eraginik.

Hau da, \mathbf{X} datu-matrizea izango da:

$$\mathbf{X} = \left[\frac{x_{ij} - \bar{x}_j}{\sqrt{m}} \right]$$

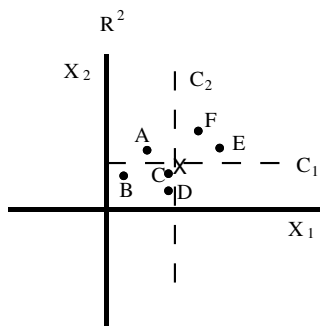
Gure helburua, puntu-hodeia dimentsio gutxi batzuetan proiektatzea da, honela, aurkezpen grafiko eskurakoiak edukiko ditugularik. Horretarako, metodo matematiko baten bidez, n aldagaien sistema batetik n ardatz berri dituen sistema batera iritsiko gara, hauek ordena behekorra jarraituz, puntu-hodeiaren informazioa maximizatzen dutelarik.

Doikuntza Ortogonalaren kasuan, datu-taularen informazioa murrizteak ez dauka zentzu askorik, kasu honetan, metodoa bi aldagai bakarrik dituen kasurako konkretatzen baitugu. Honela, gai honen helburu bakarra, pedagogikoa da, eta n aldagaien kasurako egingo dugun analisi orokorra ulergarriagoa izango zaigu.

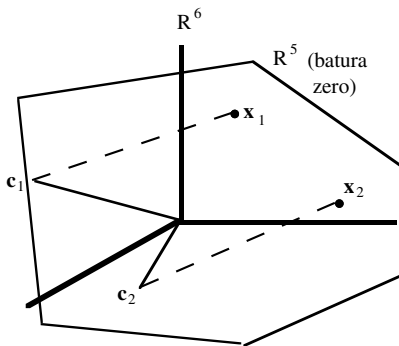
Hurrengo bost grafiken bidez, sei automobil mota direla suposatuz, prozesua laburtzen dugu, grafikoki.

JATORRIZKO ARDATZAK

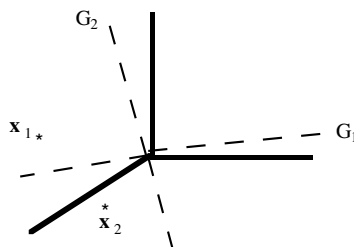
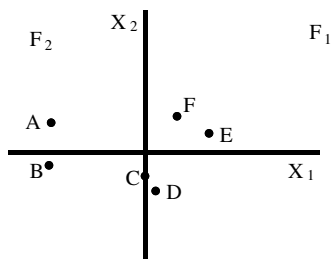
INDIBIDUOEN PUNTU-HODEIA



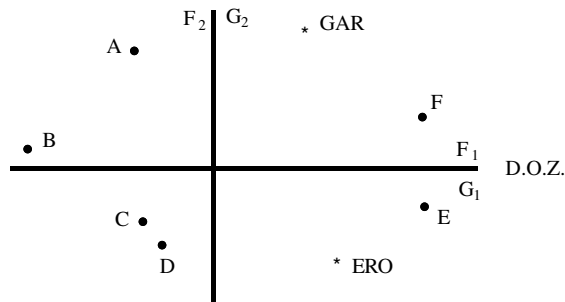
ALDAGAIEN PUNTU-HODEIA



ARDATZ BERRIAK



BATERAKO AURKEZPEN GRAFIKOA



$D.O.Z.$ lortu aurretik puntu-hodeiaren inertzia definituko dugu.

Definizioa: Puntu-hodei baten inertzia edozein punturekiko, erreferentzi puntu horretara, hodeiaren puntuen distantzi karratuen batura da.

Hodeia zentratuak dagoenez, inertzia jatorriarekiko edo grabitate-zentruarekiko kontsideratuko dugu. Hots:

$$I = \sum_{i=1}^m d^2(0, M_i) = \sum_{i=1}^m \left[\left(\frac{x_{i1} - \bar{x}_1}{\sqrt{m}} \right)^2 + \left(\frac{x_{i2} - \bar{x}_2}{\sqrt{m}} \right)^2 \right]$$

Aurrerantzean, $x_{i1} - \bar{x}_1 / \sqrt{m}$ eta $x_{i2} - \bar{x}_2 / \sqrt{m}$ balioak x_{i1} eta x_{i2} izango dira, hau da, balio zentratuak eta birreskalatuak.

Honela, $M_i(x_{i1}, x_{i2})$, $i = 1, \dots, m$ puntu-hodeiaren puntuak izango dira.

Dakusagunez:

$$I = S_1^2 + S_2^2$$

Inertzia bi bariantzen batura da.

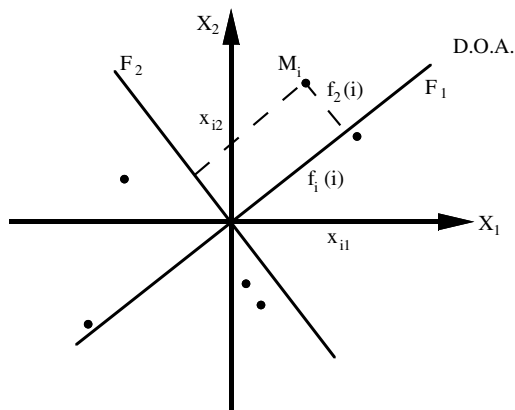
Inertzia, orduan, hodeiaren sakabanatze neurria da bere grabitate-zentruarekiko, eta puntu-hodeiaren informazioaz hitz egiten dugunean, datuen artean dagoen sakabanatzeaz ari gara.

F zuzen baten gainean edozein puntu proiektatzean, proiektzioaren $f(i)$ balioa lortzen dugu. Gure helburua, hodeian hoberen doitzen den F_1 zuzena lortzea da;

hau da, puntu guztietatik hurbilen dagoena. Datuak zentratuak daudenez, zuzen horrek jatorritik iragan behar du, hau da, grabitate-zentrutik.

Hurrengo grafikoa ikusirik, ondoko erlazioa lortzen dugu:

$$d^2(0, M_i) = f_1^2(i) + f_2^2(i)$$



Honela, inertzia adieraz daiteke:

$$I = \sum_{i=1}^m [f_1^2(i) + f_2^2(i)]$$

non $f_1(i)$, $i = 1, \dots, m$, F zuzenaren gaineko proiektzioak diren; $f_2(i)$ balioak F_2 zuzenaren gaineko proiektzioen balio berdinak dira eta F_1 zuzenaren gainean proiektatzean dauzkagun erroreak edo hondarrak bezala interpreta daitezke.

Honela:

$$I_{F_1} = \sum_{i=1}^m f_1^2(i) \quad I_{F_2} = \sum_{i=1}^m f_2^2(i)$$

Eta hodeiaren inertzia totala F_1 ardatzaren gain proiektatutako inertzian (I_{F_1}) eta F_2 ardatzaren gain proiektatutako inertzian (I_{F_2}) deskonposatzen da batukorki. Azken hori, hondar-inertzia bezala interpreta daiteke, F_1 -ek bildu ez duen informazioa biltzen baitu.

Hodeiaren inertzia (I) konstantea denez, hodeian hoberen doitzen den zuzena, hondar-inertzia minimoa edo proiektatutako inertzia maximoa egiten duen zuzena izango da.

VIII.3. DOIKUNTZA ORTOGONALAREN ZUZENAREN LORPENEA

Gure helburua proiektatutako inertziaren maximizazioa izango da.

Orokorki, zuzen baten unitate bektore zuzendaria \mathbf{u} baldin bada, edozein \mathbf{x} bektorearen proiektzioa zuzenaren gaineran $\mathbf{x}^T\mathbf{u}$ biderkadura eskalarraz lortuko da. Horrela, puntu-hodeiaren puntu guztien proiektzioak daukan bektorea \mathbf{f} , hain zuzen, lortuko da:

$$\mathbf{f} = \mathbf{X}\mathbf{u}$$

Ebatzi behar dugun problema $\left\{ \sum_{i=1}^m f^2(i) \right\}$ maximoa egitea da.

Aurreko erlazioa kontutan harturik, maximizazio problema adieraz daiteke:

$$\max_{\mathbf{u}} \mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u}$$

$$\mathbf{u}^T \mathbf{u} = 1 \text{ baldintzaz}$$

forma koadratiko baldintzatu baten maximizazioaren aurrean aurkitzen gara eta lagrangiarra eratu behar dugu.

$$\mathcal{L} = \mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} - \lambda (\mathbf{u}^T \mathbf{u} - 1) = \mathbf{u}^T \mathbf{L} \mathbf{u} - \lambda (\mathbf{u}^T \mathbf{u} - 1)$$

non, \mathcal{L} -ren maximoa lortu behar dugun.

Deribazio matriziala (V.A.7.) aplikatuz, maximizazioaren baldintza beharrezkoak lortzen ditugu:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = 2\mathbf{L}\mathbf{u} - 2\lambda\mathbf{u} = \mathbf{0} \Rightarrow \mathbf{L}\mathbf{u} = \lambda\mathbf{u}$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = -(\mathbf{u}^T \mathbf{u} - 1) = \mathbf{0} \Rightarrow \mathbf{u}^T \mathbf{u} = 1$$

Hau da, \mathbf{u}_1 soluzioa \mathbf{L} kobariantza matrizearen autobektore unitarioa da. \mathbf{L} matrize simetrikoa eta mugatu edo erdimugatu positiboa denez, $\lambda_1 \geq \lambda_2 \geq 0$ autobalio erreal eta positiboak daukala ziurtatuta dago.

Zein autobalioari dago lotuta \mathbf{u}_1 autobektore unitarioa?

Dakigunez \mathbf{u}_1 -ek inertzia proiektatua maximizatu behar du, orduan:

$$\mathbf{u}_1^T \mathbf{L} \mathbf{u}_1 = \lambda \quad \mathbf{u}_1^T \mathbf{u}_1 = \lambda = \lambda_1$$

Hots, \mathbf{L} -ren λ_1 autobalioarik handiena izango da. Horrela F_1 , D.O.Z.-ren direkzioa \mathbf{u}_1 autobektore unitarioaz finkatuta daukagu.

F_2 , \mathbf{L} -ren bigarren autodirekzioa \mathbf{u}_2 autobektore unitarioaz finkatuko da. \mathbf{u}_2 , ($\mathbf{u}_2^T \mathbf{u}_2 = 1$ eta $\mathbf{u}_2^T \mathbf{u}_1 = 0$) unitarioa eta \mathbf{u}_1 autobektorearekin ortogonalak izanik, λ_2 bigarren autobalioari dagokio.

F_1 eta F_2 ardatzek, ardatz-sistema berri bat osatzen dute.

Planoan daudenez D.O.Z. \mathbf{u}_2 bektorearekiko ortogonalak diren puntuen multzoaz adieraz daiteke. Izan bedi $\mathbf{u}_2 = (\alpha, \beta)^T$; zuzenaren adierazpena izango da:

$$\alpha \mathbf{X}_1 + \beta \mathbf{X}_2 = 0$$

$\mathbf{u}_2 = (\alpha, \beta)^T$ izanik F_1 -en autobektore unitario bat $\mathbf{u}_1 = (-\beta, \alpha)^T$ izango da eta hodeiaren m puntuen proiektzioak edo koordenatuak zuzenaren gainean izango dira:

$$f_1(i) = -\beta x_{i1} + \alpha x_{i2}$$

matrizialki, \mathbf{f}_1 proiektzio guztien bektorea izanik

$$\mathbf{f}_1 = \mathbf{X} \mathbf{u}_1$$

Bestalde, puntu bakoitzaren errorea edo hondarra, hau da, \mathbf{u}_2 bektoreaz definitutako zuzenaren gaineko proiektzioa, horrela adieraziko da:

$$f_2(i) = \alpha x_{i1} + \beta x_{i2}$$

eta matrizialki, \mathbf{f}_2 proiektzioen bektorea izanik:

$$\mathbf{f}_2 = \mathbf{X} \mathbf{u}_2$$

Propietateak:

1. \mathbf{f}_1 eta \mathbf{f}_2 proiektzioen bektoreak aldagai zentratuak dira.

Biak aldagai zentratuen konbinazio linealak izanik, zuzenki ondorioztatzen da.

2. Proiekzioen bektoreak elkarren artean koerlazio gabeko aldagaiak dira.

$$\mathbf{f}_2^T \mathbf{f}_1 = \mathbf{u}_2^T \mathbf{X}^T \mathbf{X} \mathbf{u}_1 = \mathbf{u}_2^T \lambda_1 \mathbf{u}_1 = \lambda_1 \mathbf{u}_2^T \mathbf{u}_1 = 0 \Rightarrow S_{F_2 F_1} = 0$$

3. Ardatz bakoitzaren gainean proiektatutako inertzia ardatzari dagokion autobalioa da.

$$I_{F_1} = \sum_{i=1}^m f_1^2(\mathbf{c}) = \mathbf{f}_1^T \mathbf{f}_1 = \mathbf{u}_1^T \mathbf{X}^T \mathbf{X} \mathbf{u}_1 = \mathbf{u}_1^T \lambda_1 \mathbf{u}_1 = \lambda_1$$

Era berean $I_{F_2} = \lambda_2$

Eta honela idatz dezakegu:

$$\text{tr}(\mathbf{L}) = S_1^2 + S_2^2 = \mathbf{I} = I_{F_1} + I_{F_2} = \lambda_1 + \lambda_2 = \text{tr}(\mathbf{\Lambda})$$

non, $\mathbf{\Lambda}$ matrizea, autobalioek osatzen duten matrize diagonal den.

Ardatzen fidegarritasuna. Ardatz bakoitzaren fidegarritasuna, beronen gain proiektatutako inertziaren arabera definituko dugu.

Puntu-hodeian D.O.Z. gero eta hobeago doitzen bada, inertzia totalarekiko proiektatutako inertziaren proportzioa gero eta handiagoa izango da. Doikuntza ezin hobea bada (puntuak lerrokaturik daude) inertzia guztia proiektatuta izango da eta hondar-inertziarik (edo bigarren ardatzan proiektaturik) ez daukagu.

Definizioa. τ fidegarritasunak, proiektatutako inertzia, inertzia totalarekiko bezala mugatzen du.

Hots:

$$\tau_1 = \frac{I_{F_1}}{\mathbf{I}} = \frac{\lambda_1}{S_1^2 + S_2^2} = \frac{\lambda_1}{\lambda_1 + \lambda_2}$$

$$\tau_2 = \frac{I_{F_2}}{\mathbf{I}} = \frac{\lambda_2}{S_1^2 + S_2^2} = \frac{\lambda_2}{\lambda_1 + \lambda_2}$$

Ehunekotan emanda, bi ardatzeko fidegarritasunen batura % 100 izango da.

VIII.4. HODEI DUALAREN AURKEZPEN GRAFIKOA.

Indibiduen puntu-hodeiaren analisisan, lortu ditugun ardatzak aldagaien arabera interpretatzea gustatuko litzaiguke.

Azter dezagun aldagaien puntu-hodeia edo hodei duala, indibiduen hodeiarekin egin dugun bezalaxe. Gure ikerketaren objektua hodei duala da, hau da, \mathbf{X}^T matrizearen errenkadak, aldagai-puntuak \mathbb{R}^{m-n} dauden bi puntu dira (X_1 eta X_2). Dakigunez, $\sqrt{m} S_j$, aldagai zentratuaren norma da eta zentratuz gain aldagaia berreskalatzen badugu, aldagaiaren norma desbidazio tipikoarekin bat dator, horrela hodei dualaren inertzia:

$$I = \|\mathbf{x}_1\|^2 + \|\mathbf{x}_2\|^2 = S_1^2 + S_2^2$$

Dakusagunez, indibiduo-puntuen inertzia berdina. Eraitza hori ez da harritzekoa, azken finean, bi hodeiak datu-matrize berberan daude eta informazio berbera daukate.

Berriro, hoberen doitzen den zuzena, G_1 izendatuko dugu, proiektatutako inertzia maximizatzen duena izango da. G edozein zuzenaren \mathbf{v} unitate bektore zuzendaria izanik, aldagaien proiektzioak G -ren gainean bektorialki honela lortuko ditugu:

$$\mathbf{g} = \mathbf{X}^T \mathbf{v}$$

Maximizazio problema adieraz daiteke:

$$\max_{\mathbf{v}} \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$$

$$\mathbf{v}^T \mathbf{v} = 1 \quad \text{baldintzaz}$$

Lagrangiarra izango da:

$$\mathcal{L} = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} - \lambda (\mathbf{v}^T \mathbf{v} - 1)$$

Aurreko egoera berdinean gaude, orain $\mathbf{v} \in \mathbb{R}^m$ eta ikasketaren matrizea $\mathbf{X} \mathbf{X}^T$ da. \mathbf{v}_1 soluzioa, $\mathbf{X} \mathbf{X}^T$ matrizearen autobektore unitarioa, autobalio handienari dagokiona da. Autobektore horrek G_1 zuzena mugatzen du.

Dakigunez, λ , $\mathbf{X}^T \mathbf{X}$ matrizearen autobalioa denez:

$$\mathbf{X}^T \mathbf{X} \mathbf{u} = \lambda \mathbf{u} \Rightarrow \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{u}) = \lambda (\mathbf{X} \mathbf{u})$$

Hau da, \mathbf{X} matrizeaz aurrebiderkatuz, λ , \mathbf{XX}^T matrizearen autobalioa dela ikus dezakegu eta \mathbf{XX}^T , $\mathbf{X}^T\mathbf{X}$ matrize simetrikoak, ezberdin zero diren autobalio berberak dituzte. Autobalioak $\lambda_1 \geq \lambda_2 \geq 0$ izango dira; \mathbf{XX}^T matrizeak, gainera, berdin zero diren $m-2$ autobalio ditu, bi matrizeen heina gehienez 2 baita, hots, $r(\mathbf{XX}^T) = r(\mathbf{X}^T\mathbf{X}) \leq 2$.

Hodei dualaren proiektzioak G_1 ardatzaren gainean izango dira:

$$\mathbf{g}_1 = \mathbf{X}^T \mathbf{v}_1$$

non, \mathbf{v}_1 , \mathbf{XX}^T matrizearen autobektore unitarioa, λ_1 autobalioari dagokiona den.

Izan bedi \mathbf{v}_2 , λ_2 autobalioari dagokion autobektore unitarioa. G_2 ardatzaren gaineko proiektzioak izango dira:

$$\mathbf{g}_2 = \mathbf{X}^T \mathbf{v}_2$$

\mathbf{XX}^T simetrikoa denez, lortutako ardatzak ortogonalak dira.

Propietateak:

1. Aldagaien proiektzioen bektoreen biderkadura eskalarra zero da.

$$\mathbf{g}_1^T \mathbf{g}_2 = \mathbf{v}_1^T \mathbf{XX}^T \mathbf{v}_2 = \mathbf{v}_1^T \lambda_2 \mathbf{v}_2 = \lambda_2 \mathbf{v}_1^T \mathbf{v}_2 = 0$$

2. Ardatz bakoitzaren gainean proiektatutako inertzia ardatzari dagokion autobalioa da.

$$I_{G_1} = g_1^2(1) + g_1^2(2) = \mathbf{g}_1^T \mathbf{g}_1 = \mathbf{v}_1^T \mathbf{XX}^T \mathbf{v}_1 = \mathbf{v}_1^T \lambda_1 \mathbf{v}_1 = \lambda_1$$

$$I_{G_2} = g_2^2(1) + g_2^2(2) = \mathbf{g}_2^T \mathbf{g}_2 = \mathbf{v}_2^T \mathbf{XX}^T \mathbf{v}_2 = \mathbf{v}_2^T \lambda_2 \mathbf{v}_2 = \lambda_2$$

Dakusagunez hodei dualaren ardatzen fidegarritasunak eta indibiduen puntu-hodeian lortu genituen ardatzenak berdinak dira. Hodei dualean, $m-2$ autodirekzio gehiago ditugu, baina, autobalio nuluei loturikoak direnez ez dute inertziarik jasotzen eta baztertu egingo ditugu.

VIII.5. TRANTSIZIO ERLAZIOAK

Ondoren, bi hodeien arteko proiektzioen eta bektore zuzendarien arteko erlazioak ikusiko ditugu, honela, ardatz berrien funtzioan indibiduen eta aldagaien arteko erlazioak aztertu ahal izango ditugu.

\mathbf{XX}^T -ren λ_1 autobalioari dagokion autobektore unitarioa \mathbf{v}_1 denez:

$$\mathbf{XX}^T \mathbf{v}_1 = \lambda_1 \mathbf{v}_1$$

\mathbf{X}^T matrizeaz aurrebiderkatuz eta $\mathbf{g}_1 = \mathbf{X}^T \mathbf{v}_1$ dela jakinik, honako hau daukagu:

$$\mathbf{X}^T \mathbf{X} \mathbf{g}_1 = \lambda_1 \mathbf{g}_1$$

Eta dakusagunez, \mathbf{g}_1 , $\mathbf{X}^T \mathbf{X}$ matrizearen λ_1 autobalioari lotutako autobektorea da, orduan, \mathbf{u}_1 autodirekzioan egongo den autobektore ez unitarioa, hain zuzen. $\|\mathbf{g}_1\| = \sqrt{\lambda_1}$ bere norma denez, idatz dezakegu:

$$\left. \begin{array}{l} \mathbf{g}_1 = \sqrt{\lambda_1} \mathbf{u}_1 \\ \mathbf{g}_2 = \sqrt{\lambda_2} \mathbf{u}_2 \end{array} \right\} \begin{array}{l} \text{TRANTSIZIO} \\ \text{ERLAZIOAK} \end{array}$$

non, bigarren ekuazioko, \mathbf{g}_2 eta \mathbf{u}_2 -ren arteko erlazioa, analogikoki, idatzi dugun.

Era berean, $\mathbf{X}^T \mathbf{X}$ -en λ_1 autobalioari dagokion autobektore unitarioa \mathbf{u}_1 denez:

$$\mathbf{X}^T \mathbf{X} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

\mathbf{X} matrizeaz aurrebiderkatuz eta $\mathbf{f}_1 = \mathbf{X} \mathbf{u}_1$ dela jakinik, honako hau daukagu:

$$\mathbf{XX}^T \mathbf{f}_1 = \lambda_1 \mathbf{f}_1$$

Eta dakusagunez, \mathbf{f}_1 , \mathbf{XX}^T matrizearen λ_1 autobalioari lotutako autobektorea da, orduan, \mathbf{v}_1 autodirekzioan egongo den autobektore ez unitarioa, hain zuzen. $\|\mathbf{f}_1\| = \sqrt{\lambda_1}$ bere norma denez, idatz dezakegu:

$$\left. \begin{array}{l} \mathbf{f}_1 = \sqrt{\lambda_1} \mathbf{v}_1 \\ \mathbf{f}_2 = \sqrt{\lambda_2} \mathbf{v}_2 \end{array} \right\} \begin{array}{l} \text{TRANTSIZIO} \\ \text{ERLAZIOAK} \end{array}$$

non, bigarren ekuazioko, \mathbf{f}_2 eta \mathbf{v}_2 -ren arteko erlazioa, analogikoki, idatzi dugun.

Laburki esanda, \mathbb{R}^2 edo \mathbb{R}^m espazioko autobektoreak, \mathbb{R}^m edo \mathbb{R}^2 bere hodei dualeko puntuen proiektzioen bektoreekin elkartzen dira.

$\mathbf{u}_1 = (-\beta, \alpha)^T$ eta $\mathbf{u}_2 = (\alpha, \beta)^T$ direla jakinik, aldagaien proiektzioak, honela lortuko ditugu:

$$\mathbf{g}_1 = \sqrt{\lambda_1} \begin{pmatrix} -\beta \\ \alpha \end{pmatrix}; \quad \mathbf{g}_2 = \sqrt{\lambda_2} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

VIII.6. BATERAKO AURKEZPEN GRAFIKOA ETA INTERPRETAZIOA

Bi hodeiak aztertzean lortutako emaitzak, esangura berbera daukate, nahiz eta, indibiduen funtzioan batean eta aldagaien funtzioan bestean adieraziak izan.

Adibidez, \mathbf{f}_1 indibiduen proiektzioen bektorea kontsideratuz, trantsizio erlaziotik lortzen dugu:

$$\mathbf{f}_1 = \mathbf{X}\mathbf{u}_1 = \frac{1}{\sqrt{\lambda_1}} \mathbf{X}\mathbf{g}_1$$

Hau da:

$$\begin{pmatrix} f_1(1) \\ \vdots \\ f_1(i) \\ \vdots \\ f_1(m) \end{pmatrix} = \frac{1}{\sqrt{\lambda_1}} \begin{pmatrix} x_{11} & x_{12} \\ \vdots & \vdots \\ x_{i1} & x_{i2} \\ \vdots & \vdots \\ x_{m1} & x_{m2} \end{pmatrix} \begin{pmatrix} g_1(1) \\ g_1(2) \end{pmatrix}$$

Orduan:

$$f_1(i) = \frac{1}{\sqrt{\lambda_1}} (x_{i1}g_1(1) + x_{i2}g_1(2))$$

Dakusagunez, indibiduo baten proiektzioa F_1 -en gainean, G_1 -en gaineko aldagaien proiektzioen konbinazio lineala da. Konbinazio linealaren koefizienteak indibiduo bakoitzarako aldagaien balio zentratuak dira. Indibiduoak aldagai baterako hartzen duen balioa, batezbestekoa baino handiagoa denean, koefiziente positiboa da eta txikiagoa denean, negatiboa. Hau da, grafiko berean indibiduen eta aldagaien proiektzioak adieraziz, honako hau daukagu: indibiduo bakoitza, batezbestekoa baino balio handiagoak hartzen dituen aldagaien alde berean proiektatuko da. Hauxe da, baterako aurkezpen grafikoaren arrazoi.

F_1, G_1 eta F_2, G_2 izaeraz ezberdinak badira ere (F_1 -en bektore zuzendaria $\mathbf{u}_1 \in \mathbb{R}^2$ eta G_1 -en bektore zuzendaria $\mathbf{v}_1 \in \mathbb{R}^m$) batera aurkeztuko ditugu aldagaiek indibiduoak eta indibiduoek aldagaiak interpretatzera laguntzen baitute. Indibiduen arteko distantziak ikusiz beraien arteko antzekotasunak eta ezberdintasunak ikus daitezke; halaber, aldagaiekin batera, badakigu grafikoki ere adieraztean zeintzuk diren antzekotasun edo ezberdintasun horien kausa.

VIII.7. DOIKUNTZA ORTOGONALAREN ZUZENA ALDAGAI NORMATUENTZAKO

Aldagaien eskala ezberdina izanik, analisisian eragina eduki ez dezan, analisisia aldagai normatutarako egingo da. Kasu honetan datu-matrize normatuaz abiatuz doikuntza ortogonalaren zuzenarako adierazpen orokorra daukagu, halaber, hondar-ardatzaren edo bigarren ardatzaren adierazpena, eta ondorioz, indibiduo-eta aldagai-puntuaren proiektiotarako.

Kasu honetan:

$$\mathbf{X} = \left(\frac{x_{ij} - \bar{x}_j}{S_j \sqrt{m}} \right) \quad \text{non } i = 1, \dots, m \text{ eta } j = 1, 2$$

$\mathbf{X}^T \mathbf{X}$ diagonalizatzen den matrizea, koerlazio-matrizea da eta bi aldagai besterik ez ditugunez, $|\mathbf{R} - \lambda \mathbf{I}| = 0$ ekuazio karakteristikoa, honako hau da:

$$\begin{vmatrix} 1 - \lambda & r \\ r & 1 - \lambda \end{vmatrix} = 0$$

Ekuazioa ebatziz, λ -ren bi balioak: $(1+r)$ eta $(1-r)$ dira.

Lehen kasua: $r > 0$

Autobaliorik handiena $\lambda_1 = 1+r$ izango da eta $\lambda_2 = 1-r$ txikiena.

α eta β balioak lortzeko $\mathbf{R} \mathbf{u}_2 = \lambda_2 \mathbf{u}_2$ ebatziko dugu:

$$\begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = (1-r) \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

Hortik $\alpha = -\beta$. Autobektoreak unitarioa izan behar duenez:

$$\alpha = 1 / \sqrt{2} = 0.707 \quad \text{eta} \quad \beta = -1 / \sqrt{2} = -0.707$$

Honela, D.O.Z.-ren adierazpen orokorra:

$$0.707 \mathbf{X}_1 - 0.707 \mathbf{X}_2 = 0$$

lehen eta hirugarren koadranteen erdikariari dagokion zuzena, hain zuzen ere.

Autobektore unitarioak, honako hauek dira:

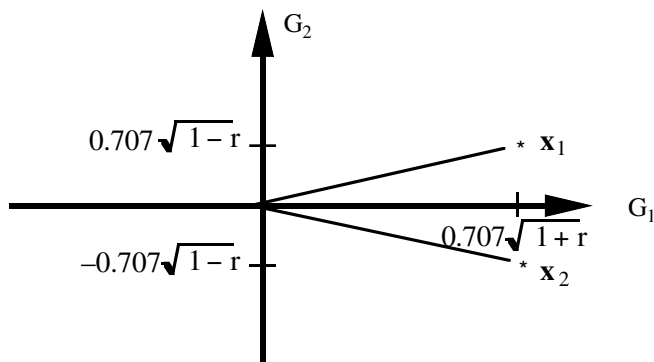
$$\mathbf{u}_1 = \begin{pmatrix} 0.707 \\ 0.707 \end{pmatrix} \quad \mathbf{u}_2 = \begin{pmatrix} 0.707 \\ -0.707 \end{pmatrix}$$

Eta indibidueno (i=1,...,m) ardatzen gaineko proiekzioak:

$$\begin{aligned} f_1(i) &= 0.707 x_{i1} + 0.707 x_{i2} \\ f_2(i) &= 0.707 x_{i1} - 0.707 x_{i2} \end{aligned}$$

Aldagaien proiekzioak berriz:

$$\begin{aligned} g_1(1) &= 0.707\sqrt{1+r} & g_2(1) &= 0.707\sqrt{1-r} \\ g_1(2) &= 0.707\sqrt{1+r} & g_2(2) &= -0.707\sqrt{1-r} \end{aligned}$$



Analitikoki nahiz grafikoki ikus dezakegunez, aldagaiak G_1 lehen ardatzarekiko simetrikoak dira. $r \sim 1$ baldin bada $g_1(j) \sim 1$ izango da eta aldagai puntuak G_1 -gaineko proiekzioetatik oso hurbil egongo dira; eta $r = 1$ baldin bada batera etorriko dira. $r \sim 0$ baldin bada bi ardatzetan aldagai bakoitzaren proiekzioak berdintsuak izango dira eta $r = 0$ baldin bada berdinak. Ondoren, $g_k(j)$ proiektzioa j. aldagaiaren eta k. ardatzaren arteko koerlazioa dela ikusiko dugu.

Trantsizio erlaziotik lortzen dugu:

$$\mathbf{g}_k = \mathbf{X}^T \frac{\mathbf{f}_k}{\sqrt{\lambda_k}} \quad k = 1, 2$$

Orduan:

$$g_k(j) = \mathbf{x}_j^T \frac{\mathbf{f}_k}{\sqrt{\lambda_k}} = r_{\mathbf{x}_j, \mathbf{f}_k} \quad k = 1, 2 \quad \text{eta} \quad j = 1, 2$$

Bigarren kasua: $r < 0$.

Kasu honetan, ondoko emaitzak lortzen ditugu:

$$\lambda_1 = 1 - r \quad , \quad \lambda_2 = 1 + r$$

D.O.Z. $0.707 \mathbf{X}_1 + 0.707 \mathbf{X}_2 = 0$ da; indibiduen proiektzioak:

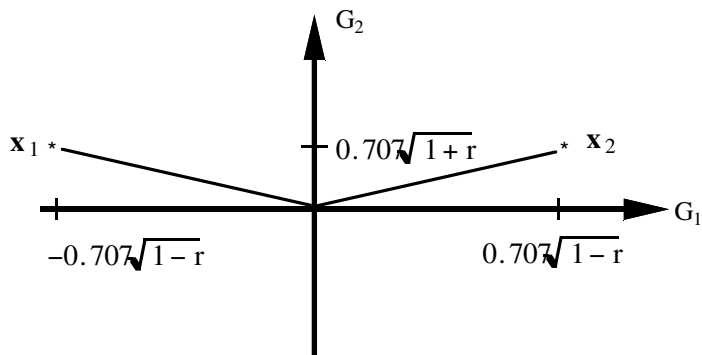
$$f_1(i) = -0.707 x_{i1} + 0.707 x_{i2}$$

$$f_2(i) = 0.707 x_{i1} + 0.707 x_{i2}$$

eta aldagaien proiektzioak:

$$g_1(1) = -0.707\sqrt{1-r} \quad g_2(1) = 0.707\sqrt{1+r}$$

$$g_1(2) = 0.707\sqrt{1-r} \quad g_2(2) = 0.707\sqrt{1+r}$$



Kasu honetan, G_2 bigarren ardatzarekiko aldagaien posizioa simetrikoa da. Koerlazioa ($r \sim -1$) -etik hurbil baldin badago $g_1(j)$ proiektzioen tamainua 1etik hurbil egongo da eta ($r \sim 0$) koerlazio ahula baldin badago proiektzioen tamainua

Otik hurbilago egongo da konkretuki $r = 0$ baldin bada bi ardatzen gaineko proiektzioak aldagai bakoitzarako, balio absolututan, berdinak izango dira.

Bi kasuetan, hau da, $r > 0$ edo $r < 0$ denean aldagaien arteko koerlazioa gero eta indartsuagoa denean, gero eta errepresentagarriagoa, izango da lehen ardatza. Lehen kasuan ($r > 0$), lehen ardatzak bi aldagaietan balio altuak (batezbestekoa baino handiagoak) hartzen dituzten indibiduoak eta balio baxuak (batezbestekoa baino txikiagoak) hartzen dituztenak kontrajartzen ditu. Hau gertatzen denean, "neurri" efektuaren aurrean gaudela esaten da, ardatzak indibiduoak neurtzen baititu. Bigarren kasuan ($r < 0$) alde berean dauden inbiduek, alde horretan kokatzen den aldagairako balio altuak eta beste aldagairako balio baxuak hartzen dituzte.

Lehen ardatzaren fidegarritasuna $\tau_1 = (1 + |r|)/2$ da eta bigarren ardatzarena $\tau_2 = (1 - |r|)/2$.

IX. DATU ANIZKOITZEN ANALISIA

IX.1. SARRERA

IX.2. ANALISI OROKORRA

IX.2.1. \mathbb{R}^n espazioaren \mathbb{R}^q azpiespazio baten bidez egindako doikuntza

IX.2.2. \mathbb{R}^m espazioaren \mathbb{R}^q azpiespazio baten bidez egindako doikuntza

IX.2.3. \mathbb{R}^n eta \mathbb{R}^m espazioen arteko erlazioa

IX.2.4. X datu taularen berreraketa

IX.3. OSAGAI NAGUSIZKO ANALISIA

IX.3.1. \mathbb{R}^n espazioan egindako analisia

IX.3.2. \mathbb{R}^m espazioan egindako analisia

IX.3.3. Baterako aurkezpen grafikoak

IX.3.4. Adibidea

IX.4. KORRESPONDENTZI ANALISI

FAKTORIALA

IX.4.1. Hodeiak, masak eta distantziak

IX.4.2. \mathbb{R}^n espazioan egindako analisia

IX.4.3. \mathbb{R}^m espazioan egindako analisia

IX.4.4. \mathbb{R}^n eta \mathbb{R}^m espazioen arteko erlazioa

IX.4.5. Maiztasun-taularen berreraketa

IX.4.6. Interpretaziorako laguntzak

IX.4.7. Korrespondentzia Anitzeko Analisi Faktoriala

IX.5. ELEMENTU GEHIGARRIAK

IX.6. SAILKAPEN-METODOAK

IX.6.1. Goranzko Sailkapen Hierarkikoa

IX.6.1.1. Bariantzaren metodoa

IX.6.2. Adibidea

IX.1. SARRERA

Azken hogeitako edo hogeita hamar urteetan, konputagailu eta kalkulu automatikoarekin batera, estatistika deskribatzaile anizkoitza asko garatu da J.P. Benzecri-ren eskola frantsesaren.

Datu-Analisiaren ezagutzen diren metodo hauek, bi motakoak dira:

- Metodo faktorialak
- Sailkapen-metodoak

Datu-multzo handi berdinei aplikatzen zaizkie, berorien ikasketa egiterakoan osagarri bezala erabiltzen direlarik.

Liburu honen lehen gaiaren sarreran esaten genuenez, estatistika deskribatzailearen helburua datu-multzo handien ikasketak egitea da, haien barne egitura gardenduz.

Ikasketa horietan, metodo faktorialak eta sailkapen-metodoak, datu-taula handien analisi bakoitzean bata besteari lagunduz, dira urrunago eramaten gaituztenak.

Zenbakizko datuen bidez, errealitate anizkoitzaz partzialki ohartu ahal badugu ere, metodo estatistiko egokiak aplikatzen dituzten kalkulu-programen bidez, errealitatean dauden estrukturak aztertu ahal ditugu.

Laburki esanez, metodo faktorialak diagonalizazio-metodoak dira, hauen bidez datutan dauden ardatz nagusi batzuei buruzko bihurtzen direlarik.

Normalki, deskribatu nahi ditugun objektuak aurkezpen grafiko jarraietan kokatzen ditugu; ondoren, taldekatzen saiatzen gara puntu-hodeian dauden konstelazioak ikusi ahal izateko, eta sailkapen metodoen bidez eskurakoi bihurtzen dira.

Sailkapen metodoak, metodo algoritmikoak dira; egokia den algoritmo bat era errepikariaz aplikatuz oharpenen edo indibiduen sailak edo klaseak lortzen ditugu.

Gai honetan, metodo faktorialen oinarrian dagoen Analisi orokorra aurkezten da, ondoren, Osagai Nagusizko Metodoa eta Korrespondentzi Analisi

Faktoriala ikusiko dira, ondoren Korrespondentzia Anitzeko Analisi Faktorialaren ideia aurkezten da eta azkenik, sailkapen metodoaz arituz Bariantzaren Metodoa ere laburki aurkeztuko da.

IX.2. ANALISI OROKORRA

Atal honetan, analisi faktorialen doikuntza metodoek duten gune teorikoa azalduko dugu.

Edozein metrika edo Analisi Kanonikoa erabiltzeak aurkezpen dotoreagoa egiteko aukera emango liguke, baina oso utilak ez diren garapen matematikotan sartuz, Datu-Analisiaren oinarrian dagoen ideia nagusia ilunduta utziko genuke ikasleentzat.

Idea hauxe da: Datu-matrize batek, bi espazio bektorial desberdinetan, hodei itxurako aurkezpen grafikoak egiteko aukera ematen du eta hodei horietan egindako doikuntzak, erlazio sinple batzuen bidez loturik daude.

Izan bedi $\mathbf{X}(m,n)$ datu-matrizea. Bere elementu orokorraz adieraziz $\mathbf{X} = [x_{ij}]$ izango da eta $m.n$ balioez osaturik dago.

$m.n$ balio kopurua oso handia izan daiteke eta honen aurrean gure galdera hauxe da: aurkez ditzakegu $m.n$ balioak balio kopuru txiki baten bidez, taulan daukagun sakabanatze adieragarrienaz jabetuz?

Demagun \mathbf{u}_1 eta \mathbf{v}_1 , hurrenez hurren, m eta n osagai dituzten zutabe-bektoreak direla eta berorien bidez \mathbf{X} datu-matrizea $(m+n)$ balio horien bidez berrerratu ahal dugula.

$$\underset{(m,n)}{\mathbf{X}} = \underset{(m,1)}{\mathbf{u}_1} \underset{(1,n)}{\mathbf{v}_1^T}$$

Kasu honetan \mathbf{X} datu-matrizearen heina 1 da. Normalki, oso zaila izango da \mathbf{X} horrela berrerratea, orduan, q heinaren hurbilketa bilatuko da.

Hau da:

$$\mathbf{X} = \mathbf{u}_1 \mathbf{v}_1^T + \mathbf{u}_2 \mathbf{v}_2^T + \dots + \mathbf{u}_q \mathbf{v}_q^T + \mathbf{E}$$

$\mathbf{E}(m,n)$ hondar-matrizea izanik.

\mathbf{u}_k \mathbf{v}_k bektoreen $q(m+n)$ osagaien bidez \mathbf{X} matrizearen $m \cdot n$ osagaiak, modu egokian berrerratu edo berreraiki daitezke, \mathbf{E} matrizearen osagaiak txikiak eta baztergarriak baitira.

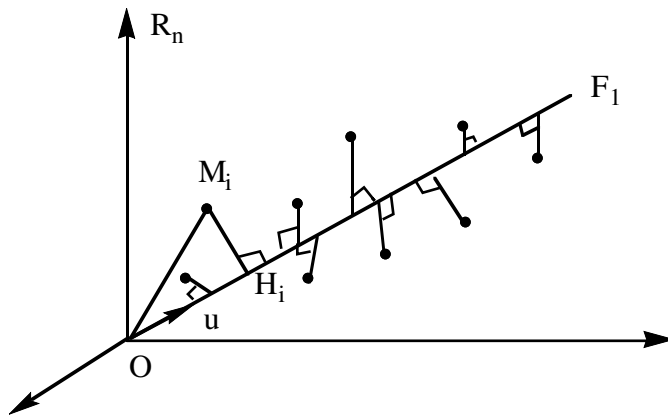
Aurreneko problema, metodo faktorialeko loturiko aurkezpen geometrikoen bidez ebaziko dugu:

\mathbf{X} datu-taulatik, dakigunez (ikus X.2) bi aurkezpen grafiko eratorzen dira, m errenkadak, \mathbb{R}^n espazio bektorialeko m punturen koordinatuak bezala kontsidera daitezke eta n zutabeak, \mathbb{R}^m espazio bektorialeko n punturen koordinatuak bezala. Hau da, errenkada-puntuen hodeia eta zutabe-puntuen hodeia \mathbb{R}^n eta \mathbb{R}^m espaziotan ditugarik.

IX.2.1. \mathbb{R}^n espazioaren \mathbb{R}^q azpiespazio baten bidez egindako doikuntza

Aurrerago ikusiko dugunez, \mathbb{R}^n espazioan ditugun m errenkada-puntuak, \mathbb{R}^q azpiespazio bateko ardatzetan dituzten koordinatuen bidez eta ardatz berri hauek dituzten osagaien bidez, berrera daitezke. Honela, taularen $m \cdot n$ balioak $q(m+n) = qm + qn$ balioen bidez berrerratuak geratuko dira.

Orduan, ohizko distantzia euklidearrak hornitutako \mathbb{R}^n espazioko \mathbb{R}^q azpiespazio baten bidez hodeia doitzea da gure helburua. Hasieran, hoberen doitzuz, jatorritik pasatzen den F_1 zuzena bilatuko dugu:



Izan bedi \mathbf{u} zuzen horren bektore unitarioa. Orduan,

$$\mathbf{u}^T \mathbf{u} = 1 \quad \text{eta} \quad \sum_{i=1}^n u_i^2 = 1$$

\mathbf{X} datu-taularen errenkadetan \mathbb{R}^n espazioko m puntuak ditugunez \mathbf{u} bektore unitarioaz atzebiderkatuz F_1 zuzenean puntu horien proiektzioak edo koordinatuak lortuko ditugu

$$\text{Hau da: } \begin{matrix} \mathbf{X} & \mathbf{u} & = & \mathbf{f} \\ (\mathbf{m}, \mathbf{n}) & (\mathbf{n}, 1) & & (\mathbf{m}, 1) \end{matrix}$$

\mathbf{f} bektorearen m osagaiak, F_1 zuzenaren gain lortzen dituzten luzerak dira.

\mathbb{R}^2 -espazioan planteatzen genuen doikuntza ortogonalaren gaian gertatzen zen bezalaxe (ikus VIII.2), hemen ere, \mathbb{R}^n espazioan, hain zuzen, puntu bakoitzetik jatorrira dagoen distantzia karratua bi distantzia karratutan deskonposatu da, hauek, F_1 zuzenera daukan e_i distantzia karratua eta F_1 gainean daukan f_i proiektzioaren karratua direlarik.

$$\text{Irudian: } \overline{OM}_i^2 = \overline{MH}_i^2 + \overline{OH}_i^2 \quad \text{edo} \quad \overline{OM}_i^2 = e_i^2 + f_i^2$$

Ondorioz, puntuek ardatzera daukaten distantzia karratuak minimizatzea (karratu txikienen erizpidea) edo zuzenaren gain daukaten proiektzio karratuak maximizatzea batera ematen da.

Puntu guztiak batera hartuz:

$$\sum_i \overline{OM}_i^2 = \sum_i \overline{MH}_i^2 + \sum_i \overline{OH}_i^2$$

$$\sum_i e_i^2 \text{ minimoa edo } \sum_i f_i^2 \text{ maximoa bilatzea baliokidea da.}$$

Arazoa, bada, honela aurkez daiteke:

Zein da, proiektzio $\mathbf{f}^T \mathbf{f} = \sum_i f_i^2$ karratuen batukaria maximoa egiten duen F_1 ardatzaren \mathbf{u}_1 bektore unitarioa?.

$$\text{Hots: } \mathbf{f}^T \mathbf{f} = \mathbf{u}_1^T \mathbf{X}^T \mathbf{X} \mathbf{u}_1 \text{ maximoa } \mathbf{u}_1^T \mathbf{u}_1 = 1 \text{ baldintzaz.}$$

Honela, hodeia hoberen doitzen duen dimentsio bateko azpiespazioa lortzen da, hau da, \mathbf{u}_1 -ek sortarazten duen azpiespazioa, $\mathbf{u}_1^T \mathbf{X}^T \mathbf{X} \mathbf{u}_1$ forma koadratikoa maximoa egiten duenlarik.

Hodeia hoberen doitzen duen bi dimentsioko azpiespazioak \mathbf{u}_1 daukala, absurdora eramanez erraz froga daiteke: \mathbf{u}_1 ez balu, \mathbf{u}_1 edukiz hobe doituko lukeen bi dimentsioko beste azpiespazioa lortuko genuke.

Honela, \mathbf{u}_2 bektore unitarioa eta \mathbf{u}_1 -ekiko ortogonal ($\mathbf{u}_2^T \mathbf{u}_2 = 1$ eta $\mathbf{u}_2^T \mathbf{u}_1 = 0$) lortuz, $\mathbf{u}_2^T \mathbf{X}^T \mathbf{X} \mathbf{u}_2$ forma koadratikoa maximoa egiten duen bi dimentsioko azpiespazioa sortaraziko da.

Analogikoki, era errepikariaz arrazonatuz zera ikusiko dugu: $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q$ bektore unitarioek, non \mathbf{u}_q , aldez aurretik lortutako $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{q-1}$ bektoreekiko ortogonal baita, $q \leq n$ dimentsioko azpiespazioa sortarazten dute, $\mathbf{u}_q^T \mathbf{u}_q = 1$ baldintzaz, $\mathbf{u}_q^T \mathbf{X}^T \mathbf{X} \mathbf{u}_q$ forma koadratikoa maximoa egiten delarik.

Maximo baldintzatuaren ebazpidea: Izan bedi λ Lagrange-ren biderkatzailea, $\mathbf{u}^T \mathbf{u} = 1$ izanik, $\mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} - \lambda (\mathbf{u}^T \mathbf{u} - 1)$ adierazpenaren maximoa lortu behar dugu.

Dakigunez, optimizazio-baldintza beharrezkoak aplikatuz, \mathbf{u} bektorearen osagai bakoitzarekiko deribatuak zero egin behar dira.

Osagai bakoitzarekiko deribatuak eginez, zero balioari berdinduz, eta gero ditugun ekuazioak matrizialki idatziz ondoko ekuazio matriziala daukagu: (berdin litzateke matrizialki deribatzea)

$$2 \mathbf{X}^T \mathbf{X} \mathbf{u} - 2 \lambda \mathbf{u} = 0$$

$$\text{Hau da:} \quad \mathbf{X}^T \mathbf{X} \mathbf{u} = \lambda \mathbf{u} \quad (1)$$

Dakusagunez, λ balioa eta \mathbf{u} bektorea $\mathbf{X}^T \mathbf{X}$ matrize simetrikoaren autobalioa eta berari loturiko autobektorea, hurrenez hurren, dira.

(1) berdintasuna, \mathbf{u}^T bektoreaz aurrebiderkatuz:

$$\mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} = \lambda$$

Hau da, bilatzen ari garen maximoa λ autobalioa da. Honela, \mathbf{u} bektorea λ_1 autobalio handienari dagokion $\mathbf{X}^T \mathbf{X}$ matrizearen \mathbf{u}_1 autobektorea izango da.

Hoberen doitzen den bi dimentsioko azpiespazioa lortu nahi badugu, bigarren zuzen bat behar dugu \mathbf{t} berorren bektore unitarioa izanik.

$\mathbf{t}^T \mathbf{X}^T \mathbf{X} \mathbf{t}$ maximoa $\mathbf{t}^T \mathbf{u}_1 = 0$ eta $\mathbf{t}^T \mathbf{t} = 1$ izanik lortu behar dugu.

Hau da: $\mathbf{t}^T \mathbf{X}^T \mathbf{X} \mathbf{t} - \lambda (\mathbf{t}^T \mathbf{t} - 1) - \mu \mathbf{t}^T \mathbf{u}_1$, bi biderkatzaile dituen lagrangearra idazten dugu.

Berriro ere, \mathbf{t} bektorearen osagaiekiko deribatu partzialak zero balioari berdinuz, ondoko ekuazio matriziala daukagu.

$$2 \mathbf{X}^T \mathbf{X} \mathbf{t} - 2 \lambda \mathbf{t} - \mu \mathbf{u}_1 = 0 \quad (2)$$

(2) berdintasuna \mathbf{u}_1^T bektoreaz aurrebiderkatuz, ondoren $\mathbf{u}_1^T \mathbf{X}^T \mathbf{X} = \lambda \mathbf{u}_1^T$ ordezkatzuz eta $\mathbf{u}_1^T \mathbf{t} = 0$, $\mathbf{u}_1^T \mathbf{u}_1 = 1$ direla jakinik $\mu = 0$ dela ikusten dugu.

Hau da: $\mathbf{X}^T \mathbf{X} \mathbf{t} = \lambda \mathbf{t}$

Dakusagunez, \mathbf{t} bektore unitarioa, $\mathbf{X}^T \mathbf{X}$ matrize simetrikoaren λ_2 bigarren autobalioari loturiko \mathbf{u}_2 bigarren autobektorea izango da. Prozesua jarraituz, $k = 1, 2, \dots, r$ autoelementuetara edo elementu propioetara zabaltzen da emaitza.

Honela, $\mathbf{X}^T \mathbf{X} \mathbf{u}_k = \lambda_k \mathbf{u}_k \quad \forall k \leq r \quad (3)$

r , $\mathbf{X}^T \mathbf{X}$ matrizearen heina izanik.

Azkenik, karratu txikiaren zentzuan hodeian hoberen doitzen den q dimentsioko azpiespazioaren oinarria, $\mathbf{X}^T \mathbf{X}$ matrize simetrikoaren q handien autobalioei loturiko q autobektoreez osaturik dago.

$\mathbf{X}^T \mathbf{X}$ matrizea, simetrikoa eta definitu edo erdidefinitu positiboa izanez, hauek dituzten propietateetatik abiatuz, beste frogapen baten bidez, aurreneko emaitzara iritsiko ginateke.

IX.2.2 \mathbb{R}^m espazioaren \mathbb{R}^q azpiespazio baten bidez egindako doikuntza

Era berean, \mathbb{R}^m espazioan ditugun n zutabe-puntuak \mathbb{R}^q azpiespazioko ardatzetan dituzten koordinatuen bidez eta ardatz berri hauek dituzten osagaien bidez berrera daitezke.

Izan bedi G_1 jatorritik pasatzen den zuzena eta \mathbf{v} beronen direkzio-kosinuen bektorea; G_1 zuzenaren gainean puntuen proiektzioen karratuen batura maximoa egitean, hodeian hoberen doitzen den zuzena, karratu txikiaren zentzuan lortzen dugu.

$$\text{Hau da: } \mathbf{X}^T \mathbf{v} = \mathbf{g}$$

(n,m)(m,1) (n,1)

\mathbf{g} zutabe-bektorearen n errenkadak proiektzioen n balioak dira, orduan, $\mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$ forma koadratikoa, maximoa egin nahi dugun balioa, izango da.

$\mathbf{v}^T \mathbf{v} = 1$ baldintzaz, $\mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}$ maximoa egiten duen \mathbf{v} bektorea aurkituz, problema ebatzia edukiko dugu.

\mathbb{R}^n espazioaren kasuan bezala, matrize simetriko baten diagonalizazioaren aurrean gaude, \mathbb{R}^m espazioan, $\mathbf{X} \mathbf{X}^T$ matrize simetrikoaren autobektorea \mathbf{v} baita.

μ_1 autobalio handienari dagokion autobektorea \mathbf{v}_1 izango da. Ondoren, $\mathbf{v}_2 \dots \mathbf{v}_q$ finkatuz, hodeian hoberen doitzen den q dimentsioko azpiespazioa sortarazten duten bektoreak lortu ditugu.

Dakigunez, $\mathbf{X}^T \mathbf{X}$ eta $\mathbf{X} \mathbf{X}^T$ matrizeak \mathbf{X} matrizearen heina berbera daukate r , hain zuzen, eta $q < r$ izango da.

IX.2.3 \mathbb{R}^n eta \mathbb{R}^m espazioen arteko erlazioa

\mathbf{v}_k eta \mathbf{u}_k bektoreen arteko erlazioak aurkitu nahi ditugu.

$$\text{Definizioz } \mathbf{X} \mathbf{X}^T \mathbf{v}_k = \mu_k \mathbf{v}_k \quad (4)$$

μ_k eta \mathbf{v}_k , $\mathbf{X} \mathbf{X}^T$ matrizearen k . autobalioa eta berari loturiko k . autobektorea, hurrenen hurren, izanik.

Diagonalizazioaren bidez, desberdin zero diren r autobalio lortuko dira eta $r \leq (m,n)$, hau da, \mathbf{X} matrizearen dimentsiorik txikiena baino txikiagoa edo berdina izango da r , \mathbf{X} matrizearen heina r delako.

(4) berdintasuna \mathbf{X}^T matrizeaz aurrebiderkatuz:

$$\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{v}_k) = \mu_k (\mathbf{X}^T \mathbf{v}_k)$$

Baina, $\mathbf{X}^T \mathbf{X}$ diagonalizazioan $\mathbf{X}^T \mathbf{X} \mathbf{u}_k = \lambda_k \mathbf{u}_k$ ((3) berdintasuna) daukagunez, $\mathbf{X} \mathbf{X}^T$ matrizearen \mathbf{v}_k bakoitzari ($k \leq r$) $\mathbf{X}^T \mathbf{X}$ matrizearen $\mathbf{X}^T \mathbf{v}_k$ dagokio μ_k autobalio berdinari lotuta.

Dakusagunez, $\mathbf{u}_k = h_k \mathbf{X}^T \mathbf{v}_k$ ($h_k =$ konstantea) idatz dezakegu bektore proportzionalak baitira.

λ_k k. autobalioa, handiena izatean, nahitaez:

$$\mu_k \leq \lambda_k \quad \text{izango da, eta } \{\mu_k\} \subset \{\lambda_k\}$$

Hots, μ_k autobalioen multzoa λ_k autobalioen multzoaren azpimultzoa da.

Era berean, (3) berdintasuna \mathbf{X} matrizeaz aurrebiderkatuz:

$$\mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{u}_k) = \lambda_k (\mathbf{X} \mathbf{u}_k)$$

Honela, $\mathbf{X}^T \mathbf{X}$ matrizearen \mathbf{u}_k bakoitzari $\mathbf{X} \mathbf{X}^T$ matrizearen $\mathbf{X} \mathbf{u}_k$ dagokio, λ_k autobalio berdinari lotuta.

Orduan, $\mathbf{v}_k = h_k' \mathbf{X} \mathbf{u}_k$ (h_k' konstantea) eta μ_k , $\mathbf{X} \mathbf{X}^T$ matrizearen k. autobalio handiena denez

$$\lambda_k \leq \mu_k \quad \text{eta } \{\lambda_k\} \subset \{\mu_k\}$$

Hots: $\lambda_k = \mu_k \quad \forall k \leq r$ idatz dezakegu.

$$\text{Dakigunez,} \quad \mathbf{u}_k^T \mathbf{u}_k = 1, \quad \mathbf{v}_k^T \mathbf{v}_k = 1$$

$$\text{Honela,} \quad \mathbf{v}_k^T \mathbf{v}_k = h_k'^2, \quad \mathbf{u}_k^T \mathbf{X}^T \mathbf{X} \mathbf{u}_k$$

eta $\mathbf{X}^T \mathbf{X} \mathbf{u}_k = \lambda_k \mathbf{u}_k$ ordezkatuz.

$$\mathbf{v}_k^T \mathbf{v}_k = h_k'^2, \quad \mathbf{u}_k^T \lambda_k \mathbf{u}_k = h_k'^2 \lambda_k = 1$$

$$h_k' = \frac{1}{\sqrt{\lambda_k}}$$

Analogikoki, $h_k = \frac{1}{\sqrt{\lambda_k}}$ dela lortzen da.

Orduan,

$$\mathbf{u}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{X}^T \mathbf{v}_k \quad (5)$$

(n,1) (n,m)(m,1)

$$\mathbf{v}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{X} \mathbf{u}_k \quad (6)$$

(m,1) (m,n)(n,1)

\mathbf{u}_k bektore unitarioaren euskarria den F_k ardatza, \mathbb{R}^n espazio bektorialaren k . ardatz faktoriala deitzen da eta \mathbf{v}_k bektore unitarioaren euskarria den G_k ardatza \mathbb{R}^m espazio bektorialaren k . ardatz faktoriala, hain zuzen ere, biek, λ_k autobalio berdinari dagokion autodirekzioa finkatzen dutelarik.

Dakusagunez, \mathbb{R}^n (edo \mathbb{R}^m)-ren k . ardatzaren gain hodeiaren puntuek dituzten koordinatuak (osagai nagusiak) $\mathbf{X} \mathbf{u}_k$ (edo $\mathbf{X}^T \mathbf{v}_k$) dira. Orduan, (5), (6) erlazioetan ikus daitekeenez, koordinatu horiek espazio batean eta ardatzaren osagai unitarioak beste espazioan proportzionalak dira.

IX.2.4 X datu taularen berreraiketa

Dakusagun, nola, \mathbb{R}^n eta \mathbb{R}^m espazioetan ditugun hodeien puntuak beraien posizioetan, gutxi gora-behera, berreratzen ditugun karratu txikiaren zentzuan hoberen doitzen den azpiespazio baten oinarrian dituzten koordinatuen bidez.

Suposa dezagun, λ_1 , lehen autobalioa, besteak baino askoz garrantzitsuagoa dela. Kasu honetan, $\lambda_1 = \mathbf{u}_1^T \mathbf{X}^T \mathbf{X} \mathbf{u}_1$ balioak proiektzioz jatorriarekiko distantzia karratuak neurtzen ditu eta lehen ardatzak puntuen posizioak ongi berreratuko ditu. Normalki, q lehen ardatz faktorialen bidez, “hurbilketa onaz”, puntuen posizioak berreratuko ditugu, horretarako $\lambda_1 + \lambda_2 + \dots + \lambda_q$ baturak, $\mathbf{X}^T \mathbf{X}$ matrizearen aztarnaren proportzio garrantzitsua aurkeztu behar du. $\mathbf{X}^T \mathbf{X}$ matrize simetrikoa eta, normalki, definitu positiboa izango denaren $n = r \leq m$ autobalioak errealak dira eta autobektoreak ortogonalak.

\mathbf{U} autobektore unitarioak osatutako matrize diagonalizatzailea eta $\mathbf{X}^T \mathbf{X}$ matrizea antzekoak dira, hau da, oinarri desberdinetan transformazio lineal berbera definitzen dute.

Hots, Λ matrizea, $\lambda_1, \lambda_2 \dots \lambda_n$ autobalioak osaturiko $\begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$ matrize diagonal

izanez, $\mathbf{X}^T \mathbf{X} \mathbf{U} = \mathbf{U} \Lambda$ eta $\mathbf{U}^T \mathbf{X}^T \mathbf{X} \mathbf{U} = \mathbf{U}^T \mathbf{U} \Lambda$ edota $\mathbf{U}^{-1} \mathbf{X}^T \mathbf{X} \mathbf{U} = \Lambda$, \mathbf{U} ortogonal baita.

Dakusagunez, bada, $\mathbf{X}^T \mathbf{X}$ matrizearen antzekotasunez transformatu diagonal da Λ matrizea eta honelako matrize antzekoen propietate bat aplikatuz¹:

$$\text{tr}(\mathbf{X}^T \mathbf{X}) = \text{tr} \Lambda$$

$$\text{tr}(\mathbf{X}^T \mathbf{X}) = \sum_i \sum_j x_{ij}^2 = \sum_k \lambda_k = \text{tr} \Lambda \quad (7)$$

Hasierako zenbakizko balioak, hau da, \mathbb{R}^n espazioan puntuak dituzten koordenatuak, q ardatzen bidez lortu nahi izanez gero, ardatz faktorialetan dituzten koordenatuak eta ardatzen osagai unitarioak (direkzio-kosinuak) batera erabiliz lortuko ditugu.

(6) erlazioa honela idaziz: $\mathbf{X} \mathbf{u}_k = \sqrt{\lambda_k} \mathbf{v}_k$

Ondoren, berdintasuna \mathbf{u}_k^T bektoreaz atzebiderkatuz:

$$\mathbf{X} \mathbf{u}_k \mathbf{u}_k^T = \sqrt{\lambda_k} \mathbf{v}_k \mathbf{u}_k^T$$

eta k balio guztientzako batukaria hartuz (autobalio nuluak balira beraiei dagozkien \mathbf{u}_k autobektoreek \mathbb{R}^n espazioaren oinarria osatuko lukete),

$$\mathbf{X} \sum_{k=1}^n \mathbf{u}_k \mathbf{u}_k^T = \sum_{k=1}^n \sqrt{\lambda_k} \mathbf{v}_k \mathbf{u}_k^T$$

$$\text{eta } \mathbf{X}_{(mn,n)} = \sum_{k=1}^n \sqrt{\lambda_k} \mathbf{v}_k \mathbf{u}_k^T_{(m,1)(1,n)} = \mathbf{u} \mathbf{u}^T = \mathbf{U}^T \mathbf{U} = \mathbf{I} \quad (\text{unitate matrizea})$$

q lehen ardatzak hartzen baditugu, lortzen dugun \mathbf{X}^* matrizea \mathbf{X} datu-matrizearen hurbilketa izango da.

1. Lau propietate dituzte, diskriminatzailea, autobalioak, aztarrena (tr) eta heina (r) berbera edukitzea, hain zuzen. Ikus FZ. DE TROCONIZ A. *Matemáticas Generales I*, 1.18.2

$$\mathbf{X} \approx \mathbf{X}^* = \sum_{k=1}^q \sqrt{\lambda_k} \mathbf{v}_k \mathbf{u}_k^T$$

$\sqrt{\lambda_{q+1}} \dots \sqrt{\lambda_n}$ balioak txikiak balira hurbilketa ona izango litzateke.

Berreraketaren kalitatea, ondoko proportzioaz neur daiteke:

$$\tau_q = \frac{\sum_{k=1}^q \lambda_k}{\sum_{k=1}^n \lambda_k}$$

\mathbf{X} datu-matrizearen $m \cdot n$ balioak q ($m+n$) = $q \cdot m + q \cdot n$ balioen bidez berrerratu ditugu $\sqrt{\lambda_k} \mathbf{v}_k$, q bektoreen m balioen bidez eta \mathbf{u}_k , q bektoreen n balioen bidez.

λ_k autobalio bakoitzak k . ardatzaren gain puntuen proiektzio karratuen batura neurtzen du, hau da, k . ardatzaren gain proiektatutako hodeiaren inertzia (ikus VIII.3 $n=2$ kasurako). Honela, inertzi-tasa deitutako τ_q balioak q lehen ardatzetan proiektatutako inertziaren proportzioa edo bariantzaren proportzioa ematen dizkigu.

Hoberen doitutako \mathbb{R}^q azpiespazioan biltzen den hodeiaren sakabanatzearen proportzioa, inertzi-tasak neurtzen digu.

IX.3. OSAGAI NAGUSIZKO ANALISIA

Osagai nagusizko analisia, aztertu dugun analisi orokorraren barrutian aurkezten dugu.

Osagai nagusizko analisia “aldagai-indibiduo” taula motatan aplikatuko da, kasu honetan, errenkada-puntuen hodeia, indibiduen hodeia eta zutabe-puntuen hodeia aldagaien hodeia izango direlarik.

Datu zentratutatik abiatu eta hodei batean nahiz bestean zentratze-eragiketarik suposatzen duena kontutan hartuz, osagai nagusizko analisia bi eratara planteatu daiteke.

Osagai Nagusizko Analisia

\mathbf{X} matrizearen gai orokorra $\frac{x_{ij} - \bar{x}_j}{\sqrt{m}}$ denean, diagonalizatzen den matrize simetrikoa $\mathbf{X}^T \mathbf{X} = \mathbf{L}$ kobariantza matrizea da.

Osagai Nagusizko Analisi Normatua.

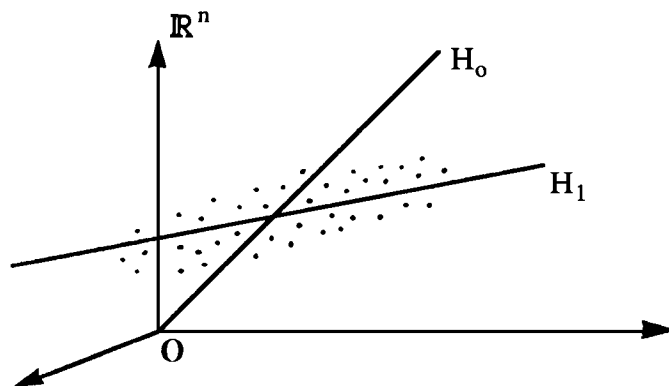
Aldagaiaren neurriak, batezbestekoen aldetik eta sakabanatzearen aldetik heterogenoak direnean \mathbf{X} matrizearen gai orokorra $\frac{x_{ij} - x_j}{S_j \sqrt{m}}$ izango da, kasu honetan diagonalizatzen den matrize simetrikoa $\mathbf{X}^T \mathbf{X} = \mathbf{R}$ koerlazio-matrizea da.

IX.3.1 \mathbb{R}^n espazioan egindako analisisia

Datu zentratutatik abiatzeak dituen abantailak hauexek dira:

- Diagonalizatzen den matrizea \mathbf{L} edo \mathbf{R} da.
- Proiekzio karratuen batura maximoa egiten duen lehen ardatza jatorritik pasatzen da; honela aztertu dugun analisi orokorraren barrutian azaldu ahal dugu osagai nagusizko analisisia.

Ondoko grafikoan ikus dezakegunez, hodeiaren grabitate-zentrua jatorrian ez balego, proiekzio karratuen batura maximoa egingo lukeen ardatza ez litzateke jatorritik pasatuko.



H_1 ardatzarentzako balio hori maximoa izango litzateke eta ez H_0 ardatzarentzako.

\mathbb{R}^n espazioan, indibiduen hodeian egindako analisisian, grabitate-zentrua jatorrira eramaten da edota ardatzen translazio paraleloa egiten da. $1/\sqrt{m}$ koefizientearen bidez, diagonalizatzen den $\mathbf{X}^T \mathbf{X}$ matrizea kobariantza matrizea izatea lortzen da.

Aldagaien sakabanatzeak oso desberdinak direnean, zentratzeaz gainera, beste aldaketa egitea beharrezkoa da, honela, osagai normatutako analisisan indibiduen arteko distantziak aldagai desberdinekiko, ondoan ikusten den bezala “eskalan berdindu ditugu”.

Hots:

$$\begin{aligned} d^2(i, i') &= \sum_{j=1}^n \left[\frac{x_{ij} - \bar{x}_j}{S_j \sqrt{m}} - \frac{x_{i'j} - \bar{x}_j}{S_j \sqrt{m}} \right]^2 = \\ &= \frac{1}{m} \sum (x_{ij} - x_{i'j})^2 / S_j^2 \end{aligned}$$

Dakusagunez, indibiduen hurbiltasunak kalkulatzean, aldagaiek ekarpen analogoa daukate, desbidazio tipikoen arabera hartu baititugu koordenatuak.

Kasu honetan, dakigunez, diagonalizatzen den $\mathbf{X}^T \mathbf{X}$ matrizea koerlazio-matrizea da.

Indibiduo-puntuak k . ardatzaren gain dituzten koordenatuak, $\mathbf{X}\mathbf{u}_k$ zutabe-bektorearen errenkadak dira, $\mathbf{X}\mathbf{u}_k = \sqrt{\lambda_k} v_k$ (6) formularen arabera.

IX.3.2 \mathbb{R}^m espazioan egindako analisisa

Espazio honetan ditugun n puntuak aldagai-puntuak dira. Dakigunez, \mathbb{R}^n espazioan egindako analisiak \mathbb{R}^m espazioaren analisisa induzitzen du. Halere, i eta j azpiindizeek betebeharrak desberdinak dituzte eta \mathbf{X} datu-matrizean egindako transformazioek interpretazio geometriko desberdinak dituzte (ikus VII. gaia).

\mathbb{R}^m espazioan, aldagaien hodeian, zentratze-eragiketaren bidez aldagai-puntuak \mathbb{R}^{m-1} azpiespazio batera proiektatzen ditugu, konkretuki, lehen erdikariarekiko proiektzio paraleloa eginez, koordenatuen batura zero ematen duen \mathbb{R}^{m-1} azpiespaziora proiektatzen ditugu.

Aldagai-puntuak jatorrirra daukan distantzia karratua aldagaiaren bariantza da:

$$d^2(j, o) = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \bar{x}_j)^2 = S_j^2$$

j, j' bi aldagaien arteko distantzia karratua:

$$\begin{aligned} d^2(j, j') &= \frac{1}{m} \sum_{i=1}^m \left[(x_{ij} - \bar{x}_j) - (x_{ij'} - \bar{x}_{j'}) \right]^2 \\ &= S_j^2 + S_{j'}^2 - 2S_{jj'} \end{aligned}$$

Aldiz, Osagai normatutako analisisan ardatzen eskalak aldatzeak, hau da, \mathbb{R}^m espazioan ditugun koordenatuak, $S_j \sqrt{m}$ aldagaiaren normaren balioaz zatituz, aldagai-puntuak normatzen ditugu eta guztiak jatorriarekiko 1 distantziaz kokatzen dira.

Hots:

$$d^2(j, o) = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \bar{x}_j)^2 / S_j^2 = 1$$

n aldagai-puntuak, koordenatuen jatorrian zentratuak erradioz bat daukan hiperesfera batean kokatzen dira.

Ondokoak, j eta j' bi aldagai puntuen arteko distantzia karratua:

$$d^2(j, j') = \frac{1}{m} \sum_{i=1}^m \left[\frac{x_{ij} - \bar{x}_j}{S_j} - \frac{x_{ij'} - \bar{x}_{j'}}{S_{j'}} \right]^2$$

Karratua eginez eta batukaria aplikatuz, hauxe da ateratzen dugun emaitza:

$$d^2(j, j') = S_j^2 / S_j^2 + S_{j'}^2 / S_{j'}^2 - 2r_{jj'} = 2(1 - r_{jj'})$$

Honela, aldagai-puntuen hurbiltasunak koerlazioen arabera interpreta daitezke. Oso hurbil dauden bi aldagai-puntuen arteko koerlazioa positiboa eta sakona izango da ($r_{jj'} \simeq 1$), eta oso urrun dauden bi aldagai-puntuen koerlazioa negatiboa eta sakona izango da ($r_{jj'} \simeq -1$). Aldagai-puntuek k ardatzaren gain dituzten koordenatuak $\mathbf{X}^T \mathbf{v}_k$ zutabe-bektorearen osagaiak dira, eta (5) formularen arabera $\mathbf{X}^T \mathbf{v}_k = \sqrt{\lambda_k} \mathbf{u}_k$.

Dakusagunez, bada, aldagai-puntuen koordenatuak k ardatzaren gain, \mathbf{R} koerlazio-matrizearen diagonalizazioan kalkulatuak λ_k , \mathbf{u}_k , kantitateen funtzioan kalkulatuak dira.

Dakigunez, (ikus VII. gaia) bi aldagai-puntuk osatzen duten angeluaren kosinua aldagaien arteko koerlazio-koefizientea da. Analisi normatuaren kasuan, aldagai-puntuak zentru-normatuak direnez, kosinua aldagaien biderkadura eskalarra da.

Aldagai-puntuek k . ardatzaren gain dituzten abzisak $\mathbf{X}^T \mathbf{v}_k$ zutabe-bektorearen osagaiak dira.

Dakusagunez, abzisa bakoitza aldagai-bektore bakoitzaren (\mathbf{X}^T matrizearen errenkada) eta \mathbf{v}_k bektorearen arteko biderkadura eskalarra da. Honela, aldagai-puntu bakoitzak ardatz baten gain daukan abzisa, aldagaia eta aldagai guztien konbinazio lineala den \mathbf{v}_k aldagai berriaren arteko koerlazio-koefizientea da.

Analisi Orokorrean ikusten genuenez, λ_k autobalio bakoitzak k . ardatzaren gain proiektzio karratuen batura neurtzen du, hau da, k . ardatzaren gain proiektatutako hodeiaren inertzia.

$\mathbf{X}^T \mathbf{X}$ matrizeak, inertzi-matrizea izanik, bere aztarrean hodeiaren jatorriarekiko inertzia biltzen du.

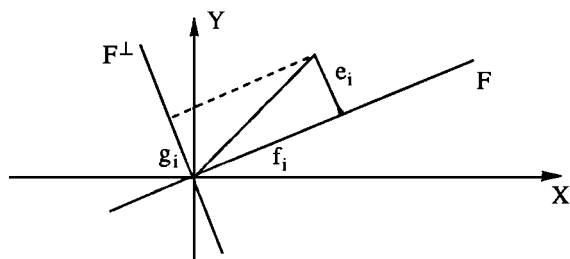
Bestalde, diagonalizazioan egindako transformazio lineal ortogonalaren bidez, jatorriarekiko baterako inertzia, n ardatz ortogonaletan proiektatutako inertzitan deskonposatzen da (ikus (7) berdintasuna).

\mathbb{R}^n , n -dimentsioko espazio bektoriala, dimentsio bateko n azpiespazio ortogonaletan deskonposatu da, n ardatz faktorialetan, hain zuzen.

\mathbb{R}^n , n -dimentsioko espazio bektoriala, bi espazio ortogonaletan deskonposa dezakegu, hauek betegarriak direlarik.

Honela, q lehen ardatzetan hodeiaren sakabanatzearen zatirik adierazgarriena biltzen denean, \mathbb{R}^n espazioa \mathbb{R}^q eta \mathbb{R}^{n-q} azpiespazio betegarritan deskonposatzen da.

$n = 2$ denean (ikus VIII gaia) oso erraz ikus daiteke hodeiaren baterako inertzia F zuzen baten gainean proiektatutako inertzian eta zuzen horrekiko daukan hondar-inertzian deskonposatzen dela, konkretuki, F zuzenaren ortogonal den F^\perp zuzenaren gainean proiektatutako inertzia eta F -rekiko hondar-inertzia berdinak dira.



Puntu bakoitzerako $g_i = e_i$ eta puntu guztientzako

$$\sum_i g_i^2 = \sum_i e_i^2$$

hodeiaren baterako inertzia batukorki deskonposatzea, Pitagoras-en teoremaren aplikatzeaz zuzenki eratortzen delarik.

$n > 2$ denean eta \mathbb{R}^n espazioa \mathbb{R}^q eta \mathbb{R}^{n-q} azpiespazio betegarritan deskonposatuta daukagunean, hodeiaren inertzia bi espaziotan proiektatutako inertzien batura da, baina \mathbb{R}^{n-q} espazioan proiektatutako inertzia eta \mathbb{R}^q espazioarekiko hondar-inertzia berdinak dira.

Era berean, hodeiaren baterako inertzia batukorki deskonposatzea, “Pitagoras-en Teorema Espazioan” aplikatzeaz zuzenki eratortzen da.

Osagai Nagusiko Analisisian, Doikuntza Ortogonalean gertatzen zen bezala, diagonalizatzen dugun matrizea \mathbf{L} kobariantza matrizea edo \mathbf{R} aldagai tipifikatuen kobariantza matrizea da. Halere, $(x_{ij} - \bar{x}_j) \sqrt{m}$ edo $(x_{ij} - \bar{x}_j) S_j \sqrt{m}$ aldagai transformatuen inertzia matrizeak bezala lortu direla gogoratu behar dugu.

Edozein kasutan, λ_k autobalio bakoitzak k . ardatzaren gain puntuen proiektzioen karratuen batura batezbestekoz neurtzen du, hau da, k . ardatzaren gain proiektatutako hodeiaren bariantza.

Hodeiaren baterako bariantza:

$$\text{tr } \mathbf{L} = \sum_{j=1}^n S_j^2 = \sum_{k=1}^n \lambda_k = \text{tr } \mathbf{\Lambda}$$

edo

$$\text{tr } \mathbf{R} = n = \sum_{k=1}^n \lambda_k = \text{tr } \mathbf{\Lambda}$$

kasu honetan, koerlazio-matrizearen elementu diagonalak 1 baitira.

Era berean, q lehen ardatzetan hodeia ondo doituta kontsideratzen badugu, $\sum_{k=1}^q \lambda_k$ proiektatutako bariantza eta $\sum_{k=q+1}^n \lambda_k$ hondar-bariantza izango dira.

Azkenik:

$$\tau_q = \frac{\sum_{k=1}^q \lambda_k}{\sum_{k=1}^n \lambda_k}$$

bariantza tasa izango da, proportzio horrek, askotan portzentaiatan emanik, q lehen faktoreetan proiektatutako bariantzaren portzentaia neurtzen du, edo q lehen faktoreetan proiektatutako inertziaren portzentaia, honela, era berean, inertzia tasa deritzogu.

IX.3.3. Baterako aurkezpen grafikoak

Osagai nagusiko analisisia mende honen hasieratik ezaguna da, baina, azken hogeita hamar urte hauetan eduki duen bultzada oso handia izan da.

Arlo askotan aplikaturik, Soziologian, Ekonomian, Biologian, Geologian, Antropologian, Linguistikan... datu-taula anizkoitzetatik laburpen deskribatzaileak ateratzeko guztiz erabilgarria da.

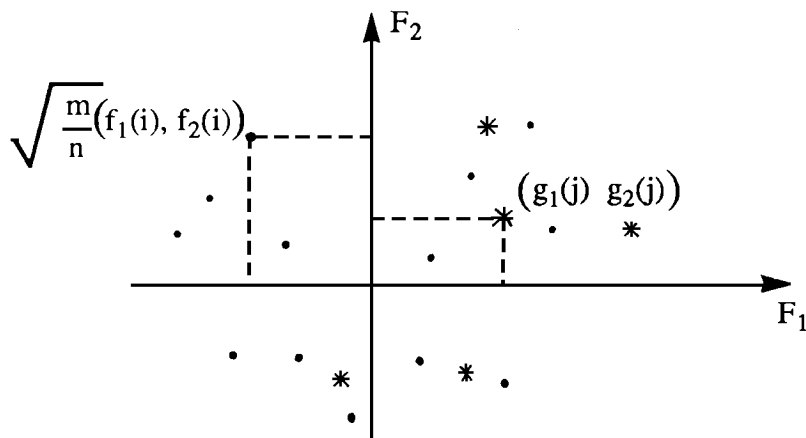
Analisi-teknika honen garapena, gaiaren hasieran esaten genuen bezala, konputagailuen eta kalkulu automatikoen garapenari loturik dago, baina praktikan, erabilpenaren garapena, baterako aurkezpen grafikoetatik eratorritako interpretagarritasunean datza.

$\mathbf{X}^T \mathbf{X}$ eta $\mathbf{X} \mathbf{X}^T$ matrize simetrikoen analisisietan lortutako emaitza paraleloek, hau da, autobalio berdinak eta autodirekzio baliokideak edukitzeak, errenkada-puntuak eta zutabe-puntuak hodeiaren ebaketa-plano berean proiektatzearen ideia sortu zuten Datu-Analisiaren eskola frantsesean.

Trantsizio-erlazioak deitzen diren (5) eta (6) erlazioetan daukagu indibiduo-puntuak eta aldagai-puntuak erlazionatzeko bidea.

(6) erlazioa ikusiz zera esan dezakegu: k . ardatzaren gain koordinatu handienak dauzkaten indibiduo-puntuek koefiziente handiagoz hartzen dute parte aldagai-puntuen hodeian k . ardatza finkatzean. Eta, (5) erlazioa ikusiz, hodeia eta puntu-mota aldatuz, ondorio berdina atera daiteke.

Bi ardatz faktorialak osatutako planoak \mathbb{R}^n edo \mathbb{R}^m espazioetan, jatorritik pasatuz egindako ebaketa da, eta lehen plano faktoriala edo plano faktorial nagusia deritzogu.



Analisiaren praktikan beste plano batzuek aurkezten dira (F_1, F_3) , (F_2, F_3) ,..., alegia, kasu bakoitzerako estatistikariak adierazgarriagoak kontsideratzen dituenak.

Analisi normatua egin bada, aldagaien koordenatuak bat baino txikiagoak izango dira, erradioz bat daukan hiperesfera batean daudela gogoratu behar dugu eta proiektzio-eragiketara murriztaile dela. \mathbb{R}^m espazioan ebaketa egitean, erradioz bat daukan zirkulua ebaki da eta puntuak zirkunferentziatik hurbil eduki ahal ditugu.

Dakusagunez, indibiduen hodeian, analisia grabitate-zentruarekiko egiten da, baina ez aldagaien hodeian, horrela aldagai-puntu denak planoaren alde berean egon daitezke; aski da horretarako beraien arteko koerlazio guztiak positiboak izatea.

Bi hodeien inertzia berdina izatean, eta normalki aldagaien kopurua indibiduen kopurua baino askoz txikiago izatean, lehenengoak planoan proiektatzean jatorritik urrunago geratuko dira. Konkreteki, ardatz bakoitzean proiektatzean indibiduo-puntuaren distantzia tipikoa jatorriarekiko $\sqrt{\lambda_k/m}$ izango da, eta $\sqrt{\lambda_k/n}$ aldagai-puntuarena. Orduan, indibiduo-puntuen koordenatuak $\sqrt{m/n}$ balioaz biderkatzen dira aldagai-puntuekin planoan bateragarriak izateko.

Ardatzen alde positiboa eta negatiboa edonolakoa da eta diagonalizazioaren algoritmoaren menpe izango da. \mathbb{R}^n espazioan k . ardatzaren bektore unitarioa aukeratzeko, bi posibilidadate daude, \mathbf{u}_k edo $-\mathbf{u}_k$, hain zuzen, baina bata aukeratuz gero \mathbb{R}^m espazioarena finkatuta dago, dakigunez $\mathbf{v}_k = \mathbf{h}_k$; $\mathbf{X} \mathbf{u}_k$ eta honen ondorioz (5), (6) trantsizio-erlazioak espazio baten elementuen abzisak k . ardatzaren gain, beste espazioan, k . ardatzaren bektore unitarioaren osagaien bidez ematen dizkigutenak, koherenteki aplikatutako dira.

Bi hodeien puntuak plano faktorialetan proiektatu ondoren eta adierazpide grafiko hauetatik ondorioak atera baino lehen, ondokoa kontutan hartu behar dugu:

Aldagai batzuren posizioek, inguruan dauden indibiduen hurbiltasunaren interpretazioaren ikuspegia emango digute, baina aldagai bat indibiduo baten ondoan egoteak ez du esangurarik, puntu baten kokapena beste hodeiaren puntu guztien kokapenarekin dago erlazionaturik eta hortik datorkio bere posizioa. Ez dugu inoiz ahaztu behar ikuspegi estatistikoa multzokoa delarik.

IX.3.4. Adibidea

Aurreko garapen teorikoen argibide gisa erabil daitekeen adibidea hauxe da:

Suposa dezagun Hego eta Ertamerikako ikasketa bat egin nahiean, hamazortzi estaturako sei aldagaien ohartutako balioak ditugula.

Aukeratu ditugun ondoko aldagaietan, munduko kontinente honen errealitatea partzialki ohartzen dugu; halere, errealitateari hurbiltzen zaion adibide honetan, ikasleak SPAD paquete-programaren azpiprograma batzuren bidez lortutako balioak jarrai ditzake kalkuladora baten bidez.

Aldagaiak

Datu-matrizeko zutabeetan ditugun aldagaiak eta identifikatzaileak hauek dira:

Nazio-Produktu Gordina / biztanle, dolarretan:	NPGB
Urteroko inflazio-tasa, 1980-1985 tarterako:	UINT
Biztanle kopurua, milaka biztanletan, 1985 urterako:	BIZK
Populazioaren hazkunde-tasa, 1973-1985 tarterako:	POHT
Hiri-populazioaren ehunekoa, 1985 urterako:	HIPO
Eskolako matrikula-tasa, 1984 urterako:	ESMT

Datu-matrizea:

Estatuak	Ident.							
VENEZUELA	VEN	3110	9.2	17323	3.3	85	74	
ARGENTINA	ARG	2130	342.8	30531	1.6	84	88	
MEXIKO	MEX	2080	62.2	78820	2.8	69	88	
PANAMA	PAN	2020	3.7	2180	2.3	50	83	
URUGUAI	URU	1660	44.6	3004	0.6	85	86	
BRASIL	BRA	1640	147.7	135539	2.3	73	78	
TXILE	TXI	1440	19.3	11990	1.7	83	95	
KOLONBIA	KOL	1320	22.5	24418	1.9	67	73	
COSTA RICA	CRI	1290	36.4	2593	2.8	45	69	
GUATEMALA	GUA	1240	7.4	7966	2.8	41	47	
EKUADOR	EKU	1160	29.7	9367	2.9	52	85	
PERU	PER	960	98.6	18653	2.4	68	85	
PARAGUAI	PAR	940	15.8	3388	2.5	41	73	
NIKARAGUA	NIK	850	33.8	3263	3.1	56	71	
DOMINIKAR ERREPUBLIKA	DOE	810	14.6	6261	2.4	56	81	
HONDURAS	HON	730	5.4	4366	3.5	39	67	
EL SALVADOR	ELS	710	11.6	5564	3.0	43	56	
BOLIVIA	BOL	470	569.1	6383	2.7	44	71	

Argiro ikusten denez, aldagaiak oso heterogenoak dira eta honek Osagai Nagusizko Analisi Normatua egitea eskatzen du.

Programaren bidez jasotako lehen estatistiko bakunetan, batezbestekoen eta desbidazio tipikoen desberdintasunak ikus daitezke, halaber, berorien bidez atera ditzakegun aldakuntza koefizienteen desberdintasunak.

Estatistiko bakunak

1 DESCRIPTION SOMMAIRE DES VARIABLES SUR 18 INDIVIDUS

	VARIABLE	MOYENNE	ECART-TYPE	MINIMUM	MAXIMUM	ABSENTS
1 /	NPGB	1364.4446	639.2376	470.0000	3110.0000	0
2 /	UINT	81.9111	142.0133	3.7000	569.1000	0
3 /	BIZK	20644.9414	32960.8437	2180.0000	oooooooo	0
4 /	POHT	2.4778	.6762	.6000	3.5000	0
5 /	HIPO	60.0556	16.4096	39.0000	85.0000	0
6 /	ESMT	76.1111	11.6232	47.0000	95.0000	0

Ondoren jasotzen dugun koerlazio-matrizea oso adierazgarria da, adibide txiki honetan, sei aldagai besterik ez baititugu.

Koerlazio-matrizea:

	NPGB	UINT	BIZK	POHT	HIPO	ESMT
NPGB	1.00	-.14	.32	-.17	.67	.35
UINT	-.14	1.00	.15	-.14	.06	.11
BIZK	.32	.15	1.00	-.06	.37	.22
POHT	-.17	-.14	-.06	1.00	-.59	-.53
HIPO	.67	.06	.37	-.59	1.00	.67
ESMT	.35	.11	.22	-.53	.67	1.00

Hiru azpimatrizen azpimarratu ditugu, 1. eta 2. aldagai ekonomikoak, 3. eta 4. populaziokoak, eta 5. eta 6. sozialak.

Bikote hauen erlazio mailak eta NPGB aldagaiak besteekin dituen erlazio mailak eta zeinuak asko esaten dute.

Koerlazio-matrizearen diagonalizazioaz egindako Osagai Nagusizko Analisisian, ondoko $\lambda_1, \dots, \lambda_6$ autobalioak ditugu.

Autobalioak

EDITION DES VALEURS-PROPRES

SOMME DES VALEURS-PROPRES

HISTOGRAMME DES PREMIERS VALEURS-PROPRES

	VALEUR-PROPRE	POURCENTAGE	POURCENTAGE CUMULE	
1	2.68648791	44.77	44.77	*****
2	1.14504230	19.08	63.86	*****
3	1.01897109	16.98	80.84	*****
4	.55648887	9.27	90.12	*****
5	.43828008	7.30	97.42	*****
6	.15473048	2.58	100.00	**

Dakusagunez:

$$\sum_k \lambda_k = 6 = \text{tr } \mathbf{R}$$

hiru lehen ardatz faktorialetan proiektatutako bariantza % 80'84 da.

Hots: $\tau_3 = 80'84$

Aldagaien koordenatuak eta interpretaziorako laguntzak.

1 EDITION DES COORDONNEES ET DES CONTRIBUTIONS DES VARIABLES

NOMS	E.TYPE *	COORDONNES						PROJECTION ANCIENS AXES UNITE *						CORRELATION VARIABLE-FACTEUR *						
		* (CARRE=CONTRIBUTION ABSOLUE) *						* (CARRE=CONTRIBUTION RELATIVE) *												
		F1	F2	F3	F4	F5	F6	F1	F2	F3	F4	F5	F6	F1	F2	F3	F4	F5	F6	*
NPGB	639.238	* -69	.53	-.13	.44	-.08	-.18	* -.42	.49	-.13	.59	-.12	-.45	* -.69	.53	-.13	.44	-.08	-.18	*
UINT	142.013	* -.12	-.82	-.43	.35	.00	-.01	* -.08	-.77	-.42	.47	.00	-.03	* -.12	-.82	-.43	.35	.00	-.01	*
BIZK	32960.844	* -.48	.11	-.76	-.42	-.07	-.04	* -.29	.10	-.75	-.56	-.11	-.10	* -.48	.11	-.76	-.42	-.07	-.04	*
POHT	.676	* .68	.38	-.44	.20	.38	.13	* .41	.35	-.44	.27	.58	.32	* .68	.38	-.44	.20	.38	.13	*
HIPO	16.410	* -.94	.08	.05	.10	-.07	.31	* -.57	.07	.05	.13	-.11	.80	* -.94	.08	.05	.10	-.07	.31	*
ESMT	11,623	* -.80	-.17	.20	-.12	.52	-.08	* -.49	-.16	.20	-.17	.79	-.21	* -.80	-.17	.20	-.12	.52	-.08	*

Aldagaien kasuan, $g_k(j)$ proiektzioak edo koordenatuak ditugu sei ardatzetarako, hau da $k \in \{1, \dots, 6\}$

Dakigunez, k . ardatz faktorialaren gain proiektatutako inertzia $\sum_j g_k^2(j) = \lambda_k$

da, kasu honetan $\mathbf{X}^T \mathbf{X} = \mathbf{R}$ diagonalizatzen denez, proiektatutako bariantza tipifikatua, hain zuzen.

Aldagai bakoitzari dagokion proiektatutako bariantzaren proportzioari aldagaiaren ekarpen absolutua deritzo. Hau da, $g_k^2(j)/\lambda_k$. aldagaiaren ekarpen absolutua da.

Dakusagunez, programak balio hauen erro karratuak ematen dizkigu edo \mathbf{u}_k bektore unitarioen osagaiak.

Hots:

$$\mathbf{u}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{X}^T \mathbf{v}_k = \frac{\mathbf{g}_k}{\sqrt{\lambda_k}}$$

Azken sei zutabeetan aldagai eta faktoreen arteko koerlazioak ditugu, hauek $\mathbf{g}_k = \mathbf{X}^T \mathbf{v}_k$ bektoreen balioak direnez, aldagai normatuen eta \mathbf{v}_k bektore unitarioaren arteko biderkadura eskalarrak, hain zuzen.

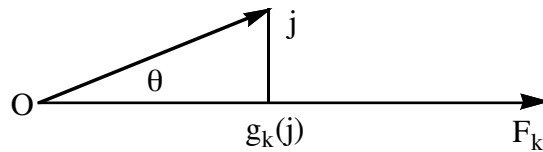
Balio hauek, $g_k^2(j)/d^2(j,0)$ balioen erro karratuak dira, analisi normatuan $d^2(j,0) = 1$ baita.

Ardatz faktorial bakoitzari dagokion aldagaiaren jatorriarekiko daukan distantzia karratuaren proportzioari, ardatzaren ekarpen erlatiboa deritzo.

Hau da, $g_k^2(j)/d^2(j,0)$ k. ardatzaren ekarpen erlatiboa da, faktoreak aldagaiari egindako ekarpen erlatiboa hain zuzen.

Dakigunez, $\sum_{k=1}^6 g_k^2(j) = d^2(j,0)$ «Pitagoras-en Teorema Espazioan» kasu honetan \mathbb{R}^{18} espazioan aplikatuz (aldagaien espazioa) lortutako sei ardatzetarako.

Ohartu behar dugu, $g_k(j) / d(j,0)$, aldagaiak ardatzarekin osatzen duen θ angeluaren kosinua dela, bien arteko koerlazioa, hain zuzen.



Hots:

$$\cos_k^2(\theta) = g_k^2(j)/d^2(j,0)$$

$$\sum_k \cos_k^2(\theta) = 1$$

Indibiduen koordinatuak eta interpretaziorako laguntzak.

 1 EDITION DES COORDONNEES ET DES CONTRIBUTIONS DES INDIVIDUS

	NOMS MASES DISTO * COORDONNES						* CONTRIBUTIONS ABSOLUES*100						* CONTRIBUTIONS RELATIVES *								
		F1	F2	F3	F4	F5	F6		F1	F2	F3	F4	F5	F6		F1	F2	F3	F4	F5	F6
VEN	1.000	1.93	* -56	.94	-23	.81	.03	.18	* 3.8325.57	1.7739.05	.05	6.90	* .16	.46	.03	.34	.00	.02	*		
ARG	1.000	1.63	* -1.06	-.53	-.13	.43	-.11	-.04	* 13.93	8.26	.5411.28	1.01	.37	* .69	.17	.01	.12	.01	.00	*	
MEX	1.000	.99	* -.65	.36	-.57	-.15	.28	-.12	* 5.19	3.81	10.36	1.35	6.12	3.19	* .42	.13	.33	.02	.08	.01	*
PAN	1.000	.41	* -.11	.26	.30	.17	.13	-.44	* .15	1.99	2.86	1.66	1.27	42.64	* .03	.17	.21	.07	.04	.48	*
URU	1.000	1.89	* -1.00	-.25	.78	-.10	-.45	.00	* 12.39	1.87	20.01	.65	15.20	.01	* .53	.03	.32	.01	.11	.00	*
BRA	1.000	2.21	* -.76	.06	-1.10	-.61	-.22	-.01	* 7.20	.11	39.64	1.95	3.72	.05	* .26	.00	.55	.17	.02	.00	*
TXI	1.000	1.03	* -.82	-.07	.52	-.16	.20	.16	* 8.30	.16	8.76	1.55	2.98	5.39	* .65	.01	.26	.03	.04	.02	*
KOL	1.000	.20	* -.18	.03	.18	-.18	-.31	.06	* .40	.03	1.06	1.90	7.22	.84	* .16	.00	.17	.16	.49	.02	*
CRI	1.000	.31	* .51	.13	.08	.08	-.01	-.14	* 3.23	.52	.19	.40	.02	3.98	* .84	.06	.02	.02	.00	.06	*
GUA	1.000	1.38	* .94	.30	-.10	-.10	-.62	-.04	* 11.03	2.67	.30	.63	29.11	.37	* .64	.07	.01	.01	.28	.28	*
EKU	1.000	.26	* .17	.06	.13	-.08	.45	-.07	* .38	.12	.53	.35	15.25	.94	* .12	.02	.06	.02	.77	.02	*
PER	1.000	.21	* -.17	-.22	.13	-.15	.23	.19	* .37	1.37	.52	1.43	4.15	8.10	* .14	.23	.08	.11	.26	.18	*
PAR	1.000	.39	* .52	-.02	.23	-.17	.03	-.20	* 3.36	.01	1.68	1.69	.07	8.75	* .69	.00	.13	.07	.00	.10	*
NIK	1.000	.36	* .51	.08	.06	.02	.15	.25	* 3.24	.17	.11	.02	1.72	13.73	* .73	.02	.01	.00	.06	.18	*
DOE	1.000	.24	* .17	-.09	.31	-.25	.18	.05	* .36	.26	3.21	3.84	2.56	.57	* .12	.04	.42	.27	.14	.01	*
HON	1.000	1.01	* .96	.18	-.06	-.07	.23	.06	* 11.32	.92	.12	.33	4.12	.74	* .90	.03	.00	.01	.05	.00	*
ELS	1.000	1.03	* .96	.12	-.02	-.09	-.26	.13	* 11.46	.43	.01	.45	5.08	3.42	* .90	.01	.00	.01	.06	.02	*
BOL	1.000	2.53	* .56	-1.33	-.50	.44	.07	-.01	* 3.85	.73	8.0211.46	.35	.02	* .12	.70	.10	.08	.00	.00	.00	*

Indibiduen kasuan, $d^2(i,0)$ jatorriarekiko daukan distantzia karratua daukagu, eta ondoren $\sqrt{m/n} f_k(i,0)$ koordinatu zuzenduak sei ardatzetarako.

Era berean, indibiduo bakoitzari dagokion proiektatutako bariantzaren proportzioari, indibiduoaren ekarpen absolutua deritzo.

Hau da, $f_k^2(i)/\lambda_k$ -i. indibiduoaren ekarpen absolutua da; kasu honetan, balio hauek ehunekotan ditugu.

Azken sei zutabeetan faktoreek indibiduoari egindako ekarpen erlatiboak ditugu.

Ekarpen erlatiboa, aldagaien kasurako definitu dugun kontzeptu berbera denez, hauxe da:

$$f_k^2(i)/d^2(i,0)$$

eta

$$\sum_{k=1}^6 f_k^2(i) = d^2(i,0)$$

Indibiduoak ardatzarekin osatzen duen angeluaren kosinua $f_k(i)/d(i,0)$ da.

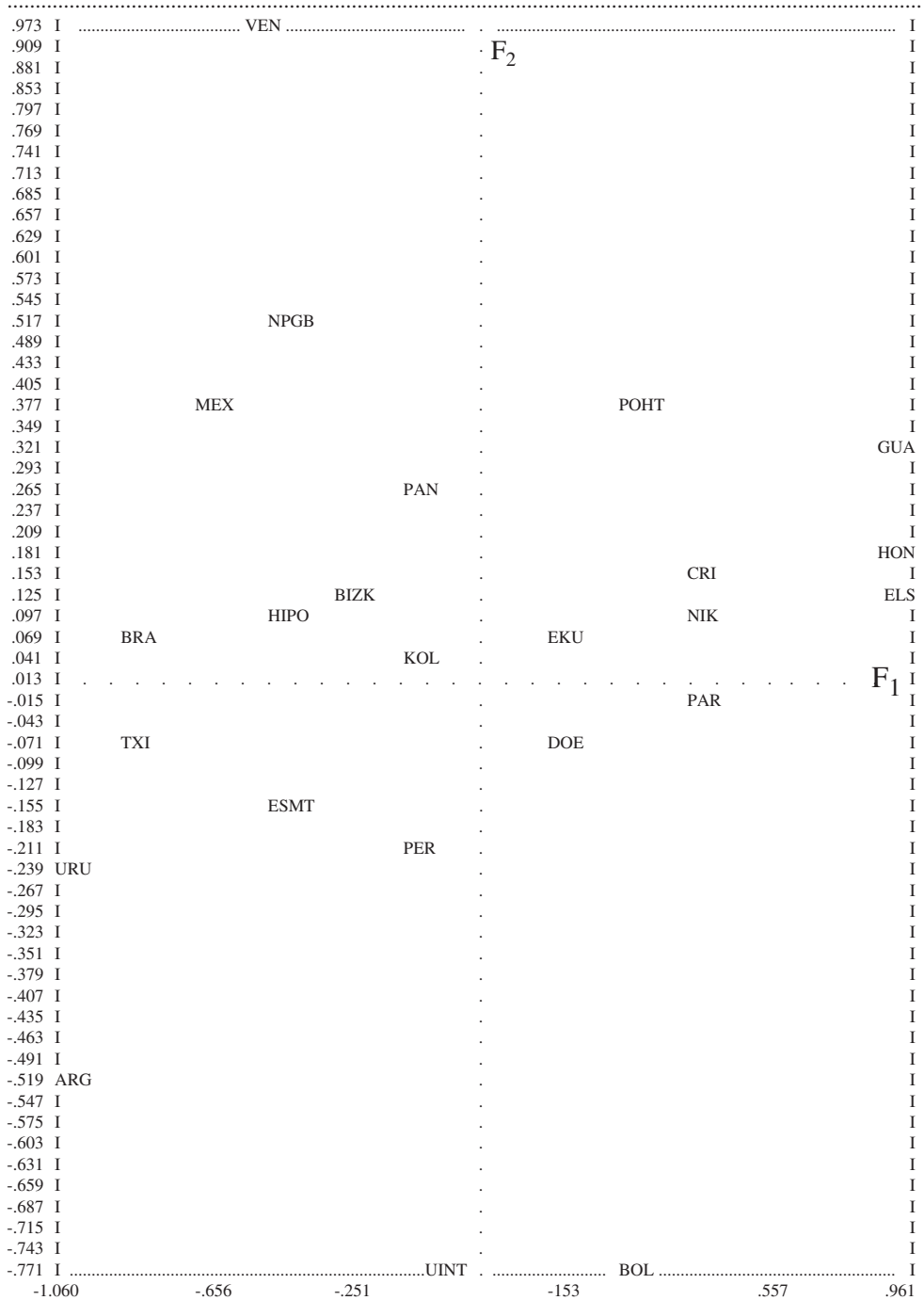
Hau da:

$$\cos_k^2(\theta) = f_k^2(i)/d^2(i,0)$$

$$\sum_k \cos_k^2(\theta) = 1$$

Baterako aurkezpen grafikoa: Plano Faktorial Nagusia

PLAN DE PROJETION DES 24 POINTS SUR LES AXES 1 ET 2



Grabitate-zentrutik gehien urruntzen diren aldagai eta indibiduoak dira koordenatu eta ekarpen absolutu eta erlatibo handienak dauzkatena ardatz bakoitzerako, hau da, puntu horiek ardatz faktoriala aldagai berri bezala definitzera eramaten gaituzte.

Lehen ardatzaren kasuan HIPO (hiri-populazioaren ehunekoa) da aldagairik nabariena alde batean, bere koordenatua edota faktorearekin daukan koerlazioa $-0'94$ da eta honen arabera faktorearen bariantzari egindako ekarpen absolutua handiena da $(-0'57)^2$ (gogoratu behar dugu ardatzaren alde batean edo bestean egotea edonolakoa dela). ESMT (eskolako matrikula-tasa) eta NPGB (Nazio-Produktu Gordina) alde berean nabariak dira.

Aldagai hauen inguruan, planoaren ezker aldean, dauden indibiduoarik nabariena ARG, URU, TXI, eta atzetik MEX, BRA, dira. Azken bi hauek ekarpenetan dituzten ezberdintasunak aipagarriak dira; Brasilek faktoreari egindako ekarpen absolutua handiagoa dauka bere koordenatua handiagoa delako, baina faktoreak Brasili egindako ekarpen erlatiboa txikiagoa da Brasilek grabitate-zentrutik daukan distantzia handiagoa delako. Honek zera esan nahi du: Brasil beste ardatz batean hobeto aurkezten dela, kasu honetan, konkretuki, hirugarren ardatzarekiko dituen balioak adierazgarriagoak dira.

Ardatzaren beste aldean POHT (populazioaren hazkunde-tasa) da ekarpen absolutu eta erlatibo nabariak daukagun aldagaia. Indibiduen artean ELS, HON, GUA eta atzetik PAR, NIK, guzti hauek planoaren eskuin aldean ikus daitezkeelarik.

Lehen ardatza, bada, garapen-ardatza bezala defini daiteke, garapen-mota batena behintzat, eta honen arabera kontinentearen hegokonoa eta Ertamerika kontrajarrita daude.

Bigarren ardatzaren kasuan, UINT (urteroko inflazio-tasa) da aldagairik nabariena eta bere inguruan BOL eta ARG daude, ardatzean proiektatzen den bariantzaren erdia baino gehiago Boliviari dagokiolarik. Puntu hauek planoaren beheko aldean ikus daitezke.

Ardatzaren beste aldean, hau da, planoaren goiko aldean, NPGB eta POHT daude, biak lehen ardatzean kontrajarrita eta koerlazio hobeagorekin genituenak eta hauen inguruan VEN, ardatzean proiektatutako bariantzaren laurdena berari dagokiona, eta MEX, azken honek ekarpen absolutu eta erlatiboak txikiak dituelarik.

Bigarren ardatza edo beste edozein beti hondar-ardatza da aurrenekoekiko; zentzu honetan bigarrenak lehenengoak azaltzen ez dituen datu-matrizearen sakabanatzearen alderdiak azaltzen ditu.

Adibide honen kasuan bigarren ardatzak, gehienbat, kontinente honetan hain berezia den inflazio-tasaren arabera dauden desberdintasunak azaltzen ditu. Halere, Nazio-Produktu Gordina eta populazioaren hazkunde-tasa aldagaien konbinazio lineala den bigarren faktorean garrantziaz sartu dira.

Hirugarren ardatzean populazioko aldagaiak dira garrantzitsuenak, nabariena BIZK (biztanle kopurua) izanik. Faktore honen arabera populazio handia edo gehikorra duten estatuak (BRA, MEX) eta populazio ez hain handia eta gehikorra dutenak (URU, TXI) kontrajartzen dira. Azkenik esan behar dugu, interpretatu nahi izan dugun hirugarren ardatza beste biekiko hondar-ardatza dela. Datu-matrizearen desberdintasun globalik garrantzitsuenak lehen ardatzean proiektatzen ziren, sakabanatzearen % 44'77 jasoz eta bigarren garrantzitsuenak eta ez hain globalak bigarrenean, sakabanatzearen % 19'08 jasoz, hain zuzen ere.

IX.4. KORRESPONDENTZI ANALISI FAKTORIALA

Datu-Analisiaren eskola frantsesaren aita bezala ezaguna den J.P. Banzécri-k Korrespondentzi Analisi Faktoriala (KAF) garatu zuen 60.eko hamarkadaren lehen urteetan.

Estatistika inferentzialetik at lehen ezer gutxi ikasiak izan ziren kontingentzi tauletan arrakasta handiz aplikatzen da orduz gero KAF-ren metodologia¹ eta honek Datu-Analisiaren eskolaren bultzada suposatu zuen.

Teknika honi buruz egin daitezkeen aurkezpen ezberdinekin loturik daude KAF-ren helburuak.

Halere, helburu orokor bat kontsidera daiteke: populazio baten gain ohartutako bi aldagaien arteko menpekotasun-erlazioen azterketa egitea, hain zuzen.

KAF taula-mota ezberdinetara aplikagarria da, kontingentzi tauletara, alegia, eta zentzu zabalean kontsideraturik kontingentzi taulak bezala har daitezkeen tauletara.

(1) K. FERNANDEZ AGIRRE. *Análisis multivariante de tablas de gran inercia. Aplicación al corpus con terminología económica en euskara*. Doktoradutza-tesia (1.988)

Kontingentzi taula

$$\begin{array}{c}
 \begin{array}{cccccc}
 & 1 & 2 & \dots & j & \dots & n \\
 \begin{array}{c} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ m \end{array} & \left[\begin{array}{cccccc}
 k_{11} & k_{12} & \dots & k_{1j} & \dots & k_{1n} \\
 k_{21} & k_{22} & \dots & k_{2j} & \dots & k_{2n} \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 k_{i1} & k_{i2} & \dots & k_{ij} & \dots & k_{in} \\
 \dots & \dots & \dots & \dots & \dots & \dots \\
 k_{m1} & k_{m2} & \dots & k_{mj} & \dots & k_{mn}
 \end{array} \right] & = & [k_{ij}]
 \end{array}
 \end{array}$$

i. errenkadaren eta j. zutabearen ebakiduran daukagun k_{ij} balioa, kontaketa baten emaitza izango da, normalki, (i,j) bikotearen maiztasun absolutua.

Halere, zentzu zabalean, kontingentzi taula definituz gero, errenkadaren baturak, k_i balioak hain zuzen, eta zutabeen baturak, k_j balioak, zentzu argia daukaten balioak izan daitezke nahiz eta berez taula horiek neurri taulak izan. Adibidez, enpresa batzuren salmenta kopuruak, esportatutako balioak... izan daitezke balio horiek.

Askotan, kontingentzi taulei menpekotasun taulak edo taula gurutzatuak deritze eta batzutan maiztasun-taulak ere. (Ikus III. ikasgaia).

Taula hauetan, errenkadak nahiz zutabeak populazio beraren partizioak dira, atribututan kategoriatan edo balio ordinaletan egindako partizioak alegia, horrela errenkadek nahiz zutabeek betebeharrak berdinarik dauzkate.

Kontingentzi taulatik errenkada-soslaiak eta zutabe-soslaiak eratortzen dira errenkada-puntuaren eta zutabe-puntuaren arteko distantziak zentzua eduki dezaten.

Honela, errenkada-soslaien taulan, k_i errenkadaren totalarekiko proportzioak eta zutabe-soslaien taulan, k_j zutabearen totalarekiko proportzioak ditugu.

Hots:

$$k = \sum_i \sum_j k_{ij} \quad \text{kontingentzi taularen totala}$$

$$f_{ij} = k_{ij} / k \quad \text{maiztasun erlatiboak}$$

$$\left. \begin{array}{l}
 f_{i.} = \sum_j f_{ij} \\
 f_{.j} = \sum_i f_{ij}
 \end{array} \right\} \quad \text{bazter-maiztasunak}$$

Errenkada-soslaien taula

	1	2	...	j	...	n	Totala
1	$\frac{f_{11}}{f_{1.}}$	$\frac{f_{12}}{f_{1.}}$...	$\frac{f_{1j}}{f_{1.}}$...	$\frac{f_{1n}}{f_{1.}}$	1
2	$\frac{f_{21}}{f_{2.}}$	$\frac{f_{22}}{f_{2.}}$...	$\frac{f_{2j}}{f_{2.}}$...	$\frac{f_{2n}}{f_{2.}}$	1
...						
i	$\frac{f_{i1}}{f_{i.}}$	$\frac{f_{i2}}{f_{i.}}$...	$\frac{f_{ij}}{f_{i.}}$...	$\frac{f_{in}}{f_{i.}}$	1
...						
m	$\frac{f_{m1}}{f_{m.}}$	$\frac{f_{m2}}{f_{m.}}$...	$\frac{f_{mj}}{f_{m.}}$...	$\frac{f_{mn}}{f_{m.}}$	1
Totala	$f_{.1}$	$f_{.2}$...	$f_{.j}$...	$f_{.n}$	

Zutabe-soslaien taula

	1	2	...	j	...	n	Totala
1	$\frac{f_{11}}{f_{.1}}$	$\frac{f_{12}}{f_{.2}}$...	$\frac{f_{1j}}{f_{.j}}$...	$\frac{f_{1n}}{f_{.n}}$	$f_{1.}$
2	$\frac{f_{21}}{f_{.1}}$	$\frac{f_{22}}{f_{.2}}$...	$\frac{f_{2j}}{f_{.j}}$...	$\frac{f_{2n}}{f_{.n}}$	$f_{2.}$
...						
i	$\frac{f_{i1}}{f_{.1}}$	$\frac{f_{i2}}{f_{.2}}$...	$\frac{f_{ij}}{f_{.j}}$...	$\frac{f_{in}}{f_{.n}}$	$f_{i.}$
...						
m	$\frac{f_{m1}}{f_{.1}}$	$\frac{f_{m2}}{f_{.2}}$...	$\frac{f_{mj}}{f_{.j}}$...	$\frac{f_{mn}}{f_{.n}}$	$f_{m.}$
Totala	1	1	...	1	...	1	

IX.4.1. Hodeiak, masak eta distantziak

Osagai Nagusizko Analisisian bezala, bi espazio desberdinetan koka gaitzke.

a) \mathbb{R}^n espazio bektorialean m errenkada-soslaien puntuak ditugu. Honela, f_i masak hornitutako i puntuaren koordinatuak ondoko multzokoak dira:

$$\left\{ \left(f_{ij}/f_i \right) ; j = 1, 2, \dots, n \right\}$$

Konkretuki, errenkada-soslaien hodeiaren m puntuak \mathbb{R}^{n-1} azpiespazio batean kokaturik daude, puntu bakoitzaren koordinatuen batura 1 baita.

Hau da:

$$\sum_j \left(f_{ij}/f_i \right) = 1 \quad \forall i = 1, 2, \dots, m$$

Honela, puntuen hurbiltasunak soslaien hurbiltasunak bezala interpreta daitezke.

b) \mathbb{R}^m espazio bektorialean n zutabe-soslaien puntuak ditugu. Honela, f_j masak hornitutako j puntuaren koordinatuak ondoko multzokoak dira:

$$\left\{ \left(f_{ij}/f_j \right) ; i = 1, 2, \dots, m \right\}$$

Konkretuki, zutabe-soslaien hodeiaren n puntuak \mathbb{R}^{m-1} azpiespazio batean kokaturik daude, puntu bakoitzaren koordinatuen batura 1 baita.

Hau da:

$$\sum_i \left(f_{ij}/f_j \right) = 1 \quad \forall j = 1, 2, \dots, n$$

Doikuntza egin aurretik puntu bakoitzaren masa kontutan hartzea guztiz zilegia da; horrela, masa txikiak dituzten puntuak ez dira gehiegi nagusitzen eta banaketa bere baitan errespetatzen da.

\mathbb{R}^n eta \mathbb{R}^m espazioetan definitzen diren distantziak hauexek dira:

$$d^2(i, i') = \sum_{j=1}^n \frac{1}{f_{.j}} \left(\frac{f_{ij}}{f_{.i}} - \frac{f_{i'j}}{f_{.i'}} \right)^2$$

$$d^2(j, j') = \sum_{i=1}^m \frac{1}{f_{i.}} \left(\frac{f_{ij}}{f_{.j}} - \frac{f_{ij'}}{f_{.j'}} \right)^2$$

Horrela definituriko distantziari χ^2 (ji karratua) deritzo.

Distantzia ponderatu honek daukan abantaila, guztiz berezkoa izateaz gain, **banaketa-baliokidetasun** izeneko propietatea betetzea da.

Propietate honen arabera i_1, i_2 errenkada-puntuak soslai berbera badira (proportzio berdinak zutabe guztietarako) berorien masen batura daukan i_0 errenkada-puntuaz ordezkatu daitezke eta \mathbb{R}^n nahiz \mathbb{R}^m espazioetan distantziak ez dira aldatzen. Berdin esan daiteke j_1, j_2 soslai berbereko bi zutabe-puntuentzako.

Propietate honek analisiaren emaitzak oso egonkor bihurtzen ditu; honela, datu-taula handi batean berdin antzeko soslaia daukaten lerroak emaitzak mantenduz batu daitezke eta honek nomenklatura baten arbitraltasunaren aurrean emaitzen garrantzia suposatzen du.

Soslaien artean definitutako χ^2 distantzia ohizko distantzia euklidear bihurtu daiteke, ponderazioak parentesien barruan sartuz.

Hau da:

$$d^2(i, i') = \sum_{j=1}^n \left(\frac{f_{ij}}{f_{i.} \sqrt{f_{.j}}} - \frac{f_{i'j}}{f_{i'.} \sqrt{f_{.j}}} \right)^2$$

Kasu honetan, ardatzen eskala aldatuz i puntuaren koordinatuak $f_{ij}/f_{i.} \sqrt{f_{.j}}$ balioak dira.

Era berean:

$$d^2(j, j') = \sum_{i=1}^m \left(\frac{f_{ij}}{f_{\cdot j} \sqrt{f_{\cdot i}}} - \frac{f_{ij'}}{f_{\cdot j'} \sqrt{f_{\cdot i}}} \right)^2$$

Kasu honetan, ardatzen eskala aldatuz j puntuaren koordinatuak $f_{ij}/f_{\cdot j} \sqrt{f_{\cdot i}}$ balioak dira.

Eskala aldatu ondoren, errenkada-soslaien grabitate-zentrua zutabeen bazter-banaketaren erro karratua da.

i puntu bakoitzaren masa edo pisua $f_{\cdot i}$ enez:

$$g_j = \sum_{i=1}^n \frac{f_{ij}}{f_{\cdot i} \sqrt{f_{\cdot j}}} f_{\cdot i} = \sqrt{f_{\cdot j}}$$

Eta zutabe-soslaien grabitate-zentrua errenkaden bazter-banaketaren erro karratua da. j puntu bakoitzaren masa edo pisua $f_{\cdot j}$ enez:

$$h_i = \sum_{j=1}^n \frac{f_{ij}}{f_{\cdot j} \sqrt{f_{\cdot j}}} f_{\cdot j} = \sqrt{f_{\cdot i}}$$

Azkenik:

$$\frac{f_{ij}}{f_{\cdot i} \sqrt{f_{\cdot j}}} - \sqrt{f_{\cdot j}} \quad \text{eta} \quad \frac{f_{ij}}{f_{\cdot j} \sqrt{f_{\cdot i}}} - \sqrt{f_{\cdot i}}$$

bi hodeien koordinatu zentratuak, hurrenez hurren, dira.

Dakusagunez, bada, puntu-hodeiaren grabitate-zentrua jatorrira eramán dugu eta ardatzen eskala-aldaketa baten bidez χ^2 distantzia ohizko distantzia euklidearra bezala idatzi dugu, honela, korrespondentzi analisi faktoriala analisi orokorraren barrutian aurkezten da.

IX.4.2. \mathbb{R}^n espazioan egindako analisia

\mathbb{R}^n espazioan diagonalizatzen den kobariantza matrizearen gai orokorra i puntuaren masa f_i dela kontutan hartuz, hauxe da:

$$S_{jj'} = \sum_{i=1}^m \left(\frac{f_{ij}}{f_i \sqrt{f_j}} - \sqrt{f_j} \right) \left(\frac{f_{ij'}}{f_i \sqrt{f_{j'}}} - \sqrt{f_{j'}} \right) f_i$$

j, j' puntuen balio zentratuen arteko biderkadura da i guztietarako eta i puntuen masaz ponderaturik.

Azkenik, biderkatzaile bakoitza izendatzaile komunean ipiniz eta $\sqrt{f_i}$ biderkatuz.

$$S_{jj'} = \sum_{i=1}^m \left(\frac{f_{ij} - f_i \cdot f_j}{\sqrt{f_i \cdot f_j}} \right) \left(\frac{f_{ij'} - f_i \cdot f_{j'}}{\sqrt{f_i \cdot f_{j'}}} \right)$$

Horrela diagonalizatzen den $\mathbf{L}(J)$ kobariantza matrizea $\mathbf{X}(m,n)$ matrizearen funtzioan idatz daiteke:

$$\mathbf{L}(J) = \mathbf{X}^T \mathbf{X}$$

(n,n) (n,m)(m,n)

non, \mathbf{X} matrizearen gai orokorra $\frac{f_{ij} - f_i \cdot f_j}{\sqrt{f_i \cdot f_j}}$ da.

Hots:

$$\mathbf{X} = \left[\frac{f_{ij} - f_i \cdot f_j}{\sqrt{f_i \cdot f_j}} \right]$$

Froga daiteke datuak zentratuz edo zentratu gabe erabiltzeak emaitza berdinetara eramaten duela.

Hau da, datuak zentratu gabe hartzen baditugu

$$\mathbf{X}^* = \left[\frac{f_{ij}}{\sqrt{f_i \cdot f_j}} \right]$$

eta kobariantza matrizea $\mathbf{L}^*(J) = \mathbf{X}^{*T} \mathbf{X}^*$ izango da.

\mathbf{u}_n bektorea, non j . osagaia $\sqrt{f_{.j}}$ baita, $\mathbf{L}(J)$ kobariantza matrizearen 0 autobalioari dagokion autobektorea da; halaber, \mathbf{u}_n (grabitate-zentrua) $\mathbf{L}^*(J)$ kobariantza matrizearen 1 autobalioari dagokion autobektorea da.

\mathbf{u}_n ez den edozein \mathbf{u}_k autobektore, $\mathbf{L}(J)$ edota $\mathbf{L}^*(J)$ matrizeen autobalio berberari dagokion autobektorea da.

k. ardatzaren gain m errenkada-puntuek dauzkaten proiektzioak, datu zentratuaz nahiz zentratu gabeaz kalkula daitezke,

$$\hat{\psi}_{(m,1)k} = \mathbf{X} \mathbf{u}_{(m,n)(n,1)k} = \mathbf{X}^* \mathbf{u}_{(m,n)(n,1)k}$$

$\mathbf{X} \mathbf{u}_n$ eta $\mathbf{X}^* \mathbf{u}_n$ bektoreen abzisak berberak baitira.

Dakusagun i puntuaren koordinatua bektore batean nahiz bestean berbera dela.

$$\begin{aligned} \hat{\psi}_{ki} &= \sum_{j=1}^n \left(\frac{f_{ij}}{f_{i.} \sqrt{f_{.j}}} - \sqrt{f_{.j}} \right) \mathbf{u}_{kj} = \\ &= \sum_{j=1}^n \frac{f_{ij}}{f_{i.} \sqrt{f_{.j}}} \mathbf{u}_{kj} \end{aligned} \quad (8)$$

Parentesia kentzean bigarren batugaia $\sum_{i=1}^n \sqrt{f_{.j}} \mathbf{u}_{kj} = 0$ da, \mathbf{u}_n eta \mathbf{u}_k autobekto-reak ortogonalak baitira.

IX.4.3. \mathbb{R}^m espazioan egindako analisia

\mathbb{R}^m espazioan diagonalizatzen den kobariantza matrizearen gai orokorra j puntuaren masa $f_{.j}$ dela kontutan hartuz, hauxe da:

$$S_{ii'} = \sum_{j=1}^n \left(\frac{f_{ij}}{f_{.j} \sqrt{f_{i.}}} - \sqrt{f_{i.}} \right) \left(\frac{f_{i'j}}{f_{.j} \sqrt{f_{i'.}}} - \sqrt{f_{i'.}} \right) f_{.j}$$

Eta biderkatzaile bakoitza izendatzaile komunean ipiniz eta $\sqrt{f_{i,j}}$ -z biderkatuz

$$S_{ii'} = \sum_{j=1}^n \left(\frac{f_{ij} - f_{i.} f_{.j}}{\sqrt{f_{i.} f_{.j}}} \right) \left(\frac{f_{i'j} - f_{i'.} f_{.j}}{\sqrt{f_{i'.} f_{.j}}} \right)$$

Hau da:

$$\mathbf{L}(\mathbf{I}) = \mathbf{X} \mathbf{X}^T$$

(n,n) (m,n)(n,m)

Dakusagunez korrespondentzian jartzen diren errenkaden eta zutabeen multzoek betebeharrak berdinak dauzkate eta espazio batetik bestera pasatzeko i, j azpiindizeen betebeharrak aldatzea aski da.

\mathbf{v}_m bektorea, non i osagaia $\sqrt{f_{i.}}$ den, $\mathbf{L}(\mathbf{I})$ kobariantza matrizearen 0 autobalioari dagokion autobektorea da; halaber, \mathbf{v}_m (grabitate-zentrua), $\mathbf{L}^*(\mathbf{I}) = \mathbf{X}^* \mathbf{X}^{*T}$ kobariantza matrizearen 1 autobalioari dagokion autobektorea da.

\mathbf{v}_m ez den edozein \mathbf{v}_k autobektore, $\mathbf{L}(\mathbf{I})$ edota $\mathbf{L}^*(\mathbf{I})$ matrizeen autobalio berberari dagokion autobektorea da.

Azkenik, k . ardatzaren gain n zutabe-puntuek dauzkaten proiektzioak:

$$\hat{\varphi}_k = \mathbf{X}^T \mathbf{v}_k = \mathbf{X}^{*T} \mathbf{v}_k$$

Eta j puntuaren koordenatua bere adierazpen errazenean:

$$\hat{\varphi}_{kj} = \sum_{i=1}^m \frac{f_{ij}}{f_{.j} \sqrt{f_{i.}}} v_{ki} \tag{9}$$

IX.4.4. \mathbf{R}^n eta \mathbf{R}^m espazioen arteko erlazioa

Analisi orokorrean ikusi genuenez $\mathbf{L}^*(\mathbf{J}) = \mathbf{X}^{*T} \mathbf{X}^*$ eta $\mathbf{L}^*(\mathbf{I}) = \mathbf{X}^* \mathbf{X}^{*T}$ matrize simetrikoek desberdin zero diren λ_k autobalio berberak dituzte eta λ_k autobalioari dagozkion $\mathbf{L}^*(\mathbf{J})$ -ren \mathbf{u}_k eta $\mathbf{L}^*(\mathbf{I})$ -ren \mathbf{v}_k autobektoreen arteko erlazioak hauek dira:

$$\mathbf{v}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{X}^* \mathbf{u}_k \tag{10}$$

$$\mathbf{u}_k = \frac{1}{\sqrt{\lambda_k}} \mathbf{X}^{*T} \mathbf{v}_k \tag{11}$$

\mathbf{v}_k bektorearen i . errenkada (10) kontutan hartuz, honela idatz daiteke:

$$v_{ki} = \frac{1}{\sqrt{\lambda_k}} \sum_{j=1}^n \frac{f_{ij}}{\sqrt{f_i \cdot f_j}} u_{kj}$$

Eta (8) erlazioan u_{kj} ordezkaturik:

$$\hat{\psi}_{ki} = \sqrt{\lambda_k} \frac{v_{ki}}{\sqrt{f_i}} \quad (12)$$

Era berean, \mathbf{u}_k bektorearen j . errenkada (11) kontutan hartuz honela idatz daiteke:

$$u_{kj} = \frac{1}{\sqrt{\lambda_k}} \sum_{i=1}^m \left(\frac{f_{ij}}{\sqrt{f_i \cdot f_j}} \right) v_{ki}$$

Eta (9) erlazioan v_{ki} ordezkaturik:

$$\hat{\varphi}_{kj} = \sqrt{\lambda_k} \frac{u_{kj}}{\sqrt{f_j}} \quad (13)$$

Azkenik, (9) erlazioan (12) erlaziotik ateratako v_{ki} ordezkaturik:

$$\hat{\varphi}_{kj} = \frac{1}{\sqrt{\lambda_k}} \sum_{i=1}^m \frac{f_{ij}}{f_j} \hat{\psi}_{ki} \quad (14)$$

eta (8) erlazioan (13) erlaziotik ateratako u_{kj} ordezkaturik:

$$\hat{\psi}_{ki} = \frac{1}{\sqrt{\lambda_k}} \sum_{j=1}^n \frac{f_{ij}}{f_i} \hat{\varphi}_{kj} \quad (15)$$

(14), (15) erlazioetan errenkada-puntuak eta zutabe-puntuak k . ardatzaren gain dituzten koordinatuen arteko erlazioak ditugu.

Dakusagunez (f_{ij} / f_i) gai orokorraz daukan matrizearen bidez, hau da, errenkada-soslaien matrizearen bidez, eta j zutabe-puntuen koordinatuen bidez ((15) erlazioa), i errenkada-puntuen koordinatuak lortzen dira.

Honela, i errenkada-puntuaren koordinatua k . ardatzaren gain ($\hat{\psi}_{ki}$), $1/\sqrt{\lambda}$ balioa salbu, j zutabe-puntuen koordinatuen barizentrua edo grabitate-zentrua da, pisuak i errenkada-soslaiaren balioak direlarik.

Era berean, j zutabe-puntuaren koordinatua k . ardatzaren gain ($\hat{\varphi}_{kj}$), $1/\sqrt{\lambda}$ balioa salbu, i errenkada-puntuen koordinatuen barizentrua da, pisuak j zutabe-soslaiaren balioak direlarik ((14) erlazioa).

Erlazio barizentrikoak bezala ezagutzen diren (14), (15) erlazioak ez dira analisi orokorraren (5), (6) erlazioen kasu bereziak, KAF-ean trantsizio-erlazioaren matrizeak ez baitira bata bestearen iraule.

Erlazio hauetan, implizituki, edozein k -rako $\lambda_k \leq 1$ dela daukagu ($1/\sqrt{\lambda_k} < 1$ ezin baita izan).

IX.4.5. Maiztasun-taularen berrerraketa

Gai honen 2.4. atalean, analisi orokorrerako ikusi genuenez, $\mathbf{X}^* = [f_{ij}/\sqrt{f_i \cdot f_j}]$ matrizearen berrerraketa hauxe da:

$$\mathbf{X}^* = \sum_{k=1}^n \sqrt{\lambda_k} \mathbf{v}_k \mathbf{u}_k^T$$

(12), (13) erlazioetatik ateratako $\mathbf{v}_k, \mathbf{u}_k$ balioak ordezkatzuz gero:

$$\frac{f_{ij}}{\sqrt{f_i \cdot f_j}} = \sum_{k=1}^n \left[\frac{1}{\sqrt{\lambda_k}} \psi_{ki} \varphi_{kj} \right] \sqrt{f_i \cdot f_j}$$

Eta dakigunez $\lambda_1 = 1$ autobalioari $\mathbf{u}_{1j} = \sqrt{f_j}$, $\mathbf{v}_{1i} = \sqrt{f_i}$ autobektoreak dagozkio, orduan, (12), (13) erlazioetan dakusagunez $\hat{\varphi}_{1j} = 1$ eta $\hat{\psi}_{1i} = 1$

$$f_{ij} = f_i \cdot f_j \left[1 + \sum_{k=2}^n \frac{1}{\sqrt{\lambda_k}} \psi_{ki} \varphi_{kj} \right]$$

Berreratze-formula honetan, $\hat{\psi}_{ki}$, $\hat{\varphi}_{ki}$ eta λ_k hirurek zerorantz jotzen dute.

Askotan, berreratze-formula, faktore normatuen bidez idazten da ere.

ψ_{ki} eta φ_{ki} faktore normatuak:

$$\psi_{ki} = \frac{1}{\sqrt{\lambda_k}} \hat{\psi}_{ki} \quad \text{eta} \quad \varphi_{ki} = \frac{1}{\sqrt{\lambda_k}} \hat{\varphi}_{ki}$$

Orduan berreratze-formula:

$$f_{ij} = f_i \cdot f_j \left[1 + \sum_{k>1} \sqrt{\lambda_k} \psi_{ki} \varphi_{kj} \right]$$

IX.4.6. Interpretaziorako laguntzak

Osagai Nagusizko Analisiaren adibidean ikusi genuen bezala, ardatzen interpretaziorako, bi koefiziente-segidak kalkulatu dira korrespondentzian ipinitako bi multzoen elementu bakoitzerako.

1) Ekarpen absolutuak

Dakigunez, elementu bakoitzari dagokion proiektatutako bariantzaren proportzioari elementuaren ekarpen absolutua deritzo.

k. ardatzaren gain m errenkada-puntuaren bariantza proiektatua, i errenkada-puntu bakoitzaren masa f_i dela jakinik, hauxe da:

$$\sum_{i=1}^m f_i \cdot \hat{\psi}_{ki}^2$$

Eta $\hat{\psi}_{ki}$ balioa ordezkatzuz:

$$\sum_{i=1}^m f_i \cdot \hat{\psi}_{ki}^2 = \lambda_k \sum_{i=1}^m v_{ki}^2 = \lambda_k \left(\mathbf{v}_k \text{ unitate bektorea baita} \right)$$

Era berean j zutabe-puntuaren bariantza proiektatua.

$$\sum_{j=1}^n f_{.j} \hat{\varphi}_{kj}^2 = \lambda_k \sum_{j=1}^n u_{kj} = \lambda_k$$

Honela, k . ardatzaren gain proiektatutako bariantzan i elementuari dagokion proportzioa:

$$Ca_k(i) = f_{.i} \hat{\varphi}_{ki}^2 / \lambda_k$$

$Ca_k(i)$ balioa, i elementuak k . ardatzean daukan ekarpen absolutua izanik.

Ohar daiteke:

$$\sum_{i=1}^m Ca_k(i) = 1 \quad \forall_k \text{ dela}$$

Analogikoki, j elementuak k . ardatzean daukan ekarpen absolutua:

$$Ca_k(j) = \frac{f_{.j} \hat{\varphi}_{kj}^2}{\lambda_k}$$

2) Ekarpen erlatiboak

Dakigunez, ardatz faktorial bakoitzari dagokion elementuaren jatorriarekiko daukan distantzia karratuaren proportzioari, ardatzaren ekarpen erlatiboa deritzo.

Ardatz faktorialek, bi espaziotan, oinarri ortonormalak osatzen dituztenez, bi espazioko G, H , grabitate-zentruetikiko elementuek dituzten distantzia karratuak ardatzetan dituzten koordinatuen karratuetan deskonposatzen dira batukorki.

Honela, \mathbb{R}^n espazioko i puntu bakoitzerako:

$$d_n^2(i, G) = \sum_{j=1}^n \left(\frac{f_{ij}}{f_{.i} \sqrt{f_{.j}}} - \sqrt{f_{.j}} \right)^2$$

Halaber, \mathbb{R}^m espazioko j puntu bakoitzerako:

$$d_m^2(j, H) = \sum_{i=1}^m \left(\frac{f_{ij}}{f_{.j} \sqrt{f_{.i}}} - \sqrt{f_{.i}} \right)^2$$

Puntu baten soslaiak batezbesteko soslaiarekin kointziditzen badu, distantzia zero izango dela ohartu behar dugu.

k. ardatzaren ekarpen erlatiboa i puntuarekiko:

$$Cr_k(i) = \frac{\hat{\Psi}_{ki}^2}{d_n^2(i, G)}$$

eta i puntuak k . ardatzarekin osatzen duen angeluaren cosinu karratua denez, koerlazio-koefiziente karratua bezala interpreta daiteke.

Ohar daiteke

$$\sum_k Cr_k(i) = 1 \quad \forall i \text{ dela}$$

Analogikoki, k . ardatzak j elementuaren posizioarekiko daukan ekarpen erlatiboa:

$$Cr_k(j) = \frac{\hat{\Phi}_{kj}^2}{d_m^2(j, H)}$$

Eta, era berean:

$$\sum_k Cr_k(j) = 1 \quad \forall j$$

IX.4.7. Korrespondentzia Anitzeko Analisi Faktoriala

Korrespondentzi Analisi Faktorialaren orokorpena bezala Korrespondentzia Anitzeko Analisi Faktoriala sortzen da.

Bi aldagai baino gehiago, kualitatiboak nahiz kuantitatiboak izan, azken kasu honetan bada kategoriatan mailakatuta, ditugunean, Korrespondentzia Anitzeko Analisi Faktorialaren metodologia analisi ikuskatzailea egiteko guztiz baliagarria da.

Teknika honen aplikazio-eremu handi bat galdera itxien inkestetako da; halaber, edozein inkesta motatan lortutako erantzunak aztertzeo, oso egokia da.

Korrespondentzia Anitzeko Analisiaren eragiketaren teknika, deskribapen logikozko taula batean Korrespondentzi Analisi Faktoriala aplikatzea besterik ez da.

Deskribapen logikozko taula

Taula hauen elementuak 0 eta 1 dira. k_{ij} taularen gai orokorra izanik, i elementuak j propietatea duenean $k_{ij} = 1$ izango da, aldiz, i elementuak j propietaterik ez duenean $k_{ij} = 0$ izango da.

Deskribapen logikozko taula batek ondoko baldintza betetzen baldin badu aukera osoko taularen erara dagoela esango dugu: J zutabeen (edo modalitateen) multzoa J_q azpimultzoen (edo galderen) Q segida da non:

$$\forall i \in I, \forall J_q \in Q, \exists j \in J_q \mid k_{ij} = 1 \wedge (j' \in J_q, j' \neq j) \Rightarrow k_{ij'} = 0$$

Honek zera esan nahi du: I indibiduen multzoko i elementu bakoitzak J_q klasean modalitate bat eta bakarrik bat edukiko duela.

Deskribapen logikozko taula batek aukera osoko taularen erara honelako itxura dauka:

	J_1	J_2	J_q	J_Q	Totala
1	010	1000	0010	010	Q
.							
.							
$Z = i$	001	0100	0100	100	Q
(m,J)							
.							
m	100	0001	0100	010	Q

Dakusagunez i indibiduoak J_q azpimultzoan bigarren modalitatea aukeratu du; inkesta baten kasua balitz, q . galderarentzako bigarren erantzuna aukeratu duela esango genuke.

Korrespondentzi Anitzeko Analisia eta Burt-en taula batetan Korrespondentzi Analisi Faktoriala aplikatzea baliokideak dira.

Burt-en taula

m indibiduen gain Q aldagai (edo galdera) ohartu baditugu, Burt-en taula (J_q, J_q') binaka gurutzatutako taula guztiez osaturik dago.

$$\mathbf{B} = \begin{array}{c} \begin{array}{cccccc} & J_1 & J_2 & \dots & J_q & \dots & J_Q \\ \begin{array}{c} J_1 \\ J_2 \\ \vdots \\ J_q \\ \vdots \\ J_Q \end{array} & \begin{array}{|c|} \hline 0 \\ \hline \end{array} & \begin{array}{|c|} \hline J_2 \\ \hline \end{array} & & \begin{array}{|c|} \hline J_q \\ \hline \end{array} & & \begin{array}{|c|} \hline J_Q \\ \hline \end{array} \\ \hline & \begin{array}{|c|} \hline J_{21} \\ \hline \end{array} & \begin{array}{|c|} \hline 0 \\ \hline \end{array} & & \begin{array}{|c|} \hline J_{2q} \\ \hline \end{array} & & \begin{array}{|c|} \hline J_{2Q} \\ \hline \end{array} \\ \hline & & & & \begin{array}{|c|} \hline 0 \\ \hline \end{array} & & \begin{array}{|c|} \hline J_{qQ} \\ \hline \end{array} \\ \hline & \begin{array}{|c|} \hline J_{Q1} \\ \hline \end{array} & \begin{array}{|c|} \hline J_{Q2} \\ \hline \end{array} & & \begin{array}{|c|} \hline J_{Qq} \\ \hline \end{array} & & \begin{array}{|c|} \hline 0 \\ \hline \end{array} \\ \hline & & & & & & \begin{array}{|c|} \hline 0 \\ \hline \end{array} \end{array} \end{array}$$

Diagonalean ditugun azpitaulak bloke diagonalak dira; honelako (J_q, J_q) blokeek, diagonalean, aukeratutako modalitateen kopuruak dituzte eta beste elementuak zero dira, bi modalitate (edo erantzun) desberdin batera ezin baitira aukeratu.

(J_q, J_q') blokeak, benetako kontingentzia taulak dira, bi galderen erantzunak gurutzatu baitira.

\mathbf{B} deitutako Burt-en taula, matrizialki, \mathbf{Z} taularen bidez honela lortuko da:

$$\mathbf{B} = \mathbf{Z}^T \mathbf{Z}$$

Korrespondentzi Analisia aplikatu ondoren lortutako faktoreak interpretatzerako, KAF-ean lortutako erlazioak eta interpretaziorako laguntzak (ekarpen absolutuak eta erlatiboak) baliagarriak dira.

IX.5. ELEMENTU GEHIGARRIAK

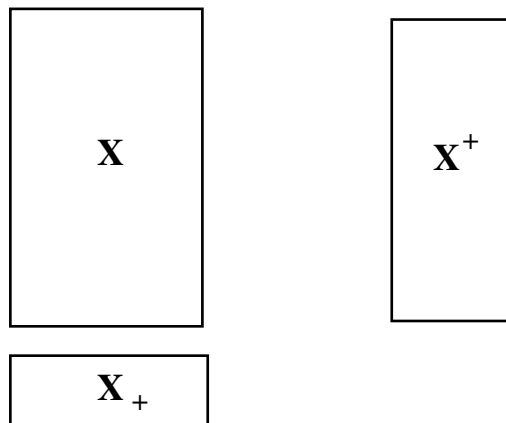
Korrespondentzi Analisi Faktorialean (bikoitza nahiz anizkoitza) eta Osagai Nagusizko Analisisan *elementu gehigarrien teknikak* asko laguntzen du analisiaren interpretazioan.

Laginetik at, n aldagaien gain ditugun indibiduo batzuen balioak ezaguturik, analisiaren lekuko bezala, hodeian, talde hori kokatzea komenigarria gerta daiteke.

Batzutan, ezaugarri interesgarri bat daukaten indibiduo-multzoak izanik, berorien grabitate-zentruak indibiduo gehigarri bezala hartuko dira. Azkenik, jatorrizko laginerako beste aldagai batzuen balioak baztertuak izanik, datu-multzoan homogenoak ez direlako gehigarri bezala aurkeztea komenigarria izatea gerta daiteke.

Arazoa, bada, hauxe da: Errenkada-puntu edota zutabe-puntu horien koordinatu faktorialak lortzea, honela ardatz eta plano faktorialetan aurkeztu ahal dira.

Eskematikoki,



X^+ nahiz X_+ tauletan ditugun elementuen koordinatu faktorialak, elementu aktiboenak bezala oso erraz kalkulatu dira eta eskola frantsesaren datu-analisiaren esparruan egiten diren estatistikaren programa-paketetan ahalbide bezala aurki daitezke.

Elementu hauek nahiz eta analisisan, aktiboki, parterik hartu ez, ardatz eta planotan proiektatzean, batzuetan, kokapen pribilegiatua edukiz asko lagundu ahal dute emaitzen azalpenetan.

IX.6. SAILKAPEN METODOAK

Gaiaren hasieran esaten genuenez estatistika deskribatzaile anizkoitzean bi analisi-mota egin daitezke eta biak datu-hodei berberan aplikagarriak dira.

Orain arte ikusi ditugun analisi faktorialen bidez hodeiaren direkzio luzatuenak bilatzen dira eta beraien gain bildutako sakabanatzerik adierazgarriena azaltzen da.

Sailkapen metodoen bidez, ordea, hodeiaren elementu (indibiduo) hurbilenak sailetan kokatzen saiatuko gara. Lortutako sailak beraien elementutan homogeenak izango dira eta sailen arteko desberdintasunak ahal den hoberen finkatuko dira.

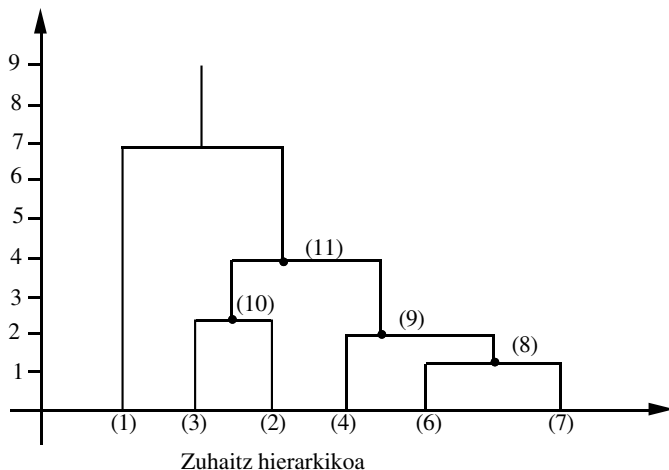
Normalki, sailkapen metodoak analisi faktorialaren emaitzak osotzeko eta aberasteko erabiltzen dira; kasu hauetan indibiduoek dituzten lehen koordinatu faktorialetatik abiatuz eta sailkapen algoritmo egoki bat aplikatuz sailak lortuko dira.

IX.6.1. Goranzko Sailkapen Hierarkikoa

Sailkapen automatikoko moten artean, zentruaren ingurukoa, beheranzko sailkapen hierarkikoa eta goranzko sailkapen hierarkikoa ditugu. Azken hau da garrantzitsuena duen koherentzia logikoagatik, hain zuzen.

Goranzko sailkapena egitean elementuak binaka elkartzen dira; honela, lehen etapan, bi hurbilenak elkartu ondoren $m-1$ elementu geratzen dira. Elkartuz lortu den elementua besteak bezala izanik, bigarren etapan, berriro bi hurbilenak elkartu ondoren $m-2$ elementu geratzen dira; honela jarraituz, azkenean, sail bat besterik ez da geratzen, honek elementu guztiak dituelarik.

Sortzen den zuhaitz hierarkikoan ebaketa bat eginez sailak lortzen dira.



Dakusagunez, lehen etapan 6, 7 elementuak elkartu dira eta elkarketa honi, (8) adabegiaz, ordenatuen ardatzean elkartze-indize txikiena dagokio; bigarren etapan 4, (6,7) elementuak elkartu dira (9) adabegian; hirugarrenean 3,2....

(8), (9), (10).... hierarkiaren adabegiak deitzen dira.

Elementuen arteko hurbilketak lehendik definitutako distantzia baten bidez finkatuko dira.

Goranzko sailkapen hierarkikoa irizpide desberdinen arabera egin daiteke, (saltu txikiena, batezbesteko distantzia, barizentrua...), hauen artean bariantzaren metodoa delarik garrantzitsuena.

IX.6.1.1. Bariantzaren metodoa

Bariantzaren metodoaren baldintza bakarra sailkatu nahi ditugun elementuak espazio euklidear batean aurkeztuak izan ahal izatea da. Honek zera esan nahi du: I multzoko m elementuak n aldagaien arabera edo n zutabeen arabera kontingentzi taulatan bezala, karakterizatuak daude. m elementuak \mathbb{R}^n espazio bektorialeko bektoreak dira, non distantzia euklidear bat definitua baita.

Bi dimentsioko espazioan hiruki zuzen batean Pitagoras-en teorema betetzen duen edozein distantziari, distantzia euklidearra deritzo.

“Aldagai-indibiduo” taulatan ohizko distantzia euklidearra definitzen da (ikus gai honen 3. atala) eta kontingentzi taulatan distantzia euklidear berezi bat χ^2 distantzia da (ikus gai honen 4. atala), hain zuzen. Ondorioz, Goranzko Sailkapen Hierarkikoa bariantzaren metodoaz edozein kasutan aplikatu daiteke.

Bariantzaren metodoa ez da bakarrik distantzien kalkuluan oinarritzen, inertzien kalkuluetan ere oinarritzen da, inertzia elementuen masa bider distantzia karratuen batura bezala definiturik. Kalkuluan sartzen diren elementuen masen batura bat baldin bada, inertzia bide batez bariantza izango da.

Köening-Huygens-en teorema

Suposa dezagun sailkatu nahi ditugun $N(I)$ hodeiaren m puntuei x_i deritzegula eta p_i berorien masak direla.

$g \in \mathbb{R}^n$ eta $y \in \mathbb{R}^n$ hodeiaren grabitate-zentrua eta hodeiaren edozein puntu izanik:

$$I_N(y) = \sum_{i=1}^m p_i(x_i, y)^2$$

Hau da, hodeiaren inertzia y puntuarekiko, puntuen masak bider puntu horrekiko distantzia karratuen batura da.

Köening-Huygens-en teoremaren bidez zera dakigu:

$$I_N(y) = d^2(g, y) + I_N(g)$$

hau da, $I_N(g)$ grabitate-zentruarekiko hodeiaren sakabanatzea minimoa dela.

Sortutako sailak beraien elementuetan homogenoak izatea eta sailen arteko desberdintasunak ahal den hoberen finkatuta izatea nahi dugunez, homogenotasun neurria bezala “sail barneko inertzia”-ren kontzeptua eta desberdintasun neurria bezala “sail arteko inertzia”-ren kontzeptua hartuko dira.

$N_1, N_2, \dots, N_k, \dots, N_h$; $N(I)$ hodeiko h sail disjuntu izanik eta $g_k, P(N_k), N_k$ sailaren grabitate-zentrua eta masa izanik:

$$I_{N_k}(g_k) = \sum_{x_i \in N_k} p_i d^2(x_i, g_k)$$

Hau da, N_k sailaren grabitate-zentruarekiko daukan “sail barneko inertzia” $I_{N_k}(g_k)$ balioa da.

$N(I)$ hodeiaren g grabitate-zentruarekiko “sail arteko inertzia” hauxe izango da:

$$\sum_{k=1}^h p(N_k) d^2(g_k, g)$$

Köening-Huygens-en teorema aplikatuz:

$$I_N(g) = \sum_{k=1}^h I_{N_k}(g_k) + \sum_{k=1}^h p(N_k) d^2(g_k, g)$$

Hau da, $N(I)$ hodeiaren inertzia, g grabitate-zentruarekiko sail barneko inertzien baturan gehi sail arteko inertzian deskonposatzen da batukorki.

I_N (g) konstante izanik, sail barneko inertzien batura minimizatzea eta sail arteko inertzia maximizatzea baliokidea da.

Algoritmoaren konkrezioa

Elkartzearen irizpidea, $N(I)$ hodeiaren inerziaren galtzea minimoa izatea da. Honela, lehen, klaseak elementu guztiak izanik, sail barneko inertzia minimoa daukaten bi klaseak elkartuko dira.

x_i, x_j lehen fasean ditugun m klaseen bi elementu izanik, sail barneko inertzia hauxe da:

$$\Delta I_{ij} = \frac{p_i p_j}{p_i + p_j} d^2 (x_i, x_j)$$

Honela, ΔI_{ij} minimoa daukaten bi elementuak elkartuko dira.

Sailkapenaren hurrengo klaseetan z_k eta x_i, x_j elkartzean sortu den beste edozein klaseren arteko sail barneko inertzia, azken honen grabitate-zentrua x izanik ondoko errepikapen-formulaz lortzen da.

$$\Delta I_{z,x} = \frac{1}{p_k + p_i + p_j} \left[(p_k + p_i) \Delta I_{ki} + (p_k + p_j) \Delta I_{kj} + p_k \Delta I_{ij} \right]$$

Era honetan, fase bakoitzean inerziaren galtze minimoa eragiten duten bi sailak elkartzean, hodeiarekiko sail barneko inertzia minimoa eta bide batez sail arteko inertzia maximoa edukitzea lortzen da.

Dakigunez, elementu guztiak sail batean elkartzean prozesua amaituko da.

Hodeiaren partizio ezberdinak zuhaitz hierarkikoaren ebaketaz lortu ahal dira.

IX.6.2. Adibidea

Hamazortzi Hego eta Ertamerikako estatuei buruz egindako Osagai Nagusien Analsiaren azterketan (3.4 atala) hiru lehen faktoreak kontutan harturik, hau da, hauen gain estatuak dituzten koordenatuen arabera sailkapen automatikoa egiteaz ikasketa osotu nahi dugu.

Goranzko sailkapen hierarkikoa bariantzaren metodoaz aplikatu ondoren lortutako emaitzak ondokoak dira:

Elkartze-indizeak salto handia eman duelako zuhaitzaren adarrak nahiko luzatzen direnean ebaketa eginez bost sail ditugu.

Dakusagunez, zuhaitzaren oinarrian sailetan ditugun indibiduoek datu-taulan daukaten posizioaren arabera ditugu; halaber, indibiduen masa 1 dela ikus daiteke.

Aplikaturako SPAD programen bidez sailen deskripzio zehatza eta lehen plano faktorialean estatu guztiak bere sailen zenbakiaz adierazita lortzen ditugu:

FIN DE L - ETAPE ** SEMIS **

CLASSE EFFECTIF POIDS CONTENU (ANCIENNES CLASSES)

1	3	3.000	1	3	6			
2	3	3.000	2	5	7			
3	5	5.000	4	8	11	12	15	
4	6	6.000	9	10	13	14	16	17
5	1	1.000	18					

VEN / 1 ARG / 2 MEX / 1 PAN / 3 URU / 2 BRA / 1 TXI / 2 KOL / 3 CRI / 4 GUA / 4 ECU / 3
PER / 3 PAR / 4 NIK / 4 DOE / 3 HON / 4 ELS / 4 BOL / 5

GRAPHIQUES DE DENSITE DES 18 INDIVIDUS

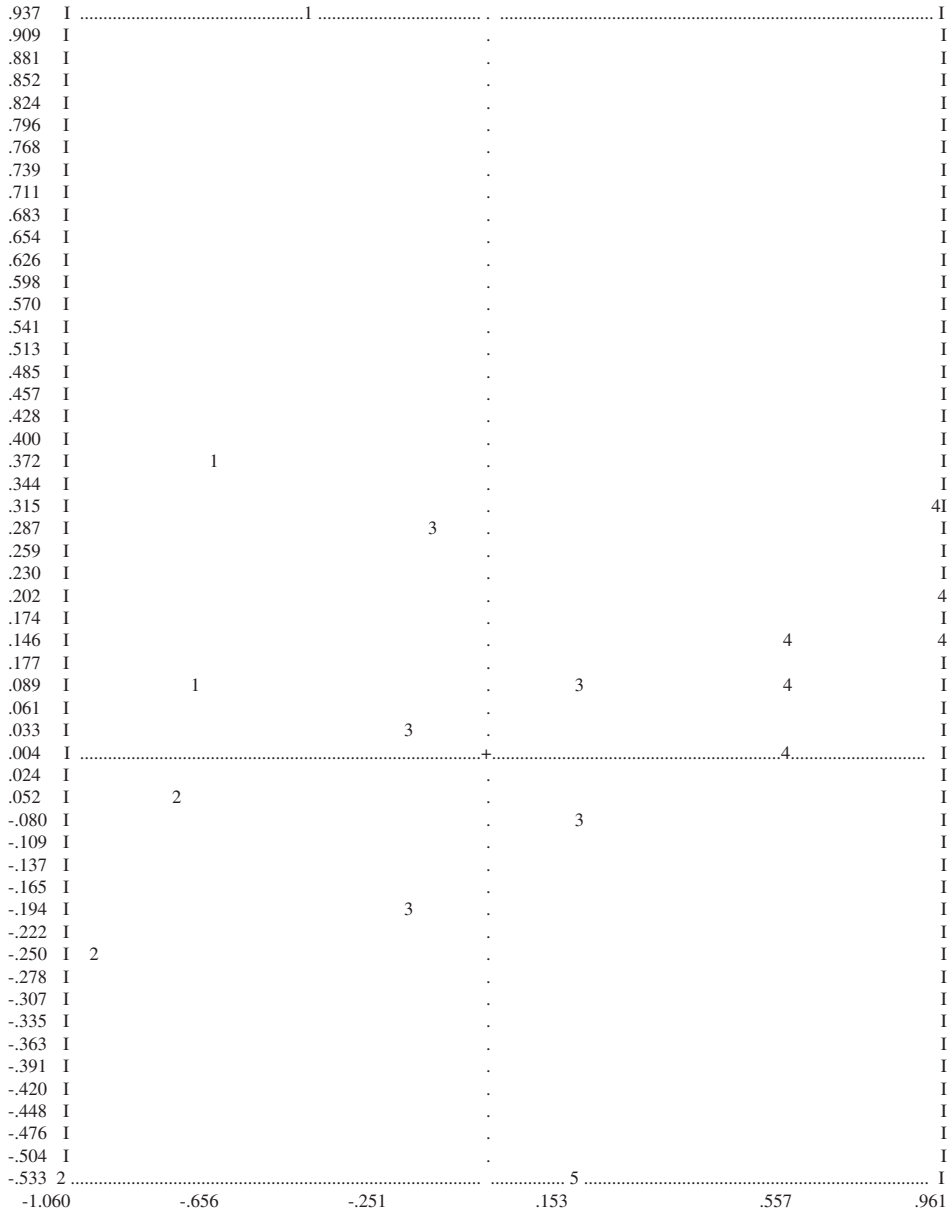
PLAN DE PROJECTION DES 18 POINTS SUR LES AXES 1 ET 2

AXE 1 / HORIZONTAL

AXE 2 / VERTICAL

1 POINTS A PLUS DE 2.3 ECARTS-TYPES DE L-ORIGINE ONT ETE RAMENES SUR LE CADRE

1



Dakusagunez, plano faktorial nagusian, sailak oso ondo bereziak agertzen dira. Lehen sailean eta goitik behera, Venezuela, Mexiko eta Brasil. Bigarrenean, Txile, Uruguay eta Argentina, hego-konoa hain zuzen. Hirugarrenean, Panama, Kolonia, Peru, Ekuador eta Dominikar Errepublika. Laugarrenean C. Rica, Nikaragua, Paraguai, Guatemala, Honduras eta El Salvador. Bosgarrenean, azkenik, duen inflazioagatik hain berezia den Bolivia.