

Oļegs Uļga-Rebrovs, Pēteris Grabusts

DATU PIRMSAPSTRĀDES METODES

DATU ANALĪZES UZDEVUMOS

Oļegs UŽGA-REBROVS, Pēteris GRABUSTS. 2024. *Datu pirmsapstrādes metodes datu analīzes uzdevumos*. Rēzekne: Rēzeknes Tehnoloģiju akadēmija.

Recenzenti:

- RTU **Sergejs PARŠUTINS**
- Mg.math., Dr.paed. **Ilmārs KANGRO** (Rēzeknes Tehnoloģiju akadēmija)

Zinātniskā monogrāfija sagatavota un izdota ar Rēzeknes Tehnoloģiju akadēmijas finansiālo atbalstu.



Publicēšanai rekomendējusi Rēzeknes Tehnoloģiju akadēmijas Studiju padome.

Monogrāfijā tiek apskatītas izplatītākās datu pirmsapstrādes metodes veiksmīgas datu analīzes nolūkos. Tiek doti datu analīzes pamatjēdzieni. Detalizēti aplūkotas darbības ar iztrūkstošiem datiem, liela uzmanība pievērsta dažāda veida anomāliju identificēšanai datos, datu atribūtu vērtību diskretizācijai un normalizācijai. Datu sagatavošanas procesā bieži ir nepieciešams pārveidot datus citos formātos ērtākai to turpmākai analīzei.

Materiāls ir papildināts ar pielikumiem, kuros iztirzāti jēdzieni, kas nepieciešami satura pilnvērtīgākai izpratnei. Datu analīzes metožu formālie apraksti papildināti ar daudzskaitlīgiem detalizētiem piemēriem, kas izskaidro katras metodes būtību un darbības principus.

Dotais materiāls paredzēts datorzinātņu un informācijas tehnoloģiju studiju programmu bakalauriem, maģistrantiem, doktorantiem un pētniekiem, kam nepieciešams veikt datu analīzi un sagatavošanu to turpmākai izmantošanai dažādos statistikas vai mākslīgā intelekta metožu pielietojumos.

Šis darbs tiek izplatīts ar internacionālo licenci



[Creative Commons Attribution 4.0 International Licence](https://creativecommons.org/licenses/by/4.0/)

© Rēzeknes Tehnoloģiju akadēmija, 2024

© Oļegs Užga-Rebrovs, Pēteris Grabusts

SATURA RĀDĪTĀJS

1. IEVADS	5
1.1. Dati, informācija un zināšanas	5
1.2. Datu veidi	7
1.3. Ievads datu analīzē	9
1.4. Sākotnējo datu pirmapstrādes procedūras	13
2. DARBĪBAS AR IZTRŪKSTOŠĀM ATRIBŪTU VĒRTĪBĀM	14
2.1. Problēmas formulējums	14
2.2. Vienkāršākās imputācijas metodes	19
2.3. K-tuvāko kaimiņu metode	22
2.4. Imputācija uz sagaidāmā rezultāta maksimizācijas pamata	30
2.5. Daudzkārtīgā imputācija	35
3. ANOMĀLIJU IDENTIFICĒŠANA DATOS	40
3.1. Kas ir anomālijas datos?	40
3.2. Statistiskās metodes anomāliju identificēšanai	43
3.3. Anomāliju identificēšanas metodes uz tuvības mēra starp objektiem pamata	50
3.4. Anomāliju identificēšana lineārajā regresijā	59
4. NEPĀRTRAUKTU ATRIBŪTU VĒRTĪBU DISKRETIZĀCIJA	72
4.1. Kas ir diskretizācija?	72
4.2. Diskretizācijas metožu klasifikācija	74
4.3. Diskretizācijas process	75
4.4. Vienkāršākās diskretizācijas metodes	77
4.5. Diskretizācija uz entropijas novērtējuma pamata	80
4.6. Uz χ^2 statistiku balstītas metodes	84
5. NORMALIZĀCIJA UN ATRIBŪTU VĒRTĪBU STANDARTIZĀCIJA	89
5.1. Definīcijas un piezīmes	89
5.2. Sākotnējo atribūtu vērtību normalizācijas metodes	89
5.3. Sākotnējo atribūtu vērtību standartizācijas metodes	94
6. ATRIBŪTU VĒRTĪBU TRANSFORMĀCIJA	97
6.1. Kad un kāpēc tiek izmantotas transformācijas?	97
6.2. Izplatītākie transformāciju veidi	100
6.3. Atribūtu vērtību transformācija normālā sadalījuma sasniegšanai	106
6.4. Atribūtu vērtību transformācija regresijas linearitātes sasniegšanai	117

IZMANTOJAMĀ LITERATŪRA.....	129
PIELIKUMI	141
P1. REGRESIJAS ATRIBŪTU VĒRTĪBU SADALĪJUMA PARAMETRU APRĒĶINI	142
P1.1. Nominālās skalas	142
P1.2. Ordinālās skalas	146
P1.3. Intervālu skalas	157
P1.4. Attiecību skalas	161
P2. METRIKAS DATU TELPĀ	162
P2.1. Kas ir metrika?.....	162
P2.2. Eiklīda attālums	163
P2.3. Manhetenas attālums.....	168
P2.4. Alternatīvi attāluma mēri starp objektiem	169
P2.5. χ^2 attālums.....	173
P2.6. Mahalanobisa distance	178
P3. NORMĀLĀ SADALĪJUMA LIKUMA PĀRBAUDE	183
P3.1. Vienkāršākās metodes sadalījumu pārbaudei uz normalitāti	183
P3.2. Sadalījuma normalitātes pārbaude pamatojoties uz sadalījuma parametriem	185
P3.3. Sadalījuma normalitātes pārbaude pamatojoties uz augstākas kārtas parametriem	191
P3.4. Sadalījuma normalitātes neparametriskā pārbaude.....	195

1. IEVADS

1.1. Dati, informācija un zināšanas

Šajā grāmatā ir aplūkotas datu pirmapstrādes metodes turpmākās analīzes nolūkos. Tāpēc ir nepieciešams ieviest formālu datu jēdziena definīciju, kā arī ar datiem cieši saistītos informācijas un zināšanu jēdzienus. Jāatzīmē, ka datu jēdziens ir saprotams vienkāršā intuitīvā līmenī. Tomēr, runājot par informācijas un zināšanu jēdzieniem, viss nav tik vienkārši. Interesentiem, kuri vēlas iegūt sīkāku informāciju par šiem svarīgajiem jēdzieniem un to savstarpējām attiecībām, ieteicams atsaukties uz darbiem [Burgin M., 2004; Zins Ch., 2007; Sanders J. D., 2016].

Mēs definēsim *datus* kā formālu objektu, faktu un notikumu attēlojumu (aprakstu) reālajā pasaulē. Šie objekti, fakti un notikumi var pastāvēt neatkarīgi vai atkarīgi no cilvēkiem. Gaisa temperatūru, vēja stiprumu un virzienu nosaka tikai dabas faktori un tie nav atkarīgi no cilvēkiem. Savukārt, preču klāsts veikalā un to izvietojums plauktos ir atkarīgs tikai un vienīgi no cilvēku rīcības.

Šīs grāmatas kontekstā ar datiem sapratīsim novērojumu, mērījumu, eksperimentu, iedzīvotāju aptauju rezultātus, dažāda veida tekstus, audio un video ierakstus un citas mūs interesējošas ārpusaules realitātes.

Datu prezentācijas formas var būt ļoti dažādas. Vienkāršākā un ļoti izplatītā forma ir datu tabula (matrica), kurā katrā rindā tiek parādīts objekts, subjekts vai cita entīcija un atribūtu vērtību (iezīmju) vektors, kas raksturo šo entīciju.

Jebkura pētījuma mērķis nav paši dati, bet gan datos ietvertā informācija.

Iespējams, zinātnes un praktiskās jomās nav neviena cita jēdziena, kam būtu tik daudz dažādu definīciju kā informācijas jēdzienam. Diemžēl šim jēdzienam nav vienotas vispārpieņemtas definīcijas. Kibernētikas tēvs Norberts Vīners informāciju definē kā to uztveres satura apzīmējumu, ko saņemam no ārējās vides. Statistikas informācijas teorijā informācija ir saistīta ar nenoteiktību, kas ietverta noteiktā ziņojumā.

Informācijas jēdziena definīciju dažādība ir saistīta ar kontekstu dažādību, kādos šis jēdziens tiek lietots.

Datu analīzes literatūra piedāvā arī daudzas informācijas jēdziena definīcijas. Darbā [Zins Ch., 2006] ir dotas daudzas šādas definīcijas.

Turpmāk mēs ņemsim par pamatu šādu vispārīgu definīciju. *Informācija liecina par struktūru, attiecībām, atkarībām un citām specifiskām iezīmēm, kas pastāv datos slēptā formā.*

Bieži vien datu vākšana ir īpaši izstrādāta, lai iegūtu atbilstošu informāciju. Informācijas iegūšana, pamatojoties uz socioloģiskām aptaujām, ir ļoti izplatīta socioloģiskajos un politiskajos pētījumos. Datu vākšanu var organizēt, lai apstiprinātu vai noraidītu hipotēzi, precizētu un padziļinātu iepriekš pabeigtu pētījumu rezultātus un daudziem citiem mērķiem.

Pēc savas būtības no datiem iegūtā informācija var būt precīza, aptuvena, neskaidra vai nepatiesa. Mēs varam runāt par precīzu informāciju, ja šāda informācija ir *apriori* ietverta datos un tā ir pareizi iegūta. Piemēram, dati satur precīzus divu mainīgo lielumu mērījumus, no kuriem viens ir funkcionāli atkarīgs no otra. Pamatojoties uz šādu datu apstrādi, mēs varam iegūt precīzu mūs interesējošo attiecību grafisku vai analītisko izteiksmi.

Lielākajā daļā gadījumu precīza informācija, kas iegūta no datiem, ir izņēmums, nevis likums. Parasti visi novērojumi un mērījumi tiek veikti ar lielākām vai mazākām kļūdām. Ja iepriekš dotajā piemērā mainīgo vērtību mērījumi tika veikti ar dažām kļūdām, tad iegūto datu apstrādes rezultātā var iegūt tikai aptuvenu vēlamās attiecības attēlojumu. Esošo informācijas precizitāti var modelēt, piemēram, izmantojot ticamības intervālus, kuros atrodas atkarīgā mainīgā faktiskās vērtības.

No datiem iegūtās informācijas nenoteiktību var izraisīt divi galvenie iemesli: (1) datu iekšējā nenoteiktība; (2) neskaidru subjektīvu datu izmantošana. Pieņemsim, ka dati atspoguļo gadījuma lieluma realizācijas. Apstrādājot šos datus, mēs varam izveidot šī nejaušā lieluma sadalījuma funkciju un/vai blīvuma funkciju un novērtēt iegūtās funkcijas parametrus. Šāda informācija atspoguļo sākotnējo datu iekšējo nenoteiktību.

Pieņemsim, ka dati atspoguļo iedzīvotāju aptaujas rezultātus par kādu konkrētu tēmu. Var sagaidīt, ka dažas atbildes uz uzdotajiem jautājumiem var būt neobjektīvas, nejaušas vai apzināti nepatiesas. Dabiski, ka šādu datu apstrādes rezultātā tiks iegūta vairāk vai mazāk neskaidra informācija. Tas ir jāpatur prātā, izmantojot šādu informāciju.

Bieži, ja nav iespējams iegūt objektīvus datus, tiek izmantoti subjektīvi ekspertu vērtējumi. Eksperti kalpo kā “mērinstrumenti”. Protams, šādu subjektīvu datu apstrādes rezultātā iegūta informācija būs vairāk vai mazāk neskaidra.

Ja ekspertiem ir grūti veikt nepārprotamu subjektīvu vērtējumu, viņi var izteikt savus spriedumus nenoteiktā formā: punktu aplēšu vietā dot intervālu vai izplūdušus skaitļus. Protams, šāda datu nenoteiktība radīs būtisku nenoteiktību informācijā, kas iegūta, apstrādājot šos datus.

Mūsdienās ir izstrādātas spēcīgas un efektīvas metodes neprecīzu un nenoteiktu datu apstrādei, taču šīs metodes netiek apskatītas šīs grāmatas kontekstā.

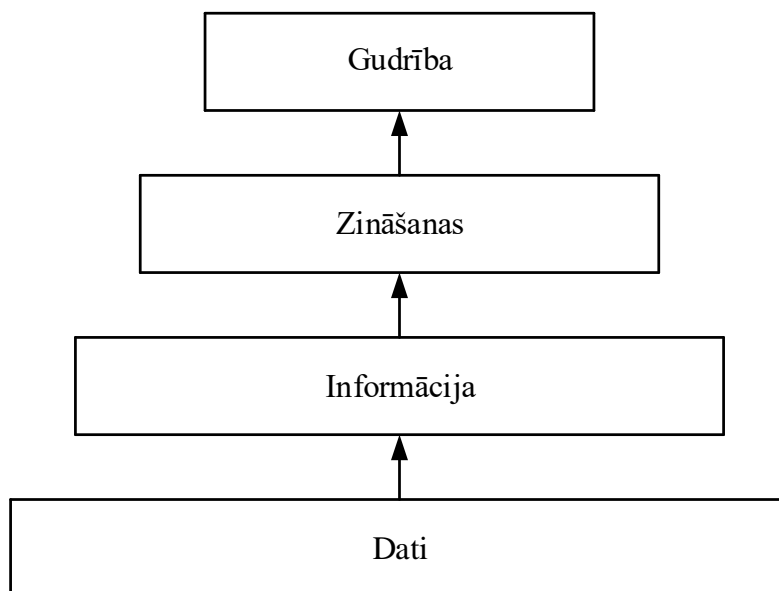
Mēs varam runāt par nepatiesu informāciju, ja šī informācija iegūta no dažāda veida viltus avotiem (datiem). Nepatiesas informācijas izmantošana var radīt ļoti nepatīkamas sekas. Uzliekot likmes uz konkrētu spēlētāju, lai uzvarētu derībās, pamatojoties uz nepatiesu informāciju, derību slēdzējs vienkārši zaudēs naudu. Tomēr neveiksmīgi rezultāti, kas iegūti nepatiesas informācijas izmantošanas dēļ, var būt daudz nopietnāki. Ja militārā operācija tiek plānota, pamatojoties uz nepatiesu informāciju par iespējamiem ienaidnieka spēkiem un darbībām, tas var radīt ievērojamus zaudējumus.

Acīmredzot attiecībā uz nepatiesu informāciju vienīgā rīcība, kas varētu būt piemērota, ir atpazīt šīs informācijas nepatiesību un nekādā gadījumā to neizmantot.

Zināšanu jēdziens ir ļoti plaša filozofiska kategorija. Datu analīzes kontekstā zināšanas definēsim kā datu apstrādes un analīzes rezultātā iegūtās informācijas izpratnes un sintēzes rezultātu.

Zināšanas var klasificēt pēc informācijas, uz kuras pamata tās iegūtas: precīzas zināšanas, aptuvenas zināšanas, neskaidras un nepatiesas zināšanas.

Ir skaidrs, ka datu analīzes kontekstā datu, informācijas un zināšanu jēdzieni ir savstarpēji saistīti jēdzieni. Darbā [Rowley J., 2007] šie jēdzieni tiek interpretēti, izmantojot šādu hierarhisko struktūru (sk. attēlu 1.1.1).



Attēls 1.1.1. Hierarhiskā struktūra, kas savieno jēdzienus dati, informācija, zināšanas un gudrība

Šo struktūru literatūrā bieži sauc par DIKW diagrammu (angļu valodas vārdu saīsinājums Data - Information - Knowledge - Wisdom). Struktūra skaidri atspoguļo hierarhiskās attiecības starp jēdzieniem: dati ir informācijas avots, zināšanas balstās uz informāciju. Gudrības jēdziens šajā

hierarhijā raksturo indivīda kognitīvās spējas, viņa rīcībā ir noteiktas zināšanas. Mēs neiedziļināsimies šajā koncepcijā, kurai ir izteikta filozofiska nozīme.

1.2. Datu veidi

Vārda “dati” izcelsme ir latīņu valodā, kas nozīmē “kaut kas, kas mums tiek dots”.

Datus var klasificēt dažādos veidos atkarībā no klasifikācijas pamatā esošajiem raksturlielumiem.

Vārda plašā nozīmē datus var klasificēt atkarībā no to iegūšanas un izmantošanas jomas. Šajā ziņā mēs varam runāt par bioloģiskiem, medicīniskiem, socioloģiskiem un citiem datiem.

Zinātniskajos pētījumos datus bieži klasificē pēc tā, kādā veidā tie iegūti:

- *Novērojumu dati*. Tie var būt vizuālo novērojumu rezultāti, dažāda veida mērījumu dati, datu ieraksti, sensoru dati un tamlīdzīgi.
- *Eksperimentālie dati*. Tie ir dati, kas iegūti iepriekš plānotu eksperimentu rezultātā.
- *Simulācijas dati*. Tādi dati tiek iegūti pētnieku interesējošā procesa datorsimulācijas rezultātā. Simulāciju bieži izmanto gadījumos, kad ir neiespējami vai ļoti dārgi ir tiešā veidā izpētīt reālu procesu.
- *Kompilētie dati*. Šāda veida dati ir datu apvienošanas rezultāts, kas iegūti no dažādiem avotiem. Šajā gadījumā bieži tiek izmantota atsevišķu datu pārveidošana. Rezultātā tiks iegūti viendabīgi dati.

Datu klasifikācija ir balstīta uz to atribūtu vērtību veidiem. Atribūtu vērtības var izmērīt dažādās skalās, kas ir datu klasifikācijas avots. Lai vienkāršotu turpmāko izklāstu, mēs izdarīsim pieņēmumu, ka mums ir darīšana ar viendimensiju datiem, t.i., ar datiem, ko raksturo viena atribūta vērtības.

Sīkāk apskatīsim atribūtu vērtību mērīšanas skalu klasifikāciju un uz tā pamata izveidoto datu klasifikāciju.

Zemāko līmeni skalu hierarhijā ieņem *nominālā skala*. Šī skala mēra atribūtus, kuru kategorijas var tikai iezīmēt vai nosaukt. Nosaukšanas skalu izmantošanas piemēri ir priekšmetu nosaukumi, cilvēku vārdi, kategoriskas pazīmes, piemēram, cilvēku dzimums, politisko partiju nosaukumi, konkursa uzvarētāju saņemto medaļu kategorijas. Dažkārt tiek lietoti nominālās skales kategoriju ciparu vai simboliskie kodējumi, piem. 1 – “vīrietis”, 2 – “sieviete”. Bet tie ir tikai citi dzimuma apzīmējumi, un nekas vairāk.

Ja atribūtu vērtības mēra nominālā skalā, ar šīm vērtībām nevar veikt nekādas darbības: vērtības nevar salīdzināt, pievienot, atņemt utt. Vienīgais, ko var izdarīt ar atribūtu vērtībām, kas izteiktas nominālā skalā, ir noteikt šo vērtību režīmu.

Dažreiz, ņemot vērā nominālās skales īpašības, to sauc arī par *nominālvērtības skalu*. Šis nosaukums labāk atspoguļo skales būtību, taču literatūrā parasti tiek lietots jēdziens *nosaukuma skala*.

Ordinālā skala ir nominālās skales paplašinājums tādā nozīmē, ka mēs varam veikt relatīvus objektu salīdzinājumus. Tipisks mērījumu piemērs ordinālajā skalā ir studentu rezultāti. Ja viens students eksāmenā ir ieguvis 7 punktus, bet otrs students iegūst 9, mēs varam teikt, ka otram studentam eksāmenā veicās labāk. Bet mēģināt novērtēt, cik labāk vai cik reizes labāk, ir bezjēdzīgi.

Citi ordinālās skales mērījumu piemēri ir subjektīvi temperatūras mērījumi: auksti, silti, karsti; subjektīvi sāpju novērtējumi: vieglas, vidēji smagas, smagas un Likerta skales mērījumi. IQ mērījumus klasificē arī pēc ordinālās skales, lai gan diezgan bieži šie mērījumi tiek kļūdaini klasificēti kā intervālu skala.

Bieži vien literatūrā nominālā un ordinālā skala tiek apvienotas zem *kategoriskās skales* vispārīgā nosaukuma.

Ja atribūtu vērtību mērījumus veic nominālā vai ordinālā skalā, tad datus sauc par *kvalitatīvajiem datiem*.

Augstāku līmeni ieņem *intervālu skala*. Tipisks šīs skalas mērījumu piemērs ir temperatūras mērīšana. Skalas raksturīga īpašība ir tāda, ka mēs varam salīdzināt mērījumu rezultātus un secināt, cik viens rezultāts ir lielāks vai mazāks par citu. Piemēram, ja mums ir divi temperatūras mērījumi: $t_1 = 15^{\circ}\text{C}$, $t_2 = 20^{\circ}\text{C}$, tad mēs varam pareizi norādīt, ka otrā vērtība ir par 5 vienībām lielāka nekā pirmā vai pirmā vērtība ir par 5 vienībām mazāka nekā otrā. Bet apgalvojumam, ka 20°C temperatūra ir divreiz lielāka par 10°C , šajā skalā nav nozīmes. Šī skalas iezīme ir saistīta ar to, ka intervālu skalai nav nosacīta sākuma punkta. Tādējādi temperatūra 0°C tiek uzskatīta par ūdens sasalšanas temperatūru. Ar tādiem pašiem panākumiem jebkuru citu temperatūru var uzskatīt par skalas nulles punktu. (Fārenheita skalā ūdens sasalšanas punkts ir 32°F).

Intervāla skalā mēs varam runāt par izmērīto vērtību intervālu attiecību. Tādējādi var apgalvot, ka temperatūras intervāls ($20^{\circ}\text{C} - 40^{\circ}\text{C}$) ir divreiz lielāks par intervālu ($10^{\circ}\text{C} - 20^{\circ}\text{C}$).

Mērījumu skalu hierarhijas augstākais līmenis ir *attiecību skala*. Tipiski mērījumu piemēri attiecību skalā ir auguma, svara vai attāluma mērījumi. Ja atribūtu vērtības mēra attiecību skalā, tad var norādīt, par cik viena vērtība ir pārāka par citu. Šādām atribūtu vērtībām var aprēķināt režīmu, mediānu un vidējo.

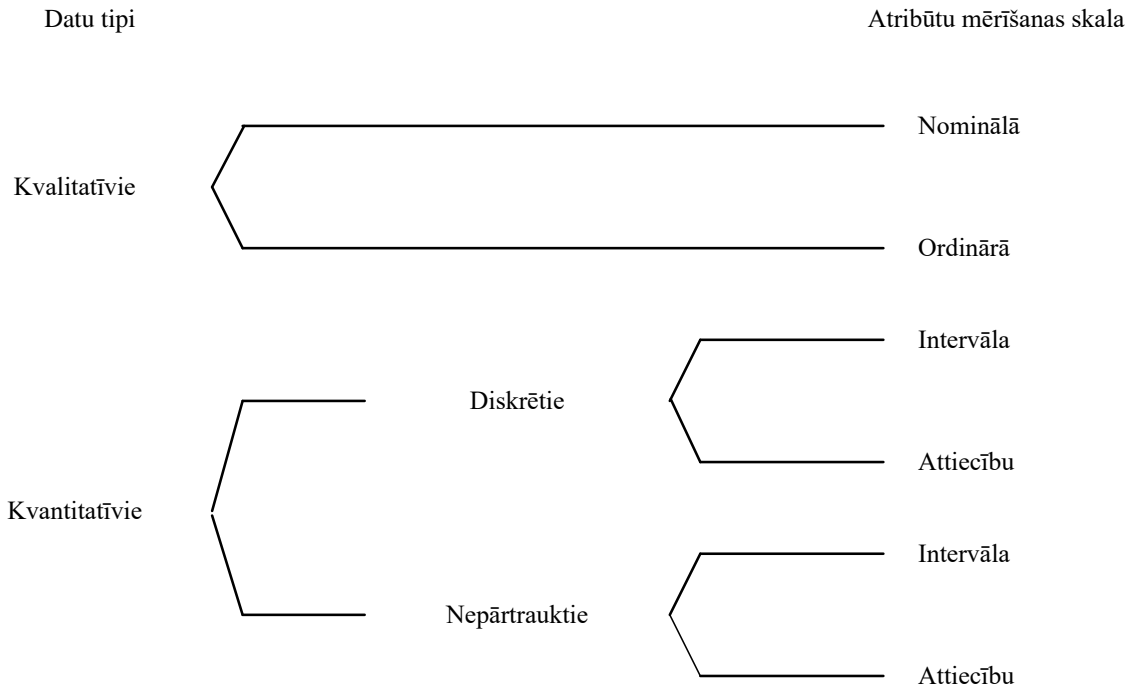
Ja atribūtu vērtības mēra intervālu skalā vai attiecību skalā, tad šādus datus sauc par *kvantitatīviem datiem*.

Vēl viena plaši izmantota vienfaktoru datu klasifikācija ir balstīta uz atribūtu vērtību būtību.

- *Nepārtraukti dati*. Atribūts var iegūt jebkuru vērtību noteiktā intervālā. Nepārtrauktu datu piemēri ir personas augums un svars, viņu mājsaimniecības mantas vērtība vai algas.

- *Diskrēti dati*. Atribūtam noteiktā intervālā var būt tikai noteiktas vērtības. Piemērs varētu būt studentu skaits lekcijā.

Vispārīgākā ieviestos jēdzienus un definīcijas. Attēlā 1.2.1 ir sniegta vispārīga vienfaktoru datu klasifikācija.



Attēls 1.2.1. Datu tipu grafiskais attēlojums

Tabulā 1.2.1 parādīti aritmētisko darbību veidi, ko var veikt ar dažāda veida datiem.

Ir arī citas pieejas datu klasifikācijai. Tā kā šīs pieejas ir ierobežoti izmantojamas, tās šajā sadaļā nav izklāstītas.

Iepriekš apskatītās datu tipu klasifikācijas tika uzrādītas vienfaktoru datu gadījumiem. Ja visu daudzfaktoru datu atribūtu vērtības tiek mērītas vienā skalā, visi iepriekš minētie jēdzieni un definīcijas ir pilnībā piemērojami šādiem datiem.

Parasti reālos daudzfaktoru datos atribūtu vērtības tiek mērītas dažādās skalās. Šādos gadījumos ir runa par *jauktiem datiem*. Šāda situācija ievērojami apgrūtina datu pirmsapstrādi un to turpmāko analīzi.

Tabula 1.2.1

Aritmētisko darbību saraksts dažādiem datu tipiem

Datu tipi un skalas	Aritmētiskās darbības			
Kvalitatīvie	Saskaitīšana	Atņemšana	Reizināšana	Dalīšana
Nominālā	nav	nav	nav	nav
Ordinārā	nav	nav	nav	nav
Kvantitatīvie	Saskaitīšana	Atņemšana	Reizināšana	Dalīšana
Intervāla	ir	ir	nav	nav
Attiecību	ir	ir	ir	ir

1.3. Ievads datu analīzē

Tā kā datu pirmapstrāde, kas ir šīs grāmatas galvenā tēma, ir paredzēta datu turpmākai analīzei, šī sadaļa iepazīstina ar datu analīzes pamatjēdzieniem. Tālāk sniegtais materiāls ir balstīts uz datiem, kas sniegti darbā [Uzhga-Rebrov O.I., 2021].

Statistika jau sen ir izstrādājusi spēcīgas metodes hipotēžu un teoriju patiesuma pārbaudei. Kā norādīts darbā [De Mast J., Kemper B. P. H., 2009], eksperimentālo metožu teorijas, hipotēžu pārbaude un modeļu veidošana ir daži no svarīgākajiem statistikas ieguldījumiem. Hipotēze tiek pārbaudīta sekojošā veidā. Tiek apkopotu eksperimenti, novērojumi vai attiecīgie dati. Šie dati tiek analizēti un, pamatojoties uz iegūtajiem rezultātiem, tiek izdarīti secinājumi par to, vai izvirzītā hipotēze ir patiesa vai noraidāma.

Šādu datu analīzi sauc par *apstiprinošo datu analīzi (ADA) (Confirmatory Data Analysis - CDA)*. Faktiski šāda analīze ir deduktīvs secinājums. Autors [Haig B.B., 2005] atzīmē, ka 150 gadus dabaszinātnēs izvēlēta metode ir bijusi hipotētiskā indukcija. Hipotētiski deduktīvā metode parasti tiek izteikta minimālos terminos: pētniekam ir jāierosina hipotēze vai teorija un tā ir vispārīgi jāpārbauda, iegūstot vienu vai vairākas prognozes. Šīs prognozes tiek pakļautas tiešai empīriskai pārbaudei. Ja prognozes apstiprina dati, tad rezultāts tiek pieņemts kā hipotēzes vai teorijas apstiprinājums. Ja prognozes nesakrīt ar datiem, tad šis fakts tiek pieņemts kā hipotēzes vai teorijas noliegums.

ADA ir ļoti nozīmīga loma zinātnisko, tehnisko, ekonomisko, sociālo un citu pētījumu rezultātā iegūto secinājumu apstiprināšanā vai noliegšanai.

Savukārt 20. gadsimta otrajā pusē, strauji attīstoties informācijas tehnoloģijām, sāka uzkrāties milzīgi dažāda veida datu apjomi. Paši dati ir pārskats par situāciju kādā jomā. Lietotājus var ļoti interesēt šajos datos slēptās zināšanas. Piemēram, uz datiem virzītajā mārketingā vēlams iegūt zināšanas par pārdošanas procesa īpatnībām, lai pieņemtu pārdomātus lēmumus par pārdošanas procesa efektivitātes uzlabošanu.

Vienkārši izsakoties, ir nepieciešama datu analīze, kas ļauj identificēt datus netiešā veidā attēlotās datu struktūras, modeļus un/vai atkarības atbilstošās iezīmes. Šāda veida analīzi sauc par *skaidrojošo datu analīzi (SDA) (Exploratory Data Analysis EDA)*.

SDA paveids ir *paredzamā datu analīze (PDA) (Predictive Data Analysis- PDA)*. Šīs analīzes būtība ir identificēt nozīmīgas iezīmes ierobežotos datu apgabalos, lai tās turpmāk varētu attiecināt uz plašāku datu apgabalu. PDA plaši izmanto medicīnā. Piemēram, izmantojot datus, kas saistīti ar noteiktas slimības diagnostiku un ārstēšanu, tiek iegūti secinājumi un prognozes, ko var veiksmīgi izmantot citos medicīniskajos pētījumos.

SDA rašanās un attīstība galvenokārt ir saistīta ar amerikāņu pētnieka Džona Tukeja un viņa kolēģu vārdu. No daudzajiem šī autora darbiem jāmin vispilnīgākais darbs [Tukey J.W, 1977]. Faktiski šis darbs nosaka SDA formālos pamatus un izskata dažādus šīs analīzes aspektus.

Kā var formāli definēt SDA? Mēs sniegsim divas alternatīvas šī jēdziena definīcijas.

Darbā [Tukey J.W, 1977] autors SDA raksturo kā (1) filozofiju vai attieksmi, nevis noteiktu formālu procedūru kopumu; (2) datu visaptverošas izpratnes process, lai no tiem iegūtu informāciju; (3) vienkāršu aprakstošu novērtējumu izmantošana datu apkopošanai un atkārtotai izteikšanai; (4) uzsverot datu grafiskos attēlojumus; (5) elastība, piemērojot analīzi datu struktūrām un attiecīgi arī neatpazītiem modeļiem; (6) koncentrēties uz eksperimentālu modeļu veidošanu un hipotēžu ģenerēšanu.

Darbā [Behrens J.T., 1997] autors sniedz savu SDA definīciju. Mēs sniegsim tā definīciju citāta veidā no šī darba.

“Šajā rakstā tiek uzskatīts, ka SDA ir īpaša datu analīzes tradīcija, kas izriet no Džona Tukeja un viņa domubiedru darba. Šo SDA tradīciju var raksturot ar: (a) neatkarīgas izpratnes saglabāšanu par datiem, kas pievēršas plašajam jautājumam “kas šeit notiek?”; b) datu grafisko attēlojumu uzturēšana; c) koncentrēšanās uz eksperimentālu modeļu veidošanu un hipotēžu ģenerēšanu, izmantojot iteratīvu modeļu specifiskāciju, modeļa atkārtotas specifiskācijas procesu; d) stabilu aplēšu izmantošana, atkārtota izteiksme un apakškopu analīze; (e) skepticisma attieksme, elastība attiecībā uz izmantotajām metodēm”.

Jebkuras skaidrojošās analīzes mērķis ir atpazīt un izskaidrot sākotnējā datu kopā slēptās pazīmes. Literatūrā nav vispārpieņemta “raksturīgo pazīmju” jēdziena. Angļu valodas literatūrā tos bieži sauc par “paraugiem” vai “parādībām”. Modeļu piemēri var būt standarta vai specifiski atribūtu vērtību sadalījumi, atsevišķu atribūtu vērtību korelācija, klasteru klātbūtne datu kopā un daudzas citas iespējas.

Kādi ir SDA galvenie mērķi? Darbā [De Mast J., Kempers B.P.H., 2009] autori apgalvo, ka “SDA mērķis nav izdarīt secinājumus par iepriekš definētiem pētījuma jautājumiem (vai tā ir modeļa konstruēšana, parametru novērtēšana vai hipotēzes apstiprināšana vai noraidīšana), patiesībā SDA problēmu risināšana bieži tiek izmantota savāktajiem datiem bez precīzi definētām hipotēzēm”.

Vissvarīgākais par SDA ir tas, ka analīzi var veikt, sākotnēji nemaz neformulējot hipotēzes, kuras varētu secināt no datiem. Tas nenozīmē, ka SDA tiek veikts “tukšā vietā”, t.i., uzsākot analīzi mums nav nojausmas, ko šīs analīzes rezultātā var iegūt. Parasti tiek formulēti pētījuma jautājumi, uz kuriem var atbildēt ar analīzi. SDA mērķis ir apstrādāt un parādīt sākotnējo datu kopu tā, lai varētu identificēt konkrētus datu modeļus. Pamatojoties uz iegūtajiem paraugiem, var izvirzīt hipotēzi par to, kādas zināšanas ir ietvertas datos. Vēl viens labi izstrādāta SDA mērķis ir identificēt iekšējos faktorus datos, kas izraisa iegūtos modeļus.

Kopumā SDA ietver šādu vispārinātu procedūru kopumu [De Mast J., Kempers B.P.H., 2009]:

1. datu attēlošana;
2. būtisku pazīmju identificēšana;
3. pamanāmu pazīmju interpretācija.

SDA izmantotie neapstrādātie dati ir atbilstoši jāapstrādā un jānorāda piemērotā (bieži vien grafiskā) formā, lai ļautu identificēt datus raksturīgās pazīmes.

SDA pēc būtības ir induktīvā pieeja, t.i., analīzes process tiek veikts saskaņā ar shēmu “no konkrētā uz vispārīgo”. Kad tiek identificēts datu modelis, tas ir jāinterpretē (jāpaskaidro). Šī ir abduktīvo secinājumu darbības joma. Mēs iepazīstināsim ar abduktīvo secinājumu definīcijām, ko ierosinājuši labi pazīstami pētnieki SDA jomā.

[De Mast J., Kempers B.P.H., 2009]: “Secinājumu veidu, kad, apvienojot jēdzienus un analogijas, daļas ir potenciāli saistītas un no tā izriet pamanāmas iezīmes, bieži sauc par abduktīvu.”

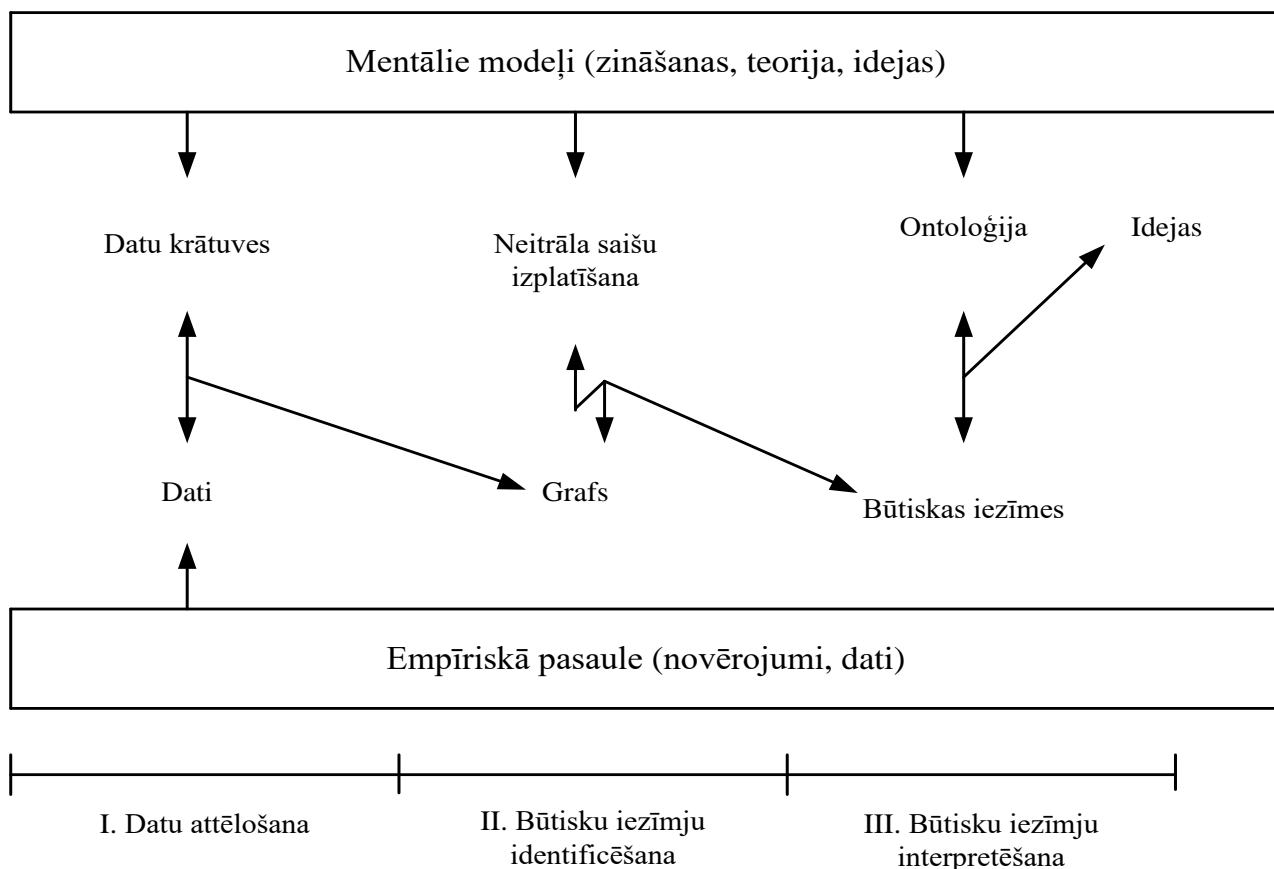
[Cheng H.Y., 1994]: “Datu analīzēm ir grūti sākt un izskaidrot datus bez jebkādas atsaucē. Tradicionāli pētnieki ir klasificējuši spriešanas veidus pētījumos kā indukciju (vadoties pēc pierādījumiem) un dedukciju (vadoties pēc teorijas vai hipotēzes). Patiesībā ir trešais veids: abdukcija. Abduktīva spriešana ne vienmēr sākas ar pilnībā izstrādātiem modeļiem vai vispār bez

modeļiem... Abduktīvajā testēšanā SDA būtība ir nevis atteikšanās no modelēšanas un sākotnējām koncepcijām, bet gan izpēte un iedrošinājums nesākt analīzi tikai ar “stingru” sākotnējo koncepciju... Abdukcija – tā nav simboliska loģika, bet gan kritiskā domāšana”.

[Haig B.B., 2005]: “Abduktīva secinājuma pamatideja ir meklējama amerikāņu filozofa un zinātnieka Pīrsa darbos. Turpmākā filozofijas un mākslīgā intelekta attīstība ir balstīta uz Pīrsa idejām, lai sniegtu pētniekiem būtisku izpratni par abduktīvo spriešanu”.

Apkopojot iepriekš minētās definīcijas, var apgalvot, ka abduktīvie secinājumi ir saistīti ar hipotēžu meklēšanu, kas vislabāk izskaidro datus atpazītos modeļus. Lai to izdarītu, pētniekam ir jāpatur prātā vai jāizmanto ārēji avoti, lai analizētu visus iespējamus identificēto modeļu skaidrojumus. Abduktīvo secinājumu procesus var attēlot diagrammas veidā, kas parādīts attēlā 1.3.1 [De Mast J., Kempers B.P.H., 2009]. Jāpiebilst, ka šīs diagrammas autori pieņem, ka visus mentālos modeļus var attēlot grafiku veidā. Realitātē potenciālo mentālo modeļu reprezentācijas iespējas ir daudz plašākas. Grafiki ir tikai viens alternatīvs veids, kā attēlot šādus modeļus.

Darbā [Haig B.B., 2005], atsaucoties uz citu avotu, ir parādīts, ka, ģenerējot hipotēzes, tiek nošķirta eksistenciālā un analogā abdukcija. *Eksistenciālā abdukcija* postulē iepriekš nezināmu paraugu esamību, bet *abdukcija pēc analogijas* balstās uz iepriekšējiem hipotēžu formulēšanas gadījumiem, lai atrastu analogu, piemērotu pašreizējam SDA uzdevumam. Kā norādīts darbā [Haig B.B., 2005], eksistenciālā abdukcija ir abdukcijas veids, ko plaši izmanto faktorus aprakstošo hipotēžu ģenerēšanai. Skaidrojošās faktoru analīzes mērķis ir atvieglot slēpto mainīgo noteikšanu. Lai gan skaidrojošās faktoru analīzes būtība un mērķis tiek diskutējams, to var saprast kā abduktīvu metodi hipotēžu ģenerēšanai.



Attēls 1.3.1. Abduktīvo secinājumu shematisks procesu attēlojums

Standarta SDA darbojas saskaņā ar šādu shēmu:

Uzdevums → Dati → Analīze → Modelis → Secinājums.

Šajā sistēmā secinājumu rezultātā tiek formulēta hipotēze, kas ticami izskaidro datu modeļus.

Pēdējos gados standarta SDA procedūras ir papildinātas ar procedūrām, kas ļauj pārbaudīt hipotēzes t.i., ir iekļauti arī ADA elementi.

Kā norādīts darbā [Cheng H.Y, 2010], SDA atbalsta mainīgos lielumus, konstrukcijas un hipotēzes, kas tiek uzskatītas par vērtīgām, bet ADA veic nākamo soli, lai to apstiprinātu. Tomēr, izmantojot atkārtotu paraugu ņemšanu (*resampling*), skaidrojošo datu analīze var vienlaikus atbalstīt un apstiprināt modeli.

Atkārtota paraugu ņemšana skaidrojošās datu analīzes kontekstā attiecas uz vispārināšanu starp paraugiem, nevis uz rezultātu, kas iegūts no viena parauga. Lai iegūtu pareizus vispārinātus rezultātus, tiek izmantota šķērsvalidācija (*crossvalidation*).

Šeit ir citāts no [Cheng HY, 2010], kurā autors parāda mūsdienu uzskatus par SDA. “Viena no problēmām ar standarta pieejām SDA ir tā, ka SDA raksturlielumi ir saistīti gan ar datu atribūtiem (sadalījums, mainīgums, linearitāte, novirzes, mērījumu skalas utt.), gan ar galīgajiem mērķiem (klasteru atpazīšana, mainīgo un slēptu modeļu un sarežģītu saišu atpazīšana). Faktiski datu atribūtu apstrāde ir vidēja posma procedūra, nevis beigu procedūra. Mūsdienu pieejas SDA pamatā ir koncentrēšanās uz mērķiem, nevis vērtībām. Tāpēc dažas modernas pieejas var apstiprināt analīzes rezultātus un tos var salīdzināt ar SDA - ADA kombināciju.” Lai ilustrētu šo apgalvojumu, tabulā 1.3.1 [Liu B. Q., 2014] ir apkopoti tradicionālās SDA un mūsdienu SDA salīdzināšanas rezultāti.

Tabula 1.3.1

Tradicionālā un mūsdienu SDA salīdzinājums

	Tradicionālais SDA	Mūsdienu SDA
Tips	Orientēts uz vērtībām	Orientēts uz mērķi
Forma	Vizualizēta raksturojoša analīze	Vispārīga raksturojoša analīze
Līdzekļi	Vienkārša aritmētika Vienkāršoti attēlojumi	Aprakstošā statistika Attīstītas datu vizualizācijas metodes Datu iegūšanas metodes
Galvenās iezīmes	Nedod noslēguma atbildes	Var nodrošināt noslēguma atbildes

Datu iegūšana (*Data Mining*) ir kļuvusi plaši izplatīta dažādās jomās. Lai iegūtu detalizētu ievadu šajā zinātniskajā un lietišķajā jomā, var atsaukties, piemēram uz [Han J., Kamber M., Pei J., 2012; Zaki M.J., Wagner M.Jr., 2014; Xu G., Zong Y., Yang Zh ., 2013]. Papildus šiem avotiem plaša informācija par datu iegūšanu ir atrodamā arī citās monogrāfijās un daudzos rakstos.

Kas ir datu iegūšana? Tas ir dažādu informācijas tehnoloģiju pielietojums lielām datu kopām, lai tos apstrādātu un identificētu interesējošos paraugus. Jāatzīmē, ka termins “datu iegūšana”, šķiet, nav pārāk labs veids, kā aprakstīt šo procesu. Tāpēc literatūrā dažreiz tiek izmantots alternatīvs nosaukums: zināšanu atklāšana no datiem (*knowledge discovery from data*). Tomēr termins “datu iegūšana” ir daudz izplatītāks nekā otrais termins, tāpēc tā lietošana ir kļuvusi par tradīciju visā pasaulē.

Kādas ir līdzības un atšķirības starp SDA un datu iegūšanu? Būtībā šīs ir ļoti līdzīgas pieejas. Abas pieejas apstrādā un parāda datus īpašos veidos, lai identificētu modeļus, t.i., nozīmīgas datu pazīmes. Abās pieejās mērķis ir izskaidrot radušos modeļus. Abas pieejas izmanto ADA, lai pārbaudītu iegūtos rezultātus.

Atšķirība starp abām pieejām ir šāda. Datu iegūšanā plaši tiek izmantotas specifiskas metodes, kas paredzētas tieši liela apjoma datu apstrādei automatiskajā režīmā. Datu iegūšanā plaši tiek izmantotas specifiskas datu pirmapstrādes metodes, piemēram, iztrūkstošo atribūtu vērtību atjaunošana, atribūtu skaita samazināšana, noņemot neefektīvās atribūtu vērtības un citas. Jebkurš datu iegūšanas process kā neatņemama sastāvdaļa ietver vienu vai otru datu analīzes veidu.

1.4. Sākotnējo datu pirmapstrādes procedūras

Kvalitatīvu analīzi var veikt tikai kvalitatīviem datiem. Datu kvalitātes sastāvdaļas ir to precizitāte, pilnīgums un saskaņotība.

Datu precizitāte ir saistīta ar atribūtu vērtību precizitāti un pareizību.

Datu pilnīgums attiecas uz datu paraugu un atribūtu vērtību pieejamību. Dati ir pilnīgi, ja datos ir visi nepieciešamie un pietiekamie modeļi un atribūtu vērtības.

Datu saskaņotība attiecas uz atbilstošo datu modeļu un atribūtu klātbūtni, pareizu kodēšanu un pareizu datu prezentācijas formātu.

Reālās situācijās avota dati var būt neprecīzi, nepilnīgi un nesakarīgi. Tam ir daudz iemeslu: neprecīzi datu ieraksti, datu vākšanas rīku bojājumi, personāla kļūdas, ievadot datus datorā, un daudzi citi, bieži subjektīvi iemesli. Tāpēc tieša neapstrādātu datu analīze var novest pie neapmierinošiem rezultātiem.

Lai veiktu kvalitatīvu datu analīzi, dati ir iepriekš jāapstrādā tā, lai tie būtu pieņemamā kvalitātē turpmākai analīzei.

Datu pirmapstrāde prasa daudz laika un pūļu. [Pyle D., 1993] atzīmēja, ka datu pirmapstrādei var būt nepieciešami līdz 60% no kopējā laika, lai veiktu pilnīgu datu analīzi. Tas arī parāda datu pirmapstrādes nozīmi.

Datu pirmapstrādi var definēt kā dažādu veidu procedūru secības veikšanu, lai oriģinālos datus iegūtu tālākai analīzei piemērotā formā. Visu attiecīgo procedūru kopumu var iedalīt četrās galvenajās kategorijās [Acuna E., 2011; Alasadi S., 2017].

1. *Datu attīrīšana*. Šis ir pirmais datu pirmapstrādes solis. Tas ietver iztrūkstošo atribūtu vērtību aizpildīšanu vai izņemšanu, trokšņainu datu izlīdzināšanu, anomāliju identificēšanu un izņemšanu no datiem, kā arī datu neatbilstību novēršanu.

2. *Datu integrācija*. Ja analīzei tiek izmantoti dati no dažādiem avotiem, šie dati ir jāintegrē vienotā veselumā. Tas var radīt nepieciešamību pielāgot atsevišķas datu sadaļas, lai tām būtu vienāds formāts, viena kodēšanas sistēma un uzdošana.

3. *Datu transformācija*. Daudzi datu analīzes algoritmi darbojas labāk, ja atribūtu vērtības tiek normalizētas vai mērogotas noteiktā veidā. Normalizācijas veikšana ir nepieciešama, lai izmantotu uz attālumu balstītus algoritmus, jo attālumu rezultāti starp atribūtiem ar lielākām vērtībām pārsvērs attālumu mērīšanas rezultātus starp atribūtiem ar mazākām vērtībām.

Citas datu pārveidošanas metodes ietver datu vispārināšanu. Šādas metodes rada jaunus atribūtus, apvienojot esošos atribūtus vai izmantojot augstāka līmeņa koncepcijas.

4. *Datu samazināšana*. Datu samazināšanas galvenā ideja ir izveidot jaunu datu izlasi, kas ir mazāka par sākotnējo datu kopu, bet saglabā visas savas pamatīpašības. Darbā [Han J., Kamber M., Pei J., 2012] tiek prezentētas šādas datu samazināšanas stratēģijas.

- *Dimensijas samazināšana*, kas saistīta ar neatbilstošu vai lieku atribūtu noņemšanu.

- *Datu saspiešana*, kad tiek izmantoti kodēšanas mehānismi, lai izveidotu samazinātu vai saspiestu sākotnējo datu attēlojumu.

- *Skaita samazināšana*, kur dati tiek aizpildīti vai tiek novērtēti, izmantojot alternatīvus mazākus to attēlojumus, piemēram, parametriskus modeļus vai neparametriskas metodes, piemēram, klasterizāciju.

- *Jēdzienu hierarhijas diskretizācija un ģenerēšana*, kur neapstrādātu datu vērtības tiek aizstātas ar rangiem (intervāliem) vai augstākiem konceptuālajiem līmeņiem.

- *Datu elementu atlase*, kur tiek atlasīta labāko datu vienību (paraugu) apakškopa no to pilnas kopas. Šī pieeja ir efektīva laika un izmaksu ziņā.

Lielākā daļa šajā sadaļā aprakstīto datu pirmapstrādes procedūru tiks detalizēti apskatītas nākamajās nodaļās.

2. DARBĪBAS AR IZTRŪKSTOŠĀM ATRIBŪTU VĒRTĪBĀM

2.1. Problēmas formulējums

Datu analīzē iztrūkstoši dati parasti ir kā likums, nevis izņēmuma gadījums. Datu iztrūkumam ir daudz iemeslu: respondentu nevēlēšanās atbildēt uz dažiem socioloģisko un politisko pētījumu jautājumiem, datu savākšanai izmantoto sensoru darbības traucējumi, datu daļas zaudējumi, personāla kļūdas datu ievades procesā utt.

Datu iztrūkums rada nopietnas problēmas datu tālākai analīzei [Kang H., 2013]:

- Datu iztrūkums samazina datu analīzes metožu statistisko nozīmīgumu.
- Datu iztrūkums var novest pie neobjektīviem vērtējumiem un secinājumiem.
- Datu izlases reprezentativitāte jūtami samazinās.
- Datu analīze kļūst daudz sarežģītāka.

Kā atzīmēts darbā [Dong Y., Peng Ch.-Y.J., 2013], iztrūkstošie dati var parādīties divos līmeņos: datu vienumu līmenī (objekti, priekšmeti utt.) un datu struktūras (atribūtu) līmenī. Datu iztrūkums vienumu līmenī rodas, ja netiek apkopota vai nav pilnīga informācija par šiem vienumiem jeb elementiem. Datu iztrūkums atribūtu līmenī nozīmē, ka datos trūkst atsevišķu atribūtu vērtību.

Kas attiecas uz datu iztrūkumu vienumu līmenī, šo problēmu var atrisināt, vispirms organizējot visu nepieciešamo datu savākšanu. Tā drīzāk ir organizatoriska, nevis konceptuāla problēma.

Daudz grūtāk ir tikt galā ar iztrūkstošajām atribūtu vērtībām. Tagad ir vispāratzīts, ka vienīgais veids, kā atrisināt šo problēmu, ir iztrūkstošo vērtību imputācija (*imputation*), t.i., aizstāšana ar citām atbilstošām vērtībām. Statistikā imputācija ir iztrūkstošo datu aizstāšanas process ar aizvietotām vērtībām. Pašlaik ir izstrādāts liels skaits metožu iztrūkstošo atribūtu vērtību imputācijai. Šo metožu analīze arī ir šīs nodaļas mērķis.

Teorētiskie un praktiskie imputācijas procesa pamati darbībām ar iztrūkstošām atribūtu vērtībām tika noteikti darbos [Rubin D.B., 1976; Little R.J.A., Rubin D.B., 2002; Schafer J.L., 1997]. Tieši šie darbi bija par pamatu efektīgu mūsdienu metožu izstrādei un izmantošanai, lai risinātu iztrūkstošo datu problēmu.

Trīs pamata faktori, kas ietekmē piemērotu datu imputācijas metodes izvēli ir [Dong Y., Peng Ch.-Y.J., 2013] sekojoši:

- iztrūkstošo datu īpatsvars;
- iztrūkstošo datu paraugi;
- iztrūkstošo datu veidošanās modelis.

Iztrūkstošo datu īpatsvars ir tieši saistīts ar šo datu analīzes kvalitāti. Diemžēl, nav pamatots iztrūkstošo datu skaita ierobežojums, pie kura analīze var tikt veikta bez būtiskiem rezultātu izkropļojumiem. Literatūrā tiek pieminēta 5% iztrūkstošo datu robeža, taču tas ir tikai empīrisks novērtējums bez jebkāda teorētiskā pamatojuma.

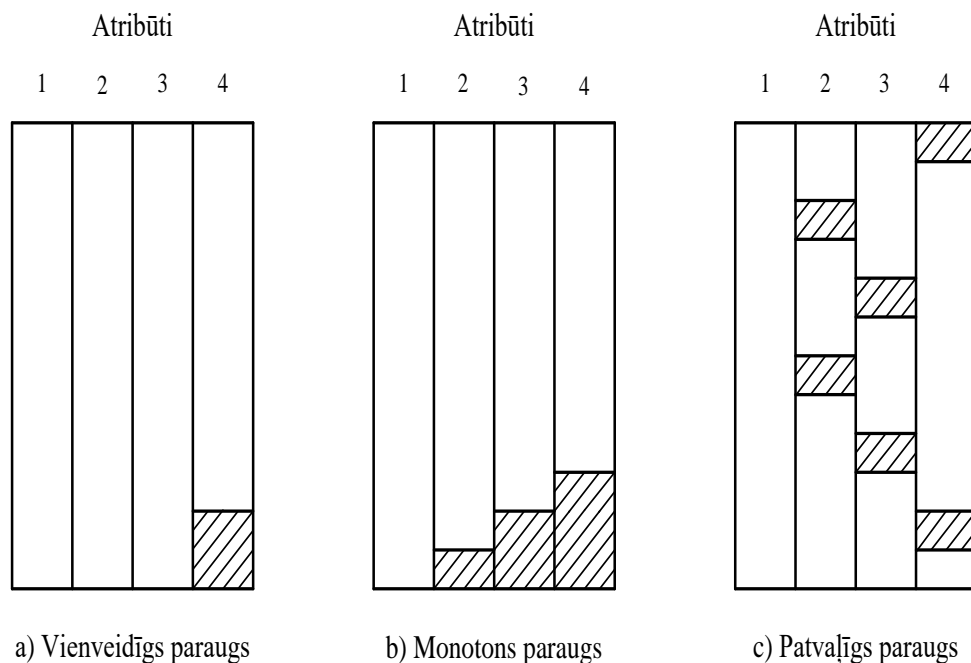
Iztrūkstošo datu paraugi tiek klasificēti šādi: vienveidīgi, monotoni un patvaļīgi.

Vienveidīgā parauga modelis rodas, ja datos trūkst vērtību tikai vienam atribūtam.

Monotonā parauga modeli raksturo fakts, ka, ja attiecīgā veidā tiks sakārtoti atribūti, tad sākot ar atribūtu *j*, visiem nākamajiem atribūtiem būs iztrūkstošas vērtības.

Patvaļīgā parauga modelī dažiem vai visiem atribūtiem iztrūkstošie dati ir patvaļīgā veidā.

Ilustrācijas nolūkā iztrūkstošo datu paraugi ir parādīti attēlā 2.1.1 [Sim J., et al., 2015].



Attēls 2.1.1. Datu paraugu grafiskais attēlojums

Vienveidīgā parauga modelī (skat. attēlu 2.1.1.a) vērtību nav tikai ceturtajā atribūtā. Monotonā parauga modelī (skat. attēlu 2.1.1.b) vērtību nav 2, 3 un 4. atribūtam. Patvaļīgā parauga modelī (skat. attēlu 2.1.1.c) vērtību nav atribūtos 2-4.

Vienveidīgā parauga modelis sastopams retos gadījumos. Par monotonā parauga modeļa piemēru var minēt aptauju rezultātus, kuros respondents neatbildot uz kādu aptaujas jautājumu, neatbild arī uz citiem jautājumiem, kas kaut kādā veidā saistīti ar neatbildēto. Patvaļīgā parauga modelis ir visbiežāk izplatītais datu analīzes uzdevumos.

Kā jau iepriekš tika atzīmēts, pastāv daudz iemeslu, kāpēc iztrūkst datu atribūtu vērtības. Ja neņem vērā konkrēto datu iztrūkuma iemeslu, ir iespējams definēt dažus vispārīgus datu iztrūkuma (atribūtu vērtības) rašanās un apstrādes paņēmienus. Šādi paņēmieni pirmo reizi tika ieviesti darbā [Rubin DB, 1976]. Šie paņēmieni tagad ir vispārpieņemti un atribūtu vērtību imputācija vienmēr balstās uz šiem paņēmieniem.

Definēsim un detalizēti apskatīsim trīs datu iztrūkuma paņēmienus. Visi šie paņēmieni balstās uz atribūtu vērtību sadalījumu. Pieņem, ka dati satur m datu elementus (datu matricas rindas) un n atribūtus (datu matricas kolonas). Tad datu matrica A būs sekojošā veidā (skat. tabulu 2.1.1).

Tabula 2.1.1

Datu matrica A

Datu elementi	Atribūti					
	a_1	a_2	a_3	a_4	a_{1n}
o_1	a_{11}	?	a_{13}	a_{14}	a_{2n}
o_2	a_{21}	a_{22}	?	a_{24}	?
o_3	?	a_{32}	a_{33}	?	a_{3n}
o_4	a_{41}	a_{42}	a_{41}	?	?
o_5	?	?	a_{51}	a_{54}	a_{5n}
....
o_m	a_{m1}	a_{m2}	?	a_{m4}	a_{mn}

Matricas rindas satur datu elementus o_1, o_2, \dots, o_m , kolonnās tiek parādīti atribūti a_1, a_2, \dots, a_n . Vērtība a_{ij} , $i = 1, \dots, m$, $j = 1, \dots, n$ i -tās rindas un j -tās kolonnas krustpunktā parāda atribūta a_j vērtību elementam o_i . Jautājuma zīmes matricas šūnās norāda, ka trūkst atbilstošo atribūtu vērtības.

Vispārīgie pieņēmumi par datiem ir tādi, ka visi datu elementi ir neatkarīgi un tiem ir identisks atribūtu vērtību sadalījums. Šis pieņēmums ir formulēts darbā [Schafer J.L., 1997].

Saskaņā ar šo pieņēmumu atribūtu vērtību daudzdimensiju sadalījuma funkciju sākotnējā datu kopā var attēlot formā

$$P(\mathbf{A} / \theta) = \prod_{j=1}^m f(\mathbf{a}_j / \theta), \quad (2.1.1)$$

kur θ - j -tā atribūta vērtību sadalījuma parametrs.

Parasti parametra θ vērtība nav zināma. Par sadalījuma funkcijas formu tiek izteikti šādi pieņēmumi:

1. Daudzfaktoru normālais sadalījums kvantitatīvo atribūtu vērtību gadījumam.
2. Multinominālais modelis kvalitatīvu atribūtu vērtību gadījumam.
3. Jauktais modelis kvalitatīvu un kvantitatīvu atribūtu vērtību kombināciju gadījumam.

Pieņēmums par kvantitatīvo atribūtu vērtību daudzfaktoru normālu sadalījumu ir diezgan ierobežojošs iztrūkstošo atribūtu vērtību imputācijas kontekstā. Bet literatūrā ir parādīts, ka šī pieņēmuma pārkāpums nopietni neietekmē imputācijas rezultātus.

Lai korekti noteiktu datu iztrūkuma paņēmienus, iepriekš definētajā matricā A , kuras šūnās vērtība 1 nosaka atribūta vērtības esamību, bet 0 – nav atribūta vērtības. Tādējādi iegūstam jaunu matricu B (skat. tabulu 2.1.2).

Tabula 2.1.2

Datu matrica B

Datu elementi	Atribūti					
	a_1	a_2	a_3	a_4	a_n
o_1	1	0	1	1	1
o_2	1	1	0	1	0
o_3	0	1	1	0	1
o_4	1	1	1	0	0
o_5	0	0	1	1	1
....
o_m	1	1	0	1	1

Mēs varam korekti attēlot mūsu sākotnējo datu kopu veidā $\mathbf{A} = (\mathbf{A}_{obs}, \mathbf{A}_{mis})$, kur \mathbf{A}_{obs} ir faktisko (novērojamo) atribūtu vērtību izlase un \mathbf{A}_{mis} ir iztrūkstošo atribūtu vērtību izlase.

Reprezentēsim iztrūkstošo atribūtu vērtību sadalījumu veidā $P(\mathbf{B} / \mathbf{A}, \xi)$, kur ξ - neesamības jeb prombūtnes parametrs. Vispārīgā veidā mēs varam paplašināt šo priekšstatu šādi:

$$P(\mathbf{B} / \mathbf{A}, \xi) = P(\mathbf{B} / \mathbf{A}_{obs}, \mathbf{A}_{mis}, \xi). \quad (2.1.2)$$

Izmantojot iepriekš sniegtos apzīmējumus un vienādojumus, mēs varam definēt šādus iztrūkstošo datu priekšstatus:

1. *Nejaušs iztrūkums (Missing at random (MAR))*. Šajā gadījumā iztrūkstošo vērtību varbūtību sadalījums ir definēts kā

$$P(\mathbf{B} / \mathbf{A}_{obs}, \mathbf{A}_{mis}, \xi) = P(\mathbf{B} / \mathbf{A}_{obs}, \xi). \quad (2.1.3)$$

Vienādojums (2.1.3) var tikt interpretēts sekojošā veidā: iztrūkstošo vērtību varbūtība ir atkarīga tikai no novērojamo atribūtu vērtību izlases \mathbf{A}_{obs} un prombūtnes parametra ξ .

Nejauša iztrūkuma gadījumā var apgalvot, ka attiecībā uz konkrētām novērotajām atribūtu vērtībām no izlases \mathbf{A}_{obs} atlikušo vērtību sadalījums starp esošajām vērtībām un iztrūkstošajām vērtībām ir vienāds.

Demonstrēsim atribūtu nejauša iztrūkuma piemēru. Pieņemsim, ka cilvēkiem, kuri vēlas mainīt profesiju, tiek organizēti karjeras pārorientēšanās kursi. Kurša sākumā dalībniekiem tika lūgts atbildēt uz testa jautājumiem, lai novērtētu viņu zināšanu un prasmju līmeni. Pieņemsim, ka daži no dalībniekiem atbildēja tikai uz daļu no jautājumiem, kas liecināja par viņu zināšanu un pieredzes zemo līmeni.

Mācību procesā daļa no dalībniekiem ar novērtētu zemu līmeni pameta kursus. Pēc kursu pabeigšanas dalībniekiem tika lūgts atbildēt uz cita testa jautājumiem, lai novērtētu viņu profesionālās pārkvalifikācijas rezultātus. Acīmredzami, ka jauniegūtajos rezultātos nebūs kursus pametušo darbinieku datu.

2. *Pilnīgi nejaušs iztrūkums (Missing completely at random (MCAR))*. Šis iztrūkuma mehānisms formāli tiek izteikts kā

$$P(\mathbf{B} / \mathbf{A}_{obs}, \mathbf{A}_{mis}, \xi) = P(\mathbf{B}, \xi). \quad (2.1.4)$$

Šis vienādojums izsaka to, ka iztrūkstošo vērtību sadalījums ir tāds pats kā pašreizējo atribūtu vērtību sadalījums. Tāpēc pieejamās vērtības varam uzskatīt par paraugu no vērtību izlases, kurā ir visas vērtības, ieskaitot tās, kas atrodas iztrūkstošo vērtību vietā. Taču iztrūkstošās vērtības nevar uzskatīt par paraugu no pilnas vērtību izlases.

Sniegsim piemēru pilnīgi nejauši iztrūkstošām atribūtu vērtībām. Organizācija nolēma veikt savu darbinieku medicīnisko pārbaudi, lai novērtētu viņu veselību. Lai to izdarītu, darbiniekiem ir jāziedo asinis analīzei, jāveic fluorogrāfija un elektrokardiogramma. Pirms asins analīzes veikšanas laborants nejauši salauza vairākas asins ampulas. Veicot elektrokardiogrammu tika bojāts aprīkojums un par dažiem darbiniekiem rezultāti netika iegūti. Par fluorogrāfijas veikšanu atbildīgais speciālists nejauši zaudēja dažus rezultātus.

3. *Nav nejaušs iztrūkums (Not missing at random (NMAR))*. Formāli tas nozīmē, ka iztrūkstošo vērtību sadalījums ir atkarīgs gan no novērotajām vērtībām, gan no iztrūkstošajām vērtībām

$$P(\mathbf{B} / \mathbf{A}_{obs}, \mathbf{A}_{mis}, \xi). \quad (2.1.5)$$

Šis ir visnelabvēlīgākais iztrūkstošo datu gadījums, jo šāda veida iztrūkstošo atribūtu vērtību imputācija nav piemērota plaši izmantotām metodēm.

Spilgts piemērs šim gadījumam var būt aptaujas rezultāti, kuras viens no jautājumiem ir personu ienākumu līmenis. Prakse rāda, ka personas ar zemu un vidēju ienākumu līmeni faktiski norāda savus reālos ienākumus. Savukārt privātpersonas ar augstu ienākumu līmeni labprātāk nenorāda savus reālos ienākumus.

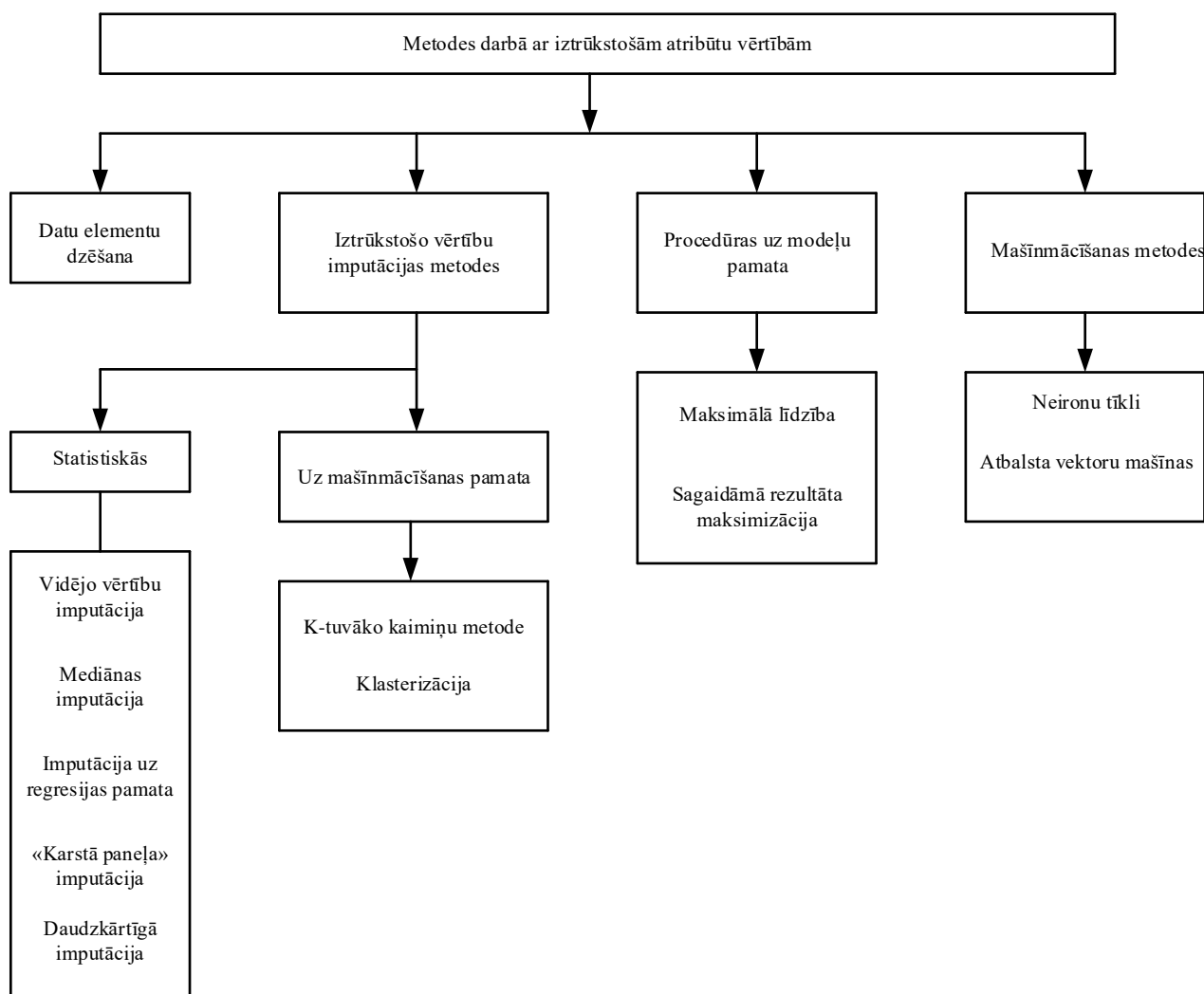
Ieviesīsim šādu svarīgu jēdzienu. Pieņem, ka MAR datu iztrūkuma mehānismā vērtību sadalījuma parametrs (θ) un iztrūkuma parametrs (ξ) ir neatkarīgi. Tad tiek uzskatīts, ka iztrūkstošo datu mehānisms ir ignorējams [Little R., J., A., Rubin D. B., 2002]. Tas pats attiecas uz NMAR mehānismu, jo NMAR ir īpašs MAR gadījums.

Lielākā daļa pašreizējo jaudīgo metožu iztrūkstošo atribūtu vērtību piedēvēšanai pieņem, ka tām trūkst MAR vai NMAR iespēju. Tomēr ļoti bieži reālās datu analīzes problēmās ir lielākas vai mazākas novirzes no šīm metodēm. Praktiskie pētījumi ir parādījuši, ka šo iespēju neizmantošana būtiski neietekmē imputācijas rezultātus. Tāpēc iztrūkstošo atribūtu vērtību imputācija dažreiz tiek veikta, nepārbaudot faktiskos iztrūkstošo datu iespēju metodes.

Šīs sadaļas noslēgumā mēs piedāvājam pamata metodes, kā rīkoties ar iztrūkstošām atribūtu vērtībām. Ir piedāvātas dažādas šo metožu klasifikācijas iespējas. Attēlā 2.1.2 sniegta klasifikācijas variants, ko var uzskatīt par dažādos literatūras avotos piedāvāto variantu vispārīnājumu.

Īsi komentēsim šīs klasifikācijas rezultātus.

Datu elementu dzēšana ietver visu elementu izdzēšanu no sākotnējiem datiem, kuriem trūkst vismaz vienas atribūta vērtības. Lai gan šī metode nav saistīta ar atribūtu vērtību imputāciju, ir lietderīgi to izmantot, ja ir liels datu apjoms un ļoti maz iztrūkstošo atribūtu vērtību (apmēram 1% - 3%).



Attēls 2.1.2. Metožu klasifikācijas grafiskais attēlojums iztrūkstošo atribūtu vērtību apstrādē

Kopumā iztrūkstošo atribūtu vērtību imputāciju var iedalīt divās lielās klasēs:

- *viena imputācija*, kad tiek aprēķinātas atsevišķas iztrūkstošās vērtības;
- *daudzkārtīgā imputācija*, kad vienlaicīgi tiek imputētas iztrūkstošo vērtību grupas.

Imputācijas metodes var iedalīt divās grupās: *statistikas metodes* un *uz mašīnmācīšanos balstītas metodes*. Pirmajā metožu grupā ietilpst vidējo vērtību imputācija, uz regresiju balstīta imputācija, mediānas imputācija, “karstā paneļa” (*Hot Deck*) imputācija un daudzkārtīgā imputācija. Visas šīs grupas metodes veic vienu imputāciju, izņemot vairākkārtēju imputāciju, kā to arī norāda pats metodes nosaukums. Šī metode dažkārt tiek klasificēta kā uz modeli balstīta procedūra, jo metode izmanto noteiktu datu modeli.

Otrā metožu grupa veic vienu imputāciju, pamatojoties vai nu uz k-tuvākā (*k-nearest*) kaimiņa metodi vai pamatojoties uz datu elementu klasterizāciju.

Uz modeļiem balstītas metodes ir vairākas imputācijas metodes un tās ietver atbilstošu skaitļošanas procedūru veikšanu. Paredzamās vērtības palielināšanas metode dod labus rezultātus, taču tai ir nepieciešami specializēti skaitļošanas rīki.

Pēdējā grupa sastāv no metodēm, kuru pamatā ir mašīnmācīšanās pieejas. Viena no metodēm izmanto mākslīgos neironu tīklus (ANN), otra metode izmanto atbalsta vektoru mašīnas (SVM) pieeju. Šīs ir visjaudīgākās mūsdienu metodes iztrūkstošo atribūtu vērtību daudzkārtējai imputācijai, taču to pielietošanai ir nepieciešami sarežģīti skaitļošanas līdzekļi.

2.2. Vienkāršākās imputācijas metodes

1. Imputācija uz izvietojuma parametru pamata

Šīs metodes ideja ir ļoti vienkārša. Ja katra atribūtu vērtību izlase tiek uzskatīta par šo vērtību sadalījumu, tad iztrūkstošo atribūtu vērtību vietā tiek izmantota sadalījuma izvietojuma parametra vērtība.

Izvietojuma parametra veids ir atkarīgs no skalas veida, kurā tiek mērītas atribūtu vērtības (sk. pielikumu P1).

Ja atribūtu vērtības mēra nominālā skalā, tad imputācijai var izmantot tikai atribūtu vērtību sadalījuma veidu.

Pieņemam, ka atribūtu vērtības a_j izmērītas nominālā skalā

$$a_j^T = (1, 1, 2, 3, ?, 1, 3, 1, 2, 2, 3, 1, 3, ?, 1, 2).$$

(Mēs izmantojam atribūtu vērtību transponēto attēlojumu, jo datu tabulā atribūtu vērtības tiek attēlotas šīs tabulas kolonās).

Šeit skaitļi 1, 2, 3 apzīmē atbilstošo nominālo atribūtu kategoriju kodus.

Iepriekš norādītajā atribūtu vērtību kopā a_j jautājuma zīmes “?” norāda uz iztrūkstošajām šī atribūta vērtībām.

Visbiežāk sastopamā atribūta vērtība tiek ņemta par sadalījuma modu. Šajā piemērā visizplatītākā kategorija ir nominālā kategorija 1, tātad tā ir sadalījuma moda un šī vērtība tiek ielikta iztrūkstošo vērtību vietā (jautājuma zīmju vietā). Rezultātā mums ir šāds pilns atribūtu vērtību kopums

$$a_j^T = (1, 1, 2, 3, 1, 3, 1, 2, 2, 3, 1, 3, 1, 1, 2).$$

Ja atribūtu vērtības mēra ordinālajā skalā, tad par šo vērtību izlases (sadalījuma) parametru var izmantot modu vai mediānu.

Uzdodam atribūtu a_j vērtību kopu, kas mērītas ordinālajā skalā

$$a_j^T = \dots \left(\begin{array}{cccc} 1 & \dots & 2 & \dots & 3 & \dots & 4 \\ 21(2) & \dots & 14(1) & \dots & 16(1) & \dots & 18(2) \end{array} \right).$$

Šeit skaitļi 1, 2, 3, 4 augšējā rindā apzīmē atbilstošo atribūta kategoriju kodus. Apakšējā rindā esošie skaitļi norāda katras kategorijas piemēru skaitu. Cipari iekavās apakšējā rindā norāda iztrūkstošo atribūta vērtību skaitu katrā kategorijā.

Dotās izlases (sadalījuma) moda ir ordinālā kategorija 1, kurai pieder lielākais piemēru skaits. Tāpēc šai kategorijai jāpievieno visi piemēri ar iztrūkstošām atribūtu vērtībām. Rezultātā mums ir šāds pilns atribūtu vērtību kopums:

$$a_j^T = \begin{pmatrix} 1 \dots\dots\dots 2 \dots\dots\dots 3 \dots\dots\dots 4 \\ 27 \dots\dots\dots 14 \dots\dots\dots 16 \dots\dots\dots 18 \end{pmatrix}.$$

Noteiksim šīs atribūtu vērtību izlases (sadalījuma) mediānu. Tā kā piemērā ordinālo kategoriju skaits ir $n=4$, sadalījuma mediāna būs kategorija, kas atrodas 2. pozīcijā, t.i., 2. kategorija. Rezultātā mums ir šāda pilna atribūtu vērtību izlase

$$a_j^T = \begin{pmatrix} 1 \dots\dots\dots 2 \dots\dots\dots 3 \dots\dots\dots 4 \\ 21 \dots\dots\dots 20 \dots\dots\dots 16 \dots\dots\dots 18 \end{pmatrix}.$$

Ja atribūtu vērtības tiek mērītas intervāla vai attiecību skalā, izvietojuma parametra varianti var būt moda, mediāna vai vidējā vērtība.

Uzdodam atribūtu vērtību kopu a_j , kas izmērītas skaitliskajā skalā (intervāla skalā vai attiecību skalā).

$$a_j^T = (10, 2, 4, ?, 6, 8, ?, 9, 4, 7, 3, 11).$$

Šīs vērtību izlases (sadalījuma) moda ir 4. Aizstājot iztrūkstošās vērtības ar modas vērtību, mums ir šāds pilns atribūtu vērtību kopums

$$a_j^T = (10, 2, 4, 4, 6, 8, 4, 9, 4, 7, 3, 11).$$

Lai noteiktu mediānu, esošās atribūtu vērtības sakārtojam augošā secībā

$$a_j^T = (2, 3, 4, 4, 6, 7, 8, 9, 10, 11).$$

Šīs virknes mediāna ir skaitlis pozīcijā $n/2 = 10/2 = 5$, t.i., skaitlis 6. Aizstājot iztrūkstošās vērtības ar mediānas vērtību, mēs iegūstam šādu pilnu vērtību kopu

$$a_j^T = (10, 2, 4, 6, 6, 8, 6, 9, 4, 7, 3, 11).$$

Aprēķinām atribūta a_j vidējo vērtību

$$\bar{a}_j = \frac{10 + 2 + 4 + 6 + 8 + 9 + 4 + 7 + 3 + 11}{10} = \frac{64}{10} = 6.4.$$

Tā kā atribūtu a_j vērtības ir dotas veselos skaitļos, noapaļosim rezultātu: $\bar{a}_j \approx 6$. Aizstājot iztrūkstošās vērtības ar vidējo vērtību, mums ir šāds pilns atribūtu vērtību kopums

$$a_j^T = (10, 2, 4, 6, 6, 8, 6, 9, 4, 7, 3, 11).$$

Aplūkotās iztrūkstošo atribūtu vērtību imputācijas metodes vienīgā priekšrocība ir tās vienkāršība un pilnīga atribūtu vērtību kopuma saglabāšana.

Metodei nav loģiska pamatojuma, aprēķinātās vērtības ir diezgan patvaļīgas. Aprēķinu kvalitāte lielā mērā ir atkarīga no tā, kuras faktiskās vērtības iztrūkst. Ja šīs vērtības ir tuvu izvietojuma parametra vērtībai, tad imputācijas rezultāti var būt veiksmīgi. Ja iztrūkstošās faktiskās vērtības ir tuvu atribūta maksimālajām vai minimālajām vērtībām, imputācijas rezultāti var būt ļoti slikti.

Vispārīgi sakot, imputācijas rezultātu pēc izvietojuma parametra vērtības nosaka atribūta sākotnējo vērtību izkliedes pakāpe un šo vērtību skaits.

Iepriekš aprakstītā metode iztrūkstošo atribūtu vērtību imputēšanai bija viena no pirmajām šāda veida metodēm. Pašlaik šī metode tiek izmantota ārkārtīgi reti.

Ja sākotnējie dati ir apmācības izlase, tad katram datu piemēram tiek piešķirta atbilstošās klases iezīme. Tad iepriekš minēto metodi var piemērot katrai piemēru klasei atsevišķi.

2. Imputācija uz "karstā paneļa metodes" pamata

Metodes pamati tika likti ASV pagājušā gadsimta 40. un 50. gados. Statistiskās izpētes nolūkos iedzīvotājiem bieži tika lūgts atbildēt uz dažāda veida testiem vai atbildēt uz aptaujās dotajiem jautājumiem. Dažādu iemeslu dēļ cilvēki bieži neatbildēja uz dažiem testa jautājumiem vai neaizpildīja dažus obligātos jautājumus aptaujās. Lai risinātu šāda veida iztrūkstošo datu problēmu, tika izmantota šāda pieeja. Piemēriem testā vai pārskatā ar iztrūkstošām vērtībām tika atlasītas citas piemēru versijas, kurās atlikušie dati atbilst piemēru datiem ar iztrūkstošo vērtību. Tad iztrūkstošā vērtība tika uzskatīta par atbilstošu vērtību no cita piemēra. Šo metodi var attiecināt arī uz gadījumiem, kad iztrūkst vairākas vērtības.

Šīs metodes acīmredzamais trūkums ir tāds, ka iztrūkstošās vērtības imputācijas pamatā ir viens līdzīgs atribūtu vērtību vektors. Visas pārējās datu izlases vispārīgās īpašības tiek ignorētas.

Cits šīs metodes variants ir "aukstā paneļa metode" (*cold deck*), kad atribūtu vērtību vektors, kas ir līdzīgs vektoram ar iztrūkstošo vērtību, tiek meklēts citā datu izlasē, nevis izmantotajā datu izlasē.

"Karstā paneļa metode" praktiski netiek izmantota mūsdienu pieejās datu pirmapstrādei. Tomēr šī metode literatūrā bieži tiek minēta kā pirmā šāda veida metode.

3. Imputācija uz regresijas atkarības pamata

Šo metodi var izmantot kvantitatīvām atribūtu vērtībām. Teorētiski to var attiecināt arī uz kvalitatīvām atribūtu vērtībām. Taču rodas šādas problēmas: (1) loģistiskās regresijas modelim (sk. pielikumu P1) ir nepieciešamas sarežģītas skaitļošanas procedūras; (2) aprēķinu rezultāti tiek interpretēti kā iztrūkstošās atribūta vērtības, kas pieder pie dažādām kvalitatīvām kategorijām, varbūtības.

Aplūkojamās imputācijas metodes būtība ir šāda. Ja atribūtu vērtību izlases a_j , a_l korelē savā starpā, tad var izveidot regresijas attiecības starp šīm vērtību kopām. Pamatojoties uz iegūto atkarību, ir iespējams noteikt otrā atribūta imputētās vērtības.

Apskatām vienkāršu ilustratīvu piemēru. Dotas divas atribūtu vērtību izlases

$$a_j^T = (5, \dots, 4, \dots, 6, \dots, 7, \dots, 6, \dots, 8, \dots, 9, \dots, 3, \dots, 7, \dots, 5, \dots, 7)$$
$$a_l^T = (7, \dots, 6, \dots, 7, \dots, ?, \dots, 13, \dots, 14, \dots, 18, \dots, 5, \dots, 12, \dots, 6, \dots, 12)$$

Izmantojot pielikumā P1.3 norādītās skaitļošanas procedūras vai izmantojot piemērotus skaitļošanas rīkus, tiek noteikts regresijas vienādojums, kas apvieno atribūtu vērtību izlases a_j , a_l .

$$a_l = -6.8 + 2.8a_j.$$

Tā kā atribūta a_l vērtība iztrūkst pie $a_j = 7$, aprēķināsim imputēto vērtību

$$a_j = -6.8 + 2.8 * 7 = 12.8 \approx 13.$$

(Mēs noapaļojām aprēķināto atribūta vērtību a_j , jo visas sākotnējās atribūta vērtības ir vesela skaitļa formā).

Rezultātā mums ir šādas pilnīgas atribūtu vērtību izlases

$$a_j^T = (5, \dots, 4, \dots, 6, \dots, 7, \dots, 6, \dots, 8, \dots, 9, \dots, 3, \dots, 7, \dots, 5, \dots, 7)$$
$$a_l^T = (7, \dots, 6, \dots, 7, \dots, 13, \dots, 13, \dots, 14, \dots, 18, \dots, 5, \dots, 12, \dots, 6, \dots, 12)$$

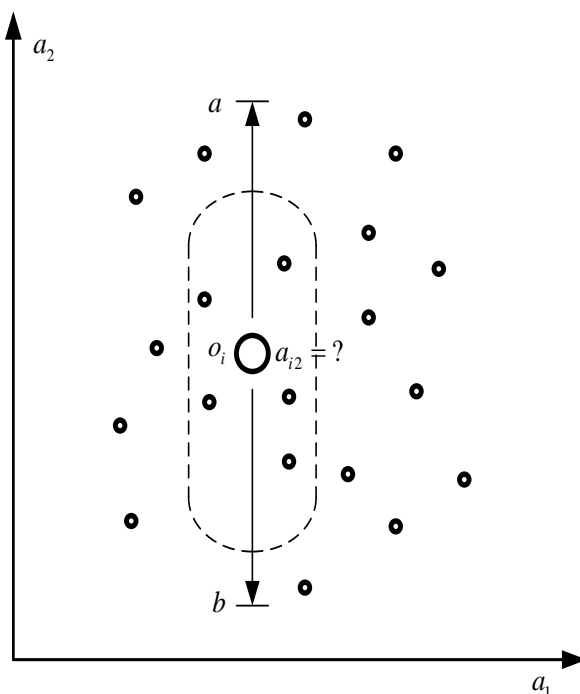
Lietojot apskatīto metodi praksē, vispirms jānovērtē atribūtu vērtību izlases korelācijas (ar iztrūkstošajām vērtībām) ar visām pārējām atribūtu vērtību kopām. Par pieņemamu datu izlasi tiek uzskatīta datu izlase, kas visvairāk korelē ar mūs interesējošo datu izlasi.

Lai gan šī imputācijas metode šķiet loģiska, tās būtiskais trūkums ir netiešais pieņēmums par lineāru saistību starp atribūtu vērtību kopām. Tā kā šis pieņēmums ne vienmēr ir spēkā, imputācija var neizdoties.

2.3. *K-tuvāko kaimiņu metode*

K-tuvāko kaimiņu metode (k-nearest neighbour method) tiek plaši izmantota iztrūkstošo atribūtu vērtību imputēšanai tās vienkāršības, lietderīguma un diezgan augstās rezultātu ticamības dēļ. Metode nav parametriska, jo tā izmanto tikai sākotnējos datus un vairāk neko citu.

Metodes pielietojums būtiski atšķiras kategoriskiem un skaitliskiem datiem. Sāksim ar to, kā izmantot šo metodi iztrūkstošo skaitlisko atribūtu vērtību imputācijai. Lai saprastu metodes būtību, pievērsīsimies attēlam 2.3.1.



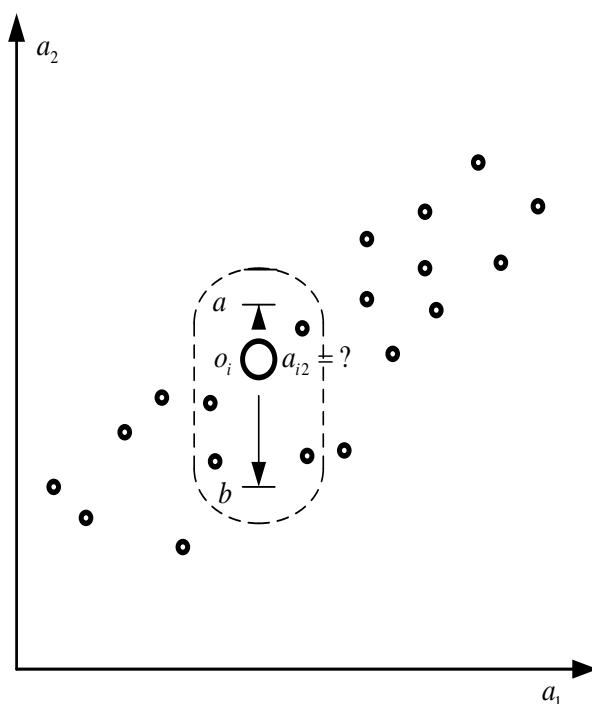
Attēls 2.3.1. Datu izlases shematisks attēlojums divu atribūtu a_1 , a_2 telpā (piemērā o_i nav atribūta vērtības a_{i2})

Kāda varētu būt iztrūkstošā atribūta a_2 vērtība piemērā o_i, a_{i2} . Atsaucoties uz attēlu 2.3.1, ir viegli pārbaudīt, vai tas var būt robežās $[a, b]$. Lai noteiktu šīs iztrūkstošās vērtības meklēšanas virzienu, par pamatu ņemam piemēra tuvākos kaimiņus o_i . Šeit tuvums tiek novērtēts pēc šo tuvāko kaimiņu atribūtu vērtībām a_1 . Pieņemsim, ka ir k tuvākie kaimiņi, tāpēc arī ir tāds metodes nosaukums. Mēs varam aprēķināt vidējo atribūta vērtību a_2 šiem kaimiņiem un ņemt to par iztrūkstošo vērtību a_{i2} .

Sprieduma loģika šeit ir ļoti vienkārša. Ja piemērs ar iztrūkstošu atribūta vērtību a_2 ir tuvs kaimiņiem pēc atribūta vērtības a_1 , tiek pieņemts, ka tā iztrūkstošā vērtība ir tuva tā tuvāko kaimiņu vidējai vērtībai a_2 .

Iztrūkstošās atribūta vērtības imputācijas rezultātu, izmantojot k-tuvāko kaimiņu metodi, var būtiski ietekmēt atribūtu vērtību a_1 un a_2 korelācija. No attēla 2.3.1 ir skaidrs, ka atribūtu vērtības šajā gadījumā nav savstarpēji saistītas (nekorelē savā starpā).

Pievērsīsimies attēlam 2.3.2.



Attēls 2.3.2. Piemēra o_i shematiskais attēlojums ar iztrūkstošo atribūta a_{i2} vērtību datu kopā ar korelētu atribūtu a_1, a_2 vērtībām

No šī attēla var izdarīt nenoliedzamu secinājumu, ka atribūtu vērtības a_1, a_2 ir pozitīvi korelētas. Tas sašaurina piemēra o_i tuvākos kaimiņu loku. Tāpēc izmantojot piemēra o_i tuvāko kaimiņu izlases vidējo atribūta vērtību a_2 kā imputācijas vērtību priekš a_{i2} , tiks iegūts ticamāks imputācijas rezultāts.

Kādas galvenās problēmas ir jāatrisina tuvāko k-kaimiņu metodes praktiskā pielietojumā? Šie uzdevumi ir:

1. Vērtības k izvēle.
2. Sākotnējo atribūtu vērtību normalizēšana.
3. Metrikas izvēle attālumu noteikšanai no piemēra ar iztrūkstošu atribūta vērtību līdz tuvākajiem kaimiņiem.

Tuvāko kaimiņu skaita izvēle ir neformalizējama problēma. Šī izvēle ir atkarīga no risināmās problēmas rakstura un ievades datu apjoma. Ļoti liela ievades datu skaita gadījumā to samazināšanai var izmantot dažādas pieejas: ierobežotu datu apgabalu definēšana, piemēru ignorēšana, kas potenciāli var nebūt potenciālie kaimiņi un tamlīdzīgi.

Lai noteiktu atbilstošo tuvāko kaimiņu skaitu, daži literatūras avoti iesaka veikt nepieciešamo imputāciju dažādām k vērtībām un par pamatu ņemt rezultātu, kas dod mazāko novērtējuma kļūdu. Tomēr šī pieeja prasa lielu skaitu aprēķinu un tāpēc to var ieteikt tikai gadījumos, kad iztrūkstošo vērtību skaits ir mazs, bet ir nepieciešama augsta imputācijas ticamība.

Atribūtu skaitlisko vērtību normalizācijas mērķis ir samazināt dažādu atribūtu vērtību amplitūdu ietekmi un nodrošināt attāluma aprēķinu pareizību, jo visu atribūtu vērtības pēc normalizācijas tiks uzrādītas vienotā veidā.

Atribūtu vērtību normalizēšanas metodes tiks detalizēti aplūkotas 5. nodaļā. Šajā nodaļā kā normalizējošos rādītājus izmantosim vidējās atribūtu vērtības.

Eiklīda attālumu (*Euclidean distance*) parasti izmanto, lai novērtētu attālumus starp piemēru ar iztrūkstošo atribūta vērtību un citiem piemēriem:

$$d(o_i, o_k) = \sqrt{\sum_j (a_{ij} - a_{kj})^2}, \quad (2.3.1)$$

kur $j = 1, 2, \dots$ - atribūtu skaits, pēc kura tiek novērtēts attālums;

a_{ij} - j -tā atribūta vērtība piemērā ar iztrūkstošu vērtību (atribūts ar iztrūkstošu vērtību nav iekļauts atribūtu sarakstā attiecīgo attālumu noteikšanai);

- a_{kj} - j -tā atribūta vērtība k -tajam piemēram.

Manhetenas attālumu (*Manhattan distance*) var izmantot kā alternatīvu attāluma mēru

$$d(o_i - o_k) = \sum_j |a_{ij} - a_{kj}|. \quad (2.3.2)$$

Tomēr šis attāluma novērtējums tiek izmantots daudz retāk.

Apskatīsim vienkāršu ilustratīvu piemēru. Tabulā 2.3.1 ir parādīta sākotnējā datu izlase.

Tabula 2.3.1

Sākotnējā datu izlase

Piemēri (objekti)	Atribūti		
	a_1	a_2	a_3
o_1	7	9	12
o_2	6	8	11
o_3	5	8	14
o_4	8	10	13
o_5	6	8	15
o_6	? (5)	7	14
o_7	7	10	? (12)
o_8	5	8	11
o_9	8	10	12
o_{10}	6	7	14
Kolonnu summas	58	85	116
Vidējā vērtība	6.444	8500	12.889

Šajos datos piemēram o_6 iztrūkst atribūta vērtības a_1 un piemēram o_7 iztrūkst atribūta a_3 vērtības. Iztrūkstošās vērtības tiek apzīmētas ar jautājuma zīmi “?”. Šo atribūtu faktiskās vērtības ir parādītas blakus iekavās. Tas tika darīts, lai salīdzinātu un analizētu imputācijas rezultātus. (Ir jāsaprot, ka šīs vērtības faktiski mums nav zināmas.)

Mēs normalizējam sākotnējās atribūtu vērtības, izmantojot vidējās atribūtu vērtības, kas norādītas tabulas 2.3.1 pēdējā rindā. Atribūtu normalizētās vērtības ir parādītas tabulā 2.3.2.

Aprēķinām Eiklīda attālumus no piemēra o_6 ar iztrūkstošu atribūta vērtību a_{61}^n līdz citiem piemēriem, izmantojot atribūtu vērtības a_2^n, a_3^n .

$$d(o_6 - o_1) = \sqrt{(a_{62}^n - a_{12}^n)^2 + (a_{63}^n - a_{13}^n)^2} = \sqrt{(0.082 - 0.106)^2 + (0.121 - 0.103)^2} =$$

$$= \sqrt{0.00058 + 0.00032} = \sqrt{0.00090} = 0.0300.$$

Tabula 2.3.2

Normalizētās atribūtu vērtības datiem no tabulas 2.3.1

Piemēri (objekti)	Normalizētās atribūtu vērtības		
	a_1^n	a_2^n	a_3^n
o_1	1.086	1.059	0.931
o_2	0.931	0.941	0.853
o_3	0.776	0.941	1.086
o_4	1.241	1.176	1.009
o_5	0.776	0.941	1.164
o_6	?	0.824	1.086
o_7	1.086	1.176	?
o_8	0.776	0.941	0.853
o_9	1.241	1.176	0.931
o_{10}	0.931	0.824	1.086

Pārējie aprēķini tiek veikti pēc analogijas. Aprēķinu rezultāti ir parādīti tabulā 2.3.3.

Tabula 2.3.3

Aprēķinātās attālumu vērtības no piemēra o_6 līdz citiem piemēriem pamatojoties uz atribūtu vērtībām o_2^n, o_3^n

	$o_6 - o_1$	$o_6 - o_2$	$o_6 - o_3$	$o_6 - o_4$	$o_6 - o_5$	$o_6 - o_7$	$o_6 - o_8$	$o_6 - o_9$	$o_6 - o_{10}$
$d(\dots)$	0.2814	0.2608	0.1170	0.3603	0.1407	-	0.2608	0.3544	0

Mēs neaprēķinām attālumu, $d(o_6 - o_7)$, jo piemērā o_7 nav atribūta vērtības a_3 .

Sakārtojam piemērus augošā secībā pēc to attāluma līdz piemēram o_6 .

$$o_{10}, o_3, o_5, o_2, o_8, o_1, o_9, o_4.$$

Pieņemam $k = 3$.

Šajā gadījumā imputētā atribūta vērtība $a_{61(i)}$ ir vienāda ar atribūtu $o_{10,1}, o_{31}, o_{51}$ vidējo vērtību.

$$a_{61(i)} = \frac{a_{10,1} + a_{31} + a_{51}}{3} = \frac{6 + 5 + 6}{3} = 5.67 \approx 6.$$

Pieņemam $k = 4$.

$$a_{61(i)} = \frac{a_{10,1} + a_{31} + a_{51} + a_{21}}{4} = \frac{6 + 5 + 6 + 6}{4} = 5.75 \approx 6.$$

Abos gadījumos mēs ieguvām labus imputācijas rezultātus. Aprēķinātā vērtība 6 būtiski neatšķiras no faktiskās vērtības 5. Tā kā reālā vērtība a_{61} nav zināma, tās aprēķinātā vērtība sniedz labu reālās vērtības tuvinājumu.

Veicam aprēķinus, lai noteiktu atribūta imputēto vērtību a_{73} . Aprēķinām Eiklīda attālumus no piemēra o_7 līdz citiem piemēriem, izmantojot atribūtu a_1^n un a_2^n vērtības. Aprēķinu rezultāti ir parādīti tabulā 2.3.4.

Tabula 2.3.4

Aprēķinātās attālumu vērtības no piemēra o_7 līdz citiem piemēriem pamatojoties uz atribūtu a_1^n, a_2^n vērtībām

	$o_7 - o_1$	$o_7 - o_2$	$o_7 - o_3$	$o_7 - o_4$	$o_7 - o_5$	$o_7 - o_6$	$o_7 - o_8$	$o_7 - o_9$	$o_7 - o_{10}$
$d(\dots)$	0.1170	0.2793	0.3895	0.1550	0.3390	-	0.3895	0.1550	0.3846

Attālums $d(o_7 - o_6)$ netika aprēķināts, jo piemērā o_6 trūkst atribūta vērtības a_1 .

Sakārtojam piemērus augošā secībā pēc to attāluma no piemēra o_7 .

$$o_1, o_4, o_9, o_2, o_5, o_{10}, o_3, o_8.$$

Pieņemam $k = 3$:

$$a_{73(i)} = \frac{a_{13} + a_{43} + a_{93}}{3} = \frac{12 + 13 + 12}{3} = 12.33 \approx 12.$$

Pieņemam $k = 4$:

$$a_{73(i)} = \frac{a_{13} + a_{43} + a_{93} + a_{23}}{4} = \frac{12 + 13 + 12 + 11}{3} = 12.$$

Abām k vērtībām tika iegūti ļoti labi imputācijas rezultāti.

Iepriekš minētajā piemērā visiem attālumiem no piemēriem ar trūkstošām atribūtu vērtībām tika piešķirta vienāda nozīme. Šajā sakarā rodas šāda ideja: attālumiem piešķirt dažādus svarus pēc principa: jo mazāks attālums, jo lielāks tā svars.

Attiecīgo attālumu svēršanai ir ierosināts ļoti liels skaits dažādu pieeju. Šajā nodaļā mēs aplūkosim vienkāršu un efektīvu pieeju sekojošā veidā - svara vērtība noteiktam attālumam tiek definēta kā šī attāluma kvadrāta apgrieztā vērtība:

$$w(\dots) = \frac{1}{(d(\dots))^2}. \quad (2.3.3)$$

Tad svērtā attāluma vērtību aprēķina kā

$$d^w(\dots) = w(\dots)d(\dots). \quad (2.3.4)$$

Aprēķinām svaru un svērto attālumu vērtības aprēķinātajām vērtībām tabulai 2.3.3. Aprēķinu rezultāti ir parādīti tabulā 2.3.5.

Tā kā mums ir darīšana ar svērtajiem attālumiem, lielāka svērtā attāluma vērtība atbilst lielākai tuvības pakāpei starp piemēriem.

Sarindojam piemērus svērtā attāluma palielināšanās attālumā no piemēra o_6 .

$$o_{10}, o_3, o_5, o_2, o_8, o_1, o_9, o_4.$$

Aprēķinātās svaru vērtības un svērtie attālumi no piemēra o_6 līdz citiem piemēriem saskaņā ar tabulu 2.3.3

	$o_6 - o_1$	$o_6 - o_2$	$o_6 - o_3$	$o_6 - o_4$	$o_6 - o_5$	$o_6 - o_7$	$o_6 - o_8$	$o_6 - o_9$	$o_6 - o_{10}$
$d(\cdot - \cdot)$	0.2814	0.2608	0.1170	0.3603	0.1407	-	0.2608	0.3544	0
$(d(\cdot - \cdot))^2$	0.0792	0.0680	0.0137	0.1298	0.0198	-	0.0680	0.1256	0
$w(\cdot - \cdot)$	12.6263	14.7016	73.0460	7.7023	50.5050	-	14.7016	7.9618	0
$d^w(\cdot - \cdot)$	3.5530	3.8342	8.5464	2.2775	7/1060	-	3.8342	2.8216	0

Mēs ieguvām tādu pašu piemēru secību, ja izmantotu nesvērtos attālumus, tāpēc tuvāko k kaimiņu izvēle un aprēķinātās imputētās atribūtu vērtības $o_{6l(i)}$ būs vienādas.

Svērto attālumu aprēķins tabulas 2.3.4 datiem tiek veikts saskaņā ar to pašu shēmu. Kopumā var būt nelielas atšķirības starp piemēru apstrādi pēc tiešā attāluma vērtībām un svērtajām attāluma vērtībām. Bet šīs iespējamās atšķirības būtiski neietekmē k tuvāko kaimiņu izvēli.

Iztrūkstošo atribūtu vērtību imputācijas vispārīgie principi pie dažādām k vērtībām, pamatojoties uz k tuvāko kaimiņu metodi, ir šādi:

1. Mazas k vērtības palielina “trokšņa” ietekmi datos un rada zemu rezultātu vispārināmību.
2. Lielākas k vērtības “nogludina” lokālos efektus un rada lielāku rezultātu vispārināmību.

Jāņem vērā: ja vienam un tam pašam atribūtam iztrūkst vairākas vērtības, attiecīgie aprēķini jāveic katrai iztrūkstošajai vērtībai atsevišķi.

Lieliem ievades datu apjomiem visi attiecīgie aprēķini tiek veikti, izmantojot piemērotus programmatūras rīkus. Lai formalizētu skaitļošanas procesus, bieži izmanto fiktīvu mainīgo (*dummy variable*). Šī mainīgā lieluma vērtības ir parādītas tāda paša izmēra tabulā (matricā) kā ieejas datu tabula. Ja datu matricas šūnā ij ir ierakstīta atribūta vērtība a_{ij} , tad fiktīvā mainīgā vērtībām attiecīgajā matricas šūnā tiek ierakstīta vērtība 1. Ja atribūta vērtības nav, tad fiktīvajam mainīgajam ieraksta vērtību 0. Piemēram, šī matrica attēlo fiktīvo mainīgo vērtības D sākotnējo datu tabulai 2.3.1.

$$D = \begin{pmatrix} 1 \dots 1 \dots 1 \\ 1 \dots 1 \dots 1 \\ 1 \dots 1 \dots 1 \\ 1 \dots 1 \dots 1 \\ 1 \dots 1 \dots 1 \\ 0 \dots 1 \dots 1 \\ 1 \dots 1 \dots 0 \\ 1 \dots 1 \dots 1 \\ 1 \dots 1 \dots 1 \\ 1 \dots 1 \dots 1 \end{pmatrix}.$$

Pēc tam tiek veikts attiecīgo attālumu aprēķins gadījumiem, kad matricas D rindās ir tikai vērtības 1 to atribūtu pozīcijās, kas ir iekļauti attāluma aprēķinos.

Dažkārt šis noteikums tiek formalizēts, ieviešot tā saucamo indikatora funkciju, kuras vērtības ir vienādas ar 1 pieejamām atribūtu vērtībām un vienādas ar 0 iztrūkstošām atribūtu vērtībām.

Tādā gadījumā tiek aprēķināts attālums starp piemēra o_i ar iztrūkstošu atribūta vērtību un piemēru o_k , fiktīvo mainīgo vērtību matricā tiek pārbaudītas paraugam atbilstošās šūnas o_k . Ja visās

šajās šūnās ir 1 (t.i., indikatora funkcijas vērtības ir vienādas ar 1), tad tiek aprēķināts attiecīgais attālums. Ja vismaz vienā no šīm šūnām ir 0 (indikatora funkcijas vērtība ir 0), tad attiecīgais attālums netiek aprēķināts.

Ja atribūtu vērtības ir norādītas kategoriskā (nominālā vai ordinārā) skalā, iztrūkstošo atribūtu vērtību imputācijas process kļūst ievērojami sarežģītāks. Tas ir saistīts ar tuvības mēra starp piemēriem noteikšanas problēmu. Skaitlisku atribūtu vērtību gadījumā attāluma aprēķins starp diviem piemēriem ir balstīts uz katra atribūta vērtību atšķirību summu šiem piemēriem. Gadījumā, ja atribūtu vērtības mēra nominālā vai ordinārā skalā, atribūtu vērtību skaitliskās atšķirības jēdzienam nav jēgas. Tāpēc ir nepieciešami citi tuvības mērījumu (attālumu) aprēķini starp attiecīgajiem piemēriem.

Diemžēl ir maz zinātnisko darbu k-tuvāko kaimiņu metodes izmantošanas kategorisko atribūtu novērtējumos. Darbā [Jönsson P., Volins C., 2004] piedāvāta metode iztrūkstošo atribūtu vērtību imputācijai gadījumam, kad visu atribūtu vērtības mēra Likerta skalā. Attālumi starp piemēriem tiek aprēķināti Likerta skalas punktus apstrādājot kā reālus skaitļus. Šī pieeja šķiet nepamatota, jo šādas darbības ar Likerta skalas punktiem nav korekta.

Darbā [Schwenger H., Iestad K., 2008] tiek piedāvāta korektāka pieeja attālumu novērtēšanai starp piemēriem, pamatojoties uz kategoriskām atribūtu vērtībām. Diemžēl šī pieeja ir piemērojama tikai gadījumos, kad visiem atribūtiem ir vienāds vērtēšanas kategoriju skaits.

Darbs [Tutz G., Ramzan Sh., 2014] ir saistīts ar jaunu piemēru klasifikāciju, novērtējot attālumus starp šo piemēru atribūtu kategoriskiem novērtējumiem un dotajām klasēm piederošo piemēru atribūtu kategoriskiem novērtējumiem. Piedāvātā attālumu novērtēšanas metode tika modificēta un izmantota, lai imputētu iztrūkstošās atribūtu vērtības [Faisal Sh., Tutz G., 2016]. Šajā sadaļā mēs piedāvājam vienkāršāko tiešo metodi, lai novērtētu attālumus starp šajā darbā piedāvātajiem attiecīgajiem piemēriem.

Pieņem, ka dati uzdoti taisnstūra izmēra tabulas veidā $m \times n$, kurā katra rinda atbilst piemēram (objekts, indivīds) o_i , $i = 1, \dots, m$ un katra kolona atbilst atribūtam a_j , $j = 1, \dots, n$. Katra atribūta a_j vērtība ir novērtēta ar nominālo vai ordinālo kategoriju r_j , $j = 1, \dots, r_j$. Atribūta a_{ij} vērtība tabulas i rindas un j kolonas krustojumā tiek parādīta kategoriskā formā kā $p_{ij}^{r_j}$. Tad attālumu starp piemēriem o_i un o_l var novērtēt ar vienkārša atbilstības koeficienta (*simple matching coefficient*) palīdzību:

$$d(o_i - o_l) = \sum_{j=1}^n I(p_{ij}^{r_j} - p_{lj}^{r_j}), \quad (2.3.5)$$

kur $p_{ij}^{r_j}$ - piemēra o_i atribūtu kategorijas vērtība a_j ;

$p_{lj}^{r_j}$ - piemēra o_l atribūtu kategorijas vērtība a_j ;

I - indikatora funkcijas vērtība, kas definēta sekojošā veidā:

$$I = \begin{cases} 0, & \text{ja } p_{ij}^{r_j} = p_{lj}^{r_j}; \\ 1, & \text{ja } p_{ij}^{r_j} \neq p_{lj}^{r_j}. \end{cases} \quad (2.3.6)$$

Citiem vārdiem sakot, par attāluma starp piemēriem o_i un o_l novērtējumu tiek ņemta nulļu un vieninieku summa, kas noteikta ar vienādojumiem (2.3.5), (2.3.6).

Ja piemērā o_i atribūtam a_{ij} ir iztrūkstošas atribūta vērtības (kategorijas). Tad, izmantojot vienādojumus (2.3.5), (2.3.6), var aprēķināt attālumus no šī piemēra līdz visiem pārējiem piemēriem, izmantojot visu atribūtu vērtības, izņemot atribūtu ar iztrūkstošo vērtību.

Pamatojoties uz iegūtajiem attāluma aprēķiniem, tiek atlasīti piemēri, kuriem ir vismazākie attālumi līdz piemēram ar iztrūkstošo vērtību. Vienkārša piemēru sakārtošana kategorisko atribūtu novērtējumos vairumā gadījumu nav iespējama, kā tas tiks parādīts nākamajā piemērā. Tāpēc arī bieži

vien nav iespējams iestatīt vēlamu tuvāko kaimiņu vērtību k . Rezultātā tuvākie kaimiņi tiek izvēlēti pēc principa - izmantot visus iespējamus.

Par imputēto vērtību tiek pieņemta atribūta vērtību mediāna ar iztrūkstošo šī atribūta vērtību, salīdzinājumā ar tuvāko kaimiņu kopumu.

Ja atlasīto atribūtu vērtību izlasei ir vairāk nekā viena mediāna (t.i., ir vienāds atribūta biežāko kategoriju skaits), tad imputētās vērtības izvēle starp tikpat bieži sastopamajām vērtībām tiek veikta nejauši vai pamatojoties uz kādu papildu nosacījumu.

Apskatām vienkāršu ilustratīvu piemēru. Ir pieejams daudz datu ar 10 piemēriem (personas). Šīs personas raksturo trīs atribūtu vērtības: a_1 - dzimums, a_2 - acu krāsa, a_3 - matu krāsa.

Atribūtam a_1 ir divas nominālās kategorijas:

p_1^1 - vīrietis;

p_1^2 - sieviete.

Atribūtam a_2 ir trīs nominālās kategorijas:

p_2^1 - brūnas acis;

p_2^2 - pelēkas acis;

p_2^3 - zilas acis.

Atribūtam a_3 ir trīs nominālās kategorijas:

p_3^1 - blondīne/blondīns;

p_3^2 - brūni mati;

p_3^3 - brunete/brunets.

Sākotnējie dati ir parādīti tabulā 2.3.7.

Tabula 2.3.7

Sākotnējā datu izlase

Piemēri	Atribūtu vērtību kategorijas		
	$(a_1) \quad p_1^i$	$(a_2) \quad p_2^j$	$(a_3) \quad p_3^k$
o_1	p_1^1	p_2^2	p_3^3
o_2	p_1^2	p_2^2	p_3^1
o_3	p_1^2	p_2^1	p_3^1
o_4	p_1^1	?	p_3^3
o_5	p_1^2	p_2^3	p_3^2
o_6	p_1^1	p_2^1	p_3^3
o_7	p_1^2	p_2^2	p_3^1
o_8	p_1^1	p_2^1	p_3^1
o_9	p_1^2	p_2^1	p_3^2
o_{10}	p_1^1	p_2^1	?

Šajā tabulā ar “?” apzīmētas iztrūkstošās atribūtu vērtības a_2 un a_3 .

Izmantojot vienādojumus (2.3.5), (2.3.6), aprēķinām attālumus no piemēra o_4 līdz citiem piemēriem, pamatojoties uz atribūtu vērtībām a_1, a_3 .

$$d(o_4 - o_1) = 0 + 0 = 0;$$

$$d(o_4 - o_2) = 1 + 1 = 2;$$

$$\begin{aligned}
d(o_4 - o_3) &= 1+1=2; \\
d(o_4 - o_5) &= 1+1=2; \\
d(o_4 - o_6) &= 0+0=0; \\
d(o_4 - o_7) &= 1+1=2; \\
d(o_4 - o_8) &= 0+1=1; \\
d(o_4 - o_9) &= 1+1=2.
\end{aligned}$$

Attālums $d(o_4 - o_{10})$ netiek aprēķināts, jo piemērā o_{10} nav atribūta vērtības $a_{10,3}$.

Izdalām piemērus ar mazākajiem attālumiem no piemēra o_4 . Tie ir piemēri o_1, o_6, o_8 . Šajos piemēros ir šādas atribūtu a_2 kategorijas: p_2^2, p_2^1, p_2^2 . Šīs vērtību izlases mediāna atbilst kategorijai p_2^2 . Tāpēc šī kategorija ir jāuzskata par iztrūkstošo atribūta vērtību a_{42} . Citiem vārdiem sakot, personai o_4 ir jābūt pelēkām acīm.

Veicam aprēķinus, lai izrēķinātu iztrūkstošo atribūta vērtību $a_{10,3}$.

$$\begin{aligned}
d(o_{10} - o_1) &= 1+1=2; \\
d(o_{10} - o_2) &= 0+1=1; \\
d(o_{10} - o_3) &= 0+1=1.
\end{aligned}$$

Attālums $d(o_{10} - o_4)$ netiek aprēķināts, jo piemērā o_4 nav atribūta vērtības a_{42} .

$$\begin{aligned}
d(o_{10} - o_5) &= 0+1=1; \\
d(o_{10} - o_6) &= 1+0=1; \\
d(o_{10} - o_7) &= 0+1=1; \\
d(o_{10} - o_8) &= 1+0=1; \\
d(o_{10} - o_9) &= 0+1=1.
\end{aligned}$$

Šajā piemērā rodas problēma ar k -tuvāko kaimiņu izvēli. Tā kā lielākā daļa attālumu ir vienādi ar 1, mēs nevaram sarindot šīs vērtības un pamatotī noteikt tuvāko kaimiņu skaitu k . Tāpēc mēs uzskatīsim visus piemērus, kuriem ir attālumi līdz piemēram o_{10} , par piemēra o_{10} tuvākajiem kaimiņiem. Šajos piemēros ir šādas atribūta a_3 vērtības: $p_3^1, p_3^1, p_3^2, p_3^3, p_3^1, p_3^1, p_3^2$. Šo vērtību izlases mediāna ir kategorija p_3^1 , tātad šī kategorija tiek uzskatīta par iztrūkstošā atribūta $a_{10,3}$ vērtību. Citiem vārdiem sakot, sievietei, kas avota datu tabulā norādīta kā o_{10} , jābūt blondīnei.

Darbā [Tutz G., Ramzan Sh., 2014] tiek piedāvāta alternatīva pieeja attālumu noteikšanai starp piemēriem, kas tiek novērtēti pēc kategorisku atribūtu vērtībām. Šī pieeja izmanto sākotnējo datu tabulu un īpašu atribūtu vērtējuma kategoriju - bināro kodējumu. Šajā nodaļā šī pieeja netiek apskatīta. Interesenti var iepazīties ar šo pieeju norādītajā darbā.

2.4. Imputācija uz sagaidāmā rezultāta maksimizācijas pamata

Šīs nodaļas iepriekšējās sadaļās apskatītās metodes iztrūkstošo atribūtu vērtību imputēšanai neņēma vērā paņēmienus, kādā veidā datos parādās iztrūkstošie dati (sk. 2.1. sadaļu). Šajā sadaļā aprakstītā metode ņem vērā iztrūkstošos datu veidošanās paņēmienus: MAR un MCAR. NMAR gadījumā šī metode nav piemērojama.

Ir daudzas praktiskas pieejas, lai noteiktu faktisko iztrūkstošo datu veidošanās paņēmieni. Lielākā daļa statistikas programmatūras pakotņu, tostarp IBM SPSS Statistica, izmanto sagaidāmā rezultāta maksimizēšanas metodes, lai pēc noklusējuma noteiktu iztrūkstošo datu veidošanās paņēmienus.

Jāņem vērā, ka turpmāk aprakstītā metode ir piemērojama tikai skaitliskiem datiem, kuros atribūtu vērtības tiek mērītas intervālu vai attiecību skalās. Lai gan ir bijuši mēģinājumi šo metodi attiecināt arī uz jauktu datu gadījumiem (kvantitatīvās un kategoriskās atribūtu vērtības), tomēr šie mēģinājumi nav guvuši plašu praktisku pielietojumu.

Sagaidāmā rezultāta maksimizēšanas princips ir formulēts dažādos veidos un izmantots dažādos pētījumos kopš 20. gadsimta 40. gadiem. Tikai darbā [Dempster A. et al., 1977] tika sniegts stingrs matemātiskais pamatojums iztrūkstošo atribūtu vērtību imputācijai un tika pierādīta rezultātu iegūšanas konverģence. Tāpat šīs metodes pielietojamība imputācijas problēmām ir izklāstīta darbā [Little P.J.A., Rubin D., 1987]. Pēdējās desmitgadēs sagaidāmās vērtības maksimizēšanas metode ir radusi plašu pielietojumu dažādās datu iegūšanas un analīzes jomās.

Sagaidāmā rezultāta vērtības maksimizēšanas metode ir balstīta uz maksimālās līdzības principa. Kāda ir šī principa būtība? Lai to izskaidrotu, mēs izmantojam vienkāršu piemēru – metam monētu un ir jānovērtē ģerboņa izkrišanas varbūtība. Monētas priekšpuse tiek saukta par “aversu”, mugurpuse - par “reversu”, bet mala par “jostu”. Monētas bieži mēdz izmantot kā vienkāršotu metamo kauliņu. Monētas aversu monētas mešanā mēdz dēvēt par “ciparu”, savukārt reversu - par “ērgli” vai “ģerboni”.

Lai novērtētu monētas izkrišanas varbūtību, tiek veikta monētas mešanas sērija. Pieņemsim, ka ir bijuši 100 metieni, ģerbonis izkritis 60 reizes un cipars 40 reizes. Apzīmēsim ar θ ģerboņa izkrišanas varbūtību: $p(\text{ģerbonis}) = \theta$, tad $p(\text{cipars}) = 1 - \theta$. Kā novērtēt vērtību θ ? Ir skaidrs, ka

$$\hat{\theta} = \frac{h}{n} = \frac{\text{Ģerboņa izkrišanas reižu skaits}}{\text{Kopējais metienu skaits}} = \frac{60}{100} = 0.60.$$

Formāli vērtības θ sadalījums ir saskaņā ar binomiālo likumu

$$p(\theta) = \binom{n}{h} \theta^h (1 - \theta)^{n-h}.$$

Šajā vienkāršajā piemērā parametra $\hat{\theta}$ novērtējums ir acīmredzams un neapšaubāms.

Pieņem, ka ir gadījumu skaitļu vērtību izlase, kuru sadalījums ir atbilstošs normālajam sadalījumam. Šajā gadījumā mēs varam viegli aprēķināt maksimālās līdzības novērtējumu vidējai vērtībai μ un dispersijai σ^2 .

Reālajā dzīvē viss nav tik vienkārši. Mums var būt tādu normālu sadalījumu izlases, kurām ir jānovērtē daudz parametru. Iespējams, ka nejaušā lieluma vērtību izlasei ir noteikts sadalījums, tāpēc ir grūti tieši novērtēt tā parametrus.

Lai sniegtu priekšstatu par šāda veida problēmām, apskatām šādu piemēru [Do Ch.B., Batzoglou S., 2008]. Pieņemsim, ka ir divas monētas A un B , un tika veiktas trīs monētas A mešanas sērijas un divas monētas B mešanas sērijas. Katra sērija sastāvēja no 10 attiecīgās monētas metieniem. Katrā metienu sērijā tika uzskaitīts izkritušo ciparu un ģerboņu skaits. Aprēķinu rezultāti ir parādīti attēlā. 2.4.1.a). Pamatojoties uz iegūtajiem rezultātiem, tika aprēķinātas varbūtības izkrist ģerbonim katras monētas mešanas laikā: $\hat{\theta}_A = 0.80$, $\hat{\theta}_B = 0.45$. Šīs aplēses tiek iegūtas, pamatojoties uz maksimālās līdzības principu.

Mainīsim nosacījumus mūsu uzdevuma risināšanai. Metienu sērijas tiek izpildītas tieši tādā pašā veidā, un tiek iegūti tie paši metienu rezultāti. Taču tagad *nav zināms*, kura monēta tika izmantota

katrā metienu sērijā. Tāpat kā iepriekš, mūsu uzdevums ir novērtēt ģerboņa izkrišanas varbūtību, nejauši metot monētas A un B .

Ar ko šis eksperiments atšķiras no iepriekšējā eksperimenta? Galvenā atšķirība ir tā, ka eksperiments tagad tiek veikts *nepilnīgu datu apstākļos* – mēs nezinām, kura monēta tiek izmantota katrā metienu sērijā. Acīmredzot tiešs veids, kā aprēķināt mūs interesējošās varbūtības, tagad nav iespējams. Ir nepieciešams cits veids, kā atrisināt šo problēmu.

Lai to izdarītu, mēs ieviešam jaunu mainīgo lielumu $z = (z_1, z_2)$, kur z_1 ir varbūtība, ka monēta A tiek izmantota konkrētajā metienu sērijā, un z_2 ir varbūtība, ka monēta B tiek izmantota konkrētajā metienu sērijā. Šāda veida mainīgo sauc par *latento mainīgo*.

Ņemot vērā sākotnējos nosacījumus, problēmu var atrisināt, izmantojot sagaidāmās vērtības maksimizācijas metodi. Aprēķinu procedūru diagramma ir parādīta attēlā 2.4.1.b.

Aplūkojamās problēmas kontekstā mēs runājam par sagaidāmo parametru aplēšu (varbūtību $\hat{\theta}_A$, $\hat{\theta}_B$) maksimizāciju, pamatojoties uz visu pieejamo informāciju. Lai inicializētu algoritmu, ir jāuzdod mūs interesējošo varbūtību sākotnējās vērtības $\hat{\theta}_A^0$, $\hat{\theta}_B^0$. Šīs vērtības var izvēlēties patvaļīgi vai pamatojoties uz dažiem papildu apsvērumiem. Mūsu piemērā ir norādītas šādas vērtības: $\hat{\theta}_A^0 = 0.60$, $\hat{\theta}_B^0 = 0.50$. Pēc tam tiek veikts iepriekš aprakstītais eksperiments un tiek iegūti tādi rezultāti, kādi parādīti tabulā attēlā 2.4.1.a, bet tagad šos rezultātus var attiecināt gan uz monētu A , gan uz monētu B .

Pamatojoties uz pašreizējām parametru vērtībām, var novērtēt varbūtības z_1 un z_2 , t.i., varbūtības, vai konkrētā metienu sērijā izmantotā monēta ir monēta A vai B . Pirmajai metienu sērijai $z_1 = 0.45$, $z_2 = 0.55$.

Katras metienu sērijas rezultāti tiek reizināti ar atbilstošajām varbūtību vērtībām z_1 , z_2 . Tādā veidā mēs iegūstam atbilstošās ģerboņu un ciparu proporcijas pie nosacījuma, ka izmestā monēta ir vai nu monēta A vai monēta B . Faktiski mēs pārdalām faktisko ģerboņu un ciparu skaitu starp abām monētām, pamatojoties uz varbūtību z_1 un z_2 vērtībām.

Visas iepriekš aprakstītās procedūras veido algoritma E soli. Šis apzīmējums cēlies no angļu vārda "*expectation*" - cerība.

Pamatojoties uz iegūtajiem rezultātiem, tiek aprēķinātas jaunās varbūtību vērtības: $\hat{\theta}_A^{(1)} = 0.71$, $\hat{\theta}_B^{(1)} = 0.58$. Šo soli apzīmē ar M — no angļu valodas vārda "*maximization*" - maksimizācija.

Izmantojot jaunās varbūtības aprēķinus $\hat{\theta}_A^{(1)}$, $\hat{\theta}_B^{(1)}$ tiek veikta nākamā algoritma soļu E un M iterācija. Rezultātā iegūstam aplēses $\hat{\theta}_A^{(2)}$, $\hat{\theta}_B^{(2)}$. Process tiek atkārtots cikliski, līdz aplēšu izmaiņas soļos (t) un $(t+1)$ kļūst mazākas par uzdoto robežvērtību. Iepriekšējā piemērā pēc desmitās algoritma iterācijas tiek iegūti šādi rezultāti: $\hat{\theta}_A^{(10)} = 0.80$, $\hat{\theta}_B^{(10)} = 0.52$. Šie rezultāti labi saskan ar rezultātiem, kas iegūti uz maksimālās līdzības pamata ar pilniem ieejas datiem, ņemot vērā, ka mūsu gadījumā rezultāti ir iegūti, pamatojoties uz nepilnīgiem datiem.

Analizējot iepriekš dotā piemēra rezultātus, var izdarīt neapstrīdamu secinājumu, ka varbūtības aplēses, kas iegūtas, pamatojoties uz sagaidāmās vērtības maksimizēšanas metodes pamata, tajā pašā laikā ir arī maksimālās līdzības aplēses, kas iegūtas ar nepietiekamiem sākotnējiem datiem.

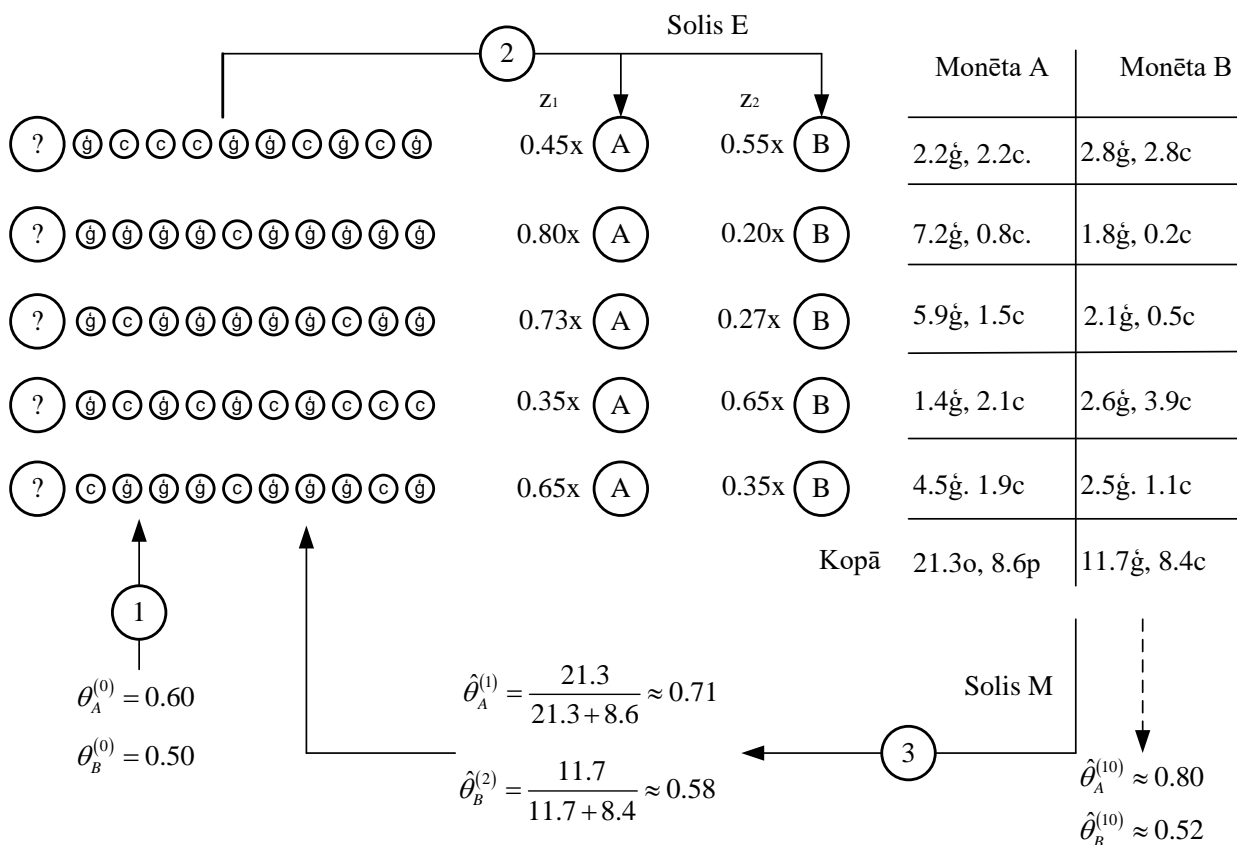
a) maksimālā līdzība

	Monēta A	Monēta B
(B) (g) (c) (c) (c) (g) (g) (c) (g) (c) (g)		5g, 5c
(A) (g) (g) (g) (g) (c) (g) (g) (g) (g) (g)	9g, 1c	
(A) (g) (c) (g) (g) (g) (g) (g) (c) (g) (g)	8g, 2c	
(B) (g) (c) (g) (c) (g) (c) (g) (c) (c) (c)		4g, 6c
(A) (c) (g) (g) (g) (c) (g) (g) (g) (c) (g)	7g, 3c	
Kopā	24g, 6c	9g, 8c

$$\hat{\theta}_A = \frac{24}{24+6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9+11} = 0.45$$

b) sagaidāmās vērtības maksimizācija



Attēls 2.4.1. Shēmas, lai iegūtu aplēses par ģerboņa un cipara izkrišanas varbūtību, metot monētas A un B

- a) pamatojoties uz maksimālās līdzības principa;
- b) pamatojoties uz sagaidāmās vērtības maksimizēšanas principa

Ieviesīsim formālas definīcijas, kas saistītas ar sagaidāmās vērtības maksimizācijas metodes izmantošanu iztrūkstošo atribūtu vērtību imputēšanai [Garsia S., et al., 2015].

Pieņemsim, ka mūsu mērķis ir modelēt atkarīgos gadījuma mainīgos kā novēroto mainīgo a un latentu mainīgo b , kas ģenerē a . Par pamatu ņemam to, ka nezināmu parametru izlase Q nosaka

varbūtību $p_\theta(a)$, $p_\theta(b)$ sadalījumus. Šo nezināmo parametru novērtēšanas process sastāv no diviem posmiem, kas tiek atkārtoti secīgi, līdz tiek panākta rezultātu konverģence.

1. Solī E - tiek aprēķināta sagaidāmā vērtība

$$Q(\theta, \theta') = \sum p_{\theta'}(b/a) \log p_\theta(b, a), \quad (2.4.1)$$

kur: $p_{\theta'}$ - nosacītā varbūtība b (atkarībā no a), pie jaunām θ' parametru vērtībām;

$p_\theta(b, a)$ - kopējā varbūtība a un b pie vecajām parametru θ vērtībām.

Kāpēc vienādojumā 2.4.1 tiek lietots termins $\log p_\theta(b, a)$? Lieta tāda, ka varbūtību vērtības, kas parādās vienādojumā 2.4.1, ir ļoti mazas un šo varbūtību reizināšanas rezultāts var nebūt precīzs. Pāreja no varbūtību vērtībām uz to logaritmiem ļauj atrisināt šo problēmu. Tā kā visas varbūtības ir mazākas par 1, mūs interesē tādas parametru izlases θ logaritmu vērtības, kas ir tuvas 0.

2. Solī M - tiek maksimizētas iepriekš aprēķinātās θ vērtības. Kā tiek panākta šī maksimizācija? Apskatām mainīgo A un B nosacītās sagaidāmās vērtības, kas iegūtas no sadalījumiem $p(a)$, $p(b)$, lai izpildītu soli M . Nosacītie sadalījumi tiek definēti kā

$$p(b/a) = \frac{p(b, a)}{p(a)}.$$

Tad gadījuma mainīgā B nosacītā sagaidāmā vērtība tiek uzdota kā

$$E(B) = \sum_b p(b) b.$$

Nosacītā sagaidāmā vērtība $\log p_\theta(b, a)$ pie dotajām a un θ' vērtībām var tikt aprēķināta sekojošā veidā:

$$\begin{aligned} E(\log p(b, a/\theta) / a, \theta') &= \sum_b p(b/a, \theta') \log p(b, a/\theta) = \\ &= \sum_b p_{\theta'}(b, a) \log p_\theta(b, a). \end{aligned} \quad (2.4.2)$$

Acīmredzot, ja $\sum_b p_{\theta'}(b/a) \log p_{\theta'}(b, a) > \sum_b p_\theta(b/a) \log p_\theta(b, a)$, tad $p_{\theta'}(a) > p_\theta(a)$.

Šī arī ir sagaidāmās vērtības maksimizācijas metodes galvenā ideja: secīgi uzlabojot varbūtību sadalījuma parametru θ vērtības, mēs varam secīgi palielināt šī sadalījuma sagaidāmo vērtību. Iepriekš šajā sadaļā mēs demonstrējām šo ideju ar vienkāršu ilustratīvu piemēru. Tagad mēs esam formāli pamatojuši šo ideju.

Kā sagaidāmās vērtības maksimizēšanas metode darbojas iztrūkstošo atribūtu vērtību imputācijas kontekstā? Lai izmantotu šo metodi, vispirms ir jānosaka dažas sākotnējās parametru vērtības atribūtu vērtību sadalījumiem. Sagaidāmās vērtības maksimizēšanas metodes algoritms labi darbojas gadījumos, ja atribūtu vērtību izlases ir sadalītas pēc normāliem vai tiem tuviem sadalījumiem.

Katrā iterācijā sagaidāmās vērtības maksimizēšanas metodes algoritms novērtē atribūtu $\hat{\mu}_j$ vidējās vērtības un variācijas/kovariācijas vērtības visām atribūtu vērtību kopām. Variācijas/kovariācijas vērtības tiek parādītas atbilstoša izmēra matricas $\hat{\Sigma}$ veidā. Pamatojoties uz iegūtajiem datiem, tiek aprēķināti regresijas parametri B , pēc kuriem tiek aprēķinātas atribūtu pašreizējās vērtības. Šīs pašreizējās vērtības tiek noteiktas pēc vienādojuma

$$A_{mis} = \mu_{mis} + (A_{obs} - \mu_{obs}) B + e, \quad (2.4.3)$$

kur A_{mis} - iztrūkstošo atribūtu vērtību izlase;

A_{obs} - esošo atribūtu vērtību izlase;

μ_{mis} - iztrūkstošo atribūtu vērtību vidējās vērtības;

μ_{obs} - esošo atribūtu vērtību vidējās vērtības;

B - regresijas parametri;

e - kļūdas rādītājs, kas tiek uzdots kā nejaušs vektors ar vidējās vērtību 0 un ar nezināmu variāciju matricu.

Iepriekš aprakstītās skaitļošanas procedūras atspoguļo sagaidāmās vērtības maksimizācijas metodes algoritma soļa E būtību iztrūkstošo atribūtu vērtību imputācijā. Vienkāršoti izsakoties, šī soļa darbības mērķis ir atkārtoti novērtēt iztrūkstošās atribūtu vērtības, pamatojoties uz novērtētajiem regresijas parametriem B .

Algoritma M solī tiek atkārtoti novērtēti atribūtu vērtību kopu sadalījumu parametri, izmantojot gan esošajā, gan iepriekšējā iterācijā aprēķinātās atribūtu vērtības.

$$\hat{\mu}^{(t+1)}(a_j) = \frac{\sum_{i=1}^m a_{ij}^{(t)}}{m}, \quad j = 1, \dots, n. \quad (2.4.4)$$

$$\hat{\Sigma}^{(t+1)} = \frac{1}{m} \sum_{i=1}^m \left(\hat{\Sigma}^{(t)} - \left(\hat{\mu}^{(t+1)} \hat{\mu}^{(t+1)} \right) \right). \quad (2.4.5)$$

Izmantojot pārvērtēšanas rezultātus $\hat{\mu}^{(t+1)}(a_j)$ un $\hat{\Sigma}^{(t+1)}$, $j = 1, \dots, n$ var pārvērtēt regresijas parametru vērtības $\hat{B}^{(t+1)}$. Pēc tam aprēķinātās atribūtu vērtības tiek atkārtoti novērtētas un process tiek atkārtots, kamēr tiks sasniegta rezultātu konverģence.

Būtisks sagaidāmās vērtības maksimizēšanas metodes trūkums ir tas, ka faktiski to izmanto, lai meklētu sadalījuma parametru maksimumu lokālās vērtības. Šie lokālie maksimumi var nebūt globālie maksimumi. Lai izvairītos no konverģences uz lokālajiem maksimumiem, ieteicams atkārtot metodes izpildi pie dažādiem sākuma nosacījumiem. Iztrūkstošo atribūtu vērtību imputācijas kontekstā tas nozīmē, ka sākotnējo atribūtu vidējo vērtību vietā iztrūkstošās vērtības tiek aizstātas ar skaitļiem, kas nav vidējās vērtības. Tādējādi metodes soļi tiek veikti pie jauniem sākuma nosacījumiem. Ja abos gadījumos iegūtie rezultāti sakrīt, tas norāda, ka ir sasniegts globālais maksimums, ja rezultāti atšķiras, tad par galarezultātu var ņemt labāko no diviem iegūtajiem.

Sagaidāmās vērtības maksimizēšanas metodes praktiskai izmantošanai ir nepieciešams ļoti liels skaitļošanas procedūru skaits. Pat ar ļoti ierobežotu ievaddatu skaitu nepieciešamo aprēķinu veikšana manuāli nav reāla, tāpēc attiecīgie aprēķini jāveic, izmantojot piemērotus programmatūras rīkus.

2.5. Daudzkārtīgā imputācija

Iepriekšējā sadaļā tika apspriesta sagaidāmās vērtības maksimizācijas metode iztrūkstošo atribūtu vērtību imputēšanai. Īsumā atkārtosim šīs metodes būtību, lai to salīdzinātu ar daudzkārtīgās imputācijas metodi, kas tiks aplūkota šajā sadaļā.

Sagaidāmās vērtības maksimizācijas metodes pamatā ir mūs interesējošo atribūtu vērtību sadalījuma parametru novērtējums, pamatojoties uz to sagaidāmo vērtību maksimizēšanu. Maksimizēšanas process tiek veikts iteratīvi. Katrā iterācijā notiek sadalījuma parametru secīga tuvināšana to faktiskajām vērtībām. Pamatojoties uz atkārtoti novērtētajām parametru vērtībām, tiek veikts nākamais iztrūkstošo atribūtu vērtību imputācijas posms. Jaunā atribūtu vērtību komplektā (novērotās + imputētās vērtības) sadalījuma parametri tiek atkārtoti novērtēti un process tiek cikliski atkārtots, līdz tiek sasniegti galīgie rezultāti.

Daudzkārtīgās imputācijas pamatā ir cits imputācijas princips. *Apriori* tiek pieņemts, ka atribūtu vērtību izlases tiek sadalītas saskaņā ar kādu sadalījuma likumu. Parasti tiek pieņemts, ka šis sadalījums ir daudzfaktoru normālais sadalījums.

Izmantojot iegūtos rezultātus, katram sadalījumam tiek atlasītas (modelētas) tā aprēķinātās vērtības no atbilstošā sadalījuma. Pēc tam, izmantojot Beijesa formulu, pamatojoties uz sākotnējo iepriekšējo sadalījumu un aprēķinātajām atribūtu vērtībām, tiek noteikts atribūtu vērtību *aposteriori* sadalījums, ņemot vērā to imputētās vērtības.

Tad tiek izveidota jauna imputēto atribūtu vērtību izlase un process tiek atkārtots, līdz rezultāti konverģē. (Rezultātu konverģences jēdziena definīcija metodē tiks dota vēlāk šajā sadaļā).

Tagad mēs varam veikt abu metožu salīdzinājumu. Abās metodēs tiek pieņemts, ka atribūtu vērtību izlases ir sadalītas saskaņā ar kādu varbūtisko likumu. Bet sagaidāmās vērtības maksimizācijas metodē, pamatojoties uz šo pieņēmumu, šo sadalījumu parametri tiek secīgi pārvērtēti. Katrā iterācijā ar aprēķinātajām parametru vērtībām tiek veikta nākamā imputēto atribūtu vērtību pārvērtēšana. Citiem vārdiem sakot, katrs nākamais aprēķins tiek veikts, pamatojoties uz iepriekšējo atribūta vērtību parametru pārvērtēšanu.

Izmantojot daudzkārtīgās imputācijas metodi, katrā iterācijā no šo vērtību *aposteriori* sadalījumiem tiek veidota nejauša imputēto atribūtu vērtību izlase.

Daudzkārtīgās imputācijas metodi pirmo reizi tika piedāvāta darbā [Rubin D.B., 1977] un tika prezentēta pilnā formā darbā [Rubin D.B., 1987]. Šīs metodes attīstību izraisīja nepieciešamība apstrādāt ļoti lielu skaitu apskatu (dažādu veidu aptauju rezultāti). Kā minēts iepriekš, šāda veida apskatos dažādu iemeslu dēļ trūkst daudz atbilžu uz uzdotajiem jautājumiem. Lai atrisinātu šo problēmu, arī tika izstrādāta daudzkārtīgās imputācijas metode.

Sagaidāmās vērtības maksimizācijas metode un citas imputācijas metodes parasti ir saistītas ar konkrētiem datu apstrādes un analīzes mērķiem un uzdevumiem. Daudzkārtīgās imputācijas metode nav vērsta uz konkrētiem mērķiem un uzdevumiem. Tā ir universāla tādā nozīmē, ka iztrūkstošo atribūtu vērtību imputācija tiek veikta datu vākšanas posmā un iegūtie rezultāti tiek izmantoti citu šo datu analīzes problēmu risināšanai. Citiem vārdiem sakot, izmantojot sagaidāmās vērtības maksimizācijas metodi, iztrūkstošo atribūtu vērtību imputāciju un turpmāko datu analīzi veic viens lietotājs konkrēta uzdevuma kontekstā. Savukārt, izmantojot daudzkārtīgās imputācijas metodi, imputācija tiek veikta neņemot vērā datu analīzes mērķa uzdevumus un iegūtos rezultātus izmanto citi lietotāji konkrētu uzdevumu kontekstā.

Daudzkārtējās imputācijas metodes īpatnība ir tāda, ka iztrūkstošo atribūtu vērtību imputācija tiek veikta nevis vienu reizi, kas ir raksturīgi visām pārējām metodēm, bet gan noteiktu skaitu m reizes. Jebkuru vai visus iegūtos rezultātus var izmantot turpmākajos pētījumos.

Daudzkārtējās imputācijas metodes izpilde ir balstīta uz šādiem trim pieņēmumiem [Sinharay S., et al., 2001]:

1. Datu modelis.

Pirmais un vissvarīgākais solis daudzkārtējās imputācijas metodes izmantošanā ir pieņemt varbūtisko modeli, kas saista visu datu izlasi A , kas ir novēroto vērtību A_{obs} un iztrūkstošo vērtību A_{mis} kombinācija, ar šo vērtību parametru sadalījumu izlasi. Izmantojot šo varbūtisko modeli un *apriori* sadalījumu pēc parametriem, var noteikt sagaidāmo sadalījumu $p(A_{mis} / A_{obs})$ iztrūkstošām vērtībām, atkarībā no novērotajām vērtībām. Izmantojot šo sagaidāmo sadalījumu, var ģenerēt imputētās atribūtu vērtības.

Vispiemērotākais modelis nepārtrauktām atribūtu vērtībām ir pieņemt daudzfaktoru normālo sadalījumu visām atribūtu vērtībām. Ja dažu atribūtu vērtības ir sadalītas saskaņā ar likumu, kas būtiski atšķiras no normālā sadalījuma, tad šīs vērtības var pārveidot tā, lai transformācijas rezultāti tiktu sadalīti saskaņā ar normālo sadalījumu. (Atribūtu vērtību transformācijas problēmas ir apskatītas 6. nodaļā). Imputētās transformētās atribūtu vērtības var tikt apgrieztas, lai iegūtu to sākotnējās vērtības.

2. Atribūtu vērtību kopu apriori sadalījums.

Tā kā uz modeli balstītā daudzkārtējās imputācijas metode pēc būtības ir balstīta uz Beijesa teorēmas pamata (*Bayes rule*), ir jānosaka vispārējs *apriori* sadalījums, no kura var secināt *aposteriori* sagaidāmo sadalījumu $p(A_{mis} / A_{obs})$. Parasti ertības labad tiek izmantots kāds neinformatīvs *apriori* sadalījums, piemēram, $p(\theta) \propto \frac{1}{\sigma^2}$.

3. Iztrūkstošo datu mehānismi.

Uz modeli balstītā daudzkārtējās imputācijas metode pieņem to, ka iztrūkstošo datu mehānisms ir vai nu pilnīgi nejaušs iztrūkums (MCAR) vai nejaušs iztrūkums (MAR).

Kopumā daudzkārtējās imputācijas metodes procesus var definēt kā Montekarlo simulāciju Markova ķēdēs. Tipisks Montekarlo simulācijas piemērs ir imputēto atribūtu vērtību nejauša atlase no to *apriori* sadalījumiem. *Apriori* paredzamo sadalījumu secīga pārvērtēšana pēc būtības veido Markova ķēdi.

Aplūkojam daudzkārtējās imputācijas metodes matemātiskos pamatus. Turpmākais izklāsts ir balstīts uz definīcijām, kas sniegtas darbā [Sinharay S., et., 2001], kas savukārt ir balstīts uz [Schafer J.L., 1997] darbu. Tiek pieņemts, ka atribūtu vērtību izlases ir nepārtraukti nejauši mainīgie. Ja atribūtu vērtību izlases ir diskreti gadījuma mainīgie, tad visas turpmākās definīcijas ir derīgas, aizstājot integrēšanas darbības ar summēšanas operācijām.

Apzīmēsim ar A sākotnējo datu izlasi. Šī izlase sastāv no divām daļām: $A = (A_{obs}, A_{mis})$, kur A_{obs} ir novēroto datu izlase (atribūtu vērtības) un A_{mis} ir iztrūkstošo datu izlase (atribūtu vērtības). Tālāk pieņemsim, ka visu atribūtu vērtības ir sadalītas saskaņā ar likumu $p(A / \theta)$, kur θ ir visu modeļa parametru izlase (vektors).

Prognozējamo sadalījumu $p(A_{mis} / A_{obs})$ var attēlot kā

$$\begin{aligned} p(A_{mis} / A_{obs}) &= \int p(A_{mis}, \theta / A_{obs}) d\theta = \\ &= \int p(A_{mis} / A_{obs}, \theta) p(\theta / A_{obs}) d\theta. \end{aligned} \quad (2.5.1)$$

Iztrūkstošo atribūtu vērtību imputēšana no izlases A_{mis} tiek veikta divos posmos. Pirmais posms ietver parametru vērtību modelēšanu no *aposteriori* sadalījuma $p(\theta / A_{obs})$. Otrais posms ir modelēt iztrūkstošo atribūtu vērtību vektoru no nosacītā *aposteriori* sadalījuma $p(A_{mis} / A_{obs}, \theta)$ izmantojot pirmajā posmā ģenerētās θ vērtības.

Iztrūkstošo atribūtu vērtību imputācijas tiek veiktas, izmantojot pirmo no vienādojumiem (2.5.1), proti

$$p(A_{mis} / A_{obs}) = \int p(A_{mis}, \theta / A_{obs}) d\theta.$$

Šajā gadījumā tiek izmantots datu pieauguma algoritms (*data augmentation algorithm*), kas tika piedāvāts darbā [Tanner M.A., Wong W.H., 1987]. Būtībā šis algoritms ir Markova ķēdes simulācija. Algoritma k iterācijā imputētās atribūtu vērtības tiek atlasītas no nosacītā paredzamā sadalījuma $A_{mis}^{(k)} \sim p(A_{mis} / A_{obs}, \theta^{(k-1)})$, bet nākamās parametru vērtības tiek atlasītas no *aposteriori* sadalījuma visai datu izlasei $\theta^{(k)} \sim p(\theta / A_{obs} / A_{mis}^{(k)})$. Sākot no noteiktām sākotnējām vērtībām, šādu procedūru secība veido Markova ķēdi $\{A_{mis}^{(k)}, \theta^{(k)}, k = 1, 2, \dots\}$.

Kā var noteikt datu pieauguma algoritma rezultātu konverģenci? Citiem vārdiem sakot, kādos apstākļos (ar kādu k vērtību) jāpārtrauc algoritma izpildes procedūru secība? Aplūkojot sagaidāmās vērtības maksimizācijas metodi iepriekšējā sadaļā, algoritma konverģence tika noteikta pēc

nosacījuma, ka $(k+1)$ algoritma solī iegūtie rezultāti maz atšķiras no rezultātiem, kas iegūti k solī. Citiem vārdiem sakot, rezultātu atšķirības nepārsniedz noteiktu sliekšni. Datu pieauguma algoritmā šī algoritma konverģence tiek noteikta sarežģītākā veidā. Fakts ir tāds, ka šajā algoritmā parametru sadalījums nemainās, veicot secīgas darbības, bet mainās nejauši izvēlētais parametru vērtības. Darbā [Schafer J.L., Olson M.K., 1998] tiek piedāvāta šāda definīcija: datu pieauguma algoritms konverģē pēc k iterācijām, ja jebkura parametra vērtība iterācijā $t \in 1, 2, \dots$ ir statistiski neatkarīga no šī parametra vērtības iterācijā $t+k$.

Pēc vienas pilnas algoritma izpildes tiek iegūta parametru vērtību izlase un imputēto atribūtu vērtību izlase. Bet, kā saka pats metodes nosaukums, lieta nebeidzas ar vienu algoritma izpildi. Algoritms tiek izpildīts noteiktu skaitu m reizu. Ņemot vērā aprēķināto atribūtu vērtību izvēles nejausību un konverģences nosacījumu, mēs iegūstam m aprēķināto atribūtu vērtību izlases. Lai veiktu nepieciešamo datu analīzi, var izmantot jebkuru no pilnajām datu izlasēm, ieskaitot novērotās un aprēķinātās atribūtu vērtības.

Kādam vajadzētu būt daudzkārtējās imputāciju skaitam m ? Darbā [Rubin D.B., 1987] tiek apgalvots, ka pietiek ar 3 līdz 5 imputācijām. Tika iegūta šāda sakarība starp daudzkārtējo imputāciju skaitu m un novērtējuma efektivitāti

$$\text{Novērtējuma efektivitāte} \sim \left(1 + \frac{\gamma}{m}\right)^{-1},$$

kur γ - iztrūkstošo atribūtu vērtību proporcija sākotnējā datu izlasē.

Ja datus iztrūkst 30% atribūtu vērtību, 94% efektivitāti var sasniegt ar 5 galīgajām aizpildītajām datu izlasēm. Ja $m=10$, tad efektivitāte ir 97%, kas dod tikai nelielu efektivitātes pieaugumu par dubultām skaitļošanas izmaksām.

Saistībā ar daudzkārtējām imputācijām rodas jautājums: kāpēc ģenerēt m imputēto atribūtu vērtību izlases, ja tikai viena no šīm izlasēm tiek izmantota turpmākai datu analīzei? Lieta tāda, ka izmantojot daudzkārtīgās imputācijas rezultātus, var tikt novērtēta kopējo imputācijas rezultātu nenoteiktības pakāpe.

Darbā [Rubin D.B., 1987] tika formulēti šādi noteikumi, lai apvienotu mūs interesējošo parametru vērtības, kas iegūtas no katras m imputācijas.

Apzīmēsim mūs interesējošo parametru ar R . Tas varētu būt atribūtu vērtību izlases vidējais lielums, regresijas koeficients starp divām atribūtu vērtību izlasēm vai jebkurš cits atbilstošs parametrs.

Ar parametra \hat{R}_i vērtību apzīmēsim mūs interesējošo sadalījumu i iterācijā. Tad šī parametra apkopoto novērtējumu visām m imputācijām var aprēķināt ar šāda vienādojuma palīdzību

$$\bar{R} = \frac{1}{m} \sum_{i=1}^m \hat{R}_i. \quad (2.5.2)$$

Novērtēšanas variācija ietver divus komponentus: variācijas katrā iegūtajā datu izlasē un variācijas starp iegūtajām pilnajām datu izlasēm.

Iekšējā variācija tiek novērtēta pēc vienādojuma

$$\bar{U} = \frac{1}{m} \sum_{i=1}^m U_i. \quad (2.5.3)$$

Acīmredzami šī variācija skaitliski ir vienāda ar vidējo variāciju visās pilnajās datu izlasēs.

Variācija starp iegūtajām pilnajām datu izlasēm aprēķina, izmantojot vienādojumu

$$B = \frac{1}{m-1} \sum_{i=1}^m (\hat{R}_i - \bar{R})^2. \quad (2.5.4)$$

Kopējā variācija T ir šo divu komponentu koriģētā summa

$$T = \bar{U} + \left(1 + \frac{1}{m}\right) B. \quad (2.5.5)$$

Kvadrātsakne no T ir standarta kļūdas aplēse, kas saistīta ar vidējo novērtējumu \bar{R} .

Viena no daudzkārtējās imputācijas metodes priekšrocībām ir spēja novērtēt nenoteiktības parametru aplēsēs, kas saistītas ar iztrūkstošo atribūtu vērtību vairākkārtēju imputāciju. Neviena cita imputācijas metode nevar novērtēt rezultātu nenoteiktību.

Novērtējam daudzkārtējās imputācijas metodes galvenās priekšrocības un trūkumus. Šī analīze ir balstīta uz apkopotajiem datiem, kas sniegti darbā [Allison P.D., 2000]. Ērtības labad mēs piedāvājam vispārinātu daudzkārtējās imputācijas procedūru secību:

1. Imputējiet iztrūkstošās atribūtu vērtības, izmantojot piemērotu modeli, kas ietver nejaušas variācijas.
2. Atkārtojiet imputācijas procesu m reizes (parasti 3 – 5 reizes).
3. Veiciet vēlamo analīzi katram (vai jebkuram) sadalījumam no iegūtās datu izlases, izmantojot piemērotas statistikas metodes.
4. Izskaitļojiet vidējās parametru vērtības visās m izlasēs.
5. Aprēķiniet standartkļūdas: (a) aprēķinot standarta vidējo vērtību kļūdas; b) m parametru aplēšu variāciju aprēķināšana visās datu izlasēs; c) abu iegūto aplēšu apvienošana.

Dauzkārtējai imputācijai ir šādas vēlamās īpašības:

1. Piemērotu nejaušu novērtējumu ieviešana imputācijas procesā padara par iespējamu iegūt objektīvus visu vēlamo parametru aprēķinus.
2. Atkārtota imputācija nodrošina ticamus standarta kļūdas aprēķinus. Vienas imputācijas metodes nevar apstrādāt papildu kļūdas imputācijas procesa laikā.
3. Dauzkārtējo imputāciju var izmantot jebkura veida datiem, izmantojot piemērotu programmatūru.

Praktiski pielietojot daudzkārtējās imputācijas metodi, ir jāievēro noteiktas prasības. Šeit mēs iepazīstinām ar trim šāda veida pamatprasībām [Rubin D. B., 1987, 1996]:

1. Iztrūkstošo datu mehānismam ir jābūt MCAR vai MAR.
2. Modelim, kas izmantots imputēto vērtību ģenerēšanai, kaut kādā ziņā jābūt “korektam”.
3. Datu analīzei izmantotajam modelim ir jāatbilst modelim, ko izmanto atribūtu vērtību imputēšanai.

Iespējamie daudzkārtējās imputācijas metodes neveiksmīgi lietojumi ir saistīti tikai ar vienas vai vairāku šo prasību neievērošanas.

Dauzkārtējās imputācijas metodes praktiskā izmantošana prasa veikt daudzas sarežģītas skaitļošanas procedūras, tāpēc jāizmanto specializētus skaitļošanas rīkus.

Pastāv arī citas pieejas iztrūkstošo atribūtu vērtību imputēšanai. Starp daudzajām šādām pieejām var minēt imputāciju, kas balstīta uz komponentu analīzi, imputāciju, kuras pamatā ir neironu tīkli.

3. ANOMĀLIJU IDENTIFICĒŠANA DATOS

Šajā un nākamajās divās nodaļās ir aplūkota datu attīrīšanas (*cleaning*) problēmas. Datu attīrīšana ir saistīta ar kļūdu un neatbilstību identificēšanu un novēršanu datos, lai uzlabotu šo datu kvalitāti. Kvalitātes problēmas datu kopās rodas nepareizas objektu vai atribūtu definīcijas, iztrūkstošu atribūtu vērtību, neinformatīvu atribūtu izmantošanas un citu iemeslu dēļ. Nepieciešamība pēc datu attīrīšanas ievērojami palielinās, ja tiek integrēti dati no dažādiem avotiem.

Šajā nodaļā mēs aplūkojam vienu datu attīrīšanas aspektu, proti, anomāliju identificēšanu datos.

3.1. Kas ir anomālijas datos?

Zinātniskajā literatūrā par datu analīzi ir vairāki termini, ar kuru palīdzību tiek raksturotas anomālijas datos: “*anomaly detection*”, “*outlier detection*” vai “*novelty detection*”. Turpmākajā izklāstā mēs izmantosim terminu “anomāliju noteikšana jeb identificēšana”.

Darbā [Hawkins D.M., 1980] dota šāda definīcija: anomālija datos ir “novērojums, kas tik ļoti atšķiras no citiem novērojumiem, ka ir aizdomas, ka to ģenerēja cits mehānisms”.

Autori [Moore Y.D., McCabe G., 1999] anomāliju definē kā datus, kas neatbilst nekādiem nosacītajiem sadalījumiem.

Darbs [Barnett V., Lewis T., 1984] sniedz šādu definīciju: “Anomālija datu kopā ir novērojums, kas nav saderīgs ar doto datu kopu”.

Apkopojot visas šīs definīcijas, šajā nodaļā anomāli dati tiek definēti kā datu paraugi (objekti), kas būtiski atšķiras no citiem datu paraugiem atribūtu vērtībās.

Dažkārt literatūrā par datu iegūšanu un analīzi par anomālijām sauc par “novirzēm”, “pretrunām” vai “novirzēm no normas”.

Sekas, ko rada anomāliju ignorēšana datos, var izraisīt nopietnas negatīvas sekas [Salgado C. et al., 2016]: (1) kļūdu dispersijas palielināšanās un statistiskās nozīmes samazināšanās; (2) samazināta normalitāte gadījumos, kad anomālija neatbilst normāla sadalījuma likumiem.

Anomāliju jēdzienu nevajadzētu jaukt ar trokšņa jēdzienu datos. Trokšnis datos ir saistīts ar nelielām atribūtu vērtību novirzēm no to patiesajām vērtībām. Tajā pašā laikā anomālijām ir raksturīgas ievērojamas novirzes atribūtu vērtībās, salīdzinot ar šo atribūtu vērtībām galvenajā datu kopā.

Dodam ieskatu pielietojumos, kas izmanto anomāliju atpazīšanu un analīzi [Hedge V.J., Austin J., 2004]:

- Krāpšanas atpazīšana kredītkaršu izmantošanā.
- Problemātisku klientu noteikšana kredīta izlietošanā.
- Nesankcionētu darbību datortīklos identificēšana.
- Aktivitāšu monitorings: krāpšanās identificēšana ar mobilajiem tālruniem.
- Tīkla veiktspējas uzraudzība, piemēram, lai identificētu tīkla vājās vietas.
- Bojājumu diagnostika - procesa uzraudzība, lai noteiktu bojājumus iekārtās, piemēram, kosmosa kuģos.
- Strukturālo bojājumu atpazīšana - ražošanas līniju monitorings.
- Satelītattēlu analīze - jaunu vai vēl neklasificētu pazīmju noteikšana.
- Jaunumu atpazīšana paraugos - jaunas robotu vai novērošanas sistēmu reakcijas.
- Kustību segmentācija - attēla iezīmju atpazīšana, kas kustas neatkarīgi no pamata attēla.
- Medicīnisko stāvokļu uzraudzība, piemēram, sirdsdarbības uzraudzība.
- Farmaceutiskie pētījumi - jaunu molekulāro struktūru noteikšana.
- Neparedzētu šūnu atpazīšana datu bāzēs, lai identificētu kļūdas, viltojumus, negaidītas šūnas.
- Nemarakētu datu atpazīšana apmācības datu kopās.

Iemesli anomāliju parādīšanās datus var būt cilvēka kļūdas, mērinstrumentu darbības traucējumi, dabiskas novirzes datu unikālo objektu atribūtu vērtībās, krāpnieciska rīcība, izmaiņas tehniskās sistēmas darbībā, iekārtu elementu bojājumi u.c.

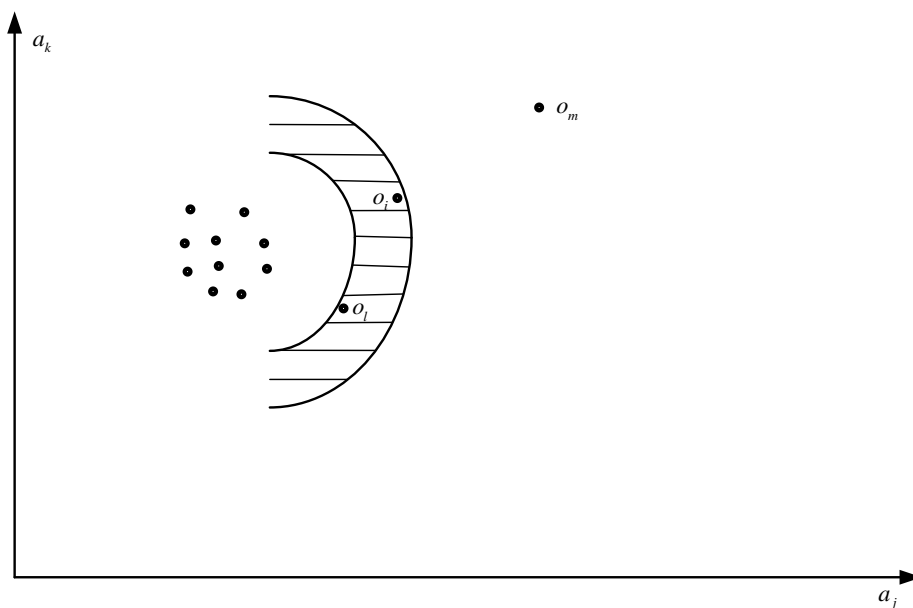
Ir jāsaprot, kādi ir galvenie mērķi anomāliju identificēšanai datus. Apskatām vienkāršu piemēru. Formatējot tekstu datorā, pareizrakstības programma automātiski atpazīst kļūdas, signalizē par to un piedāvā pareizu attiecīgā vārda pareizrakstību. Mērķis šeit ir identificēt anomāliju un savlaicīgi to novērst. Citiem vārdiem sakot, šajā kontekstā mērķis ir pilnībā attīrīt datus no anomālijām un iegūt normālus datus, t.i., tekstu bez gramatiskām kļūdām. Var teikt, ka mūs interesē tikai normāli dati un anomālijas nemaz neinteresē.

Cits piemērs. Lai identificētu pacienta iespējamo slimību, ārsts aicina pacientu veikt nepieciešamās pārbaudes un veikt izmeklējumus. Analizējot šos izmeklējumu rezultātus, ārsts izmeklē tos pacienta ķermeņa darbības parametrus, kas pārsniedz normu, t.i., meklē anomālijas. Tādējādi ārsts izmanto anomālijas, lai identificētu pacienta iespējamo slimību. Šajā gadījumā ārstu neinteresē normālie rādījumi, bet anomālijas kā svarīgas informācijas avotus.

Neatkarīgi no anomāliju identificēšanas mērķa, to noteikšanas procesi ir līdzīgi.

Anomāliju identifikāciju datus būtiski sarežģī trokšņa klātbūtne šajos datus un nepārprotamas robežas starp troksni un anomālijām neiespējamība. Apskatām attēlu 3.1.1. Šeit objekti ir attēloti ar punktiem divu atribūtu telpā: a_j, a_k . Lielākā daļa datu punktu veido kompakto kopu un tos var korekti interpretēt kā normālus datus. Datu objekta punkts o_m ļoti atšķiras no visiem citiem datiem, tāpēc objektu o_m var identificēt kā anomāliju.

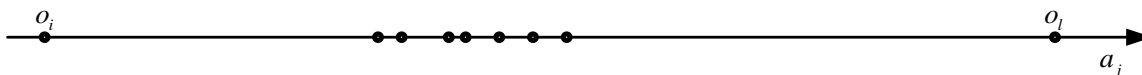
Problēmas rodas ar datu punktiem o_i un o_l . Šie punkti atrodas tā sauktajā “pelēkajā zonā”, kas ir izplūdusi robeža starp normāliem datiem un anomālijām.



Attēls 3.1.1. Trokšņu un anomāliju grafisks attēlojums

Objektus o_i un o_l datus var interpretēt kā troksni. Bet objekts o_l ir tuvāk normālu datu apgabalam, bet objekts o_i ir tuvāk anomāliju apgabalam. Tas, vai objekts o_l pieder normālu datu kopai un vai objekts o_i ir anomāls, ir jautājums, uz kuriem var atbildēt tikai rūpīgi analizējot šos datu punktus. Objektus, kuru punkti atrodas pelēkajā zonā, sauc par iespējamām anomālijām, bet objektus, kuru punkti atrodas ārpus pelēkās zonas augšējās robežas, par varbūtējām anomālijām [Tukey J., 1977]. Acīmredzot galveno problēmu rada objekti, kas iespējams ir anomālija.

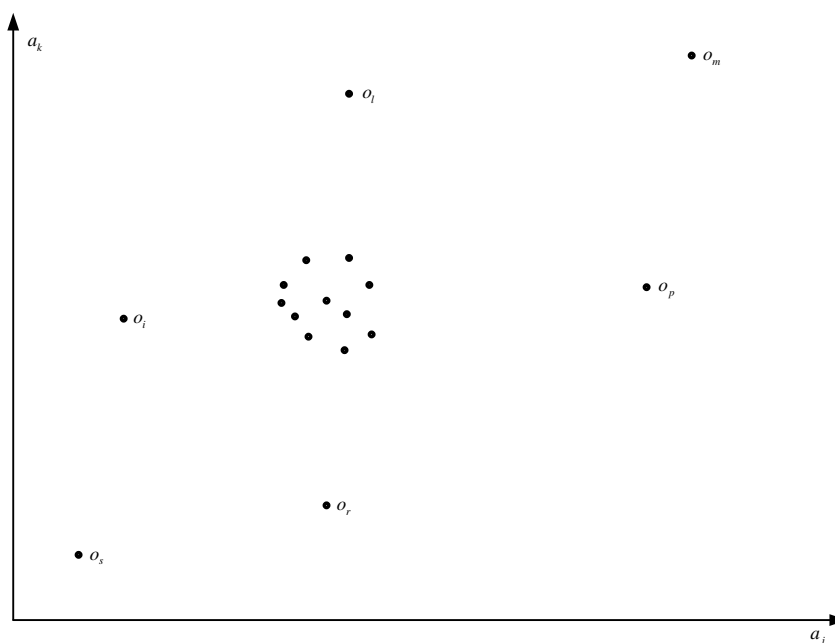
Apskatīsim vēl vienu problēmu, kas saistīta ar anomāliju identificēšanu datus. Attēlā 3.1.2 punkti parāda objektus, kurus raksturo tikai viena atribūta vērtības a_j .



Attēls 3.1.2. Objektu grafiskais attēlojums atribūtu vērtību telpā a_j

Lielākā daļa objektu punktu atrodas ierobežotā atribūtu vērtību diapazonā a_j . Acīmredzot tie ir normāli dati. Objektus o_i un o_l var uzskatīt par potenciālām anomālijām. Objekta o_i atribūta vērtība a_j ir ievērojami mazāka nekā šī atribūta vērtības lielākajai daļai datu, objekta o_l atribūta vērtība a_j ir ievērojami lielāka nekā šī atribūta vērtības lielākajai daļai datu.

Attēlā 3.1.3 punkti apzīmē datu objektus divu atribūtu telpā a_j, a_k .



Attēls 3.1.3. Objektu grafiskais attēlojums atribūtu telpā a_j, a_k

Lielākā daļa datu tiek parādīta skaidri noteiktas datu kopas veidā.

Apskatīsim objektus o_i, o_p . Šie objekti ir potenciālas anomālijas pēc atribūta a_j vērtībām, bet nav anomālija pēc atribūta a_k vērtībām, jo šī atribūta vērtības šiem objektiem atbilst atribūtu vērtībām lielākajai daļai datu. Tāpat objekti o_l, o_r ir potenciālas anomālijas pēc atribūta a_k vērtībām, bet ne pēc atribūta a_j vērtībām. Tikai objektus o_m, o_s var uzskatīt par potenciālām anomālijām, pamatojoties uz abu atribūtu vērtībām.

Situācija var būt vēl sarežģītāka, ja objekti tiek novērtēti pēc atribūtu vērtībām $n > 2$. Šādos apstākļos dažus objektus var uzskatīt par potenciālām anomālijām, pamatojoties uz viena vai vairāku atribūtu vērtībām, bet, pamatojoties uz citu atribūtu vērtībām, tie var atbilst lielākajai daļai datu. Tas noved pie acīmredzama secinājuma: lai praktiski identificētu anomālijas, kas novērtētas pēc noteiktas atribūtu kopas vērtībām, ir nepieciešams identificēt potenciālās anomālijas pēc katra atribūta vērtību kopas.

Vispārīgā gadījumā anomālija var piederēt kādam no šiem veidiem.

1. Globālā anomālija.

Objektu var uzskatīt par globālu anomāliju, ja tas būtiski atšķiras no citiem objektiem datu kopā pēc visām atribūtu vērtībām. Lai identificētu globālo anomāliju, ir jāizmanto piemērota dispersijas metrika.

2. Kontekstuālā anomālija.

Spilgts piemērs atribūtam, kas var izraisīt kontekstuālas anomālijas, ir temperatūra. Pieņem, ka objekti atrodas dažādās apsildāmās telpās un tiek izmērīta šo objektu temperatūra. Ja telpas temperatūra ir robežās no 18°C – 22°C , šī temperatūra atbilst objekta normālajam stāvoklim. Ja temperatūra ir 12°C vai 29°C , tad objekts neatrodas normālos apstākļos un to var identificēt kā kontekstuālu anomāliju.

3. Kolektīvā anomālija.

Ja noteikta objektu apakškopa būtiski atšķiras no objektu lielākās daļas, tad šādus objektus var raksturot kā kolektīvas anomālijas.

Ir trīs pamata pieejas anomāliju identificēšanai [Hodge V.J, Austin J., 2004]:

1. *veids*. Nosaka anomālijas, iepriekš nezinot informāciju par datiem. Būtībā šī ir apmācības pieeja, kas ir analogiska neuzraudzītai klasterizācijai (*unsupervised clusterization*). Algoritms klasterizē datus, nosaka attālākos punktus un atzīmē tos kā potenciālo anomāliju. Šāda veida identifikācija paredz, ka kļūdainie dati ir atdalīti no lielākās daļas datu un tāpēc tiek pieņemti kā anomālijas.

Praksē tiek izmantotas šīs pieejas dažādas variācijas. Diagnostikas pieeja identificē potenciālās anomālijas un izslēdz tās no turpmākās datu apstrādes. Daudzas diagnostikas pieejas iteratīvi izslēdz anomālijas. Process tiek atkārtots, līdz pamatdatos vairs nav anomāliju.

Vēl viena pieeja realizē veidu, kas ģenerētajā datu modelī iekļauj anomālijas un izmanto dažādas metodes šo datu tālākai apstrādei. Šādas metodes var dod iespēju ignorēt anomālijas un nosaka robežu ap lielāko normālo daļu datu.

2. *veids*. Nosaka modeļa “normalitāti”. Šī pieeja ir līdzīga uzraudzītajai klasifikācijai (*supervised classification*) un prasa iepriekšēju normālu un anomālu datu iezīmēšanu. Visi dati ārpus normālo datu apgabala tiek traktēti kā anomāli.

Šo pieeju var izmantot kā klasifikatoru, kur klasifikators tiek apmācīts par klasifikācijas modeli un pēc tam klasificē jaunus piemērus (objektus). Ja jauns piemērs atrodas normalitātes apgabalā, tas tiek klasificēts kā normāls, pretējā gadījumā piemērs tiek atzīmēts kā anomālija.

3. *veids*. Modelē tikai “normalitāti”. Šī pieeja parasti tiek izmantota, lai atpazītu jaunus papildu datus. Tā ir līdzīga daļēji uzraudzītai atpazīšanai (*semi-supervised recognition*). Tiek pieņemts, ka ņemta normālu datu kopa, un algoritms meklē anomālijas, kas varētu interesēt pētnieku.

Lielāko daļu anomāliju identificēšanas metožu var iedalīt šādās lielās klasēs:

- statistiskās metodes;
- metodes, kuru pamatā ir objektu tuvības mērs;
- metodes anomāliju identificēšanai regresijas atkarībās.

Nākamajās sadaļās tiks apspriestas visbiežāk izmantojamās metodes katrā grupā.

3.2. Statistiskās metodes anomāliju identificēšanai

Visās šāda veida metodēs tiek pieņemts, ka atribūtu vērtības tiek sadalītas saskaņā ar noteiktu (visbiežāk normālu) sadalījuma likumu. Tomēr daudzās reālās datu kopās esošie sadalījumi var būt nezināmi vai sarežģīti. Tāpēc, lai izmantotu statistikas metodes, ir rūpīgi jānosaka sadalījuma robežas, kurās atribūtu vērtības var attiecināt uz anomālijām.

Apskatīsim dažas plaši izmantotas statistikas metodes anomāliju identificēšanai.

1. z-iezīmju metode.

Šīs metodes būtība ir tāda, ka atribūtu faktiskās vērtības tiek aizstātas ar tā sauktajām z-iezīmēm. Atribūta a_{ij} vērtībai tās z-iezīmi aprēķina šādi:

$$z_{ij} = \frac{a_{ij} - \bar{a}_j}{s_j}, \quad i = 1, \dots, m, \quad j = 1, \dots, n, \quad (3.2.1)$$

kur \bar{a}_j - j-tā atribūta vidējā vērtība;

s_j - j-tā atribūta standartnovirzes vērtība.

Atribūtu vērtību pārveidošana atbilstošās z-iezīmēs ļauj novērtēt attālumu no katras atribūta vērtības līdz vidējai atribūtu vērtībai. Pastāv problēma ar iezīmes robežvērtības z_{jg} definēšanu, kas atdala normālo atribūtu vērtības un šīm vērtībām atbilstošu troksni no atribūtu vērtībām, kas atbilst anomālīgām. Ja atribūtu a_j vērtības tiek sadalītas saskaņā ar normālo sadalījuma likumu, tad atribūtu iezīmes $|z_{ij}| \geq 3$ tiek uzskatītas par anomālīgām. Citiem atribūtu vērtību sadalījumiem var iestatīt citas robežvērtības, piemēram, $|z_{ij}| \geq 2.5$ vai $|z_{ij}| \geq 2.7$, un tamlīdzīgi.

Apskatīsim vienkāršu ilustratīvu piemēru. Tabulā 3.2.1 ir parādītas divu atribūtu vērtības: a_1 un a_2 . Tabulā atribūtu vērtības uzdotas transponētā veidā. Mūsu mērķis ir pārbaudīt tabulas datus attiecībā uz anomāliju klātbūtni. Šim nolūkam mēs pārveidojam atribūtu vērtības atbilstošās z-iezīmēs. Aprēķinu rezultāti parādīti tabulās 3.2.2 un 3.2.3.

Tabula 3.2.1

Sākotnējās atribūtu a_1, a_2 vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}
a_1	7	6	8	6	20	7	7	6	5	8
a_2	9	11	10	11	9	8	10	30	12	10

Tabula 3.2.2

Atribūta a_1 vērtību transformācijas z-iezīmēs rezultāti datiem no tabulas 3.2.1

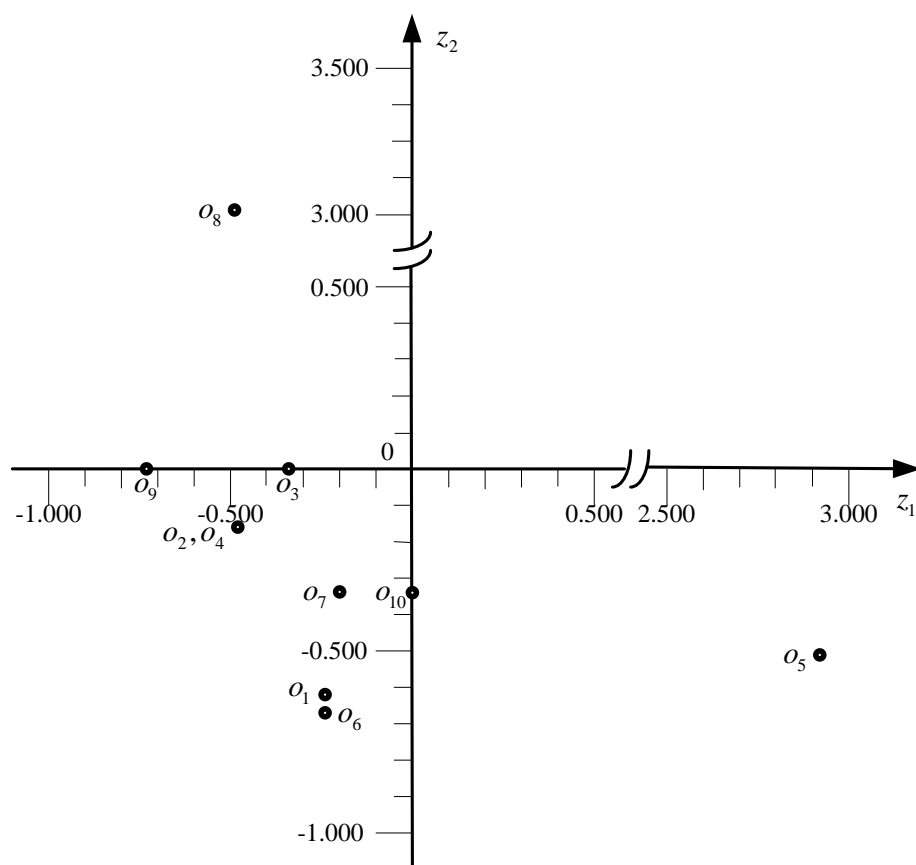
Objekti	a_{i1}	$a_{i1} - \bar{a}_1$	$(a_{i1} - \bar{a}_1)^2$	z_{i1}
o_1	7	-1	1	-0.244
o_2	6	-2	4	-0.488
o_3	8	0	0	0
o_4	6	-2	4	-0.488
o_5	20	12	144	2.930
o_6	7	-1	1	-0.244
o_7	7	-1	1	-0.244
o_8	6	-2	4	-0.488
o_9	5	-3	9	-0.732
o_{10}	8	0	0	0
	$\sum = 80$		$\sum = 168$	$\max z_{i1} = 2.930$
	$\bar{a}_1 = 8.000$		$s_1 \sqrt{\frac{168}{10}} = 4.100$	

Tabula 3.2.3

Atribūta a_2 vērtību transformācijas z -iezīmēs rezultāti datiem no tabulas 3.2.1

Objekti	a_{i2}	$a_{i2} - \bar{a}_2$	$(a_{i2} - \bar{a}_2)^2$	z_{i2}
o_1	9	-3	9	-0.512
o_2	11	-1	1	-0.171
o_3	10	-2	4	-0.342
o_4	11	-1	1	-0.171
o_5	9	-3	4	-0.512
o_6	8	-4	16	-0.683
o_7	10	-2	4	-0.342
o_8	30	18	324	3.073
o_9	12	0	0	0
o_{10}	10	-2	4	-0.342
	$\sum = 120$		$\sum = 343$	$\max z_{i2} = 3.073$
	$\bar{a}_2 = 12.000$		$s_2 = \sqrt{\frac{343}{10}} = 5.857$	

Uzskatāmības labad attēlā 3.2.1 objektu punkti $o_1 - o_{10}$ tiek grafiski attēloti z_1, z_2 iezīmju telpā.

Attēls 3.2.1. Objektu punktu grafiskais attēlojums datiem no tabulas 3.2.1 z_1, z_2 iezīmju telpā

Ja mēs ņemam $z_{1g} = z_{2g} = 3.000$ iezīmju vērtības kā robežvērtības, tad objektu a_1 ar iezīmes vērtību $z_{51} = 2.930$, kas tikai nedaudz atšķiras no robežvērtības, var traktēt kā anomāliju pēc atribūta o_5 vērtības. Atbilstoši iezīmes z_2 vērtībai anomālija noteikti ir objekts o_8 , kuram $z_{82} = 3.073 > 3.000$.

Būtisks šīs metodes trūkums ir novērtējumu \bar{a}_j un s_j atkarība no anomālijām datos. Jāatzīmē, ka šis trūkums ir raksturīgs visām pieejām, kur vidējās vērtības tiek izmantotas kā sadalījuma parametri un standarta novirzes tiek izmantotas kā diapazona parametri.

Šajā sakarā tiek piedāvāta šāda modificētu z -iezīmju metode.

2. Modificēto z -iezīmju metode.

Šo metodi piedāvāja [Iglowitz B., Hoaglin D.C., 1993]. Atribūtu vērtību iezīmes tiek aprēķinātas, izmantojot vienādojumu

$$z_{mij} = \frac{0.6745(a_{ij} - med(a_j))}{mad(a_j)}, \quad i = 1, \dots, m, \quad j = 1, \dots, n, \quad (3.2.2)$$

kur $med(a_j)$ - j -tā atribūta mediānas vērtība;

$$mad(a_j) = med\{|a_{ij} - med(a_j)|\} - j\text{-tā atribūta noviržu mediāna no vērtību mediānas.}$$

Vērtības $med(a_j)$ un $mad(a_j)$ neietekmē anomāliju klātbūtne datos.

Ja uzdota vērtību $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$ kopa, tad, lai aprēķinātu mediānu $med(X)$, vērtības $x_i \in X$ jāsakārto nedilstošā secībā

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n-1)} \leq x_{(n)},$$

kur apakšindeksi iekavās apzīmē vērtību $x_{(i)}$ numurus sakārtotajā secībā.

Mediānas vērtību aprēķina šādi:

$$med(\mathbf{X}) = x_{(n+1)/2}, \quad n - \text{nepāra}; \quad (3.2.3.a)$$

$$med(\mathbf{X}) = \frac{1}{2}(x_{(k)} + x_{(k+1)}), \quad n = 2k - \text{pāra}. \quad (3.2.3.b)$$

Vērtību $z_{mijg} > 3.5$ ieteicams izmantot kā robežvērtību starp normālo datu kopu un anomāliju kopu.

Apskatīsim vienkāršu ilustratīvu piemēru. Ņemsim par pamatu atribūtu a_1 un a_2 sākotnējās vērtības no tabulas 3.1.1. Aprēķināsim katra atribūta vērtību mediānu.

Sakārtosim atribūta a_1 vērtības nedilstošā secībā:

$$\{a_{i1}\} = 5, 6, 6, 6, 7, 7, 7, 8, 8, 20.$$

Tā kā šajā sarakstā ir 10 vērtības ($n=10$), mēs aprēķinām mediānas vērtību, izmantojot vienādojumu (3.2.3.b).

$$med(a_1) = \frac{1}{2}(7 + 7) = 7.$$

Sakārtosim atribūta a_{i2} vērtības nedilstošā secībā:

$$\{a_{i2}\} = 8, 9, 9, 10, 10, 10, 11, 11, 12, 30.$$

$$\text{med}(a_2) = \frac{1}{2}(10+10) = 10.$$

Aprēķināsim absolūtās novirzes no atribūtu mediānas. Aprēķinu rezultāti parādīti tabulā 3.2.4.

Tabula 3.2.4

Atribūtu a_1, a_2 vērtību absolūtās novirzes no mediānas $\text{med}(a_1), \text{med}(a_2)$

Objekti	a_{i1}	$ a_{i1} - \text{med}(a_1) $	a_{i2}	$ a_{i2} - \text{med}(a_2) $
o_1	7	0	9	1
o_2	6	1	11	1
o_3	8	1	10	0
o_4	6	1	11	1
o_5	20	13	9	1
o_6	7	0	8	2
o_7	7	0	10	0
o_8	6	1	30	20
o_9	5	2	12	2
o_{10}	8	1	10	0

Sakārtosim noviržu absolūtās vērtības no mediānas $\text{med}(a_1)$ nedilstošā secībā:

$$\{|a_{i1} - \text{med}(a_1)|\} = 0, 0, 0, 1, 1, 1, 1, 1, 2, 13.$$

Aprēķināsim vērtību : $\text{mad}(a_1)$.

$$\text{mad}(a_1) = \frac{1}{2}(1+1) = 1.$$

Sakārtosim noviržu absolūtās vērtības no mediānas $\text{med}(a_2)$ nedilstošā secībā:

$$\{|a_{i2} - \text{med}(a_2)|\} = 0, 0, 0, 1, 1, 1, 1, 2, 2, 20.$$

Aprēķināsim vērtību $\text{mad}(a_2)$:

$$\text{mad}(a_2) = \frac{1}{2}(1+1) = 1.$$

Ar vienādojuma (3.2.2) palīdzību aprēķināsim modificēto iezīmju z_{mi1}, z_{mi2} vērtības.

Aprēķinu rezultāti parādīti tabulā 3.2.5.

Tabula 3.2.5

Modificēto iezīmju z_{mi1}, z_{mi2} aprēķinātās vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}
z_{mi1}	0	0.6745	0.6745	0.6745	8.7685	0	0	0.6745	1.3490	0.6745
z_{mi2}	0.6745	0.645	0	0.6745	0.6745	1.3490	0	13.4900	1.3490	0

Saskaņā ar tabulu 3.2.5 mēs varam izdarīt pamatotu secinājumu, ka objekts o_5 ir anomālija pēc atribūta a_1 vērtības un objekts o_8 ir anomālija pēc atribūta a_2 vērtības.

Acīmredzami, ka anomāliju identificēšana datos, pamatojoties uz modificētām z -iezīmēm šķiet saprātīgāka un pareizāka nekā identificēšana pamatojoties uz z -iezīmēm. Šīs metodes galvenā priekšrocība ir tā, ka modificēto iezīmju aprēķinu rezultātus neietekmē ekstremālas atribūtu vērtības, kas var gadīties z -iezīmju gadījumā.

3. Metode, kas balstīta uz starpkvartiļu rangiem.

Lai izskaidrotu šīs metodes būtību, aplūkosim dažus statistikas jēdzienus. Ņemsim par pamatu sakārtoto atribūtu a_i vērtības no tabulas 3.2.1. Šīs vērtības sauc par kārtas statistikām.

Šeit ir kopsavilkums par dažiem svarīgiem jēdzieniem darbā ar kārtas statistikām.

Rangs augšup – tas ir i -tās kārtas statistikas vērtības skaitlis, ja numerācija sākas no mazākās vērtības.

A	Kārtas statistikas numurs	1	2	3	4	5	6	7	8	9	10
B	Kārtas statistikas	5	6	6	6	7	7	7	8	8	20
C	Rangs augšup	1	2	3	4	5	6	7	8	9	10
D	Rangs lejup	10	9	8	7	6	5	4	3	2	1
E	Dziļums	1	2	3	4	5	5	4	3	2	1

Ranga augšup vērtības parādītas C rindā.

Rangs lejup – tas ir i -tās kārtas statistikas vērtības skaitlis, ja numerācija sākas no lielākās vērtības (D rinda).

Apzīmēsim ar $r \uparrow (a_{i1})$ kārtas statistikas a_{i1} rangu augšup un ar $r \downarrow (a_{i1})$ – šīs kārtas statistikas rangu lejup. Acīmredzot jebkurai kārtas statistikai a_{i1} ir spēkā attiecība:

$$r \uparrow (a_{i1}) + r \downarrow (a_{i1}) = n + 1.$$

kur n ir kopējais atribūtu a_i vērtību skaits.

Ja mēs saskaitām divas vērtības no katra vērtību pāra rindās C un D, mēs iegūstam vērtību 11, kas apstiprina iepriekš minēto attiecību.

Sakārtoto vērtību mediānu aprēķina pēc vienādojumiem (3.2.3.a,b). Mediānu var definēt kā kārtas statistiku, kurai ir vislielākais dziļums. Iepriekš dotajos datos kārtas statistikām 5 un 6 vislielākais dziļums ir 5. Tāpēc mediānas vērtība tiek aprēķināta kā šīs kārtas statistikas vidējā vērtība.

$$med(a_i) = \frac{1}{2}(7 + 7) = 7.$$

Formāli definēsim sakārtotas atribūtu a_i vērtību kvartiles jēdzienu. Kvartile ir kārtas statistika, kuras dziļumu nosaka vienādojums

$$g(\text{quartile}) = (\lfloor g(\text{median}) \rfloor + 1) / 2. \quad (3.2.4)$$

Šajā vienādojumā $g(\text{median})$ – mediānas dziļums. Apzīmējums $\lfloor x \rfloor$ nozīmē lielāko veselo skaitli, kas nepārsniedz x . Piemēram, $\lfloor 4.6 \rfloor = 4$.

Aprēķināsim kvartiles dziļumu sakārtotām atribūtu a_i vērtībām, izmantojot vienādojumu (3.2.4).

$$g(\text{quartile}) = (5 + 1) / 2 = 3.$$

Tas nozīmē, ka apakšējās kvartiles Q_1 vērtība ir kārtas statistika dziļumā 3, t.i., $Q_1 = 6$. Augšējās kvartiles vērtība ir kārtas statistika dziļumā 3, t.i., $Q_3 = 8$.

Starpkvartiļu rangs ir vienāds ar starpību

$$IQR = Q_3 - Q_1 = 8 - 6 = 2.$$

Noteiksim kvartiles un starpkvartiļu rangus sakārtotām atribūtu a_2 vērtībām no tabulas 3.2.1.

A	Kārtas statistikas numurs	1	2	3	4	5	6	7	8	9	10
B	Kārtas statistikas	8	9	9	10	10	10	11	11	12	30
C	Rangs augšup	1	2	3	4	5	6	7	8	9	10
D	Rangs lejup	10	9	8	7	6	5	4	3	2	1
E	Dziļums	1	2	3	4	5	5	4	3	2	1

Kārtas statistikām 5 un 6 ir vienāds dziļums, kas vienāds ar 5. Tāpēc

$$\text{med}(a_2) = \frac{1}{2}(10 + 10) = 10.$$

Kvartiļu dziļumi ir vienādi ar

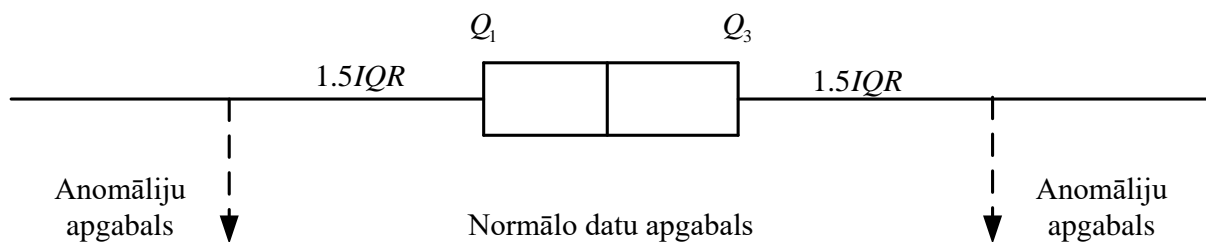
$$g(\text{quartile}) = (5 + 1) / 2 = 3.$$

Iegūstam $Q_1 = 9$, $Q_2 = 11$.

Starpkvartiļu rangs ir

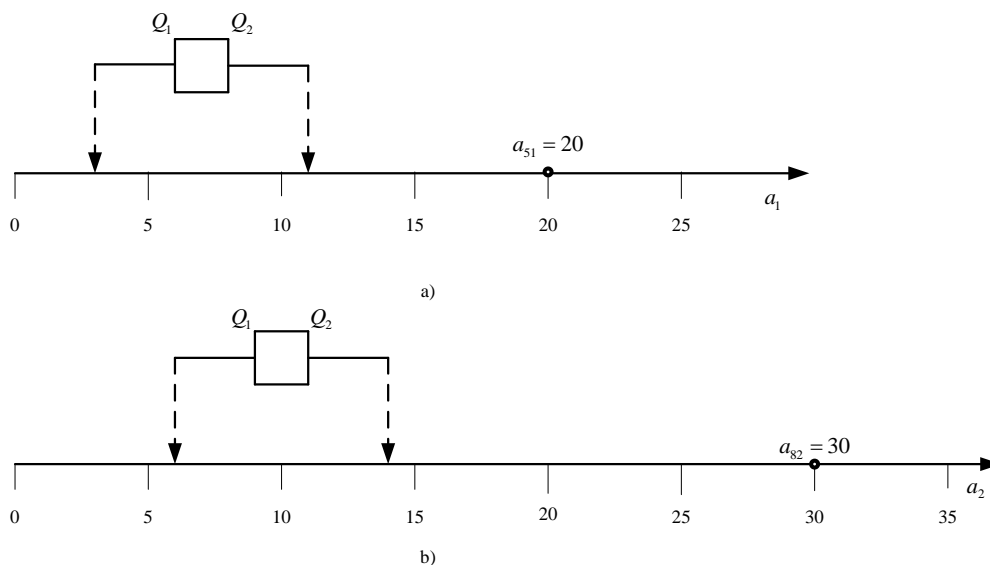
$$IQR = Q_3 - Q_1 = 11 - 9 = 2.$$

Kā starpkvartiļu rangus var izmantot anomāliju identificēšanai datos? Robežas starp normāliem datiem un anomālijām nosaka atņemot vērtību $1.5IQR$ no apakšējās kvartiles Q_1 vērtības un pievienojot vērtību $1.5IQR$ augšējai kvartilei Q_3 . Uz šīs metodes balstītā anomāliju identifikācijas shēma ir parādīta attēlā 3.2.2.



Attēls 3.2.2. Shēma robežu noteikšanai starp normāliem datiem un anomālijām izmantojot starpkvartiļu rangus

Noteiksim normālo atribūtu a_1, a_2 vērtību apgabalu robežas datiem no tabulas 3.2.1 (sk. attēlu 3.2.3).



Attēls 3.2.3. Normālo atribūtu a_1, a_2 robežu grafiskais attēlojums

Normālām atribūtu a_1 vērtībām jāatrodas intervālā [3,11]. Atribūta $a_{51} = 20$ vērtība ir acīmredzama anomālija. Normālām atribūtu vērtībām a_2 jāatrodas intervālā [7,14]. Vērtība $a_{82} = 30$ ir acīmredzama anomālija.

Šīs metodes priekšrocība ir tās spēcīgā ietekmē uz atribūtu ekstrēmajām vērtībām. Mediānas un kvartiles vērtības tiek noteiktas pēc atribūta atrašanās vietas sakārtotā atribūtu secībā, tāpēc šī atribūta ekstremālās vērtības neietekmē šo parametru noteikšanu.

3.3. Anomāliju identificēšanas metodes uz tuvības mēra starp objektiem pamata

Statistiskās metodes anomāliju identificēšanai, kas aplūkotas 3.2. sadaļā, pieņem, ka atribūtu vērtības tiek sadalītas saskaņā ar noteiktiem (visbiežāk normāliem) sadalījuma likumiem. Statistikas metožu būtība ir novērtēt atribūtu vērtību sadalījumu parametrus. Anomāliju identifikācija tiek veikta pamatojoties uz robežu noteikšanu starp normāliem datiem un anomālijām. Šīs robežas tiek noteiktas atkarībā no datu izvietojuma parametriem.

Reālās datu pirmapstrādes un analīzes problēmās atribūtu vērtību sadalījumi un šo sadalījumu parametri ir zināmi tikai diezgan retos gadījumos. Tas nozīmē, ka ir jāizmanto metodes, lai identificētu anomālijas, kas strādā ar visu datu kopumu, neņemot vērā faktisko atribūtu vērtību sadalījumu. Tā kā šie sadalījumi nav zināmi, tad nav zināms, kā to parametri tiek izmantoti par atskaites punktiem attāluma līdz objektiem novērtēšanai. Ņemot to vērā, ir nepieciešami citi principi tuvības mēru noteikšanai starp objektiem, lai šos tuvības mērus varētu izmantot kā pamatu anomāliju identificēšanai.

Ir piedāvātas daudzas metodes tuvības noteikšanai objektu kopai. Labi pazīstama šāda veida metode ir DB(p, D) - attāluma metode, kas tika piedāvāta darbā [Knorr E.M., Ng R.T., 1998, 1999]. Tā ir balstīta uz attāluma mērīšanu starp visiem objektiem sākotnējā datu kopā. Pamatojoties uz konkrētu šo attālumu analīzi, datus var identificēt anomālijas.

Vēl viena izplatīta metode ir k-tuvāko (*k-nearest*) kaimiņu metode, kuras pamatā ir attālumi starp objektiem.

Anomāliju identificēšanas metodes, kuru pamatā ir objektu grupēšana datu kopā, ir kļūvušas ļoti izplatītas. Klasterizācija attiecas uz objektu grupēšanu klasteros, pamatojoties uz šo objektu līdzību. Jāņem vērā viens svarīgs apstāklis: grupējot mēs runājam par objektiem, kas nav iezīmēti jeb marķēti. Pretējā gadījumā visi objekti nav atšķirami tādā nozīmē, ka tie nepieder noteiktām klasēm vai kaut kādā veidā atšķiras.

Jebkuras klasterizācijas metodes galvenā īpašība ir spēja apmācīties, izmantojot sākotnējo datu kopu. Šajā ziņā klasterizāciju mašīnmācības jomā sauc par neuzraudzīto apmācību.

Lai novērtētu objektu līdzības pakāpi, atribūtu vērtību kopai ir nepieciešama piemērota metrika. Parasti šim nolūkam izmanto Eiklīda attālumu jeb distanci. Dažās īpašās situācijās tiek izmantots Mahalanobisa attālums. (Dažādu attālumu metrikas datu telpā sniegta pielikumā P2).

Ideja par anomāliju identificēšanu, pamatojoties uz objektu klasterizāciju, ir šāda. Vispārīgā gadījumā starp atsevišķiem objektiem pastāv lielākas vai mazākas līdzības. Uz tā pamata objektus var apvienot atsevišķās datu kopās. No otras puses, potenciālās anomālijas var būt atsevišķi objekti, kas atrodas lielos attālos no visām izveidotajām kopām vai arī anomālijas veido mazas kopas, kas atrodas attālu no lielajiem klasteriem. Klasterizācija ļauj strukturēt sākotnējos datus tā, lai lielas normālu objektu kopas krasi atšķirtos no atsevišķām anomālijām vai ļoti mazām anomāliju kopām. Šī normālu datu klasteru nošķiršana no potenciālajām anomālijām ir pamats šo anomāliju identificēšanai.

Šajā sadaļā tiks parādītas divas galvenās datu klasterizācijas metodes anomāliju identificēšanai:

- k-vidējo metode (*k-means*);
- k-medoīdu metode (*k-medoids*).

Apskatīsim katru no šīm metodēm.

1. K-vidējo metode.

K-vidējo metode ir paredzēta tikai skaitliskām atribūtu vērtībām. Tā kā atribūtu vērtībām ir dažādas mērvienības, lai piemērotu k-vidējā metodi, šīs vērtības ir jāpārveido bezdimensiju formā. Parasti pārveidošana notiek pārvēršot sākotnējās a_{ij} atribūtu vērtības par z -iezīmēm z_{ij} pēc vienādojuma

$$z_{ij} = \frac{a_{ij} - \bar{a}_j}{s_j} \quad . \quad i = 1, \dots, m \quad . \quad j = 1, \dots, n \quad . \quad (3.3.1)$$

kur \bar{a}_j - atribūta a_j vidējā vērtība;

s_j - atribūta a_j vērtības standarta novirze.

(Atzīmēsim, ka mēs jau izmantojām sākotnējo atribūtu vērtību z transformāciju sadaļā 3.2 (vienādojums (3.2.1))).

Šajā sadaļā objektu klasterizācijas, kuras pamatā ir k-vidējie, teorētiskie pamati ir izklāstīti darbā [Bishop Ch., 2006]. Detalizētu informāciju par dažādām klasterizācijas metodēm var atrast darbos [Wu J., 2012, Alekseeva L., Borisov A., Uzhga-Rebrov O., 2015].

Pieņemsim, ka mums ir m objekti, piemēri vai citas entītijas (datu vienības), kas novērtētas pēc n atribūtu vērtībām. Citiem vārdiem sakot, datu punkti atrodas n -dimensiju atribūtu vērtību telpā.

Pirmā k-means algoritma procedūra ir klasteru skaita noteikšana. Šo uzdevumu var atrisināt, pamatojoties uz iepriekšēju datu analīzi vai izmantojot šādas palīgmetodes:

- elkoņa metode (*elbow method*) [Torndike R.L., 1953].
- silueta metode (*silhouette method*) [Rousseeuw P.J., 1987].

Pieņemsim, ka ir noteikts klasteru skaits k . Mēs varam izvēles veidā katram klasterim piešķirt tā centrālo vērtību (centroīdu) μ_k . $k = 1, \dots, k$. Dažreiz vērtību μ_k sauc par prototipu, kas saistīts ar k -to klasteri, $k = 1, \dots, k$. Var pieņemt, ka vektori μ_k ir atbilstošo klasteru centri. Tad klasterizācijas problēmu, kas balstīta uz k -vidējā metodi, var formulēt šādi: jānosaka vektoru μ_k kopa un visus datu objektus (punktus) piešķirt klasteriem tā, lai attālumu kvadrātu summa no jebkura objekta k -klasterī līdz tā centram μ_k būtu minimāla.

Lai formāli attēlotu objekta o_i saistību ar k -to klasteri, mēs ieviešam mainīgo $r_{ik} \in \{0, 1\} \dots$ $i \in \{1, \dots, m\}$, $k \in \{1, \dots, k\}$ kā bināro indikatoru kopu. Ja objekts o_i pieder klasterim k , tad $r_{ik} = 1$ un $r_{ik} = 0$ visām pārējām k vērtībām. Pamatojoties uz šo formalizāciju, var definēt mērķa funkciju sekojošā veidā

$$J = \sum_{i=1}^m \sum_{k=1}^k r_{ik} \|\mathbf{a}_i - \mu_k\|^2 \quad . \quad (3.3.2)$$

kur \mathbf{a}_i - atribūtu vērtību vektors i -tajam objektam;

μ_k - k -tā klastera centrs;

r_{ik} - iepriekš definētais binārais indikatora mainīgais.

Vienādojums (3.3.2) attēlo attālumu kvadrātu summu no objektiem, kas pieder k -tajam klasterim līdz šī klastera centram μ_k .

Mūsu uzdevums ir atrast r_{ik} vērtības $i = 1, \dots, m$, $k = 1, \dots, k$ un $\{\mu_k\}$ tā, lai samazinātu mērķa funkciju (3.3.2). To var izdarīt ar iteratīvas procedūras palīdzību, kurā katra iterācija ietver divus secīgus soļus, kas atbilst r_{ik} un μ_k vērtību secīgai optimizācijai.

Sākotnēji tiek atlasītas μ_k , $k=1, \dots, k$ vērtības. Pēc tam mērķa funkcija J tiek minimizēta attiecībā pret r_{ik} vērtībām, pie tam μ_k vērtības ir fiksētas. Pēc tam μ_k vērtības tiek pārrēķinātas, pamatojoties uz klasterim k attiecināmo objektu atribūtu vērtībām.

Tad abas darbības tiek atkārtotas, līdz rezultāti konverģē (μ_k vērtības ir vienādas ar iepriekšējā iterācijā iegūtajām vērtībām).

Var apgalvot, ka vērtību r_{ik} izmaiņas atbilst paredzamās vērtības maksimizēšanas algoritma solim E . (Šis algoritms tika parādīts 2.4. sadaļā iztrūkstošo atribūtu vērtību imputācijas kontekstā).

Vērtību r_{ik} optimizēšana nesastāda būtiskas problēmas. Tā kā mērķa funkcijas J vērtība (vienādojums (3.3.2)) ir lineāri atkarīga no r_{ik} vērtībām un r_{ik} vērtības pie dažādiem i ir neatkarīgas, mēs varam veikt optimizāciju katrai atsevišķai i vērtībai, izvēloties vērtību $r_{ik}=1$ jebkuram k , kas dod minimālo vērtību $\|\mathbf{a}_i - \mu_k\|^2$. Faktiski mēs vienkārši attiecinām objektu o_i tuvākajam klasterim. Formāli šo apgalvojumu var izteikt kā

$$\begin{aligned} & 1. \text{ Ja } k = \operatorname{argmin} \|\mathbf{a}_i - \mu_k\|^2; \\ & 0. \text{ citādi.} \end{aligned} \quad (3.3.3)$$

Tagad ieviesīsim μ_k optimizācijas principu pie fiksētām r_{ik} vērtībām. Mērķa funkcija J ir kvadrātiskā funkcija no μ_k vērtībām un to var minimizēt, iestatot tās pirmo atvasinājumu vienādu ar nulli attiecībā pret μ_k :

$$2 \sum_{i=1}^m r_{ik} (\mathbf{a}_i - \mu_k) = 0.$$

Atrisinot šo vienādojumu attiecībā pret μ_k iegūstam

$$\mu_k = \frac{\sum_i r_{ik} \mathbf{a}_i}{\sum_i r_{ik}}. \quad k=1, \dots, k. \quad (3.3.4)$$

Šis izteiksmes saucējs parāda klasterim attiecināmo datu punktu (objektu) skaitu k . Var secināt, ka vērtība μ_k , $k=1, \dots, k$ katrā iterācijā tiek aprēķināta kā klasterim attiecināto objektu atribūtu vektoru vidējā vērtība k . Līdz ar to arī ir šāds metodes nosaukums: klasterizācija uz k -vidējo pamata.

Dodam k -vidējā algoritma procedūru sarakstu:

1. Jānosaka klasteru skaits k .
2. Jāatlasa k punktus atribūtu vērtību telpā un jāuzdod tos kā klasteru sākotnējos centroīdus.
3. Jāizskaitļo attālumus no katra datu punkta (objekta) līdz katram no dotajiem centroīdiem, izmantojot Eiklīda attālumu.
4. Jāattiecina katru objektu klasterim, kuram Eiklīda attālums starp objektu un dotā klastera centroīdu ir minimāls.
5. Pārrēķina jaunās centroīdu koordinātes un izdara 4. punktā dotās darbības.
6. Jāatkārto 3.–5. darbības, līdz centroīdu koordinātes vairs nemainās.

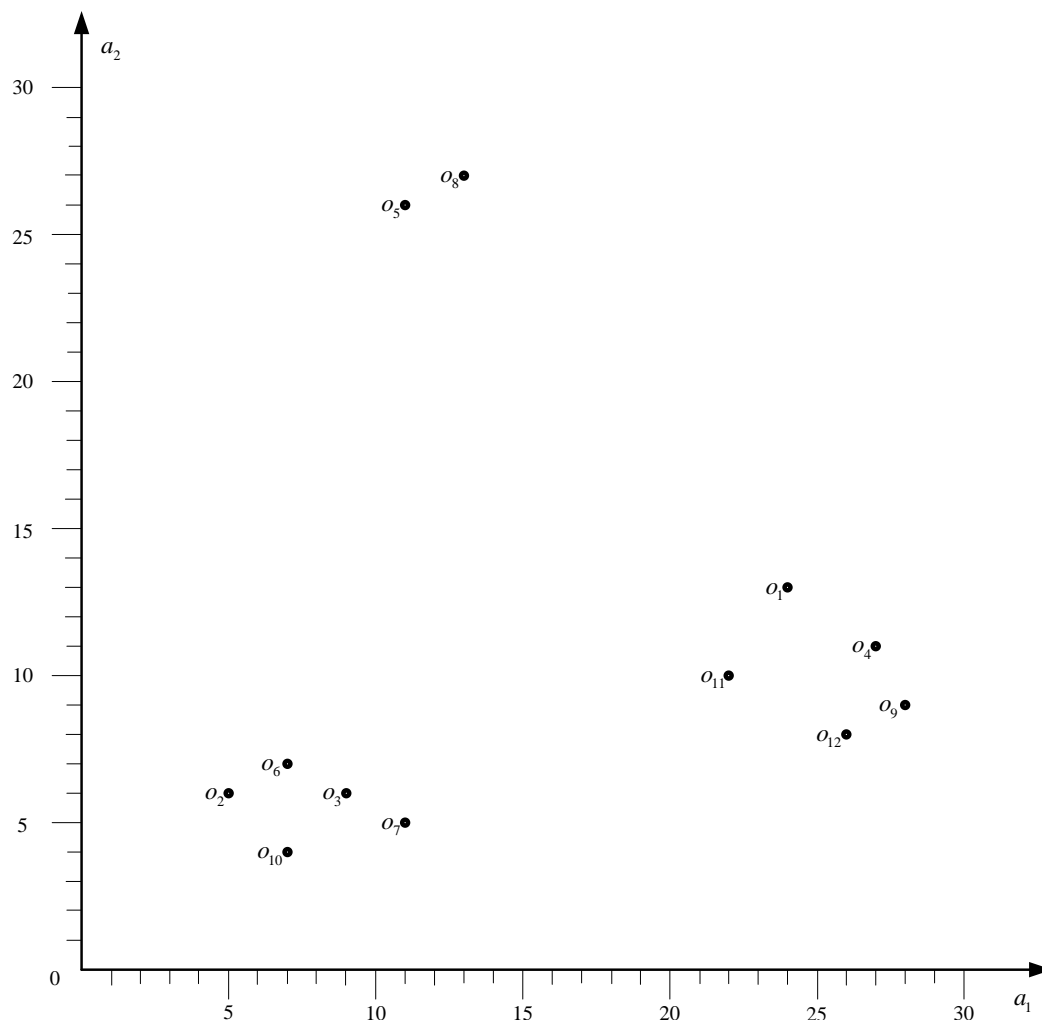
Apskatīsim vienkāršu ilustratīvu piemēru. Objektus $o_1 - o_{12}$ raksturo skaitlisko atribūtu vērtības a_1 , a_2 (sk. tabulu 3.3.1).

Tabula 3.3.1

Objektu $o_1 - o_{12}$ sākotnējās atribūtu a_1, a_2 vērtības

	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	o_{11}	o_{12}
a_1	24	5	9	27	11	7	11	13	28	7	22	26
a_2	13	6	6	11	26	7	5	27	9	4	10	8

Uzskatāmības labad dati no tabulas 3.3.1 grafiski parādīti attēlā 3.3.1.



Attēls 3.3.1. Sākotnējo datu no tabulas 3.3.1 grafiskais attēlojums

Mūsu mērķis ir identificēt anomālijas sākotnējā datu kopā ar k -vidējā klasterizācijas algoritma palīdzību.

Pārveidojam sākotnējās atribūtu vērtības z -iezīmju formā. Aprēķināsim atribūtu a_1 un a_2 vidējās vērtības un standartnovirzes. Iegūstam

$$\bar{a}_1 = 15.83. \quad s_1 = 8.44;$$

$$\bar{a}_2 = 11.00. \quad s_2 = 7.36.$$

Izmantojot vienādojumu (3.3.1) mēs pārveidojam atribūtu a_1, a_2 sākotnējās vērtības. Pārveidošanas rezultāti ir parādīti tabulā 3.3.2.

Tabula 3.3.2

Pārveidotās atribūtu vērtības datiem no tabulas 3.3.1

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	o_{11}	o_{12}
z_1	0.97	-1.28	-0.81	1.32	-0.57	-1.05	-0.57	-0.33	1.44	-1.05	0.73	1.20
z_2	0.27	-0.68	-0.68	0	2.04	-0.54	-0.82	2.17	-0.27	-0.95	-0.14	-0.41

Uzskatāmības labad pārveidotās atribūtu z_1, z_2 vērtības ir grafiski parādītas attēlā 3.3.2. Pamatojoties uz pārveidoto datu struktūru mēs iestatām vērtību $k = 3$.

Uzdosim klasteru centroīdu sākotnējās vērtības. Tā kā šajā piemērā datu struktūra ir tāda, ka objekti veido trīs grupas novērtēšanas telpā z_1, z_2 , mēs izvēlamies centroīdu sākotnējās koordinātes aptuveni katras grupas centrā.

$$c_1^{(0)} : (-0.90, -0.80);$$

$$c_2^{(0)} : (1.10, -0.20);$$

$$c_3^{(0)} : (-0.50, 2.10).$$

Attēlā 3.3.2 sākotnējie centroīdi ir parādīti ar krustiņiem.

Aprēķināsim Eiklīda attālumus no katra objekta līdz katram centroīdam. Aprēķinu rezultāti ir parādīti tabulā 3.3.3.

Tabula 3.3.3

Eiklīda attāluma vērtības no objektiem līdz sākotnējiem klasteru centroīdiem

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	o_{11}	o_{12}
c_1	4.67	0.15	0.02	5.57	1.65	0.09	0.11	9.14	5.75	0.04	3.10	4.55
c_2	0.24	5.89	3.88	0.09	7.81	4.74	3.17	7.66	0.13	5.28	0.17	0.05
c_3	5.51	2.63	7.83	7.72	0.01	7.27	8.53	0.03	9.38	9.60	6.53	9.19
Klasteris priekš objekta	c_2	c_1	c_1	c_2	c_3	c_1	c_1	c_3	c_2	c_1	c_2	c_2

Rezultātā mums ir šādas kopas:

$$c_1 : \{o_2, o_3, o_6, o_7, o_{10}\};$$

$$c_2 : \{o_1, o_4, o_9, o_{11}, o_{12}\};$$

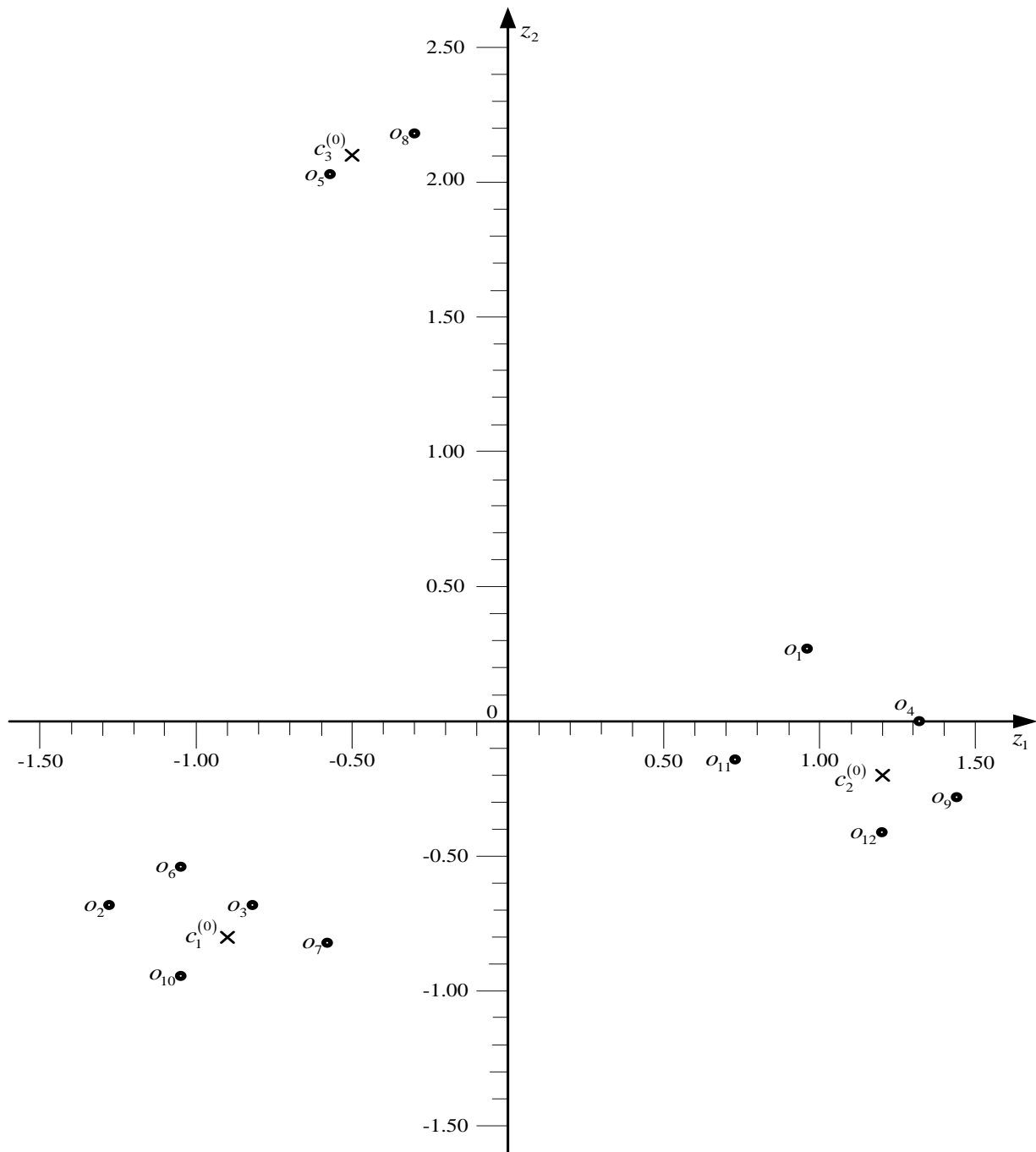
$$c_3 : \{o_5, o_8\}.$$

Kā arī bija sagaidāms, klasteros nokļuva objekti, kas tika sagrupēti attiecībā pret uzdotajiem centroīdiem. Šajā gadījumā tas izskaidrojams ar vienkāršo datu struktūru šajā piemērā. Reālās datu priekšapstrādes problēmās objekti var pārvietoties no viena klastera citā, kad centroīdu koordinātes izmainās katrā algoritma iterācijā.

Atkārtoti pārvērtēsim klasteru centroīdu koordinātes kā atbilstošajos klasteros iekļauto objektu atribūtu transformēto vērtību vidējās vērtības:

- klasteris c_1 :

$$\bar{z}_1(c_1^{(1)}) = \frac{-1.28 + (-0.81) + (-1.05) + (-0.57) + (-1.05)}{5} = \frac{-4.76}{5} = -0.95;$$



Attēls 3.3.2. Datu kopas grafiskais attēlojums pārveidoto iezīmju telpā $z_1 \cdot z_2$

$$\bar{z}_2(c_1^{(1)}) = \frac{-0.68 + (-0.68) + (-0.54) + (-0.82) + (-0.95)}{5} = \frac{-3.67}{5} = -0.73.$$

Mums ir $c_1^{(1)} = (-0.95, -0.73)$.

- klasteris c_2 :

$$\bar{z}_1(c_2^{(1)}) = \frac{0.97 + 1.32 + 1.44 + 0.73 + 1.20}{5} = \frac{5.66}{5} = 1.13;$$

$$\bar{z}_2(c_2^{(1)}) = \frac{0.27 + 0 + (-0.27) + (-0.14) + (-0.41)}{5} = \frac{-0.55}{5} = -0.11.$$

Mums ir $c_2^{(1)} = (1.13, -0.11)$.

- klasteris c_3 :

$$\bar{z}_1(c_3^{(1)}) = \frac{-0.57 + (-0.33)}{2} = \frac{-0.90}{2} = -0.45;$$

$$\bar{z}_2(c_3^{(1)}) = \frac{2.04 + 2.17}{2} = \frac{4.21}{2} = 2.10.$$

Mums ir $c_3^{(1)} = (-0.45, 2.10)$.

Aprēķinātās jaunās $c_k^{(1)}$ centroīdu koordinātes maz atšķiras no sākotnējām $c_k^{(0)}$ koordinātēm.

Nākamā algoritma iterācija netiks veikta tālāk norādītā iemesla dēļ. Ja mēs aprēķināsim attālumus no objektiem līdz jaunajiem centroīdiem $c_k^{(2)}$, tad šīs koordinātes sakrītīs ar koordinātēm $c_k^{(1)}$. Tas nozīmē, ka ir panākta rezultātu konverģence. Reālās problēmās, lai panāktu konverģenci, var būt nepieciešams diezgan liels iterāciju skaits.

Paanalizēsim iegūtos rezultātus. Ir skaidrs, ka klasteros c_1 un c_2 ir normāli dati. Klasterī c_3 ir tikai divi objekti, kas atrodas attālu no citiem objektiem. Tāpēc ir nopietni iemesli apgalvot, ka objekti o_5 , o_8 šajā datu kopā ir anomālijas.

Principā šādu secinājumu var izdarīt, pamatojoties uz attēla 3.3.1 grafiskā attēlojuma vizuālu analīzi. Bet mēs formāli atrisinājam problēmu ar anomāliju identificēšanu šajos datos, lai sniegtu priekšstatu par klasterizācijas algoritma k -vidējais darbības principu.

Jāatzīmē šādas aplūkotā algoritma priekšrocības:

- algoritms ir vienkārši izpildāms;
- tas labi darbojas lielās datu kopās;
- labi pielāgojas jauniem objektiem.

Algoritma trūkumi ir:

- jūtība pret anomālijām;
- neveiksmīga klasteru skaita izvēle var novest līdz neapmierinošiem rezultātiem.

2. k -medoīdu metode.

k -medoīdu metode ir klasiska klasterizācijas metode, kuras pamatā ir sākotnējās objektu kopas sadalīšana *a priori* noteiktā skaitā grupu (klasteru). Tas ir līdzīgs k -vidējais metodei ar dažām īpašām iezīmēm. k -vidējais metodē jauši izvēlēti punkti atribūtu vērtību telpā tiek ņemti kā klasteru sākotnējie centroīdi. k -medoīdu metodē kā sākotnējie klasteru centri tiek atlasīti konkrēti objekti datu kopā (medoīdi). Klasteru skaits k ir noteikts *a priori*. Objekti tiek attiecināti klasteriem, pamatojoties uz aprēķiniem par attālumiem no medoīdiem līdz visiem objektiem datu kopā. Šī procedūra tiek veikta tieši tādā pašā veidā kā k -vidējais metode.

Būtiskākā atšķirība starp abām metodēm ir centroīdu un medoīdu pārvērtēšana. Izmantojot k -vidējais metodi jaunās centra koordinātes tiek aprēķinātas kā visu atribūtu vidējā vērtība visiem objektiem, kas iekļauti katrā klasterī. Izmantojot k -medoīdu metodi katrā iterācijā par pašreizējo medoīdu tiek izvēlēts cits objekts un tiek aprēķināti attālumi no visiem objektiem klasterī līdz jaunajam medoīdam. Ja attālumu summa līdz jaunajam medoīdam ir mazāka nekā attālumu summa līdz iepriekšējam medoīdam, jaunais medoīds tiek uzskatīts par pašreizējo medoīdu. Ja attālumu summa līdz iepriekšējai medoīdai ir mazāka nekā attālumu summa līdz jaunajam medoīdam, iepriekšējo medoīdu uzskata par pašreizējo darba medoīdu. Process tiek atkārtots iteratīvi, līdz tiek

identificēts medoīds, kuram ir vismazākā attālumu summa līdz objektiem noteiktā klasterī salīdzinājumā ar visiem citiem potenciālajiem medoīdiem.

k -medoīdu metode pirmo reizi tika piedāvāta darbā [Kaufman L., Rousseeuw P., 1987] kā k -vidējais metodes uzlabota versija PAM (*Partitioning Around Medoid*). Vēlāk tika piedāvātas uzlabotas šī algoritma versijas, piemēram, CLARA algoritms.

Formāli k -medoīdu algoritma mērķa funkcija tiek izteikta formā

$$s_k = \arg \min \sum_{a_{ij} \in c_k} (a_{ij} - m_k)^2, \quad (3.3.5)$$

kur c_k , $k = 1, \dots, k$, - k -tais klasteris;

m_k - k -tā klastera medoīda.

k -medoīdu algoritms ietver šādu procedūru secību:

1. Jānosaka klasteru skaits k .
2. Jāatlasa punktu atribūtu vērtību telpā un jāuzdod tos kā klasteru sākotnējos centroīdus.
3. Jāizskaitļo attālumus no katra datu punkta (objekta) līdz katram no dotajiem centroīdiem, izmantojot Eiklīda attālumu.
4. Jāattiecina katru objektu klasterim, kuram Eiklīda attālums starp objektu un dotā klastera centroīdu ir minimāls.
5. Pārrēķina jaunās centroīdu koordinātes un izdara 4. punktā dotās darbības.
6. Jāatkārto 3.–5. darbības, līdz centroīdu koordinātes vairs nemainās.

Apskatīsim vienkāršu piemēru. Ņemsim par pamatu sākotnējos normalizētos datus no tabulas 3.3.2. Izvēlēsimies klasteru skaitu $k = 3$. Izvēlēsimies sākotnējos medoīdus klasteriem c_1, c_2, c_3 :

$$m_1^{(0)} = o_3, \quad m_2^{(0)} = o_4, \quad m_3^{(0)} = o_5.$$

Mēs nenoteiksim, vai objekti pieder klasteriem. Šī procedūra tika veikta iepriekšējā piemērā. Datu struktūra šajā piemērā ir tāda, ka nosakot objektu piederību klasteriem pamatojoties uz jebkuriem sākotnējiem medoīdiem objektu grupās (attēls 3.3.2), mēs iegūsim tieši tādu pašu objektu sadalījumu klasteros, kas tika iegūti iepriekšējā piemērā. (Šis apgalvojums, protams, attiecas tikai uz šo piemēru).

Aprēķināsim attālumus no objektiem kopās līdz sākotnējiem medoīdiem un aprēķināsim šo attālumu summas. Aprēķinu rezultāti ir parādīti tabulā 3.3.4.

Tabula 3.3.4

Eiklīda attālumi no objektiem līdz sākotnējiem medoīdiem $m_1^{(0)}$ un $m_2^{(0)}$ un šo attālumu summas

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	o_{11}	o_{12}	s_k
$m_1^{(0)} = o_3$		0.22	0			0.08	0.08			0.13			0.51
$m_2^{(0)} = o_4$	0.20			0					0.08		0.37	0.19	0.84
$m_3^{(0)} = o_5$					0			0.08					0.08

Tabulas pēdējā kolonnā ir norādītas mērķa funkciju vērtības kā attālumu summas pa attiecīgajām tabulas rindām.

Uzdosim medoīdus šādi:

$$c_1 : m_1^{(1)} = o_6 \quad (-1.05, -0.54);$$

$$c_2 : m_2^{(1)} = o_1 \quad (0.97, 0.27).$$

Klasterim c_3 nav jēgas veikt medoīdu maiņu. Šajā klasterī ir tikai divi objekti. Par sākotnējo medoīdu tika izvēlēts objekts o_5 . Ja tagad izvēlēsimies par pašreizējo medoīdu objektu o_8 , mēs iegūsim tādu pašu attālumu no objekta o_5 līdz medoīdai o_8 , kas tika iegūts iepriekš.

Formāli izpildot algoritmu pēc jaunu medoīdu piešķiršanas, ir nepieciešams no jauna noteikt attālumus no visiem objektiem līdz šiem medoīdiem. Šī procedūra ir nepieciešama, jo medoīdu maiņa var izraisīt dažu objektu attiecināšanu uz citiem klasteriem. Mēs neveicam šo procedūru, jo datu struktūra šajā piemērā ir tāda, ka ņemot jebkuru objektu klasteros c_1 , c_2 par medoīdu, iepriekš definētie klasteri netiks izmainīti.

Attālumu aprēķinu rezultāti no objektiem klasteros c_1 , c_2 līdz jaunajiem medoīdiem $m_1^{(1)}$, $m_2^{(1)}$ ir parādīti tabulā 3.3.5.

Tabula 3.3.5

Eiklīda attāluma vērtības no objektiem klasteros c_1, c_2 līdz pašreizējiem medoīdiem $m_1^{(1)}, m_2^{(1)}$

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	o_{11}	o_{12}	s_k
$m_1^{(1)} = o_6$		0.07	0.08		-	0	0.31	-		0.17			0.63
$m_2^{(1)} = o_1$	0			0.19	-			-	0.22		0.23	0.51	1.15

Varam secināt, ka iepriekšējie rezultāti (tabulā 3.3.4) ir labāki par pašreizējiem rezultātiem (tabulā 3.3.5) tādā ziņā, ka tiem ir mazākas mērķa funkciju $s_k, k = 1, 2$, vērtības. Klasterī c_3 izvēloties kādu no objektiem o_5 vai o_8 par medoīdu, tiks iegūti tādi paši rezultāti.

Mēs neatkārtosim iterācijas c_1 . No attēla 3.3.2 vizuālas analīzes var secināt, ka jebkura objekta o_2 , o_7 vai o_{10} izvēle par medoīdu dod lielus attālumus no objektiem klasterī līdz potenciālajam medoīdam. Tāpēc šīs iterācijas kļūst bezjēdzīgas.

Veiksim nākamo iterāciju klasterim c_2 . Par pašreizējo medoīdu izvēlēsimies objektu o_9 . Aprēķinu rezultāti ir parādīti tabulā 3.3.6.

Tabula 3.3.6

Eiklīda attālumu vērtības no objektiem klasterī c_2 līdz pašreizējai medoīdai $m_2^{(2)} = o_9$

Objekti	o_1	o_4	o_9	o_{11}	o_{12}	s_2
$m_2^{(2)} = o_9$	0.22	0.08	0	0.52	0.08	0.90

Mēs neveiksim turpmākas iterācijas klasterī c_2 . Attēla 3.3.2 vizuālā analīze parāda, ka izvēloties objektus o_{11} vai o_{12} par pašreizējo medoīdu, mēs iegūstam attālumu summu, kas ir lielāka nekā medoīdai $m_2^{(0)}$.

Apkopojot iegūstam šādus rezultātus:

1. Klasteris c_1 ar objektu o_2 kā medoīdu.
2. Klasteris c_2 ar objektu o_4 kā medoīdu.
3. Klasteris c_3 ar objektiem o_5, o_5 vai o_8 kā medoīdu.

Anomāliju identifikācija, pamatojoties uz k -medoīdu klasterizācijas metodi tiek veikta tieši tāpat kā anomāliju identifikācija ar k -vidējais klasterizācijas metodi. Mūsu piemērā anomālijas ir objekti o_5, o_8 .

3.4. Anomāliju identificēšana lineārajā regresijā

Pieņemsim, ka mums ir divas atribūtu a_j, a_k vērtību kopas un mūsu uzdevums ir izveidot regresijas līkni, kas parāda attiecības starp atribūtu vērtībām. Tas ir triviāls uzdevums normālu datu gadījumā, t.i., datiem bez anomālijām. Bet, ja datos ir anomālijas, tās var būtiski ietekmēt regresijas attiecību raksturu.

Attēlā 3.4.1.a ir parādīta regresijas līkne, kas konstruēta uz datiem bez anomālijām. Attēlā 3.4.1.b ar tiem pašiem datiem papildus ir punkts O_i , kuram ir liela atribūta a_k vērtība. Šis punkts var būt potenciāla anomālija. Šī punkta dēļ ir mainījies regresijas līkne, tā atrodas virs datu punktiem un tāpēc neatbilst normālo datu daļai. Attēlā 3.4.1.c pie tiem pašiem datiem papildus ir punkts O_i , ko literatūrā mēdz dēvēt par sviras vai līdzsvara punktu. Šim punktam ir augsta atribūta a_j vērtība un tas arī varētu būt potenciāla anomālija. Šī iemesla dēļ regresijas līkne atrodas zem galvenajiem datu punktiem un īsti neatbilst normāliem datiem. Attēlā 3.4.1.d anomālija ir objekts O_i , kuram ir lielas abu atribūtu vērtības. Neņemot vērā to, ka objekta punkts O_i atrodas uz regresijas līknes, šis objekts ir anomālija, jo tas atrodas tālu no lielākās datu daļas.

Ir piedāvātas daudzas metodes anomāliju identificēšanai lineārajā regresijā. Tās var iedalīt statistiskajās un robustajās (*robust*) metodēs. Apskatīsim dažas metodes no abām grupām.

1. Statistikas metodes.

Lai lineārajā regresijā varētu identificēt anomālijas, ir nepieciešams noteikt regresijas vienādojumu priekš sākotnējiem datiem.

Detalizēts regresijas līknes noteikšanas un konstruēšanas procedūru apraksts ir sniegts pielikumā P1.3. Regresijas vienādojums ir dots formā

$$y = b_0 + b_1x, \quad (3.4.1)$$

kur x - faktoriālā pazīme (faktoriālais mainīgais);

y - rezultatīvā pazīme (rezultatīvais mainīgais).

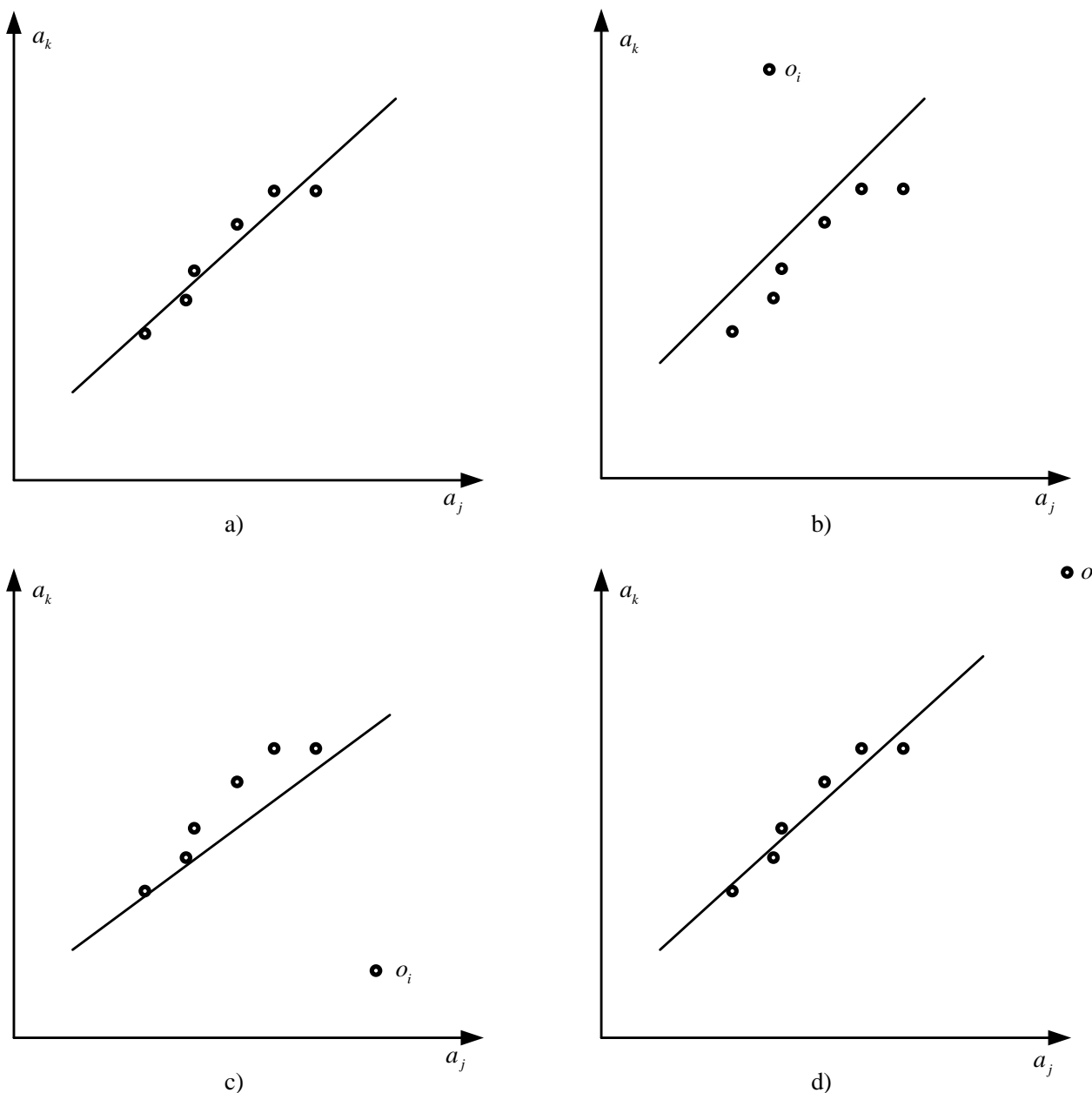
Regresijas koeficientu vērtības aprēķina, izmantojot šādus vienādojumus:

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (3.4.2)$$

kur \bar{x} - faktoriālā mainīgā vidējā vērtība x ;

\bar{y} - rezultatīvā mainīgā vidējā vērtība y .

$$b_0 = \bar{y} - b_1\bar{x}. \quad (3.4.3)$$



Attēls 3.4.1. Anomāliju ietekmes uz regresijas atkarību grafisks attēlojums

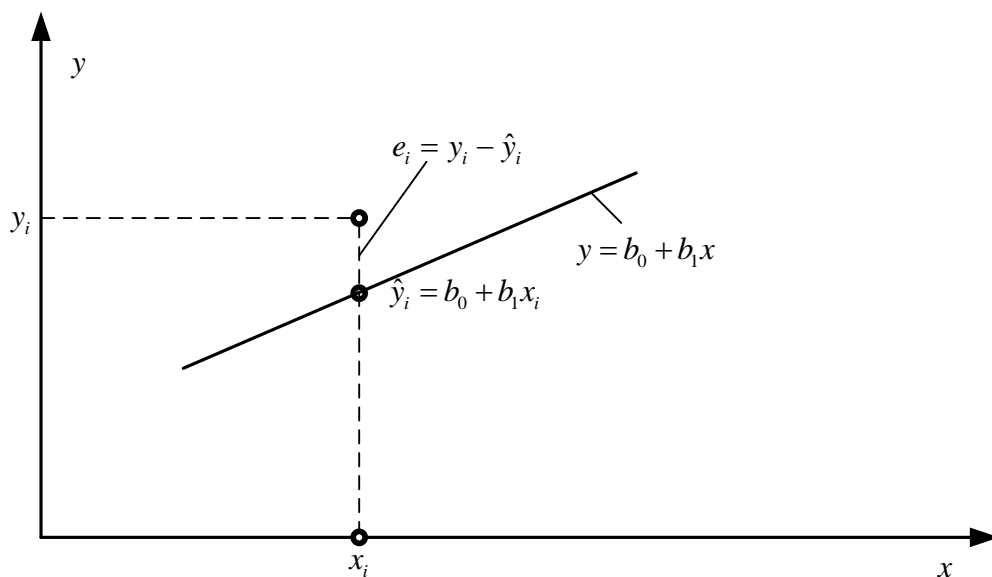
Koeficienta b_1 vērtība nosaka regresijas līknes slīpumu attiecībā pret horizontālo asi x . Parametra b_0 vērtība nosaka regresijas līknes krustošanās punktu ar vertikālo y asi.

Uzdodot divas mainīgā x vērtības, tiek aprēķinātas atbilstošās rezultatīvā mainīgā y vērtības, izmantojot vienādojumu (3.4.1), tādējādi var izveidot regresijas līkni.

Attēlā 3.4.2 shematiski attēloti datu punkti (x_i, y_i) un konstruētā regresijas līkne $y = b_0 + b_1 x$. Mainīgajiem x un y būs lielākas vai mazākas novirzes no regresijas līknes. Punkts \hat{y}_i šajā attēlā atbilst vērtībai y , kas aprēķināta ar regresijas vienādojumu pie vērtības $x = x_i$. Atšķirība

$$e_i = (y_i - \hat{y}_i), \quad i = 1, \dots, n \quad (3.4.4)$$

parāda datu punktu y_i novirzes no regresijas līknes (kļūdas).



Attēls 3.4.2. Datu punktu un regresijas līknes shematisks attēlojums

Izmantojot sākotnējo datu punktu novirzes no regresijas līknes, pilnu regresijas modeli var attēlot formā

$$y = b_0 + b_1x + \varepsilon, \quad (3.4.5)$$

kur ε - kļūdas lielums, kas parāda datu punktu novirzes no regresijas līknes.

Standarta statistika izdara šādus pieņēmumus par kļūdas lielumu:

1. ε ir nulles vidējais rādītājs.
2. ε ir nemainīga variācija σ^2 .
3. Individuālās novirzes (kļūdas) ir nekorelēti gadījuma lielumi.
4. Novirzes (kļūdas) ε_i atbilst normālajam sadalījumam.

Lai regresijā identificētu anomālijas, ir jādefinē regresijas vienādojums, jākonstruē regresijas līkne un jānosaka datu punktu novirzes ε_i no regresijas līknes.

Viena no statistiskajām metodēm anomāliju identificēšanai regresijā ir šāda. Tiek aprēķinātas standartnoviržu vērtības

$$d_i = \frac{e_i}{\sqrt{MS_{RES}}}. \quad i = 1, \dots, n, \quad (3.4.6)$$

kur e_i aprēķina pēc vienādojuma (3.4.4);

MS_{RES} - aprēķinātā noviržu kvadrātu vidējā vērtība pēc vienādojuma

$$MS_{RES} = \frac{\sum_{i=1}^n e_i^2}{n}. \quad (3.4.7)$$

Ja kādam objektam o_i $d_i \geq 3$, šo objektu var identificēt kā potenciālu anomāliju.

Apskatīsim šādu ilustratīvu piemēru. Tabulā 3.4.1 ir parādītas objektu $o_1 - o_{10}$ atribūtu vērtības a_1, a_2 .

Tabula 3.4.1

Objektu $o_1 - o_{10}$ sākotnējās atribūtu a_1, a_2 vērtības

Objekti	a_1	a_2
o_1	4.00	4.50
o_2	7.00	6.50
o_3	10.00	8.50
o_4	11.00	20.00
o_5	15.00	10.00
o_6	6.00	5.50
o_7	14.00	10.00
o_8	8.00	6.50
o_9	12.00	8.00
o_{10}	13.00	8.50
	$\sum = 100.00, \bar{a}_1 = 10.00$	$\sum = 88.00, \bar{a}_2 = 8.80$

Skaidrības labad dati no tabulas 3.3.1 ir grafiski parādīti attēlā 3.4.3. Šis datu kopums ir jāanalizē, lai tajā identificētu anomālijas.

Tabulas 3.4.1 pēdējā rindā ir norādītas atribūtu a_1, a_2 vērtību summas attiecīgajās kolonās un vidējās atribūtu vērtības: $\bar{a}_1 = 10.00$, $\bar{a}_2 = 8.80$.

Veiksim starppaprēķinus, lai noteiktu regresijas koeficientus b_0 un b_1 . Aprēķinu rezultāti parādīti tabulā 3.4.2.

Izmantojot vienādojumu (3.4.2), aprēķinām

$$b_1 = \frac{70.00}{120.00} = 0.58.$$

Pēc vienādojuma (3.3.3):

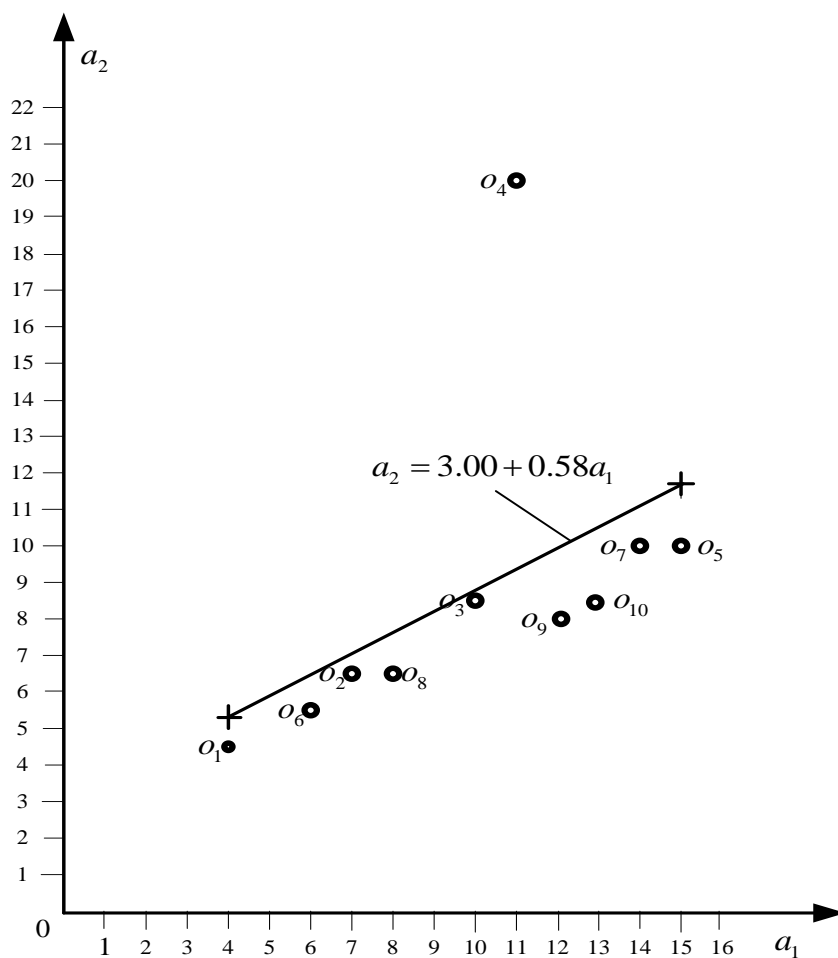
$$b_0 = 8.80 - 0.58 * 10 = 8.80 - 5.80 = 3.00.$$

Pieņemsim $a_{11} = 4.00$.

$$\hat{a}_{12} = 3.00 + 0.58 * 4.00 = 5.32.$$

Pieņemsim $a_{51} = 15.00$.

$$\hat{a}_{52} = 3.00 + 0.58 * 15.00 = 11.70.$$



Attēls 3.4.3. Tabulas 3.4.1 datu grafiskais attēlojums

Tabula 3.4.2

Regressijas koeficientu b_0, b_1 noteikšanas rezultāti

Objekti	$a_{i1} - \bar{a}_1$	$(a_{i1} - \bar{a}_1)^2$	$a_{i2} - \bar{a}_2$	$(a_{i1} - \bar{a}_1) * (a_{i2} - \bar{a}_2)$
o_1	-6.00	36.00	-4.30	25.00
o_2	-3.00	9.00	-2.30	6.90
o_3	0.00	0.00	-0.30	0.00
o_4	1.00	1.00	11.20	11.20
o_5	5.00	25.00	1.20	6.00
o_6	-4.00	16.00	-3.30	13.20
o_7	4.00	16.00	1.20	4.80
o_8	-2.00	4.00	-2.30	4.60
o_9	2.00	4.00	-0.80	-1.60
o_{10}	3.00	9.00	-0.30	-0.90
		$\Sigma = 120.00$		$\Sigma = 70.00$

Izmantojot izskaitļotās punktu vērtības, konstruējam regresijas līkni. Šī līkne ir parādīta attēlā 3.4.3.

Redzams, ka šī līkne slikti atbilst sākotnējiem datiem, jo tā atrodas virs lielākās datu daļas.

Aprēķināsim datu punktu standartizētās novirzes no konstruētās regresijas līknes. Starppaprēķinu rezultāti ir parādīti tabulā 3.4.3.

Tabula 3.4.3

Starpprezultātu aprēķinu vērtības standartizētās novirzes d_{i2} noteikšanai

Objekti	a_{i2}	\hat{a}_{i2}	$e_{i2} = a_{i2} - \hat{a}_{i2}$	e_{i2}^2
o_1	4.50	5.32	-0.82	0.672
o_2	6.50	7.06	-0.56	0.314
o_3	8.50	8.80	-0.30	0.090
o_4	20.00	9.38	10.62	112.784
o_5	10.00	11.70	-1.70	2.890
o_6	5.50	6.48	-0.98	0.960
o_7	10.00	11.12	-1.12	1.254
o_8	6.50	7.64	-1.14	1.300
o_9	8.00	9.96	-1.96	3.842
o_{10}	8.50	10.54	-2.04	4.162
				$\Sigma = 128.268$

Pēc vienādojuma (3.3.7):

$$MS_{RES} = \frac{128.268}{10} = 12.827.$$

$$\sqrt{MS_{RES}} = \sqrt{12.827} = 3.581.$$

Izmantojot vienādojumu (3.4.6) nosakām standartizētās novirzes vērtības. Aprēķinu rezultāti parādīti tabulā 3.4.4.

Tabula 3.4.4

Standartizētās novirzes vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}
d_{i2}	-0.230	-0.156	-0.083	2.966	-0.475	-0.274	-0.313	-0.318	-0.577	-0.570

Vērtība $d_{i2} = 2.966 \approx 3$ ļauj secināt, ka objekts o_4 ir anomālija.

Lai parādītu anomālijas ietekmi uz regresijas kvalitāti, izņemsim šo anomāliju (objektu o_4) no sākotnējās datu kopas. Samazinātā datu kopa ir parādīta tabulā 3.4.5.

Tabula 3.4.5

Samazināta datu kopa no tabulas 3.4.1

Objekti	a_1	a_2
o_1	4.00	4.50
o_2	7.00	6.50
o_3	10.00	8.50
o_5	15.00	10.00
o_6	6.00	5.50
o_7	14.00	10.00
o_8	8.00	6.50
o_9	12.00	8.00
o_{10}	13.00	8.50
	$\sum = 89.00, \bar{a}_1 = 9.89$	$\sum = 68.00, \bar{a}_2 = 7.56$

Tabulā 3.4.6 parādīti starppaprēķinu rezultāti regresijas koeficientu b_0, b_1 noteikšanai.

Tabula 3.4.6

Regresijas koeficientu b_0, b_1 noteikšanas starppaprēķinu rezultāti

Objekti	$a_{i1} - \bar{a}_1$	$(a_{i1} - \bar{a}_1)^2$	$a_{i2} - \bar{a}_2$	$(a_{i1} - \bar{a}_1) * (a_{i2} - \bar{a}_2)$
o_1	-5.89	34.69	-3.05	18.02
o_2	-2.89	8.35	-1.06	3.06
o_3	0.11	0.01	0.94	0.10
o_5	5.11	26.11	2.44	12.47
o_6	-3.89	15.13	-2.06	8.01
o_7	4.11	16.89	2.41	10.03
o_8	-1.89	3.57	-1.06	2.00
o_9	2.11	4.45	0.44	0.93
o_{10}	3.11	9.67	0.94	2.02
		$\sum = 118.87$		$\sum = 57.54$

Pēc vienādojuma (3.4.2):

$$b_1 = \frac{57.44}{118.87} = 0.48.$$

Pēc vienādojuma (3.4.3):

$$b_0 = 7.56 - 0.48 * 9.89 = 2.81.$$

Tādējādi mums ir šāds regresijas vienādojums:

$$a_2 = 2.81 + 0.48 * a_1.$$

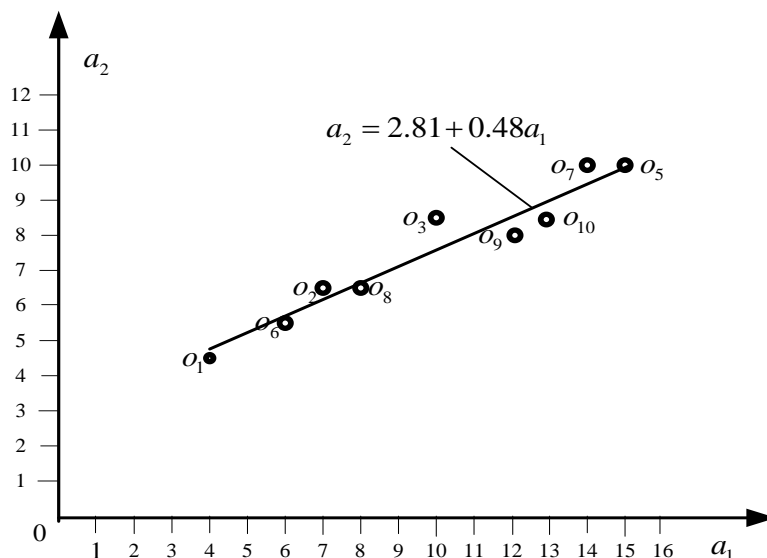
Pieņemsim $a_{11} = 4.00$.

$$a_{12} = 2.81 + 0.48 * 4.00 = 4.73.$$

Pieņemsim $a_{32} = 15.00$.

$$a_{32} = 2.81 + 0.48 * 15.00 = 10.01.$$

Attēls 3.4.4 grafiski attēlo sākotnējo samazināto datu kopu un regresijas līkni, kas izveidota uz šīs kopas pamata.



Attēls 3.4.4. Samazinātās datu kopas grafiskais attēlojums un konstruētā regresijas līkne

Vienkārša vizuālā analīze parāda, ka regresijas līkne attēlā 3.4.4 labi atbilst datiem, atšķirībā no regresijas līknes attēlā 3.4.3, kas slikti atbilst datiem, jo tad datos bija anomālija.

Var redzēt, ka regresijas līknes attēlos 3.4.3 un 3.4.4 ir līdzīgas regresijas līknēm attēlos 3.4.1.a,b, kas shematiski attēlo regresijas līknes bez anomālijām datos un ar anomālijas punktu.

Alternatīva anomāliju identificēšanai regresijā ir izmantot studentizētās novirzes vērtības (*studentized residual*). Pieņemsim, tāpat kā iepriekš, $e_i = y_i - \hat{y}_i$, kur y_i ir rezultātīvā vērtība priekš mainīgā x_i i -tās vērtības, \hat{y}_i ir rezultātīvā vērtība, kas noteikta no regresijas vienādojuma.

Ieviesīsim mainīgo h_i , ko aprēķina pēc vienādojuma

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}. \quad i = 1, \dots, n, \quad (3.4.8)$$

kur n - faktoriālā mainīgā vērtību skaits;

x_i - faktoriālā mainīgā i -tā vērtība;

\bar{x} - faktoriālā mainīgā vidējā vērtība.

Faktiski vērtība h_i ir attāluma novērtējums no vērtības x_i līdz vidējai vērtībai \bar{x} .

Ieviesīsim vēl vienu mainīgo S_E :

$$S_E = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n - k - 1}}, \quad (3.4.9)$$

kur $e_i = y_i - \hat{y}_i$ - punkta y_i novirze no regresijas līknes;

n - rezultatīvā mainīgā y vērtību skaits;

k - faktoriālo mainīgo skaits (šajā gadījumā $k = 1$).

Faktiski S_E ir summētā novirze visā noviržu kopumā $\{e_i\}$ no regresijas līknes.

Studentizētās i -tās novirzes vērtību aprēķina pēc vienādojuma

$$r_i = \frac{e_i}{S_E \sqrt{1-h_i}}, \quad (3.4.10)$$

kur h_i aprēķina pēc vienādojuma (3.4.8);

S_E aprēķina, izmantojot vienādojumu (3.4.9).

Rodas jautājums par robežvērtības r noteikšanu. Dažādi literatūras avoti piedāvā dažādas robežvērtības. Šajā sadaļā mēs noteiksim robežvērtību $r = 2$. Tas nozīmē, ka objektu o_i , kuram $r_i \geq 2$ var uzskatīt par potenciālu anomāliju datus.

Aprēķināsim h_i iepriekšējā piemēra datiem. Ņemsim par pamatu $(a_{i1} - \bar{a})^2$ un $\sum_{i=1}^{10} (a_{i1} - \bar{a}_1)^2$ vērtības tabulas 3.4.2 trešās kolonas. Šīs vērtības ir parādītas tabulas 3.4.7 2. rindā.

Izmantojot vienādojumu (3.4.8), mēs aprēķinām h_i vērtības. Šīs vērtības ir parādītas tabulas 3.4.7 3. rindā.

Tabula 3.4.7

$(a_{i1} - \bar{a})^2$ un $\sum_{i=1}^{10} (a_{i1} - \bar{a}_1)^2$ vērtības no tabulas 3.4.2 un aprēķinātās h_i vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	$\sum_{i=1}^{10} (a_{i1} - \bar{a}_1)^2$
$a_{i1} - \bar{a}_1$	36.00	9.00	0.00	1.00	25.00	16.00	16.00	4.00	4.00	9.00	120.00
h_i	0.400	0.175	0.00	0.108	0.308	0.233	0.233	0.133	0.133	0.175	

Lai aprēķinātu S_E vērtības, mēs izmantojam starpības $a_{i2} - \bar{a}_2$ no tabulas 3.4.2 4. kolonas. Šo starpību vērtības ir parādītas tabulas 3.4.8 2. rindā. Trešajā rindā ir parādītas šo starpību kvadrātā vērtības. Tabulas apakšējā labajā šūnā ir parādīta šo starpību summa.

Tabula 3.4.8

Aprēķinu rezultāti S_E vērtības noteikšanai

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	$\sum_{i=1}^{10} (a_{i2} - \bar{a}_2)^2$
$a_{i2} - \bar{a}_2$	-4.30	-2.30	-0.30	11.20	1.20	-3.30	1.20	-2.30	-0.80	-0.30	
$(a_{i2} - \bar{a}_2)^2$	18.49	5.29	0.09	125.44	1.44	10.89	1.44	5.29	0.64	0.09	169.10

Izmantojot vienādojumu (3.4.9), aprēķinām

$$S_E = \sqrt{\frac{169.100}{10-1-1}} = \sqrt{\frac{169.100}{8}} = \sqrt{21.124} = 4.598.$$

Tabulas 3.4.9 2. rindā ir parādītas novirzes e_{i2} vērtības no tabulas 3.4.3 4. kolonas. Tabulas 3. rindā parādītas studentizētās noviržu vērtības, kas aprēķinātas, izmantojot vienādojumu (3.4.10).

Tabula 3.4.9

Noviržu e_{i2} un studentizēto noviržu r_i vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}
e_{i2}	-0.82	-0.56	-0.30	10.62	-1.70	-0.98	-1.12	-1.14	-1.96	-2.04
r_i	-0.224	-0.314	-0.065	2.445	-0.444	-0.243	-0.278	-0.266	-0.458	-0.489

Būtiskākā ir studentizētā novirze $r_4 = 2.445$, tāpēc objektu o_4 var pamatoti uzskatīt par anomāliju dotajā datu kopā.

Studentizēto noviržu vērtību svarīga iezīme ir tā, ka to aprēķinos tiek izmantotas faktoriālā mainīgā lieluma attālums no tā vidējās vērtības.

Lai novērtētu rezultatīvā mainīgā vērtību noviržu no regresijas līknes ietekmes pakāpi uz pašu regresiju, tiek izmantoti dažādi skaitliskie rādītāji. Piedāvājam divus plaši izmantojamus šāda veida rādītājus.

1. DFFITS.

Šis rādītājs i -tajam rezultatīvajam mainīgajam tiek aprēķināts šādā veidā:

$$DFFITS_i = r_i \sqrt{\frac{h_i}{1-h_i}}, \quad (3.4.11)$$

kur r_i - studentizētās novirzes vērtība faktoriālā mainīgā i -tajai vērtībai (pēc vienādojuma (3.4.10));

h_i - aprēķina, izmantojot vienādojumu (3.4.8).

Tabulas 3.3.10 2. rindā attēlotas studentizēto noviržu r_i vērtības, kas ņemtas no tabulas 3.4.8 3. rindas. Tabulas 3. rindā attēlotas h_i vērtības, kas ņemtas no tabulas 3.4.7 3. rindas. Tabulas 4. rindā ir parādītas DFFITS vērtības, kas aprēķinātas pēc vienādojuma (3.4.11).

Tabula 3.4.10

 r_i , h_i un aprēķinātās DFITS vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}
r_i	-0.224	-0.314	-0.065	2.445	-0.444	-0.243	-0.278	-0.266	-0.458	-0.489
h_i	0.400	0.175	0.000	0.108	0.308	0.233	0.233	0.123	0.133	0.175
DFFITS _i	-0.099	-0.067	0.000	1.085	-0.245	-0.134	-0.153	-0.032	-0.179	-0.225

Izmantojot DFFITS rādītāju, tiek pieņemts, ka objekts o_i , kam $DFFITS_i > 1$, ir anomālija. Pamatojoties uz tabulas 3.4.10 pēdējā rindā dotajiem datiem, var izdarīt acīmredzamu secinājumu, ka objekts o_4 ir anomālija. Šis rezultāts sakrīt arī ar iepriekš iegūto rezultātu.

2. Kuka attālums (Cook's distance).

Šis rādītājs atspoguļo regresijas koeficienta vērtību atšķirību, kas aprēķināta datu kopā esošā i -tā objekta klātbūtnē vai bez tā. Kuka attālumu aprēķina, izmantojot vienādojumu

$$D_i = \frac{e_i^2}{S_{EP}^2} \left(\frac{h_i}{(1-h_i)^2} \right), \quad (3.4.12)$$

kur $e_i = y_i - \hat{y}_i$ - punkta y_i novirze no regresijas līknes;

n - rezultatīvā mainīgā y vērtību skaits;

k - faktoriālo mainīgo skaits (šajā gadījumā $k = 1$);

h_i - aprēķina, izmantojot vienādojumu (3.4.8);

$$S_E^2 = \frac{\sum_{i=1}^n e_i^2}{n - k - 1}. \quad (3.4.13)$$

Aprēķināsim Kuka attāluma vērtību iepriekš dotajiem datiem. Tabulas 3.4.11 2. rindā ir parādītas noviržu e_{i2} vērtības no tabulas 3.4.3 4. kolonas. Tabulas 3. rindā ir norādītas h_i vērtības no tabulas 3.4.10 3. rindas. Tabulas 4. rindā ir norādītas aprēķinātās D_i vērtības.

Vērtības D_i nenosaka robežvērtību starp normālo datu apgabalu un anomālijām. Tikai objekts ar vislielāko D_i vērtību var tikt identificēts kā potenciāla anomālija. Tabulā 3.4.11 objektam o_4 ir vislielākā D_i vērtība: $D_4 = 0.034$, tātad šis objekts ir potenciāla anomālija.

Tabula 3.4.11

Vērtības e_i , h_i un aprēķinātās D_i vērtības

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}
e_i	-0.82	-0.56	-0.30	10.62	-1.70	-0.98	-1.12	-1.14	-1.96	-2.04
h_i	0.400	0.175	0000	0.108	0.308	0.233	0.233	0.133	0.133	0.175
D_i	-0.021	-0.003	0000	0.034	-0.026	-0.090	-0.10	-0.005	-0.008	-0.012

Šajā sadaļā tika aplūkotas statistiskās metodes anomāliju identificēšanai lineārajā regresijā. Iegūtos rezultātus var attiecināt arī uz daudzfaktoru lineārās regresijas gadījumiem, kad ievērojami palielinās aprēķinu sarežģītība. Tā vietā, lai darbotos ar diskrētām skaitliskām vērtībām, daudzfaktoru lineārā regresija veic vektoru un matricu aprēķinus. Šāda veida aprēķinus var veikt tikai izmantojot piemērotus skaitļošanas rīkus.

2. Robustas metodes anomāliju identificēšanai regresijā.

Statistiskajām metodēm anomāliju noteikšanai regresijā ir viens būtisks trūkums: regresijas atkarības definīcija ir balstīta uz visiem sākotnējiem datiem, ieskaitot anomālijas. Tā rezultātā var rasties regresijas līkne, kas neatbilst lielākajai datu daļai. Spilgts piemērs tam ir regresijas līknes attēlos 3.3.3 un 3.3.4. Pirmā līkne ir izveidota izmantojot visu sākotnējo datu kopu, tostarp anomāliju objektu o_4 . Otrā regresijas līkne tiek veidota uz datu kopas pamata, no kuras ir izslēgts objekts o_4 . Acīmredzami otrā līkne labāk atbilst lielākajai datu daļai.

Lai izslēgtu anomāliju ietekmi uz konstruētās regresijas līkni, tiek izmantotas robustas metodes regresijas atkarības noteikšanai. Ir daudz spēcīgu šāda veida metožu. Šajā sadaļā aplūkosim vienkāršāko un visbiežāk izmantoto metodi: M - novērtēšanu. Regresijas M - novērtēšanas metodes teorētiskie pamati tika likti darbā [Huber P.J., 1981]. Līdzīgi kā iepriekš, x ir faktoriālais mainīgais un y ir rezultātīvais mainīgais. Standarta lineārajā regresijā, lai noteiktu regresijas sakarību starp faktoriālo mainīgo un rezultātīvo mainīgo, mērķis ir minimizēt rezultātīvā mainīgā vērtību kvadrātu noviržu summu no regresijas līknes:

$$\min \sum_{i=1}^n e_i^2 = \min \sum_{i=1}^n (y_i - x_i' b)^2, \quad (3.4.14)$$

kur b - korelācijas koeficients, kas atspoguļo regresijas līknes slīpumu attiecībā pret horizontālo asi x .

Robustajā M - novērtēšanā vienādojuma (3.3.14) kvadrātiskās novirzes tiek aizstātas ar funkcijas $\rho(u_i)$ vērtībām, kur

$$u_i = \frac{e_i}{S_m}. \quad (3.4.15)$$

Vienādojumā (3.4.15) e_i ir rezultatīvā i -tā mainīgā vērtības novirze no regresijas līknes.

$$S_m = \frac{\text{med} |e_i - \text{med} \{e_i\}|}{0.6475} = \frac{\text{mad} \{e_i\}}{0.6475}, \quad (3.4.16)$$

kur $\text{mad}(e_i)$ - vērtību e_i noviržu mediāna no kopas mediānas vērtībām $\{e_i\}$.

Robustās M - novērtēšanas mērķa funkcija uzdota sekojošā veidā:

$$\sum_{i=1}^n \rho(u_i) = \sum_{i=1}^n \left(\frac{e_i}{S_m} \right) = \sum_{i=1}^n \rho \left(\frac{y_i - x_i' \beta}{S_m} \right). \quad (3.4.17)$$

Funkcijai $\rho(u_i)$ jāatbilst šādām prasībām:

1. $\rho(u_i) \geq 0$;
2. $\rho(0) = 0$;
3. $\rho(u_i) = \rho(-u_i)$;
4. $\rho(u_j) \geq \rho(u_i)$ priekš $|u_j| \geq |u_i|$.

Robustās regresijas atkarības M - novērtēšanas mērķis ir minimizēt mērķa funkciju (3.4.17):

$$\min \sum_{i=1}^n \rho(u_i) = \min \sum_{i=1}^n \left(\frac{e_i}{S_m} \right) = \min \sum_{i=1}^n \rho \left(\frac{y_i - x_i' \beta}{S_m} \right). \quad (3.4.18)$$

Robustā M - novērtēšana ir balstīta uz maksimālās līdzības principa.

Ir piedāvātas dažādas funkcijas $\rho(u_i)$ definēšanas iespējas. Izplatīta pieeja ir izteikt funkciju $\rho(u_i)$ ar svērtās Tukey (*Tukey method*) izteiksmes palīdzību:

$$\frac{u_i^2}{2} - \frac{u_i^4}{2c^2} + \frac{u_i^6}{6c^4} \cdot |u_i| \geq c, \quad \rho(u_i) = \frac{c^2}{6}, \quad |u_i| > c, \quad (3.4.19)$$

kur $c = 4.685$.

Lai noteiktu β vērtību, ir jāatrisina vienādojums (3.4.17). Lai to izdarītu, ir jāizmanto funkcijas $\rho(u_i)$ atvasinājums $\Psi = \rho'(u_i)$ un vērtības x_{ij} - i -tā novērojuma j -tais parametrs (atribūts). Tad

$$\sum_{i=1}^n x_{ij} \Psi \left(\frac{y_i - \sum_{i=1}^n x_{ij} \beta_j}{S_m} \right) = 0. \quad (3.4.20)$$

Lai nodrošinātu vienādojuma (3.4.20) atrisinājumu, svara funkciju $w(e_i)$ definē šādi:

$$w(e_i) = \frac{\Psi\left(\frac{y_i - x_i\beta}{S_m}\right)}{\left(\frac{y_i - x_i\beta}{S_m}\right)}. \quad (3.4.21)$$

Tā kā $u_i = \frac{e_i}{S_m}$, tad vienādojumu (3.4.21) var attēlot veidā

$$\left(1 - \left(\frac{u_i}{c}\right)^2\right)^2, |u_i| \leq c, \quad w_i = w(e_i) = 0, \quad |u_i| > c, \quad (3.4.22)$$

kur $c = 4.685$.

Tagad vienādojums (3.4.20) iegūst šādu veidu

$$\sum_{i=1}^n x_{ij} w_i (y_i - x_i' \beta) = 0. \quad (3.4.23)$$

Matricas formā vienādojumu (3.4.23) var attēlot sekojošā veidā:

$$\beta_{l+1} = \mathbf{X}^T \mathbf{W}_0 \mathbf{X} \beta = \mathbf{X}^T \mathbf{W}_0 \mathbf{Y}, \quad (3.4.24)$$

kur \mathbf{W}_0 - diagonālā svara matrica izmērā $n \times n$;

\mathbf{X} - faktoriālo vērtību matrica ($n \times (p+1)$). (Priekš viendimensiju lineārās regresijas $p = 1$).

Vērtības β aprēķināšanas procedūras tiek iteratīvi atkārtotas atbilstoši vienādojumam

$$\beta_{l+1} = (\mathbf{X}^T \mathbf{W}_l \mathbf{X}) (\mathbf{X}^T \mathbf{W}_l \mathbf{Y}). \quad (3.4.25)$$

Vienādojumā (3.4.24) \mathbf{W}_0 ir svara matrica, kurā i -tās diagonāles elements ir vienāds ar w_{i0} . Šo vienādojumu sauc par svēro mazāko kvadrātu vienādojumu.

Rodas jautājums: kā var noteikt regresijas koeficienta sākotnējo vērtību $\hat{\beta}^{(0)}$? Parasti pirmais solis ir regresijas koeficienta noteikšana, izmantojot standarta lineārās regresijas metodi:

1. Novērtēt regresijas koeficienta $\hat{\beta}^{(0)}$ sākotnējo vērtību, izmantojot standarta lineārās regresijas metodi.
2. Noteikt novirzes vērtības $e_i = y_i - \hat{y}_i$.
3. Aprēķināt S_m vērtību, izmantojot vienādojumu (3.4.16).
4. Aprēķināt $u_i = \frac{e_i}{S_m}$ vērtības, izmantojot vienādojumu (3.4.15).
5. Aprēķināt svaru $w(e_i)$ vērtības, izmantojot vienādojumu (3.4.21).
6. Aprēķināt $\hat{\beta}^{(l)}$ vērtību, izmantojot mazāko kvadrātu svēro metodi, izmantojot svarus $w(e_i)$, kas aprēķināti pēc iepriekšējā procedūras.
7. Atkārti 2. – 6. procedūras, izmantojot iepriekšējā iterācijā iegūto $\hat{\beta}^{(l)}$ vērtību.

Par algoritma apstāšanās kritēriju kalpo starpība starp secīgām $\hat{\beta}^{(l+1)} - \hat{\beta}^{(l)}$ vērtībām. Ja šī starpība nepārsniedz *a priori* noteikto robežvērtību, algoritms pārstāj darboties.

4. NEPĀRTRAKTU ATRIBŪTU VĒRTĪBU DISKRETIZĀCIJA

4.1. Kas ir diskretizācija?

Lai izprastu nepārtrauktu atribūtu vērtību diskretizācijas nozīmi, mērķus un uzdevumus, sadaļas sākumā par šo tēmu sniedzam citātus no diviem avotiem.

[Liu H., et. al., 2002]. "Dati parasti parādās jauktā formātā: nominālā, diskrētā un/vai nepārtrauktā formātā. Diskrētie un nepārtrauktie dati ir kārtas datu tipi ar attiecībām starp to vērtībām, un nominālās vērtības tiem neuzliek nekādas atkarības. Diskrētās vērtības ir intervāli nepārtrauktā vērtību spektrā. Lai gan atribūta nepārtraukto vērtību skaits var būt bezgalīgi liels, diskrēto vērtību skaits bieži ir ierobežots un mazs. Divu veidu vērtības rada atšķirības apmācības un klasifikācijas kokos. Viens lēmumu koka izvades piemērs var vēl vairāk ilustrēt atšķirību starp diviem datu veidiem. Kad tiek veidots lēmumu koks, tiek atlasīta viena pazīme, lai izveidotu šķautni tās vērtībām. Ja ir saderīgas nepārtrauktas un atsevišķas funkcijas, ir normāli izvēlēties nepārtrauktu pazīmi, jo tai ir vairāk iezīmju... Izvēloties nepārtrauktu atribūtu nākamajā koka līmenī, var ātri sasniegt "tīru stāvokli" ar visiem piemēriem, kas pieder vienai klasei... Tas noved pie sliktas klasifikatora veikspējas. Tāpēc nav prātīgi izmantot nepārtrauktas vērtības, lai sadalītu lēmumu koka mezglu. Nepārtrauktas pazīmes ir nepieciešams diskretizēt vai nu pirms lēmumu koka izveides vai koka konstruēšanas procesā... Diskrētu vērtību izmantošanai salīdzinājumā ar nepārtrauktām vērtībām ir daudz citu priekšrocību. Diskrētās pazīmes ir tuvāk zināšanu reprezentācijas līmenim nekā nepārtrauktas. Datus var arī samazināt un vienkāršot, izmantojot sākotnējo izlasi. Gan lietotājiem, gan ekspertiem atsevišķas funkcijas ir vieglāk saprotamas, lietojamas un izskaidrojamas. Veiksmīga paraugu ņemšana var ievērojami paplašināt daudzu apmācības algoritmu robežas".

[Hemada B., Lakshmi K.S.V., 2013]. "...lielu skaitu statistikas un mašīnmācīšanās metožu var izmantot datu kopām, kas pilnībā sastāv no nominālvērtībām. Tomēr ļoti liela daļa reālo datu kopu ietver nepārtrauktus mainīgos... Viens problēmas risinājums ir sadalīt skaitliskos mainīgos vairākos intervālos un uzskatīt katru intervālu kā kategoriju. Šo nepārtraukto mainīgo sadalīšanas procesu kategorijās parasti sauc par diskretizāciju.

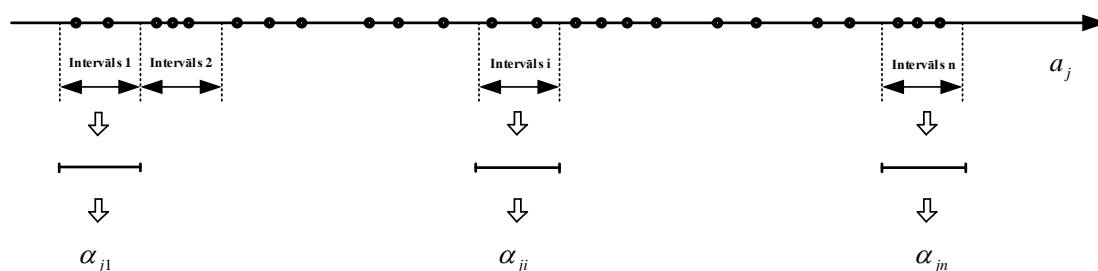
Nepārtraukta atribūtu izlase ir svarīga datu iegūšanas un mašīnmācīšanās algoritmu priekšapstrādes pieeja. Efektīva izlases metode var ne tikai samazināt prasības sistēmas modelim un uzlabot mašīnmācīšanās algoritmu efektivitāti, bet arī padarīt no diskretizētās datu kopas iegūtās zināšanas kompaktākas, vieglāk saprotamas un lietojamas... Diskretizācijas rezultāts ir ne tikai saistīts ar pašu algoritmu, bet arī ar datu sadalījumu un sadalīto punktu skaitu. Ja dažādām datu kopām tiek lietots viens un tas pats algoritms, mēs varam iegūt atšķirīgus rezultātus. Paraugu ņemšanas metodes efektivitāti varam uzzināt no pēcapstrādes rezultātiem. Tas, vai izlases metode ir laba, ir saistīts arī ar vēlāk izmantoto secinājumu algoritmu.

Diskrētu vērtību izmantošanai ir daudz priekšrocību salīdzinājumā ar nepārtrauktajām vērtībām: (1) Diskretizācija samazina nepārtraukto atribūtu vērtību skaitu, kā rezultātā samazinās sistēmas atmiņas prasības. (2) Diskrētās pazīmes ir tuvāk zināšanu līmeņa atspoguļošanai nekā nepārtrauktas. (3) Gan lietotājiem, gan ekspertiem atsevišķas funkcijas ir vieglāk saprotamas, lietojamas un izskaidrojamas. (4) Diskretizācija padara apmācību precīzāku un ātrāku".

Šie citāti parāda nepārtrauktu atribūtu vērtību diskretizācijas metožu būtību. Nosacītā atribūta a_j nepārtraukto vērtību diskretizācijas process shematiski parādīts attēlā 4.1.1.

Viss nepārtrauktu vērtību diapazons, atribūtu a_j vērtības, ir sadalīts noteiktā intervālu skaitā n . Ir piedāvāti dažādi principi nepārtraukta lieluma vērtību sadalīšanai intervālos. Daži no pamatprincipiem tiks apspriesti vēlāk šajā nodaļā. Katrs intervāls tiek kodēts noteiktā veidā. Mašīnmācīšanās problēmās parasti ir ierasts kodēt intervālus, izmantojot dažus nominālos

nosaukumus (kategorijas). Šīs kategorijas ir tikai atbilstošo intervālu apzīmējumi un nesatur nekādu citu informāciju.



Attēls 4.1.1. Nepārtrauktas atribūtu a_j vērtības diskretizācijas procesa shematiskais attēlojums

Attēlā rezultējošie intervāli ir apzīmēti ar α_{ji} , $i = 1, \dots, n$. Jebkura nepārtraukta lieluma vērtība, kas ietilpst intervālā α_{ji} , tiek attēlota ar šo iezīmi un visas attiecīgās turpmākās datu apstrādes procedūras ir attēlotas šo intervālu kontekstā.

Statistiskās analīzes problēmās arī tiek izmantots nepārtrauktu lielumu dalījums intervālos. Piemēram, nepārtraukts mainīgais attēlo izlasē iekļauto indivīdu vecumu un pētījuma mērķis ir iegūt sadalījuma funkciju pa indivīdiem atbilstoši viņu vecumam. Lai to izdarītu, vecuma vērtību diapazons tiek sadalīts piemērotā intervālu skaitā un šajos intervālos tiek veiktas turpmākas datu analīzes procedūras. Bet jāpatur prātā, ka intervālu kopai ir šāda secība: jebkuram intervālam to sakārtotajā secībā katra nepārtraukta lieluma vērtība intervālos, kas atrodas pa kreisi no šī intervāla, būs mazāka par šī intervāla kreisās puses robežas vērtību. Tāpat katra nepārtraukta lieluma vērtība intervālos pa labi no noteiktā intervāla būs lielāka par šī intervāla labo robežvērtību.

Statistiskās izpētes kontekstā iegūto intervālu secībai var būt liela nozīme. Mašīnmācīšanās uzdevumu kontekstā (klasifikācija, lēmumu izvade utt.) iegūto intervālu secībai nav nozīmes. Piemēram, mēs zinām no apmācības kopas, ka septiņi objekti, kuriem atribūtu vērtības α_j atrodas intervālā α_{ji} , pieder klasei c_k un viens objekts ar atribūta vērtību α_j no šī intervāla pieder klasei c_l . Pieņemsim, ka klasificētajam objektam ir atribūta vērtība α_j , kas atbilst intervālam α_{ji} . No šiem datiem mēs varam secināt, ka šis objekts, iespējams, ar varbūtību $7/8$ pieder klasei c_k un, iespējams, ar varbūtību $1/8$ pieder klasei c_l .

No tā izriet svarīgs secinājums: lai izveidotu klasifikatora modeli, mūs neinteresē atribūtu vērtības α_j kā tādas, bet gan attiecības starp atribūtu vērtībām noteiktā intervālā un objektu piederību klasēm. Citiem vārdiem sakot, ir svarīga sakarība: “intervāla iezīme – klases iezīme”. Tāpēc intervālu kodēšana klasifikācijas problēmās un citos mašīnmācīšanās uzdevumos tiek veikta nominālajās skalās.

Var dot šādu formālu definīciju: diskretizācija ir atribūta nepārtraukto vērtību diapazona sadalīšana noteiktā intervālu skaitā. Visas atribūtu vērtības šajā intervālā tiek traktētas vienā šī intervāla nominālajā kategorijā. Turpmākās datu apstrādes un analīzes darbības tiek veiktas, izmantojot iegūtās nominālo intervālu iezīmes.

Nepārtrauktu atribūtu vērtību diapazona sadalīšana intervālos nav triviāls uzdevums. Kā norādīts darbā [Kotsiansis S., Kanellopoulos D., 2006], ir jāatrod kompromiss starp informācijas kvalitāti (viendabīgi intervāli attiecībā uz prognozēto atribūtu) un statistisko kvalitāti (pietiekams izlases lielums katrā intervālā, lai nodrošinātu vispārinājumus). Hī-kvadrāta tests koncentrējas uz statistisko priekšstatu, bet uz entropiju balstīts tests koncentrējas uz informācijas priekšstatu. Citi kritēriji mēģina atrast kompromisu starp informācijas un statistikas īpašībām.

4.2. Diskretizācijas metožu klasifikācija

Diskretizācijas metodes var klasificēt pēc daudzām pazīmēm. Pazīme šajā kontekstā ir principi, uz kuriem balstās diskretizācijas metodes. Pašlaik ir vispārpieņemtas šādas pazīmes, kas ir dažādu diskretizācijas metožu pamatā.

1. *Neuzraudzīts - uzraudzīts.*

Neuzraudzītās metodes diskretizācijas veikšanai neizmanto informāciju par klasēm, kurām objekti pieder. Diskretizācija tiek veikta pamatojoties tikai uz dotajām atribūtu vērtību kopām. Tipiski šīs metožu grupas pārstāvji ir diskretizācija pēc vienāda platuma intervāliem un vienāda biežuma intervāliem. *Uzraudzītās* metodes nepārtrauktu atribūtu vērtību diapazonu sadalīšanai izmanto informāciju par objektu piederību noteiktām klasēm.

2. *Globāli - lokāli.*

Globālās metodes veic diskretizāciju pilnām attiecīgo atribūtu nepārtrauktu vērtību kopām. Citiem vārdiem sakot, diskretizācija tiek veikta pilnam objektu kopumam (piemēriem).

Lokālās metodes veic diskretizāciju atsevišķām attiecīgo atribūtu apakškopām, t.i., diskretizācija tiek veikta noteiktai objektu apakškopai (piemēriem).

3. *Diskretizācijas virziens.*

Pēc šīs pazīmes diskretizācijas metodes iedalās *lejupejošās* metodēs un *augšupejošās* metodēs.

Lejupejošās metodes sākas ar pilnu nepārtrauktu atribūtu vērtību kopu. Pirmais diskretizācijas algoritma solis ir sadalīt visu atribūtu vērtību diapazonu divos intervālos. Intervālu dalīšanas process turpinās secīgi, līdz tiek sasniegts algoritma apstāšanās kritērijs.

Šīs grupas metodēm problēma ir griezuma punktu noteikšana, kas nosaka robežu starp diviem jaunizveidotiem intervāliem. Mēs pievērsīsimies šim jautājumam vēlāk šajā nodaļā.

Augšupejošās metodes sākas ar pieņēmumu, ka katrai nepārtrauktai atribūta vērtībai ir atsevišķs intervāls. Algoritma sākotnējā solī tiek apvienotas atsevišķas vērtības un tiek veidoti daži sākotnējie intervāli. Intervālu sapludināšanas process turpinās, līdz tiek sasniegts algoritma apstāšanās kritērijs.

Šīs grupas metodēm problēma ir izvēlēties piemērotu nosacījumu, lai atlasītu pašreizējos intervālus kā sapludināšanas kandidātus. Mēs pievērsīsimies šim jautājumam vēlāk šajā nodaļā.

4. *Statisks - dinamisks.*

Statiskās metodes veic nepārtrauktu atribūtu vērtību diskretizāciju pirms datu analīzes algoritma piemērošanas. Citiem vārdiem sakot, uz šīm metodēm balstīta datu izlase ir neatņemama datu pirmapstrādes sastāvdaļa.

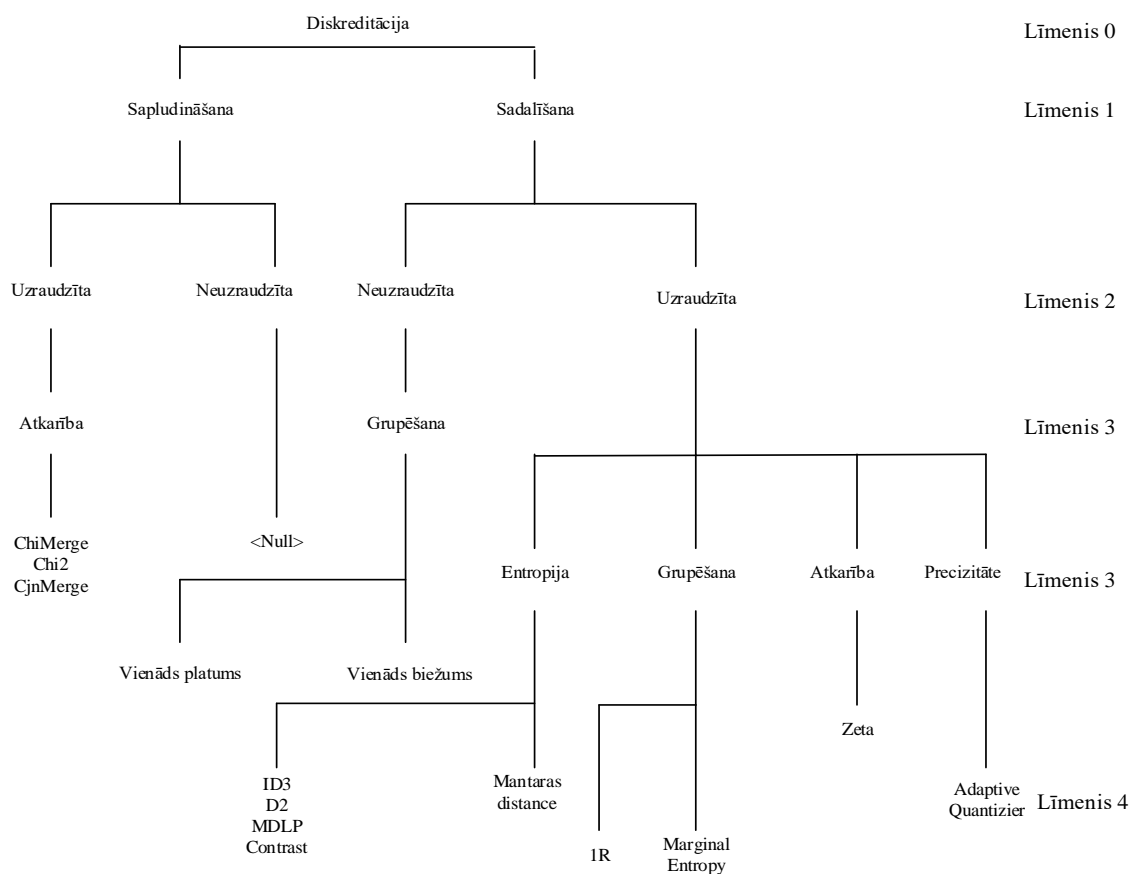
Dinamiskās metodes veic nepārtrauktu atribūtu vērtību diskretizāciju datu apstrādes laikā. Piemēram, izmantojot populāro C4.5 klasifikācijas algoritmu, attiecīgo atribūtu nepārtrauktās vērtības tiek diskretizētas klasifikatora modeļa veidošanas laikā.

5. *Tieša - pieaugoša.*

Tiešās metodes sadala nepārtrauktu atribūtu vērtību diapazonus *apriori* noteiktā intervālu skaitā. Diskretizācija notiek, pamatojoties uz vienāda platuma vai vienāda biežuma intervāliem.

Pieaugošajām metodēm ir raksturīgs tas, ka diskretizācija sākas no kāda attiecīgā atribūta vērtību sākuma stāvokļa, piemēram, no viena intervāla ar *lejupejošām* metodēm vai no vairākiem intervāliem ar vienu atribūta vērtībām *augšupejošās* metodēs. Diskretizācijas procesa laikā pašreizējie intervāli tiek sadalīti ar *lejupejošām* metodēm vai apvienoti intervālos ar *augšupejošām* metodēm. Abos gadījumos ir nepieciešams apstāšanās kritērijs, lai noteiktu, kad diskretizācijas procesam jābeidzas.

Pamatojoties uz iepriekš dotajām pazīmēm, var izveidot vispārīgu klasifikācijas shēmu konkrētām diskretizācijas metodēm. Pamatā ir klasifikācijas shēma, kas piedāvāta darbā [Liu H. et al., 2002]. Šī diagramma ar dažām izmaiņām ir parādīta attēlā 4.2.1.



Attēls 4.2.1. Diskretizācijas metožu hierarhiskā klasifikācijas shēma

Šī hierarhiskā diagramma satur 4 līmeņus. 1. līmenī diskretizācijas metodes tiek iedalītas divās lielās grupās: *Sapludināšana* un *Sadalīšana*, atkarībā no tā, kurš no intervālu veidošanas principiem ir konkrētās metodes pamatā.

2. līmenī gan saplušināšanas, gan sadalīšanas metodes tiek iedalītas uzraudzītajās un neuzraudzītajās.

3. līmenī ir norādīti konkrētie intervālu veidošanas principi.

Struktūras galīgā pozīcija (4. līmenis) atspoguļo konkrētās izmantotās diskretizācijas metodes.

Lai noteiktu klasifikācijas pazīmes, kas raksturo konkrēto diskretizācijas metodi, ir nepieciešams detalizētāk apskatīt dotās hierarhiskās struktūras elementus.

4.3. Diskretizācijas process

Nepārtraukto atribūtu vērtību diskretizācijas procedūru secība shematiski parādīta attēlā 4.3.1 [Liu H. et al., 2002].

Mēs sniegsim attiecīgo procedūru aprakstu un raksturojumu.

1. Šķirošana.

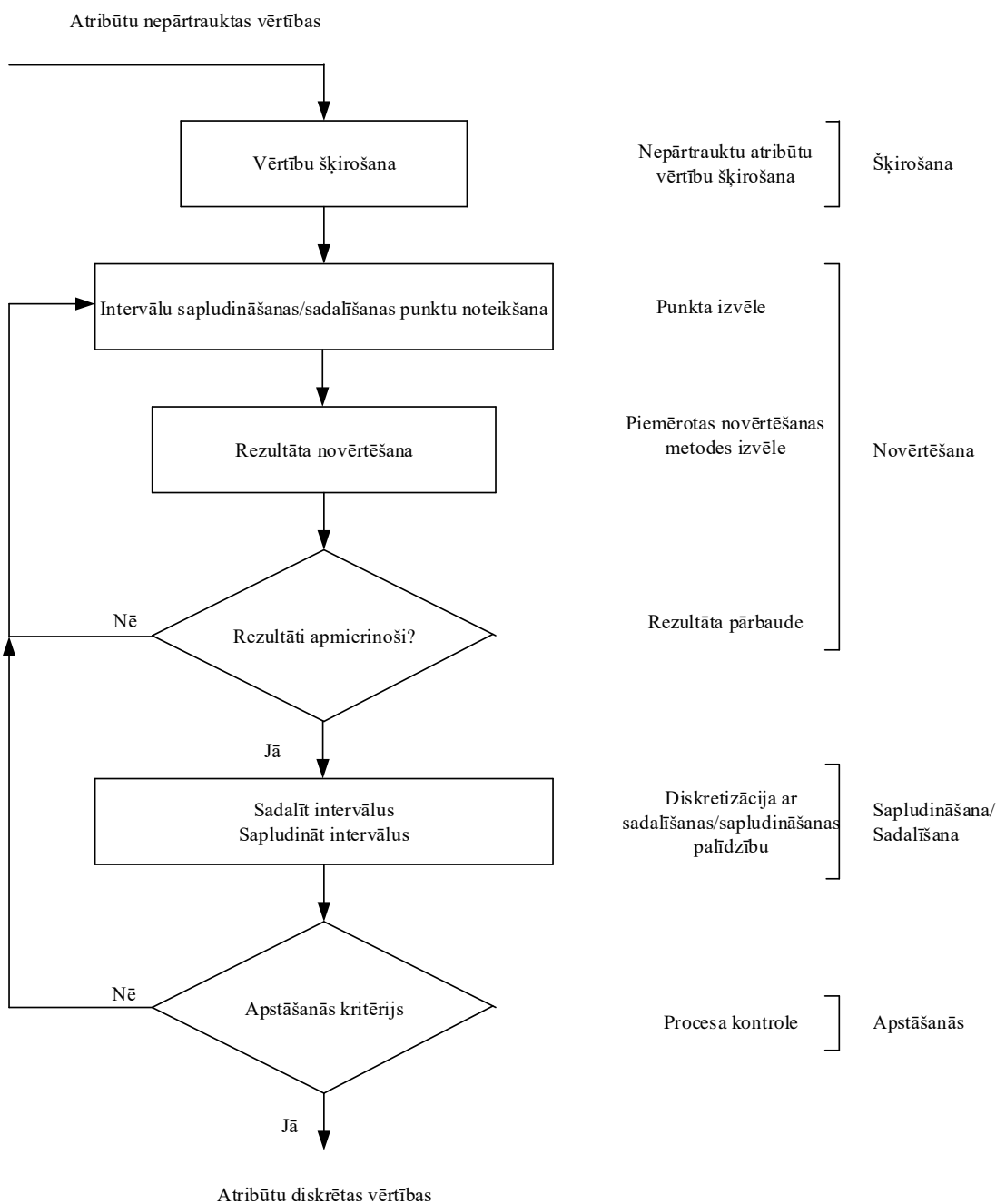
Pirms diskretizācijas procedūru veikšanas nepārtrauktas atribūtu vērtības jāsakārto augošā vai dilstošā secībā.

2. Intervālu sapludināšanas vai sadalīšanas punktu noteikšana.

Ja tiek izmantotas lejupejošas diskretizācijas metodes, tad nepārtrauktu atribūtu vērtību diapazonā vai vienā no iepriekš izveidotajiem intervāliem ir jānosaka griezuma punkta kandidāts. Ja tiek izmantota augšupejoša diskretizācijas metode, ir jāidentificē sapludināšanai blakus esošu intervālu kandidātu pāris.

3. Rezultātu novērtēšana.

Sapludināšanas/sadalīšanas rezultāts ir jāpārbauda saskaņā ar piemērotu kritēriju. (Daži biežāk izmantojamie intervālu sapludināšanas/sadalīšanas kritēriji tiks parādīti turpmākajās nodaļas sadaļās.) Ja pašreizējais intervālu sadalījuma novērtēšanas rezultāts ir pozitīvs, griezuma punkts tiek fiksēts un iegūtie intervāli tiek pakļauti turpmākai sadalīšanai. Ja sadalīšanas rezultāts neatbilst izmantotajam kritērijam, tiek izvēlēts cits griezuma punkts un process tiek atkārtots, līdz tiek sasniegts pozitīvs rezultāts.



Attēls 4.3.1. Nepārtrauktu atribūtu diskretizācijas procedūru vispārīgā shēma

Sapludinot intervālus, tiek pārbaudīta arī jaunizveidotā intervāla atbilstība kvalitātes kritērijam. Ja pārbaudes rezultāts ir pozitīvs, sapludināšanas process turpinās. Ja sapludināšanas rezultāts ir neapmierinošs, apvienošanai tiek atlasīts cits blakus esošu intervālu pāris. Process tiek atkārtots, līdz tiek identificēts piemērots intervālu pāris sapudināšanai.

4. Procesu apturēšana.

Intervālu sapludināšanas/sadalīšanas process apstājas, kad tiek sasniegts apstāšanās kritērijs. Kopumā apstāšanās kritērija formulēšana ir neformāls uzdevums. Šis kritērijs ir noteikts konkrētā problēmā, pamatojoties uz kompromisu starp nepārtraukto vērtību skaitu intervālos un intervāla iezīmes un klases iezīmes atbilstību iegūtajiem intervāliem.

Atbilstoša apstāšanās kritērija noteikšana, it īpaši, ja katram atribūtam ir liels skaits nepārtrauktu vērtību, ir sarežģīts uzdevums. Šīs problēmas risināšana ir iespējama tikai pamatojoties uz esošo datu rūpīgu analīzi. Vienīgais, ko var droši apgalvot, ka problēmai nav globāli optimāla risinājuma. No tā izriet, ka ir jāpanāk piemērots kompromiss starp pretrunīgajām prasībām gan attiecībā uz diskretizācijas rezultātiem, gan attiecībā uz pašas problēmas risināšanas rezultātiem.

Iepriekš aprakstītās diskretizācijas procedūras tiek veiktas katram atribūtam ar nepārtrauktām vērtībām. Jāatzīmē, ka patlaban ir izstrādātas metodes, kas ļauj vienlaicīgi veikt visu nepārtrauktu atribūtu vērtību diskretizāciju. Bet šīs metodes ir diezgan sarežģītas un paliek ārpus šīs nodaļas tēmas.

4.4. Vienkāršākās diskretizācijas metodes

Šajā sadaļā ir apskatītas populārākās uzraudzītās/neuzraudzītās metodes nepārtrauktu atribūtu vērtību diskretizācijai.

1. Diskretizācija ar vienāda platuma intervāliem.

Šīs metodes ideja ir sadalīt nepārtraukto atribūtu a_j vērtību diapazonu *a priori* uzdotā vienāda platuma intervālu skaitā. Uzdotam k intervālu skaitam katra intervāla platums tiek aprēķināts kā

$$\text{Intervāla platums} = \frac{a_{j\max} - a_{j\min}}{k}, \quad (4.4.1)$$

kur $a_{j\max}$ - atribūta a_j maksimālā vērtība;

$a_{j\min}$ - atribūta a_j minimālā vērtība.

Sākotnējā atribūtu vērtību diapazona griezum punkti ir definēti kā

$$\text{Griezuma punkts } i = a_{i\min} + i \cdot (\text{Intervāla platums}). \quad (4.4.2)$$

Apskatīsim vienkāršu ilustratīvu piemēru. Tabulā 4.4.1 ir parādīta sākotnējā nepārtraukto atribūtu a_j vērtību kopa.

Tabula 4.4.1

Nepārtraukto atribūtu a_j vērtību sākotnējā kopa

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8
a_{ij}	22	26	10	16	18	28	12	24

Šīs nepārtrauktās vērtības ir nepieciešams diskretizēt, izmantojot četrus ($k=4$) vienāda platuma intervālus.

Vispirms sakārtosim atribūtu a_j vērtības augošā secībā. Sakārtotās vērtības ir parādītas tabulā 4.4.2.

Tabula 4.4.2

Sakārtotas atribūtu a_j vērtības datiem no tabulas 4.4.1

Objekti	o_3	o_8	o_4	o_5	o_1	o_8	o_2	o_6
a_{ij}	10	12	16	18	22	24	26	28

No šīs tabulas mums ir: $a_{j\min} = 10$, $a_{j\max} = 28$. Aprēķināsim intervālu platumu, izmantojot vienādojumu (4.4.1).

$$\text{Intervāla platums} = \frac{28-10}{4} = \frac{18}{4} = 4.5.$$

Tā kā diskretizācijai mēs izmantojam četrus intervālus, atribūtu a_j vērtību diapazonā būs trīs griezuma punkti. Mēs nosakām šo punktu vērtības, izmantojot vienādojumu (4.4.2).

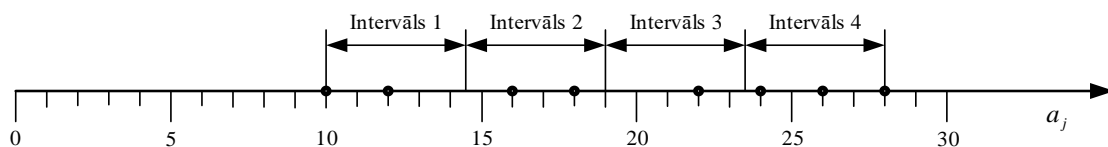
$$\text{Griezuma punkts 1} = 10+1*4,5=14,5;$$

$$\text{Griezuma punkts 2} = 10+2*4,5=19;$$

$$\text{Griezuma punkts 3} = 10+3*4,5=23,5.$$

Šai metodei nav nepieciešams īpašs apstāšanās kritērijs, jo diskretizācijai intervālu skaits ir noteikts *a priori*.

Skaidrības labad iegūtie diskretizācijas rezultāti ir grafiski parādīti attēlā 4.4.1.



Attēls 4.4.1. Diskretizācijas rezultātu grafisks attēlojums nepārtrauktu atribūtu a_j vērtībām no tabulas 4.4.1

Vienāda platuma intervālos var nebūt vienāds datu punktu skaits, kas ietilpst katrā intervālā. Attēlā 4.4.1 ir redzams, ka 1. un 2. intervālā ir divi punkti, 3. intervālā ir viens punkts, 4. intervālā ir trīs punkti.

Būtisks vienāda platuma intervālu diskretizācijas metodes trūkums ir tas, ka intervālu veidošanas procesu var negatīvi ietekmēt anomāliju klātbūtne datos. Lai parādītu šī apgalvojuma pareizību, pieņemsim, ka atribūta vērtība $a_{6j} = 28$ ir kļūdaini uzrakstīta - $a_{6j} = 88$. Ar šiem datiem mums ir:

$$\text{Intervāla platums} = \frac{88-10}{4} = 19.5.$$

$$\text{Griezuma punkts 1} = 10 + 1 * 19,5 = 29,5;$$

$$\text{Griezuma punkts 2} = 10 + 2 * 19,5 = 49.$$

$$\text{Griezuma punkts 3} = 10 + 3 * 19,5 = 68,5.$$

Tagad visi datu punkti, izņemot punktu $a_{6j} = 88$ ietilpst 1. intervālā, 2. un 3. intervālā nav punktu, bet 4. intervālā ir viens datu punkts. Acīmredzot šādai diskretizācijai nav praktiskas nozīmes.

2. Diskretizācija ar vienāda biežuma intervāliem.

Pieņemam, ka mums ir n nepārtrauktu atribūtu a_j vērtības un intervālu skaits tiek norādīts *a priori*. Tad katrā intervālā ir jāietver aptuveni n/k atribūtu vērtības. Šeit jēdziens “aptuveni” ir saistīts ar to, ka pie dotajām n un k vērtībām ne vienmēr ir iespējams iegūt veselu dalījuma vērtību.

Ņemsim par pamatu sakārtotās atribūtu a_j vērtības no tabulas 4.4.2. Mēs diskretizējam šo vērtību kopu, izmantojot četrus vienāda biežuma intervālus. Katrā intervālā ir jābūt $\frac{8}{4} = 2$ vērtībām. Tādējādi mums ir šādi intervāli: [10,12], [16,18], [22,24], [26,28].

Atribūtu vērtības katrā no intervāliem ir diezgan tuvu viena otrai. Tagad pieņemsim, ka vērtības vietā $a_{6j} = 28$ kļūdains ir uzrakstīta vērtība $a_{6j} = 88$. Pie kā novedīs šī kļūda? Acīmredzot, ka pirmajos trīs intervālos nekādas izmaiņas nenotiks. Bet ceturtais intervāls tagad satur vērtības [26,88]. Šķiet, ka tik liela atšķirība starp vērtībām neatbilst citiem datu punktiem. Tāpēc anomālijas datus rada arī nepamatotus diskretizācijas rezultātus. Lai gan anomāliju ietekmes sekas šajā gadījumā nešķiet tik katastrofālas kā izlases gadījumā ar vienāda platuma intervāliem, tomēr tās var ļoti negatīvi ietekmēt turpmākās datu analīzes rezultātus.

Jāņem vērā, ka vienāds datu punktu skaits katrā intervālā nenozīmē vienādu intervālu platumu. Šai metodei nav nepieciešams arī īpašs apstāšanās kritērijs, jo intervālu skaits ir iestatīts *a priori*.

2.1. 1R klasifikācijas metode.

Vienāda platuma un vienāda biežuma intervālu metodes ir tipiskas neuzraudzītas diskretizācijas metožu pārstāvji. Pat ja objekti sākotnējā datu kopā ir saistīti ar dotajām klasēm, izmantojot šīs metodes informācija par objektu klasēm tiek pilnībā ignorēta.

Lai novērstu šo trūkumu [Holte R.C., 1993] ierosināja iepriekš aprakstīto metožu modifikāciju, kas tika nosaukta par 1R metodi. Šīs metodes būtība ir šāda. Ignorējot informāciju par objektu klasēm, tiek diskretizētas nepārtrauktu atribūtu a_j sākotnējās vērtības ar kādu no minētajām divām metodēm. Tad tiek pielāgotas robežas starp intervāliem (griezuma punkti), lai izvairītos no situācijām, kad objekti abos intervālos, kas atrodas blakus robežai starp šiem intervāliem, pieder vienai un tai pašai klasei. Minimālajam datu punktu skaitam jebkurā no intervāliem ir jābūt vismaz 6, izņemot galējo labo intervālu, kurā var būt patvaļīgs datu punktu skaits. Šī prasība nosaka arī algoritma apstāšanās kritēriju.

Apskatīsim vienkāršu ilustratīvu piemēru. Tabulā 4.4.3 ir parādītas nepārtrauktu atribūtu a_j vērtības. Tabulas pēdējā rindā ir norādītas to klašu iezīmes, kurām objekti pieder. Šīs nepārtrauktās vērtības ir nepieciešams diskretizēt ar 1R metodi.

Tabula 4.4.3

Sākotnējā atribūtu a_j vērtību kopa

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8	o_9	o_{10}	o_{11}	o_{12}	o_{13}	o_{14}	o_{15}	o_{16}
a_{ij}	13	5	24	8	17	7	25	20	11	15	19	22	10	12	14	18
Klases	c_2	c_1	c_2	c_1	c_2	c_2	c_2	c_2	c_2	c_2	c_1	c_1	c_1	c_1	c_2	c_1

Sakārtosim atribūtu a_j vērtības augošā secībā, rezultāti ir doti tabulā 4.4.4.

Tabula 4.4.4

Sakārtotas atribūtu a_j vērtības no tabulas 4.4.3

Objekti	o_2	o_6	o_4	o_{13}	o_9	o_{14}	o_1	o_{15}	o_{10}	o_5	o_{16}	o_{11}	o_8	o_{12}	o_3	o_7
a_{ij}	5	7	8	10	11	12	13	14	15	17	18	19	20	22	24	25
Klases	c_1	c_2	c_1	c_1	c_2	c_1	c_2	c_2	c_2	c_2	c_1	c_1	c_2	c_1	c_2	c_2

Izmantojot tabulas 4.4.4 datus, mēs sadalām datu punktus (atribūtu a_j vērtības un to atbilstošos objektus) divos vienāda biežuma intervālos. Kreisajā intervālā ir objekti $o_2, o_6, o_4, o_{13}, o_9, o_{14}, o_1, o_{15}$ ar atribūtu a_j vērtībām no 5 līdz 14, labajā intervālā ir objekti $o_{10}, o_5, o_{16}, o_{11}, o_8, o_{12}, o_3, o_7$ ar atribūtu a_j vērtībām no 15 līdz 25.

Analizēsim robežu starp intervāliem. Kreisā intervāla labajā pusē ir divi objekti (o_1 un o_{15}), kas pieder klasei c_2 , labā intervāla kreisajā pusē ir objekti o_{10} un o_5 , kas arī pieder klasei c_2 . Saskaņā ar 1R metodes prasībām, ir jāpielāgo robeža starp intervāliem. Ņemsim vērā, ka lielākā daļa objektu kreisajā intervālā pieder klasei c_1 un lielākā daļa objektu labajā intervālā pieder klasei c_2 . Tāpēc šķiet loģiski pārvietot robežu starp intervāliem divas pozīcijas pa kreisi. Pēc tam jauninātajā kreisajā intervālā būs objekti $o_2, o_6, o_4, o_{13}, o_9, o_{14}$ ar atribūtu a_j vērtībām no 5 līdz 12, bet labajā intervālā būs objekti $o_1, o_{15}, o_{10}, o_5, o_{16}, o_{11}, o_8, o_{12}, o_3, o_7$ ar atribūtu a_j vērtībām no 13 līdz 25.

Mēs neveiksim nekādus citus robežu pielāgojumus starp intervāliem, jo labajā pusē tagad ir tikai seši objekti ar atbilstošām atribūtu a_j vērtībām. Apstāšanās kritērijs ir sasniegts.

4.5. Diskretizācija uz entropijas novērtējuma pamata

Datu klasifikācijas kontekstā entropija ir informācijas apjoma novērtējums, uz kura pamata tiek konstruēta klasifikatora shēma. Plaši zināmajā ID3 datu klasifikācijas algoritmā [Quinlan J.R., 1986] entropijas aplēses tiek izmantotas, lai noteiktu zarus konstruējamā klasifikācijas koka mezglos.

Nepārtraukto atribūtu vērtību diskretizācijas sakarā autori [Fayyad U., Irani K., 1993] ierosināja izmantot entropijas aplēses, lai veidotu diskretizējamus intervālus.

Parasti datu kopas S entropijas aprēķins tiek definēts kā

$$H(S) = \sum_{i=1}^m p_i \log_2 p_i, \quad (4.5.1)$$

kur p_i - varbūtība, ka dotais objekts no kopas S pieder klasei c_i , $i = 1, \dots, m$.

Apzīmēsim ar n objektu skaitu kopā S un ar n_i , $i = 1, \dots, m$, objektu apakškopas no S , kas pieder klasei c_i . Tad varbūtības p_i vērtību vienādojumā (4.5.1) aprēķina kā

$$p_i = \frac{n_i}{n}, \quad i = 1, \dots, m. \quad (4.5.2)$$

Vienādojumā (4.5.1) \log_2 ir logaritms ar bāzi 2. Šādu logaritmu izmantošana ir saistīta ar faktu, ka informācija tiek kodēta bits.

Entropijas aplēses var atrasties intervālā $[0, \log_2 m]$, kur m ir klašu skaits datu kopā.

Minimālā vērtība $H = 0$ tiek sasniegta, ja visi objekti sākotnējā datu kopā pieder vienai un tai pašai klasei i . Tad $p_i = 1$ un $\log_2 1 = 0$. Šādā hipotētiskā situācijā nav nepieciešams izveidot klasifikācijas modeli, jo nav vajadzīga objektu klasifikācija.

Maksimālā entropijas vērtība $H = \log_2 m$ tiek sasniegta hipotētiskā situācijā, kad visi sākotnējās datu kopas objekti ir vienmērīgi sadalīti pa visām m klasēm. Šī ir datu kopas maksimālā iespējamā nenoteiktība, pie kuras nav iespējams pareizs klasifikatora modelis.

Apskatāmā diskretizācijas algoritma ideja ir šāda. Nepārtraukto atribūtu a_j vērtības tiek sakārtotas augošā vai dilstošā secībā. Tādējādi tiek izveidota jauna datu kopas objektu secība, kas atbilst sakārtotajām atribūtu vērtībām. Iegūtajā objektu secībā to sākotnējā sadalīšana tiek veikta tā,

lai kreisajā intervālā būtu viens objekts, kas atbilst mazākajai atribūta a_j vērtībai, bet labajā intervālā ir visi pārējie objekti.

Pēc tam kreisajam intervālam secīgi tiek pievienots viens objekts. Objektu skaits attiecīgajā labajā intervālā tiek samazināts par vienu. Galīgajā sadalījumā būs viens objekts ar vismazāko atribūta a_j vērtību labajā intervālā un visi pārējie objekti kreisajā intervālā.

Faktiski katrā iterācijā griezuma punkts nobīdās pa labi par vienu datu punktu.

Katram sadalījuma variantam tās entropijas novērtējums tiek aprēķināts, izmantojot vienādojumu

$$H(l) = -p_{left} \sum_{i=1}^m p_{i,left} \log_2 p_{i,left} - p_{right} \sum_{i=1}^m p_{i,right} \log_2 p_{i,right} . \quad (4.5.3)$$

Šajā vienādojumā:

$l = 1, \dots, n - 1$ - pašreizējā griezuma punkta numurs (iterācijas numurs);

p_{left} - varbūtība, ka l -tās iterācijas reizē objekts nokļūst kreisajā intervālā. Praksē vērtību p_{left} aprēķina kā objektu skaita kreisajā intervālā attiecību pret kopējo objektu skaitu;

p_{right} - varbūtība, ka l -tās iterācijas reizē objekts nokļūst labajā intervālā. Praksē vērtību p_{right} aprēķina kā objektu skaita labajā intervālā attiecību pret kopējo objektu skaitu;

$p_{i,left}$ - varbūtība, ka objekts kreisajā intervālā pieder klasei c_i . Šīs varbūtības vērtību aprēķina kā objektu skaita piederību klasei c_i kreisajā intervālā attiecībā pret kopējo objektu skaitu šajā intervālā;

$p_{i,right}$ - varbūtība, ka objekts labajā intervālā pieder klasei c_i . Šīs varbūtības vērtību aprēķina kā objektu skaita piederību klasei c_i labajā intervālā attiecībā pret kopējo objektu skaitu šajā intervālā;

Aprēķini, izmantojot vienādojumu (4.5.3), tiek veikti visiem alternatīvajiem griezuma punktiem. Efektīvais griezuma punkts ir punkts l^* , kas dod viszemāko entropijas aprēķinu $H(l^*)$ l^* -tajā iterācijā.

Algoritma pirmā pilnā soļa izpildes rezultātā sākotnējā objektu kopa tiek sadalīta divos intervālos. Ja nepieciešams, diskretizācija turpinās iegūtajos intervālos un intervālos, kas iegūti algoritma otrās un turpmākās darbības rezultātā.

Šī algoritma apstāšanās kritēriju var norādīt divos alternatīvos veidos:

1. *A priori* ir norādīts galīgs intervālu skaits.
2. Kā apstāšanās kritērijs tiek izmantots minimālā apraksta garuma princips (*minimum description length principle*).

Mēs neizmantosim minimālā apraksta garuma principu tā matemātiskās sarežģītības dēļ. Sīkāka informācija par šo principu atrodama darbos [Grünwald P.D., 2007; Rissanens J., 2007].

Apskatīsim vienkāršu ilustratīvu piemēru. Tabulā 4.5.1 ir parādīta sākotnējā datu kopa: objektu kopa un šo objektu atribūtu a_j vērtību kopa. Tabulas 3. rindā ir norādītas klases, kurām pieder attiecīgie objekti.

Tabula 4.5.1

Sākotnējo datu kopa

Objekti	o_1	o_2	o_3	o_4	o_5	o_6	o_7	o_8
a_{ij}	7	21	10	12	8	14	20	18
Objektu klases	c_1	c_2	c_1	c_1	c_2	c_1	c_2	c_1

Nepārtraukto atribūtu a_j vērtības ir nepieciešams diskretizēt divos intervālos, izmantojot entropijas aplēses, lai noteiktu efektīvo griezuma punktu.

Sakārtosim atribūtu a_j vērtības augošā secībā. Datu kopa ar sakārtotām atribūtu vērtībām ir parādīta tabulā 4.5.2.

Tabula 4.5.2

Datu kopa no tabulas 4.5.1 ar sakārtotām atribūtu a_j vērtībām

Objekti	o_1	o_5	o_3	o_4	o_6	o_8	o_7	o_2
a_{ij}	7	8	10	12	14	18	20	21
Objektu klases	c_1	c_2	c_1	c_1	c_1	c_1	c_2	c_2

Izpildīsim diskretizācijas algoritma darbības. Sadalīsim sākotnējo objektu kopu divos intervālos.

1. iterācija.

- kreisais intervāls: objekts o_1 ;
- labais intervāls: objekti $o_5, o_3, o_4, o_6, o_8, o_7, o_2$.

(Intervālos iekļāvām nevis atribūtu a_j vērtības, bet objektus ar atbilstošām atribūtu vērtībām. Lieta tāda, ka diskretizējot nepārtrauktas atribūtu vērtības ar šeit piedāvāto metodi, mūs interesē tikai objekti un to piederība attiecīgajām klasēm. Atribūtu a_j vērtības kalpo tikai datu sakārtošanai. Visi turpmākie aprēķini ietver tikai objektu skaitu iegūtajos intervālos un objektu skaitu katrā no klasēm. Pēc diskretizācijas procesa pabeigšanas vienkārši jāpāriet no objektiem uz atbilstošo atribūtu vērtībām).

Apskatīsim kreiso intervālu. Tā kā tajā ir tikai viens objekts o_1 , ir skaidrs, ka $p_{left} = \frac{1}{8} = 0.125$.

Tā kā objekts o_1 pieder klasei c_1 , $p_{1,left} = 1$, $\log_2 1 = 0$, un izteiksmes (4.5.3.) labās puses kreisā summa ir vienāda ar 0.

Apskatīsim labo intervālu. Tajā ir 7 objekti, tāpēc $p_{right} = \frac{7}{8} = 0.875$. Četri objekti šajā intervālā pieder klasei c_1 un 3 objekti pieder klasei c_2 . No šejienes $p_{1,right} = \frac{4}{7} = 0.571$.

$$p_{2,right} = \frac{3}{7} = 0.429.$$

Veiksim pilnu aprēķinu, izmantojot vienādojumu (4.5.3).

$$\begin{aligned} H(1) &= -0 - 0.875 * (0.571 * \log_2 0.571 + 0.429 * \log_2 0.429) = \\ &= -0.875 * (0.571 * (-0.8084) + 0.429 * (-1.2209)) = \\ &= -0.875 * (-0.4616 - 0.5238) = -0.875 * (-0.9854) = 0.8622. \end{aligned}$$

2. iterācija.

- kreisais intervāls: objekti o_1, o_5 ;
- labais intervāls: objekti $o_3, o_4, o_6, o_8, o_7, o_2$.

Tagad kreisajā intervālā ir 2 objekti, bet labajā - 6 objekti. Tāpēc

$$p_{left} = \frac{2}{8} = 0.250; p_{right} = \frac{6}{8} = 0.750.$$

Viens objekts kreisajā intervālā pieder klasei c_1 un viens objekts pieder klasei c_2 . Tāpēc

$$p_{1,left} = \frac{1}{2} = 0.500; p_{2,left} = \frac{1}{2} = 0.500.$$

Četri objekti labajā intervālā pieder klasei c_1 un divi objekti pieder klasei c_2 . Tāpēc

$$p_{1,right} = \frac{4}{6} = 0.667; p_{2,right} = \frac{2}{6} = 0.333.$$

Aprēķināsim sadalījuma entropijas novērtējumu, izmantojot vienādojumu (4.5.3).

$$\begin{aligned} H(2) &= -0.250 * (0.500 \log_2 0.500 + 0.500 \log_2 0.500) - \\ &\quad - 0.750 * (0.667 \log_2 0.667 + 0.333 \log_2 0.333) = \\ &= -0.250 * (0.500 * (-1) + 0.500 * (-1)) - 0.750 * (0.667 * (-0.5842) + 0.333 * (-1.5864)) = \\ &= -0.250 * (-0.500 - 0.500) - 0.750 * (-0.3896 - 0.5827) = \\ &= -0.250 * (-1) - 0.750 * (-0.9723) = 0.2500 + 0.7299 = 0.9799. \end{aligned}$$

Pārējie aprēķini tiek veikti pēc analogijas. Aprēķinu rezultāti ir apkopoti tabulā 4.5.3.

Tabula 4.5.3

Entropijas aprēķinu rezultāti potenciālajiem griezumam punktiem

Iterācijas numurs (l)	1	2	3	4	5	6	7
H(l)	0.8622	0.9799	0.9502	0.9056	0.7946	0.9371	0.8622

Kā izriet no šīs tabulas datiem, zemākais entropijas novērtējums tika sasniegts algoritma 5. iterācijā. Tāpēc ir iegūts šāds objektu sadalījums:

- kreisais intervāls: objekti o_1, o_5, o_3, o_4, o_6 ;
- labais intervāls: objekti o_8, o_7, o_2 .

Šī diskretizācija jāuzskata par veiksmīgu, jo kreisajā intervālā tikai objekts o_5 pieder klasei c_2 , visi pārējie objekti pieder klasei c_1 . Labajā intervālā tikai objekts o_8 pieder klasei c_1 , visi pārējie objekti pieder klasei c_2 .

Mūsu galvenais mērķis šajā piemērā ir diskretizēt nepārtraukto atribūtu a_j vērtības. Izmantojot iepriekš minēto objektu sadalījumu intervālos un datus no tabulas 4.5.2, katrā no intervāliem ir šādas atribūtu a_j vērtības:

- kreisais intervāls: 7, 8, 10, 12, 14 $\Rightarrow \alpha_{j_1} = [7, 14]$;
- labais intervāls: 18, 20, 21 $\Rightarrow \alpha_{j_2} = [18, 21]$.

Ja mums ir jauns objekts ar atribūta vērtību $a_j = 11$, tad ar lielu varbūtības pakāpi varam teikt, ka šis objekts pieder klasei c_1 .

4.6. Uz χ^2 statistiku balstītas metodes

χ^2 - tas ir statistiskais novērtējums, ar kura palīdzību tiek veikts nozīmības tests attiecībā "atribūtu vērtības - objektu piederība klasēm". Uz χ^2 statistiku balstīto diskretizācijas metožu būtība ir tāda, ka blakus intervālos, kas veidojas kādā algoritma iterācijā, tiek noteikts, vai divas kategorijas var tikt apvienotas vai nē, pamatojoties uz to biežumu sadalījumu datu kopā. Ja atribūtu vērtības divos blakus intervālos ir neatkarīgas no pazīmju klases iezīmēm, tad šādus intervālus var apvienot vienā kopējā intervālā. Pretējā gadījumā intervālu apvienošana nav iespējama.

Kā piemēru nepārtrauktu atribūtu vērtību diskretizācijai, pamatojoties uz χ^2 statistiku, apskatīsim ChiMerge metodi [Kerber R., 1992]. To var raksturot kā augšupejošu uzraudzītu metodi.

Kā šī metode darbojas? Dota datu kopa, kas sastāv no n objektiem, no kuriem katrs tiek novērtēts ar atbilstošo atribūta a_j vērtību. Turklāt visiem objektiem ir norādītas klašu iezīmes.

Sākotnēji tiek pieņemts, ka katrs datu punkts veido vienu intervālu, t.i., sākotnējais intervālu skaits ir vienāds ar objektu skaitu datu kopā.

Vispirms tiek veidoti primārie intervāli ar objektiem, kas pieder tikai vienai klasei.

Tad pirmajā iterācijā katram blakus esošu intervālu pārim apēķina χ^2 vērtību:

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^p \frac{(A_{ij} - E_{ij})^2}{E_{ij}}, \quad (4.6.1)$$

kur p - klašu skaits;

A_{ij} - dažādu klašu iezīmju i -tajam intervālam un j -tajai klasei;

Lai noteiktu vērtību E_{ij} vienādojumā (4.6.1), mēs ieviešam šādus apzīmējumus:

R_i - objektu skaits i -tajā intervālā, $R_i = \sum_{i=1}^p A_{ij}$;

C_j - objektu skaits j -tajā klasē, $C_j = \sum_{j=1}^m A_{ij}$;

N - kopējais objektu skaits, $N = \sum_{i=1}^p C_j = \sum_{i=1}^m R_i$.

Tad E_{ij} vērtība tiek interpretēta kā sagaidāmais biežums

$$E_{ij} = \frac{R_i C_j}{N}. \quad (4.6.2)$$

χ^2 vērtības tiek aprēķinātas katram blakus esošu intervālu pārim, kas iegūti sākotnējā algoritma iterācijā. Pēc tam iegūtās aplēses tiek salīdzinātas savā starpā. Blakus esošu objektu pāris, kuriem ir viszemākais χ^2 vērtējums, tiek sapludināts. Tādējādi tiek izveidots jauns intervālu kopums.

Procesu iteratīvi atkārto, līdz tiek sasniegts apstāšanās kritērijs. Apstāšanās kritēriju var norādīt divos alternatīvos veidos:

1. Diskretizācijas intervālu skaits ir *apriori* zināms.
2. Tiek noteikta χ^2 robežvērtība. (Ieteicams uzdot šo robežvērtību intervālā $[0.90, 0.99]$).

Apskatīsim ilustratīvu piemēru. Par pamatu ņemsim tabulas 4.5.1 datus. Tabulā 4.5.2 sakārtotie dati ir atkārtoti tabulā 4.6.1.

Tabula 4.6.1

Sakārtotie dati no tabulas 4.5.2

Objekti	o_1	o_5	o_3	o_4	o_6	o_8	o_7	o_2
a_{ij}	7	8	10	12	14	18	20	21
Klases	c_1	c_2	c_1	c_1	c_1	c_1	c_2	c_2

Izmantojot ChiMerge metodi, nepārtraukto atribūtu a_j vērtības ir jādiskretizē divos intervālos.

Sākumā pieņemam, ka katrs intervāls ietver vienu datu objektu. Tādējādi mums ir sākotnējais 8 intervālu kopums. Veiksim objektu sadalīšanu. Definēsim primāros intervālus tā, lai katrā intervālā būtu objekti, kas pieder vienai klasei. Mums ir

Intervāls 1: objekts o_1 ;

Intervāls 2: objekts o_5 ;

Intervāls 3: objekti o_3, o_4, o_6 ;

Intervāls 4: objekti o_8, o_7 ;

Intervāls 5: objekts o_2 .

(Tāpat kā iepriekšējā sadaļā, veicot attiecīgos aprēķinus, darbosimies ar intervāliem, kas satur objektus, nevis atbilstošās atribūtu a_j vērtības. Tas tiek darīts, lai uzskatāmāk atspoguļotu aprēķinu procesus un jaunu intervālu veidošanos. Pāreja uz atribūtu vērtību intervāliem var viegli tikt izpildīta pēc visu algoritma iterāciju pabeigšanas).

Lai noteiktu intervālu pāri - sapludināšanas kandidātus, mēs aprēķinām χ^2 katram intervālu pārim.

1. iterācija.

- intervālu pāris $Int.1, Int.2$.

Lai atvieglotu aprēķinus, mēs izveidojam tabulu 4.6.2.

Tabula 4.6.2

Darba tabula aprēķinu veikšanai

	1. klase ($j = 1$)	2. klase ($j = 2$)	Σ pa rindām
$Int.1(i = 1)$	A_{11}	A_{12}	R_1
$Int.2(i = 2)$	A_{21}	A_{22}	R_2
Σ pa kolonnām	C_1	C_2	N

Tabulā ir divas kolonas, jo sākotnējā datu kopā ir divas objektu klases. Ja ir lielāks klašu skaits, kolonu skaitam jāatbilst objektu klašu skaitam.

Šajā tabulā:

A_{11} - c_1 klasei piederošo objektu skaits 1. intervālā;

A_{12} - c_2 klasei piederošo objektu skaits 1. intervālā;

A_{21} - c_1 klasei piederošo objektu skaits 2. intervālā;

A_{22} - c_2 klasei piederošo objektu skaits 2. intervālā;

R_1, R_2 - skaitļu summa tabulas rindās;

C_1, C_2 - skaitļu summas tabulas ailēs;

N - visu skaitļu summa tabulas rindās vai kolonnās.

Aizpildīsim šo tabulu.

Tabula 4.6.3

Darba tabula aprēķinu veikšanai

	Klase c_1 ($j=1$)	Klase c_2 ($j=2$)	Σ pa rindām
$Int.1(i=1)$	$A_{11} = 1$	$A_{12} = 0$	$R_1 = 1$
$Int.2(i=2)$	$A_{21} = 0$	$A_{22} = 1$	$R_2 = 1$
Σ pa kolonām	$C_1 = 1$	$C_2 = 1$	$N = 2$

Izmantojot vienādojumu (4.6.2), mēs aprēķinām E_{ij} vērtības:

$$E_{11} = \frac{R_1 C_1}{N} = \frac{1 \cdot 1}{2} = 0.500; \quad E_{12} = \frac{R_1 C_2}{N} = \frac{1 \cdot 1}{2} = 0.500;$$

$$E_{21} = \frac{R_2 C_1}{N} = \frac{1 \cdot 1}{2} = 0.500; \quad E_{22} = \frac{R_2 C_2}{N} = \frac{1 \cdot 1}{2} = 0.500.$$

Izmantojot vienādojumu (4.6.1), mēs aprēķinām $\chi^2_{1,2}(1)$ vērtību. Šeit apakšindekss 1,2 nozīmē, ka mēs runājam par 1. un 2. intervālu sākotnējā intervālu secībā. Skaitlis (1) nozīmē 1. iterāciju.

$$\begin{aligned} \chi^2_{1,2}(1) &= \frac{(A_{11} - E_{11})^2}{E_{11}} + \frac{(A_{12} - E_{12})^2}{E_{12}} + \frac{(A_{21} - E_{21})^2}{E_{21}} + \frac{(A_{22} - E_{22})^2}{E_{22}} = \\ &= \frac{(1 - 0.500)^2}{0.500} + \frac{(0 - 0.500)^2}{0.500} + \frac{(0 - 0.500)^2}{0.500} + \frac{(1 - 0.500)^2}{0.500} = \\ &= 0.500 + 0.500 + 0.500 + 0.500 = 2.000. \end{aligned}$$

- intervālu pāris $Int.2, Int.3$.

Tabula 4.6.4

Darba tabula aprēķinu veikšanai

	c_1 ($j=1$)	c_2 ($j=2$)	Σ
$Int.2(i=1)$	0	1	1
$Int.3(i=2)$	3	0	3
Σ	3	1	4

(Mēs izmantojam indeksus $i=1$ priekš $Int.2$ un $i=2$ priekš $Int.3$, lai veiktu aprēķinus līdzīgā veidā.)

$$E_{11} = \frac{1 \cdot 3}{4} = 0.750; \quad E_{12} = \frac{1 \cdot 1}{4} = 0.250; \quad E_{21} = \frac{3 \cdot 3}{4} = 2.250; \quad E_{22} = \frac{3 \cdot 1}{4} = 0.750.$$

$$\begin{aligned} \chi^2_{2,3}(1) &= \frac{(0 - 0.750)^2}{0.750} + \frac{(1 - 0.250)^2}{0.250} + \frac{(3 - 2.250)^2}{2.250} + \frac{(0 - 0.750)^2}{0.750} = \\ &= 0.750 + 0.250 + 2.250 + 0.750 = 4.000. \end{aligned}$$

- intervālu pāris $Int.3, Int.4$.

Tabula 4.6.5

Darba tabula aprēķinu veikšanai

	c_1 ($j=1$)	c_2 ($j=2$)	Σ
$Int.3(i=1)$	3	0	3
$Int.4(i=2)$	0	2	2
Σ	3	2	5

$$E_{11} = \frac{3*3}{5} = 1.800; E_{12} = \frac{3*2}{5} = 1.200; E_{21} = \frac{2*3}{5} = 1.200; E_{22} = \frac{2*2}{5} = 0.800.$$

$$\chi_{3,4}^2(1) = \frac{(3-1.800)^2}{1.800} + \frac{(0-1.200)^2}{1.200} + \frac{(0-1.200)^2}{1.200} + \frac{(2-0.800)^2}{0.800} = 0.800 + 1.200 + 1.200 + 1.800 = 5.000.$$

- intervālu pāris *Int.4, Int.5*.

Tabula 4.6.6

Darba tabula aprēķinu veikšanai

	$c_1 (j=1)$	$c_2 (j=2)$	Σ
<i>Int.4(i=1)</i>	0	2	2
<i>Int.5(i=2)</i>	1	0	1
Σ	1	2	4

$$E_{11} = \frac{2*1}{3} = 0.667; E_{12} = \frac{2*2}{3} = 1.333; E_{21} = \frac{1*1}{3} = 0.333; E_{22} = \frac{1*2}{3} = 0.667.$$

$$\chi_{4,5}^2(1) = \frac{(0-0.667)^2}{0.667} + \frac{(2-1.333)^2}{1.333} + \frac{(1-0.333)^2}{0.333} + \frac{(0-0.667)^2}{0.667} = 0.667 + 0.333 + 1.336 + 0.667 = 3.000.$$

Vismazākā no visām aprēķinātajām aplēsēm ir $\chi_{1,2}^2(1) = 2.000$. Tāpēc ir nepieciešams sapludināt pašreizējos intervālus *Int.1, Int.2*. Rezultātā ir šādi starpintervāli:

2. iterācija.

Int.1(2): objekti o_1, o_5 ;

Int.2(2): objekti o_3, o_4, o_5 ;

Int.3(2): objekti o_8, o_7 ;

Int.4(2): objekts o_2 .

Šo intervālu pāru χ^2 aprēķini tiek veikti pēc analogijas ar iepriekš sniegtajiem aprēķiniem. Aprēķinu rezultāti:

$$\textit{Int.1(2), Int.2(2)}: \chi_{1,2}^2(2) = 1.875;$$

$$\textit{Int.2(2), Int.3(2)}: \chi_{2,3}^2(2) = 5.000;$$

$$\textit{Int.3(2), Int.4(2)}: \chi_{3,4}^2(2) = 2.940.$$

Mazākā no šīm vērtībām ir $\chi_{1,2}^2(2) = 1.875$. Tāpēc ir nepieciešams sapludināt intervālus. Rezultātā mums ir šādi trīs nākošie starpintervāli:

3. iterācija.

Int.1(3): objekti o_1, o_5, o_3, o_4, o_6 ;

Int.2(3): objekti o_8, o_7 ;

Int.3(3): objekts o_2 .

$$Int.1(3), Int.2(3): \chi_{1,2}^2(3) = 3.750;$$

$$Int.2(3), Int.3(3): \chi_{2,3}^2(2) = 5.000.$$

No tā izriet, ka $Int.1(3), Int.2(3)$ viņiem ir jābūt sapludinātiem. Rezultātā mums ir šādi divi intervāli:

$Int.1$: objekti o_1, o_5, o_3, o_4, o_6 ;

$Int.2$: objekti o_8, o_7, o_2 .

Tā kā tika iegūti divi rezultējošie intervāli, tika sasniegts apstāšanās kritērijs un algoritma darbība beidzas.

Pārejot no objektiem uz atbilstošajām nepārtraukto atribūtu a_j vērtībām, mums ir

$Int.1$ (kreisais intervāls): 7, 8, 10, 12, 14 $\Rightarrow \alpha_1 = [7, 14]$;

$Int.2$ (labais intervāls): 18, 20, 21 $\Rightarrow \alpha_2 = [18, 21]$.

Jāatzīmē, ka šie intervāli ir tādi paši kā intervāli, kas iegūti no tiem pašiem sākuma datiem iepriekšējās sadaļas piemērā, izmantojot pilnīgi citu diskretizācijas metodi. Šī sakritība norāda uz abu metožu efektivitāti.

5. NORMALIZĀCIJA UN ATRIBŪTU VĒRTĪBU STANDARTIZĀCIJA

5.1. Definīcijas un piezīmes

Jebkuras datu pirmapstrādes procedūras mērķis ir sagatavot datus un pārveidot tos formā, kas ir vispiemērotākā šo datu tālākai analīzei.

Ieejas datus atribūtu vērtības var izmērīt dažādās skalās un to vērtības dažādiem atribūtiem var būt ļoti atšķirīgas. Pieņemsim, ka starp dažādām iedzīvotāju grupām tiek veikti socioloģiski pētījumi. Izmantotie atribūti ietver personu vecumu un gada ienākumus. Individu vecums grupās ir ierobežots intervālā [25–75] gadi. Savukārt individu gada ienākumi var svārstīties ļoti plašās robežās gan katras grupas ietvaros, gan starp grupām.

Turpmāk pieņemsim, ka pētījuma nolūkos tiek mēģināts identificēt individu apakšgrupas (klasterus), pamatojoties uz abu atribūtu vērtībām. Klasterizāciju var veikt ar k-tuvāko kaimiņu metodi vai k-vidējais metodi. Izmantojot katru no šīm metodēm, ir jānovērtē attālumi starp punktiem (individu) datu telpā. Visbiežāk šim nolūkam tiek izmantots Eiklīda attālums. Ja tiek izmantotas sākotnējās atribūtu vērtības, tad attiecīgo attālumu aprēķinos lielā mērā dominēs gada ienākumu vērtības, jo šīs vērtības vienmēr ir daudz lielākas nekā atbilstošās individu vecuma vērtības.

Arī daudzās citās datu analīzes metodēs, mašīnmācīšanās un neironu tīklos nevar tiešā veidā izmantot datu sākotnējās atribūtu vērtības.

Kritiski svarīga ir kritēriju vērtību normalizācija lēmumu pieņemšanas problēmās. Fakts ir tāds, ka daudzas metodes daudzkritēriju lēmumu izvēlei izmanto lineāras kritēriju vērtību kombinācijas katram no alternatīvajiem risinājumiem. Plaši pazīstamā TOPSIS metode izmanto aprēķinus par attālumiem no katra alternatīvā risinājuma līdz “ideāli pozitīvajam” un “ideāli negatīvajam”. Abos gadījumos iegūtajos attāluma aprēķinos dominēs lielas kritēriju vērtības.

Lai risinātu problēmas, kas saistītas ar neapstrādātām atribūtu vai kritēriju vērtībām, šīs vērtības ir jāpārvērš piemērotā formā.

Izdarīsim ļoti svarīgu piezīmi. Ir trīs galvenās pieejas oriģinālo atribūtu vērtību pārveidošanai kompaktā bezdimensiju formā, proti, normalizēšana, standartizācija un transformācija. Diemžēl, šie jēdzieni literatūrā bieži tiek lietoti kā sinonīmi. Tas rada neskaidrības un šo svarīgo jēdzienu nepareizu lietošanu. Šajā darbā mēs izmantojam šādas definīcijas.

Normalizēšana ir atribūta sākotnējo vērtību pārveidošana bezdimensiju formā, izmantojot īpašas paša atribūta vērtības vai šo vērtību kombinācijas, vai kādu no attiecīgā atribūta vērtību sadalījuma parametriem normalizēšanai.

Standartizācija ir sākotnējo atribūtu vērtību pārveidošana bezdimensiju standarta formā, izmantojot tikai attiecīgo atribūtu vērtību sadalījuma parametrus.

Transformācija ir sākotnējo atribūtu vērtību pārveidošana citā formā, kuras pamatā ir kāda matemātiska funkcija.

Dažreiz literatūrā tiek izmantots vispārīgāks atribūtu vērtību *mērogošanas jēdziens*, kas tiek saprasts kā atribūtu vērtību pārveidošana vēlamajā bezdimensiju formā.

Šajā nodaļā tiks parādītas labi zināmās metodes sākotnējo atribūtu vērtību normalizācijai un standartizēšanai. 6. nodaļā tiks apskatītas metodes datu atribūtu vērtību transformācijai.

5.2. Sākotnējo atribūtu vērtību normalizācijas metodes

Sagrupēsim normalizācijas metodes, ņemot vērā normalizējošā faktora noteikšanas principus.

1. *Metozu grupa ar vienu normalizējošā faktora vērtību.*

- *normalizācija attiecībā pret maksimālo atribūta vērtību*

$$a_{ij}^n = \frac{a_{ij}}{a_{j\max}}, \quad (5.2.1)$$

kur a_{ij} - objekta o_i sākotnējā atribūta a_j vērtība;

$a_{j\max}$ - maksimālā atribūta a_j vērtība.

Normalizācija ar vienādojuma (5.2.1) palīdzību ir iespējama, ja atribūtu a_j vērtības mēra attiecību skalā. Normalizācijas rezultātus izsaka arī attiecību skalā.

Piemērs. Tabulā 5.2.1 ir parādītas atribūtu a_j sākotnējās vērtības.

Tabula 5.2.1

Sākotnējās atribūtu vērtības a_j

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9

Atribūta maksimālā vērtība $a_{j\max} = 9$. Normalizējot visas sākotnējās vērtības pret šo vērtību, mēs iegūstam atribūtu vērtības, kas parādītas tabulā 5.2.2.

Tabula 5.2.2

Normalizētas atribūtu vērtības a_j no 5.2.1 tabulas

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}^n	0.333	0.556	0.778	0.444	0.889	1.000

Šo normalizācijas metodi raksturo fakts, ka visas normalizētās vērtības ir mazākas par 1. Maksimālajai atribūta vērtībai tā normalizētā vērtība ir 1.

- *normalizācija pret vidējo vērtību*

$$a_{ij}^n = \frac{a_{ij}}{\bar{a}_j}, \quad (5.2.2)$$

kur a_{ij} - objekta o_i sākotnējā atribūta a_j vērtība;

\bar{a}_j - atribūta a_j vidējā vērtība.

Normalizācija ar vienādojuma (5.2.2) palīdzību ir iespējama, ja atribūtu a_j vērtības mēra attiecību skalā. Normalizācijas rezultātus izsaka arī attiecību skalā.

Piemērs. Ņemsim par pamatu sākotnējās atribūtu a_j vērtības no tabulas 5.2.1. Vidējā atribūta a_j vērtība ir

$$\bar{a}_j = \frac{3+5+7+4+8+9}{6} = \frac{36}{6} = 6.$$

Sākotnējās atribūtu a_j vērtības, kas normalizētas ar vienādojumu (5.2.2.), ir parādītas tabulā 5.2.3.

Tabula 5.2.3

Atribūtu vērtību a_j normalizācijas rezultāti tabulas 5.2.1 datiem pēc vienādojuma (5.2.2)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
a_{ij}^n	0.500	0.833	1.167	0.667	1.333	1500

Jāatzīmē, ka, normalizējot ar vienādojumu (5.2.2), normalizētās vērtības sākotnējām vērtībām, kas ir mazākas par tās vidējo vērtību, ir mazākas par 1 un sākotnējām vērtībām, kas lielākas par vidējo vērtību, ir lielākas par 1.

- normalizācija pret atribūtu vērtību standarta novirzi

$$a_{ij}^n = \frac{a_{ij}}{s_j}, \quad (5.2.3)$$

kur a_{ij} - objekta o_i sākotnējā atribūta a_j vērtība;

s_j - atribūtu vērtību standarta novirze a_j no tās vidējās vērtības

$$s_j = \sqrt{\frac{\sum_{i=1}^n (a_{ij} - \bar{a}_j)^2}{n-1}}.$$

Normalizācija ar vienādojuma (5.2.3) palīdzību ir iespējama, ja atribūtu a_j vērtības mēra attiecību skalā. Normalizācijas rezultātus izsaka arī attiecību skalā.

Piemērs. Ņemsim par pamatu sākotnējās atribūtu vērtības no tabulas 5.2.1. Šī atribūta vidējā vērtība ir $\bar{a}_j = 6$. Aprēķināsim standartnovirzes vērtību.

$$\begin{aligned} s_j &= \sqrt{\frac{(3-6)^2 + (5-6)^2 + (7-6)^2 + (4-6)^2 + (8-6)^2 + (9-6)^2}{5}} = \\ &= \sqrt{\frac{9+1+1+4+4+9}{5}} = \sqrt{\frac{28}{5}} = \sqrt{5.600} = 2.366. \end{aligned}$$

Atribūta a_j sākotnējās vērtības un tā normalizētās vērtības ir parādītas tabulā 5.2.4.

Tabula 5.2.4

Atribūtu a_j vērtību normalizācijas rezultāti tabulas 5.2.1 datiem pēc vienādojuma (5.2.3)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
a_{ij}^n	1.268	2.113	2.958	1.691	3.381	3.804

2. *Metožu grupa, kas izmanto normalizēšanu attiecībā pret atribūtu vērtību starpībām.*

- normalizēšana attiecībā pret vērtību starpību $a_{j\max} - a_{j\min}$.

Šai metodei ir divas variācijas:

$$a_{ij}^n = \frac{a_{ij}}{a_{j\max} - a_{j\min}}; \quad (5.2.4)$$

$$a_{ij}^n = \frac{(a_{ij} - a_{j\min})}{(a_{j\max} - a_{j\min})}; \quad (5.2.5)$$

Vienādojumos (5.2.4) un (5.2.5):

a_{ij} - objekta o_i atribūta a_j sākotnējā vērtība;

$a_{j\min}$ - minimālā atribūta a_j vērtība;

$a_{j\max}$ - maksimālā atribūta a_j vērtība.

Normalizācija ar vienādojuma (5.2.4) palīdzību ir iespējama, ja atribūtu a_j vērtības mēra attiecību skalā. Normalizācijas rezultātus izsaka arī attiecību skalā.

Normalizācija ar vienādojuma (5.2.5) palīdzību ir iespējama, ja atribūtu a_j vērtības mēra attiecību skalā vai intervālu skalā. Normalizācijas rezultāti tiek izteikti intervālu skalā.

Piemērs. Ņemsim par pamatu sākotnējās atribūtu vērtības no tabulas 5.2.1.

No šīs tabulas mums ir $a_{j\min} = 3$, $a_{j\max} = 9$, $a_{j\max} - a_{j\min} = 9 - 3 = 6$. Sākotnējo atribūtu vērtību normalizēšanas rezultāti pēc vienādojuma (5.2.4) ir parādīti tabulā 5.2.5.

Tabula 5.2.5

Atribūtu a_j vērtību normalizācijas rezultāti tabulas 5.2.1 datiem pēc vienādojuma (5.2.4)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
a_{ij}^n	0.500	0.833	1.167	0.667	1.333	1.500

Jāatzīmē, ka normalizēšanas rezultāti tabulā 5.2.5 sakrīt ar normalizēšanas rezultātiem attiecībā pret vidējo vērtību tabulā 5.2.3. Šajā gadījumā tā ir sakritība, jo vērtību atšķirība $a_{j\max} - a_{j\min} = 6$ sakrīt ar atribūta vidējo vērtību $\bar{a}_j = 6$.

Normalizēsim atribūta sākotnējās vērtības, izmantojot vienādojumu (5.2.5). Normalizācijas rezultāti ir parādīti tabulā 5.2.6.

Tabula 5.2.6

Atribūtu a_j vērtību normalizācijas rezultāti tabulas 5.2.1 datiem pēc vienādojuma (5.2.5)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
$a_{ij} - a_{j\min}$	0	2	4	1	5	6
a_{ij}^n	0	0.333	0.667	0.167	0.833	1.000

Normalizācijas vienādojuma (5.2.5) raksturīga iezīme ir tāda, ka atribūta normalizētās vērtības vienmēr atrodas intervālā no 0 līdz 1. Pie vērtības $a_{ij} = a_{j\min}$ vienādojuma (5.2.5) skaitītājs ir vienāds ar 0 un normalizācijas rezultāts arī ir vienāds ar 0. Pie vērtības $a_{ij} = a_{j\max}$ vienādojuma (5.2.5) skaitītājs un saucējs ir vienādi un normalizācijas rezultāts ir 1.

3. Normalizācijas metožu grupa attiecībā pret atribūtu vērtību summām.

Pirmā metode šajā grupā ietver sākotnējo atribūtu a_j vērtību normalizēšanu attiecībā pret visu atribūtu vērtību summu.

$$a_{ij}^n = \frac{a_{ij}}{\sum_{i=1}^m a_{ij}}. \quad (5.2.6)$$

Izmantojot šo metodi, atribūtu a_j vērtības jāmēra attiecību skalā. Normalizācijas rezultātus izsaka arī attiecību skalā.

Piemērs. Ņemsim par pamatu sākotnējās atribūtu a_j vērtības no tabulas 5.2.1. No šīs tabulas mums ir

$$\sum_{i=1}^6 a_{ij} = 3 + 5 + 7 + 4 + 8 + 9 = 36.$$

Normalizācijas rezultāti ir parādīti tabulā 5.2.7.

Tabula 5.2.7

Atribūtu vērtību a_j normalizācijas rezultāti tabulas 5.2.1 datiem pēc vienādojuma (5.2.6)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
a_{ij}^n	0.083	0.139	0.194	0.111	0.222	0.250

Otrā metode no šīs grupas ietver sākotnējo atribūtu vērtību normalizēšanu attiecībā pret kvadrātsakni no atribūtu vērtību kvadrātu summas.

$$a_{ij}^n = \frac{a_{ij}}{\sqrt{\sum_{i=1}^m a_{ij}^2}}. \quad (5.2.7)$$

Normalizācija ar vienādojumu (5.2.7) ir iespējama, ja atribūtu a_j vērtības mēra attiecību skalā. Normalizācijas rezultātus izsaka arī attiecību skalā.

Piemērs. Ņemsim par pamatu sākotnējās atribūtu a_j vērtības no tabulas 5.2.1. No šīs tabulas mums ir

$$\sum_{i=1}^6 a_{ij}^2 = 3^2 + 5^2 + 7^2 + 4^2 + 8^2 + 9^2 = 9 + 25 + 49 + 16 + 64 + 81 = 244.$$

$$\sqrt{\sum_{i=1}^6 a_{ij}^2} = \sqrt{244} = 15.620.$$

Normalizācijas rezultāti ir parādīti tabulā 5.2.8.

Tabula 5.2.8

Atribūtu vērtību a_j normalizācijas rezultāti tabulas 5.2.1 datiem pēc vienādojuma (5.2.7)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
a_{ij}^n	0.192	0.320	0.448	0.256	0.512	0.576

Pastāv arī citas metodes atribūtu vērtību normalizēšanai. Normalizējot kritēriju vērtības lēmumu pieņemšanas problēmās, ir jāņem vērā tas, ka daži kritēriji atspoguļo ieguvumus, bet citi - izmaksas (zaudējumus), kas saistītas ar alternatīvu lēmumu rezultātiem.

Šajā nodaļā šīs citas metodes netiek aplūkotas, jo galvenā uzmanība tiek pievērsta datu pirmapstrādes metodēm. Interesentiem ieteicams meklēt darbu [Vafaei N. et al., 2018], lai iegūtu detalizētāku informāciju par vērtēšanas kritēriju normalizēšanas metodēm lēmumu pieņemšanas problēmās.

5.3. Sākotnējo atribūtu vērtību standartizācijas metodes

Visizplatītākā ir standartizācijas metode, kuras pamatā ir atribūtu vērtību sadalījuma parametri: vidējā vērtība un standartnovirze:

$$a_{ij}^s = \frac{a_{ij} - \bar{a}_j}{s_j}, \quad (5.3.1)$$

kur a_{ij} - objekta o_i sākotnējā atribūta a_j vērtība;

\bar{a}_j - atribūta a_j vidējā vērtība;

s_j - atribūtu vērtību a_j standarta novirze no tās vidējās vērtības.

Standartizāciju ar vienādojumu (5.3.1) var veikt, ja atribūtu a_j vērtības mēra attiecību skalā vai intervālu skalā. Standartizācijas rezultāti tiek izteikti intervālu skalā.

Standartizācijas ar vienādojumu (5.3.1) raksturīga iezīme ir tāda, ka standartizēto atribūtu vērtību sagaidāmā vērtība ir 0 un to standartnovirze ir 1.

Piemērs. Ņemsim par pamatu sākotnējās atribūtu a_j vērtības no tabulas 5.2.1. Iepriekšējā sadaļā tika aprēķināti šādi tā vērtību sadalījuma parametru aprēķini:

$$\bar{a}_j = 6; s_j = 2.366.$$

Standartizēsim atribūtu a_j vērtības pēc vienādojuma (5.3.1). Standartizācijas rezultāti ir parādīti tabulā 5.3.1.

Tabula 5.3.1

Atribūtu vērtību a_j standartizācijas rezultāti tabulas 5.2.1 datiem tabulas pēc vienādojuma (5.3.1)

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
$a_{ij} - \bar{a}_j$	-3	-1	1	-2	2	3
a_{ij}^s	-1.268	-0.423	0.423	-0.845	0.845	1.268

Alternatīva metode atribūtu vērtību standartizēšanai ir standartizācija, kuras pamatā ir atribūtu vērtību sadalījuma mediāna un novirze no mediānas (Vēbera standartizācija):

$$a_{ij}^s = \frac{(a_{ij} - \text{med}(a_j))}{1.4826 \text{mad}(a_j)}, \quad (5.3.2)$$

kur a_{ij} - objekta o_i sākotnējā atribūta a_j vērtība;

$\text{med}(a_j)$ - atribūtu vērtību a_j sadalījuma mediāna;

$\text{mad}(a_j)$ - atribūtu vērtību a_j noviržu kopas mediāna no mediānas $\text{med}(a_j)$.

Standartizāciju pēc vienādojuma (5.3.2) var veikt, ja atribūtu a_j vērtības mēra attiecību skalā vai intervālu skalā. Standartizācijas rezultāti tiek izteikti intervālu skalā.

Vērtību $\text{med}(a_j)$ aprēķina šādi. Vispirms n atribūtu vērtības a_{ij} jāsakārto un jāindeksē augošā secībā: $a_{(1)j}, a_{(2)j}, \dots, a_{(n)j}$. Vērtība $\text{med}(a_j)$ tiek aprēķināta, izmantojot vienādojumus

$$\text{med}(a_j) = a_{(n+1)j/2}, \quad n - \text{nepāra}; \quad (5.3.2.a)$$

$$\text{med}(a_j) = \frac{1}{2}(a_{(k)j} + a_{(k+1)j}), \quad n = 2k - \text{pāra}. \quad (5.3.2.b)$$

Ja n ir nepāra skaitlis, tad par $\text{med}(a_j)$ vērtību tiek ņemta atribūta a_j vērtība, kas atrodas $(n+1)/2$ pozīcijā tā vērtību sakārtotajā secībā. Ja n ir pāra skaitlis, tad vērtība $\text{med}(a_j)$ tiek ņemta par vidējo vērtību no atribūta a_j vērtībām, kas atrodas k -tajā un $(k+1)$ pozīcijā vērtību sakārtotajā secībā.

Pēc vērtības $\text{med}(a_j)$ aprēķināšanas, izmantojot vienādojumus (5.3.2.a) vai (5.3.2.b), tiek aprēķinātas atribūtu vērtību a_j novirzes no vērtības $\text{med}(a_j)$. Vērtība $\text{mad}(a_j)$ tiek aprēķināta, izmantojot vienādojumus (5.3.2.a) vai (5.3.2.b), bet atribūta vērtību a_j vietā tagad tiek izmantotas šī atribūta novirzes vērtības no mediānas $\text{med}(a_j)$.

Piemērs. Par pamatu ņemsim atribūtu a_j vērtības no tabulas 5.2.1. Lai aprēķinātu $\text{med}(a_j)$, mēs sakārtojam šīs vērtības augošā secībā:

3, 4, 5, 7, 8, 9.

Tā kā mums ir pāra skaits atribūtu vērtību - 6, mēs aprēķinām vērtību $\text{med}(a_j)$, izmantojot vienādojumu (5.3.2.b).

$$\text{med}(a_j) = \frac{1}{2}(a_{(3)j} + a_{(4)j}) = \frac{1}{2}(5 + 7) = 6.$$

Aprēķini par atribūtu a_j vērtību novirzēm no mediānas $\text{med}(a_j)$ ir parādīti tabulā 5.3.2.

Tabula 5.3.2

Atribūtu vērtību a_j noviržu no mediānas $\text{med}(a_j)$ aprēķini tabulas 5.2.1 datiem

a_{ij}	3	4	5	7	8	9
$a_{ij} - \text{med}(a_j)$	3	2	1	1	2	3

Sakārtosim novirzes vērtības augošā secībā: 1, 1, 2, 2, 3, 3.

Mēs aprēķinām vērtību $mad(a_j)$, izmantojot vienādojumu (5.3.2.b).

$$mad(a_j) = \frac{1}{2}(a_{(3)j} - med(a_j) + a_{(4)j} - med(a_j)) = \frac{1}{2}(2 + 2) = 2.$$

Vienādojuma (5.3.2) saucējs ir vienāds ar $1.4826 * 2 = 2.9652$.

Standartizēsim atribūtu a_j vērtības pēc vienādojuma (5.3.2). Starppaprēķinu rezultāti un standartizētās atribūtu vērtības ir parādītas tabulā 5.3.3.

Tabula 5.3.3

Starppaprēķinu rezultāti un standartizētās atribūtu a_{ij} vērtības tabulas 5.2.1 datiem

Objekti	o_1	o_2	o_3	o_4	o_5	o_6
a_{ij}	3	5	7	4	8	9
$a_{ij} - med(a_{ij})$	-3	-1	1	-2	2	3
a_{ij}^s	-1.012	-0.337	0.337	-0.674	0.674	1.012

6. ATRIBŪTU VĒRTĪBU TRANSFORMĀCIJA

6.1. Kad un kāpēc tiek izmantotas transformācijas?

Datu analīzes veikšana, izmantojot sākotnējās atribūtu vērtības, dažkārt var būt problemātiska dažādās īpašās situācijās. Pieņemsim, ka atribūtu a_j vērtības atrodas intervālā $[10, 10000]$. Acīmredzot darboties ar šādām vērtībām ir ļoti neērti. Risinājums var būt šī atribūta vērtību transformēšana citā piemērotā mērogā. Ja sākotnējo atribūtu a_{ij} vērtību vietā izmantosim vērtības $\lg(a_{ij})$, tad jaunajā skalā atribūtu vērtības atradīsies intervālā $[1, 4]$. Ja decimāllogaritmā vietā mēs izmantojam naturālos logaritmus, tad jaunās atribūtu vērtības atradīsies intervālā $[2.303, 9.210]$. Jebkurā gadījumā darbība ar transformētām atribūtu vērtībām būs vienkāršāka un ērtāka nekā darbība ar sākotnējām vērtībām.

Šāda veida atribūtu vērtību transformēšana bieži tiek izmantota praksē. Piemēram, ūdeņraža jonu koncentrācija (pH) vienmēr tiek izteikta logaritmiskā skalā tieši šī parametra sākotnējo vērtību lielās izkliedes dēļ.

Bieži ikdienā mēs izmantojam dažādas transformācijas, pat nedomājot par to. Ja mēs mainām noteiktu eiro summu uz Polijas zlotiem, mēs formāli veicam vienā skalā mērītas naudas summas (eiro) lineāru pārveidošanu līdzvērtīgā citā skalā (zlotos).

Datu pirmapstrādes un analīzes procesos tiek izmantotas sarežģītākas atribūtu vērtību nelineāras transformācijas. Kādi ir šādu pārveidojumu mērķi? Kā problēmas ieskatu apskatīsim dažus ilustratīvus piemērus.

Statistisko parametru datu analīze balstās uz *a priori* pieņēmumiem par datu struktūru un īpašībām. Viens no šādiem pieņēmumiem ir tāds, ka attiecīgo atribūtu vērtības ir normāli sadalītas. Ja tas tā patiešām ir, tad daudzas statistiskās analīzes procedūras var veikt standarta veidā: analizēt atribūtu vidējās vērtības, novērtēt atribūtu vērtību izplatību attiecībā pret to vidējām vērtībām, konstruēt ticamības intervālus, analizēt variācijas starp paraugiem utt.

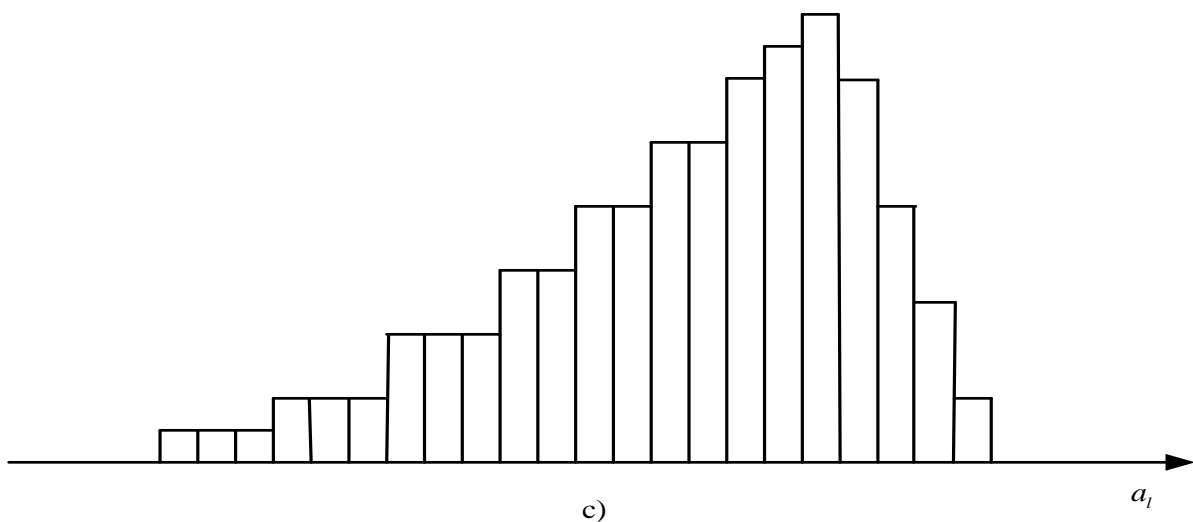
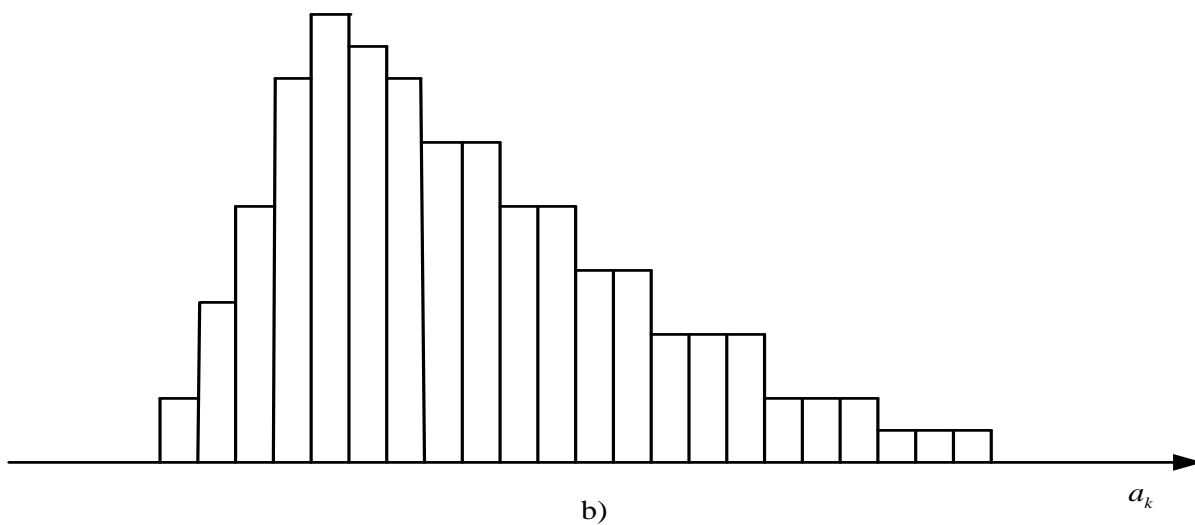
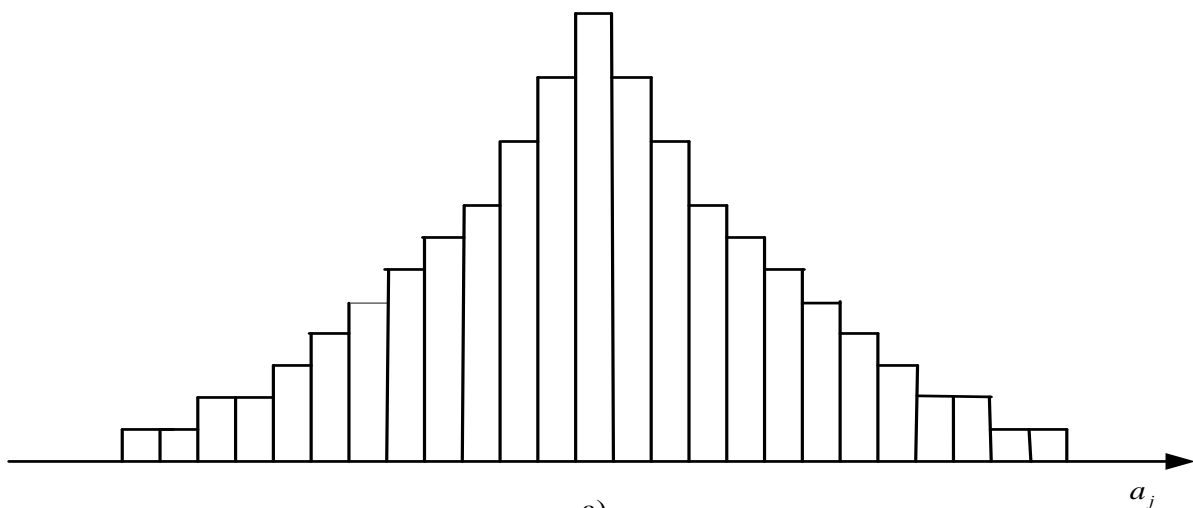
Attēlā 6.1.1 parādītas trīs nosacītās atribūtu a_j , a_k , a_l vērtību sadalījuma histogrammas.

Histogramma attēlā 6.1.1.a attēlo to, ka atribūtu a_j vērtības ir sadalītas saskaņā ar likumu, kas ir tuvs normālam sadalījumam. Šīm vērtībām var piemērot jebkuru mūs interesējošo parametru statistiskās analīzes procedūru.

Histogramma attēlā 6.1.1.b norāda uz to, ka palielinoties atribūtu a_k vērtībām, to parādīšanās biežums sākotnēji strauji palielinās. Pēc maksimālās vērtības sasniegšanas notiek pakāpeniska frekvenču samazināšanās. Rezultātā mums ir asimetrisks slīpums, uz labo pusi vērsts atribūtu a_k vērtību sadalījums.

Grafiks attēlā 6.1.1.c norāda uz to, ka sākumā atribūtu a_l vērtību parādīšanās biežums lēnām palielinās. Pēc maksimālās vērtības sasniegšanas notiek strauja frekvenču samazināšanās. Rezultātā mums ir asimetrisks slīpums, uz kreiso pusi vērsts sadalījums.

Ir acīmredzams, ka atribūtu a_k , a_l vērtību sadalījums ir tālu no normālā sadalījuma un tiem nevar piemērot statistiskās analīzes procedūras, kas ir derīgas normāliem sadalījumiem.



Attēls 6.1.1. Atribūtu a_j , a_k , a_l sadalījumu histogrammas

Ja mums ir histogrammas, kas līdzīgas tām, kas parādītas attēlā 6.1.1, tad veicot šo histogrammu vizuālo analīzi, var izdarīt aptuvenus secinājumus par šo sadalījumu būtību. Reālās situācijās attiecīgajiem atribūtiem var būt ļoti liela vērtību kopa, kas attēlota datu tabulas (matricas) kolonnās. Kā šādus atribūtu vērtību sadalījumus var pārbaudīt, vai tie ir normāli? Pielikumā P3 ir parādītas izmantojamās parametriskās un neparametriskās metodes, lai pārbaudītu sadalījumu atbilstību normālajam sadalījumam.

Kādi rīcības varianti ir iespējami gadījumos, kad attiecīgā atribūtu vērtību sadalījums ir līdzīgs attēlos 6.1.1.b vai 6.1.1.c redzajiem? Darbā [Sakia R., 1992] tiek piedāvātas šādas iespējamās darbības iespējas:

- (i) Ignorējiet pieņēmumu pārkāpumus un veiciet analīzi tā, it kā visi pieņēmumi būtu izpildīti.
- (ii) Izlemiet, kāds ir pareizais pieņēmums pārkāptā vietā un izmantojiet procedūru, kurā tiek ņemts vērā jaunais pieņēmums.
- (iii) Izstrādājiet jaunu modeli, kam ir svarīgi sākotnējā modeļa aspekti un kas atbilst pieņēmumiem, piemēram, izmantojot efektīvu datu transformāciju vai filtrējiet dažus aizdomīgus punktus, kurus var uzskatīt par anomālijām.

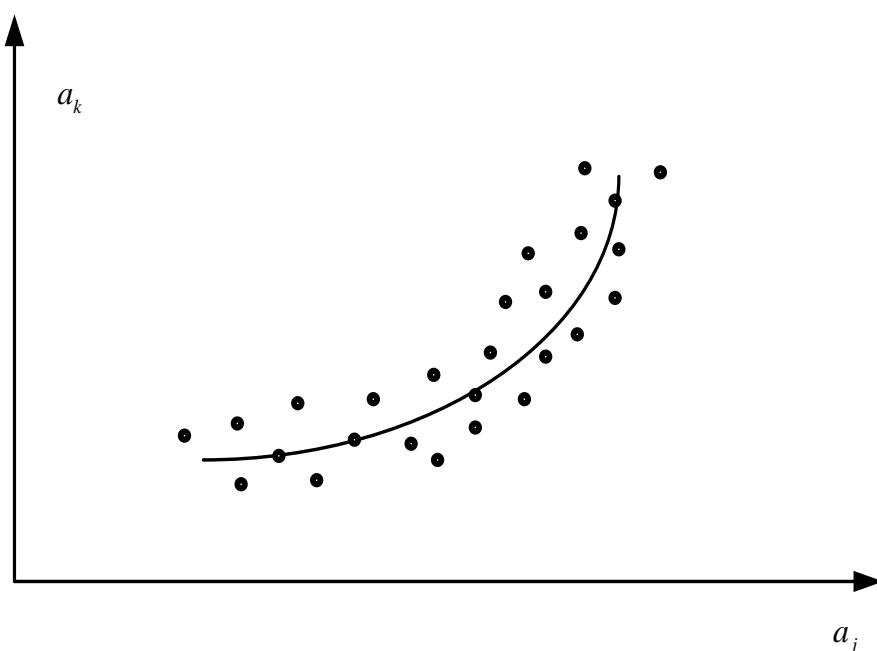
Acīmredzot i) iespēja jebkurā gadījumā nav pieņemama. Atribūtu vērtību statistiskās analīzes veikšana, pieņemot, ka tās ir normāli sadalītas, lai gan patiesībā tās nav, var novest pie kļūdainiem analīzes rezultātiem.

Variants (ii) ir diezgan pieņemama, ja attiecīgā atribūta vērtību sadalījums precīzi vai aptuveni atbilst kādam likumam, kas nav normāls sadalījums.

Šķiet, ka variants (iii) ir visdaudzsološākais. Ja ir iespējams pārveidot atribūtu a_k , a_l vērtības no attēliem 6.1.1.b,c tādā veidā, ka to transformētās vērtības tiks aptuveni sadalītas saskaņā ar normālo sadalījuma likumu, tad problēmu var uzskatīt par atrisinātu. Nepieciešamo analīzi var veikt transformētajām atribūtu vērtībām. Ja nepieciešams, analīzes rezultātus var apgriezt pārveidot sākotnējo atribūtu vērtību telpā. Šī vajadzība bieži ir saistīta ar analītisko rezultātu interpretāciju. Lietotājam var būt ļoti grūti interpretēt analīzes rezultātus transformēto atribūtu vērtību telpā. Apgrieztā transformācija ļauj mums atrisināt šo problēmu.

Tādējādi ir iezīmēts viens atribūtu vērtību transformācijas mērķis. Atribūtu vērtību sadalījumiem, kas ir vērsts pa kreisi vai pa labi, vajag pārveidot šīs vērtības tā, lai transformētajām vērtībām sadalījums būtu pēc iespējas tuvāks normālajam sadalījumam.

Apskatīsim citu virzienu, kurā plaši tiek izmantotas atribūtu vērtību transformācijas. Pievērsīsimies attēlam 6.1.2.



Attēls 6.1.2. Datu punktu grafiskais attēlojums atribūtu a_j un a_k vērtību telpā un regresijas sakarība starp šiem punktiem

Kā izriet no šī attēla, atribūtu a_j , a_k vērtības ir savstarpēji saistītas (starp to vērtību kopām pastāv korelācija). Lai analizētu šādus datus, ir iespējamas šādas iespējas:

- (i) Mēģiniet modelēt atkarības starp atribūtu vērtībām, izmantojot regresijas vienādojumus.
- (ii) Transformējiet sākotnējās atribūtu vērtības tā, lai atkarības starp to pārveidotajām vērtībām varētu modelēt, izmantojot lineāro regresiju.

Daudzās reālās datu apstrādes un analīzes problēmās priekšroka dodama variantam ii).

Piezīmēsim, ka gan sākotnējo sadalījumu normalizēšanas gadījumā, gan regresijas atkarības linearizācijas gadījumā sākotnējo atribūtu vērtību transformācija nav brīnumlīdzeklis, kas ļauj sasniegt vēlamos rezultātus visos gadījumos. Sākotnējo atribūtu vērtību transformācijas var un ir jāveic, ja iegūtie rezultāti ir piemēroti turpmākai datu analīzei. Ja plaši pazīstamu datu transformācijas metožu izmantošana neļauj iegūt vēlamos rezultātus, jāizmanto citas (t.sk. no parametriem neatkarīgas) datu apstrādes un analīzes metodes.

Ieviesīsim formālu definīciju [Osborne J.W, 2007]. “Datu transformācijas ir matemātisko modifikāciju pielietošana mainīgo vērtībām. Ir daudz dažādu datu transformāciju, sākot no konstantu pievienošanas līdz reizināšanai, celšana kvadrātā vai citā pakāpē, konvertēšana logaritmiskajās skalās, invertēšanai un attēlošanai, kā arī vērtību kvadrātsaknes vilkšanu”.

Pieminēšanas vērti ir divi citāti no darba [Lee D.K., 2020], kas akcentē atribūtu vērtību transformācijas vietu un lomu datu pirmapstrādes kontekstā.

“Dažādi pieņēmumi, tādi kā normalitāte, lineāras attiecības, bieži ir nepieciešami parametriskās analīzes metodēs. Dati, kas savākti no eksperimentiem, bieži pārkāpj šos pieņēmumus. Mainīgo lielumu transformācija sniedz iespēju bez sistemātiskām kļūdām padarīt datus pieņemamus parametriskai statistiskai analīzei. Mainīgo transformācijas mērķis ir ļaut veikt parametrisku statistisko analīzi un tās galvenais mērķis ir laba rezultātu interpretācija ar transformētajiem mainīgajiem. Mainīgo lielumu transformācija parasti maina to vērtību sākotnējos raksturlielumus. Apgrieztā transformācija ir būtiska rezultātu interpretācijai.

... Lielākajai daļai parametriskās statistiskās analīzes metožu ir nepieciešami pieņēmumi par normālo sadalījumu. Ja tie tiek pārkāpti, statistikas rezultāti no nenormāliem sadalījumiem var izraisīt nopietnas kļūdas... Lai gan centrālā robežas teorēma var aptvert nenormalitāti, ja izlases lielums ir pietiekams, daudzi eksperimentālie dati neatbilst pieņemumam par normālo sadalījumu, neskatoties uz to, ka ir salīdzinoši liels datu izlases lielums. Par laimi, vienkārša statistikas pieeja, mainīgo transformācija, nodrošina metodi datu sadalījuma pārveidošanai no nenormāla uz normālu. Turklāt mainīgo lielumu transformācija var veidot lineāras attiecības starp mainīgajiem no nelineārās attiecības un var stabilizēt aprēķinātās lineārās modelēšanas variācijas”.

“Atbilstoši datu sadalījuma raksturlielumiem var izmantot dažādas transformācijas metodes, lai panāktu normalitātes testa izpildi. Šāda veida transformācijas var padarīt datus simetriski sadalītus. Šīs transformācijas metodes var izmantot, lai nodrošinātu lineāras attiecības starp mainīgajiem. Ir labi zināms, ka daudzas statistiskās modelēšanas metodes balstās uz lineārām sakarībām un lineāras attiecības veidošana starp transformētajiem mainīgajiem, kas atvieglo modeļa statistisko novērtēšanu. Tipisks piemērs ir *logit* transformācija, ko izmanto binomiālajai loģistikas regresijai. Logit transformācija pārvērš notikumu varbūtības koeficientu logaritmos, ļaujot veikt regresijas analīzi starp rezultāta mainīgo un neatkarīgo mainīgo, kas darbojas kā lineārs prognozētājs”.

Darbā [McCune B., Grace J. 2002] tiek apgalvots, ka lielākajai daļai datu to transformācijām ir priekšrocības - tiek uzlaboti pieņēmumi par normālo sadalījumu, linearitāti, variācijām utt.; atribūtu vērtības tiek padarītas salīdzināmas, kad tās tiek mērītas dažādos mērogos.

6.2. Izplatītākie transformāciju veidi

1. *Logaritmisko transformāciju saime.*

Jebkuras logaritmiskās transformācijas būtība ir aizstāt sākotnējās atribūtu a_j vērtības ar šo vērtību logaritmiem:

$$a'_{ij} = \log(a_{ij}). \quad (6.2.1)$$

Šādas transformācijas ļoti bieži izmanto dažādos bioloģiskos un medicīnas pētījumos. Kādi ir tam iemesli? Pirmkārt, komponentu koncentrācijas dažādos bioloģiskajos paraugos var atšķirties daudzkārtīgi. Darbība ar šādiem datiem ir ļoti sarežģīta un neērta. Šādu datu pārveidošana logaritmiskā formā ļauj ievērojami samazināt transformēto vērtību izmaiņu diapazonu, kas ievērojami vienkāršo datu apstrādi un analīzi. Otrkārt, daudzu bioloģisko rādītāju sākotnējo vērtību variācijas ir proporcionālas to paredzamo vērtību pakāpēm. Šādu rādītāju logaritmiskās transformācijas ir optimālas tādā ziņā, ka ļauj iegūt transformēto vērtību normālus sadalījumus vai tuvu normālam sadalījumam.

Mēs runājam par logaritmisko pārveidojumu saimi, jo šādas transformācijas var veikt, izmantojot logaritmus ar dažādām bāzēm: logaritmus ar bāzi 10 (decimāllogaritmi), logaritmus ar bāzi $e \approx 2.718$ (naturālie logaritmi) un logaritmus ar bāzi 2 (binārie logaritmi).

Lai uzskatāmi nodemonstrētu pārveidoto vērtību atkarību no izmantotajām logaritmu bāzēm, tabulā 6.2.1 ir parādītas sākotnējās atribūtu a_j vērtības diapazonā $[1,10]$ un to transformētās vērtības, izmantojot logaritmus ar dažādām bāzēm. Attēlā 6.2.1 transformāciju rezultāti tiek attēloti grafiskā veidā.

Kādus secinājumus var izdarīt, analizējot attēlu 6.2.1? Acīmredzot visas transformācijas saglabā tādu pašu atribūta a_j transformēto vērtību secību, kāda pastāv tā sākotnējo vērtību kopā. Ja $a_{ij} > a_{kj}$, tad $\lg(a_{ij}) > \lg(a_{kj})$, $\ln(a_{ij}) > \ln(a_{kj})$ un $\log_2(a_{ij}) > \log_2(a_{kj})$. Tomēr sakarības starp sākotnējo un transformēto atribūtu vērtībām nav lineāras.

Tabula 6.2.1

Sākotnējās atribūtu a_j vērtības un to transformētās vērtības ar dažādām logaritmu bāzēm

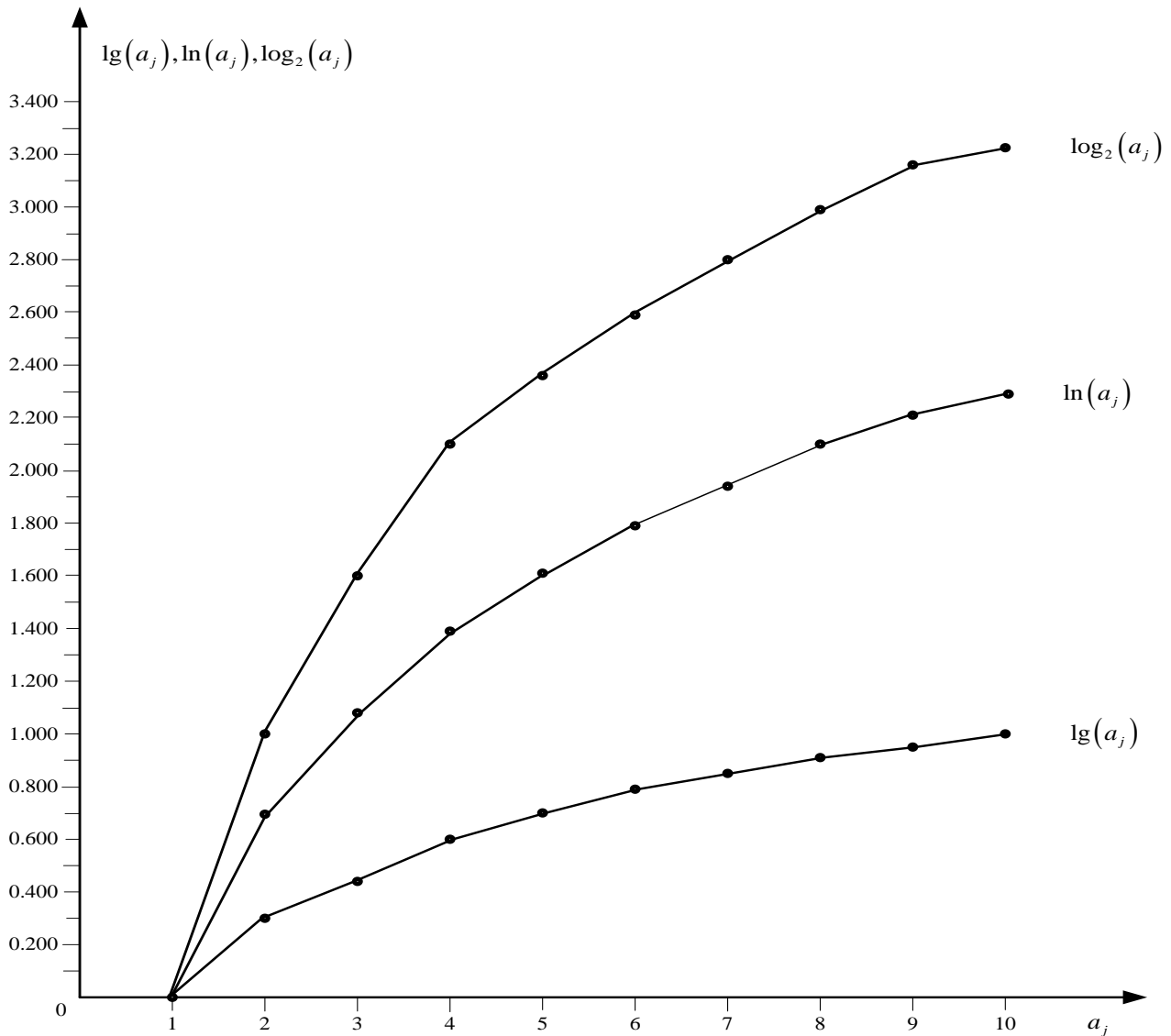
a_j	1	2	3	4	5	6	7	8	9	10
$\lg(a_{ij})$	0	0.301	0.477	0.602	0.699	0.778	0.845	0.903	0.954	1.000
$\ln(a_{ij})$	0	0.693	1.097	1.386	1.609	1.792	1.946	2.079	2.197	2.302
$\log_2(a_{ij})$	0	1.000	1.585	2.000	2.322	2.585	2.807	3.000	3.169	3.222

Atšķirības starp atribūta transformētajām vērtībām nav tieši proporcionālas atšķirībām starp šī atribūta atbilstošajām sākotnējām vērtībām. Palielinoties atribūta sākotnējām vērtībām a_j , atšķirības starp tā transformētajām vērtībām mēdz samazināties. Tieši šī logaritmisko transformāciju īpašība tiek izmantota, lai pārveidotu sākotnējos atribūtu vērtību sadalījumu, kas ir asimetrisks, uz labo pusi vērsts sadalījums, formā, kas ir tuvu normālajam sadalījumam.

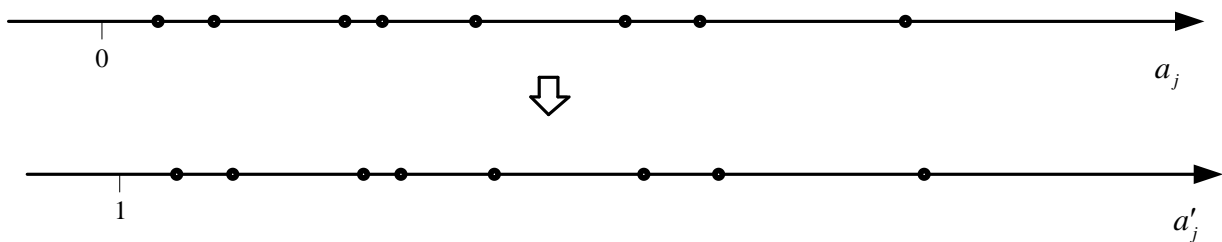
Ja visas atribūtu a_j vērtības ir lielākas par 1, tad šo vērtību transformācija logaritmiskā formā nav problēma. Ja visas atribūtu a_j vērtības ir stingri lielākas par 0, tad logaritmiski transformētās vērtības sākotnējām atribūtu vērtībām, kas ir mazākas par 1, būs negatīvas. Lai izvairītos no negatīvām transformētajām vērtībām, var piemērot šādu sākotnējo atribūtu vērtību mērogošanu:

$$a'_{ij} = a_{ij} + 1. \quad (6.2.2)$$

Šī mērogošanas maiņa shematiski parādīta attēlā 6.2.2.



Attēls 6.2.1. Sākotnējo un logaritmiski transformēto atribūtu a_j vērtību datiem no tabulas 6.2.1 grafisks attēlojums



Attēls 6.2.2. Sākotnējās atribūtu a_j vērtības mērogošanas procesa shematiskais attēlojums pēc vienādojuma (6.2.2)

Ja nav vēlamas ne tikai negatīvas, bet arī nulles transformētās atribūtu a_j vērtības, tad var izmantot šādu sākotnējo atribūtu vērtību mērogošanu:

$$a''_{ij} = a_{ij} + 1 + \varepsilon, \quad (6.2.3)$$

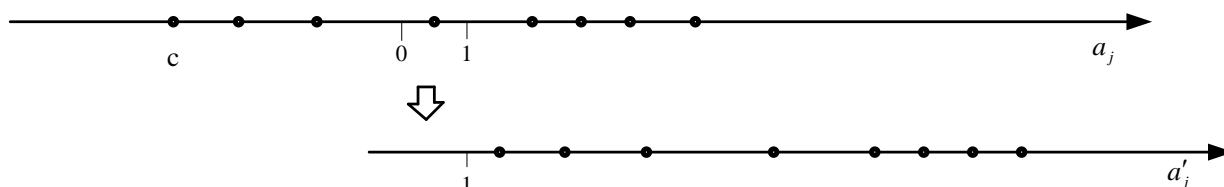
kur ε - kļūdas lielums (mazs pozitīvs skaitlis, piemēram, $\varepsilon = 0.01$ vai $\varepsilon = 0.001$).

Problēmas rodas gadījumos, ja dažas no atribūta a_j sākotnējām vērtībām ir negatīvas. Lai nodrošinātu to, ka visas transformētās atribūtu vērtības nav negatīvas, var izmantot šādu sākotnējo atribūtu vērtību mērogošanu:

$$a'_{ij} = a_{ij} + c + 1, \quad (6.2.4)$$

kur c - atribūta a_j lielākā negatīvā vērtība.

Šādas mērogošanas maiņas shēma ir parādīta attēlā 6.2.3.



Attēls 6.2.3. Sākotnējās atribūtu a_j vērtības mērogošanas procesa shematiskais attēlojums pēc vienādojuma (6.2.4)

Šajā attēlā simbols c apzīmē atribūta a_j lielāko negatīvo vērtību. Pievienojot summu $c + 1$ šai vērtībai, mēs panākam to, ka jaunajā skalā vērtība c atbildīs 1 un visas atribūta vērtības būs pozitīvas.

Ja transformētās atribūtu a_j nulles vērtības nav vēlamas, var izmantot šādu mērogošanu:

$$a''_{ij} = a_{ij} + c + 1 + \varepsilon, \quad (6.2.5)$$

kur parametra ε nozīme ir izskaidrota vienādojumā (6.2.3).

Jebkuras logaritmiskās transformācijas svarīga iezīme ir ievērojams transformēto atribūtu vērtību diapazona samazinājums, salīdzinot ar tā sākotnējo vērtību diapazonu.

Transformāciju, izmantojot decimāllogaritmus, var pielietot atribūtiem ar ekstremāliem vērtību diapazoniem. Tad transformēto vērtību diapazons ļoti būtiski mainās. Tomēr jāņem vērā, ka arī transformēto vērtību kopas mērogs ievērojami samazinās. Piemēram, mums ir divas sākotnējo atribūtu a_j vērtības: 1500 un 1600. Transformējot šīs vērtības ar decimāllogaritmiem, mēs iegūstam vērtības: 3.176 un 3.204. Atšķirība starp šīm vērtībām ir ļoti maza, kas var radīt problēmas datu analīzes laikā.

Nelieliem sākotnējo atribūtu vērtību diapazoniem ir vēlama transformācija, izmantojot naturālos logaritmus vai logaritmus ar bāzi 2. Lietojot šādus logaritmus, mēs iegūstam arī ievērojamu samazinājumu transformēto atribūtu vērtību diapazonos, savukārt relatīvi tuvu transformēto vērtību mērogs nesamazinās tik daudz kā decimāllogaritmu lietošanas gadījumā.

2. Pakāpes transformāciju saime.

Vispārīgā gadījumā transformētās atribūtu a_j vērtības var iegūt formā

$$a^t_{ij} = (a_{ij})^\alpha, \quad \alpha \neq 0. \quad (6.2.6)$$

Pie dažādām α vērtībām mēs iegūsim dažādus transformācijas rezultātus:

$\alpha = 3$ - kubiskā transformācija;

$\alpha = 2$ - kvadrātiskā transformācija;

$\alpha = 1$ - transformētās atribūtu vērtības sakrīt ar sākotnējām vērtībām (acīmredzot šādi transformācijai nav jēgas.)

$\alpha = 1/2$ - kvadrātsaknes transformācija.

Pie vērtībām $\alpha < 1$, notiek lielo atribūtu vērtību “saspiešana”. Šo pakāpes transformācijas īpašību var izmantot, lai normalizētu asimetriskus, uz labo pusi vērstus sadalījumus.

Tabulā 6.2.2 parādīti sākotnējo atribūtu a_j vērtību pakāpes transformācijas pie dažādām α vērtībām.

Tabula 6.2.2

Transformētās atribūtu a_j vērtības pie dažādiem α

a_j	1	2	3	4	5	6	7	8	9	10
$\alpha = 2$	1	4000	9000	16 000	25 000	36 000	49 000	64 000	81 000	100 000
$\alpha = 1.5$	1	2.828	5.196	8000	11.180	14.697	18.520	22.627	27 000	31.623
$\alpha = 0.5$	1	1.414	1.732	2000	2.236	2.449	2.645	2.829	3000	3.162
$\alpha = 0.25$	1	1.189	1.136	1.414	1.495	1.565	1.625	1.682	1.732	1.778

Tabulā var redzēt ka, pie $\alpha > 1$ palielinoties atribūta a_j sākotnējām vērtībām, to transformētās vērtības strauji palielinās. Šajā gadījumā atribūtu vērtībām $a_{ij} > 1$ transformētās atribūtu vērtības pārsniedz sākotnējās vērtības.

Pie $\alpha < 1$ palielinoties atribūta a_j sākotnējām vērtībām, palielinās arī to transformētās vērtības. Bet šis pieaugums notiek lēni un visas transformētās vērtības (izņemot $a_{ij} = 1$) ir mazākas par atbilstošām sākotnējām atribūtu vērtībām.

3. Box-Cox transformāciju saime.

Šis pakāpes transformācijas variants pirmo reizi tika piedāvāts darbā [Box G.E.P., Cox D.R., 1964] tikai pozitīvām sākotnējo atribūtu vērtībām. Vēlāk to attiecināja arī uz negatīvām atribūtu vērtībām. Sākotnējo atribūtu vērtību transformācijas var veikt ar sekojošu vienādojumu palīdzību.

- visas atribūtu a_j vērtības ir pozitīvas:

$$a_{ij}^t(\lambda) = \frac{(a_{ij})^\lambda - 1}{\lambda}, \lambda \neq 0;$$

$$= \log(a_{ij}), \lambda = 0,$$
(6.2.7)

kur λ - transformācijas parametrs. Vērtības λ var variēt diapazonā $[-5, 5]$.

- atribūtu a_j vērtības ir negatīvas:

$$a_{ij}^t(\lambda) = \frac{(a_{ij} + \lambda_2)^{\lambda_1} - 1}{\lambda_1}, \lambda_1 \neq 0;$$

$$\log(a_{ij} + \lambda_2),$$
(6.2.8)

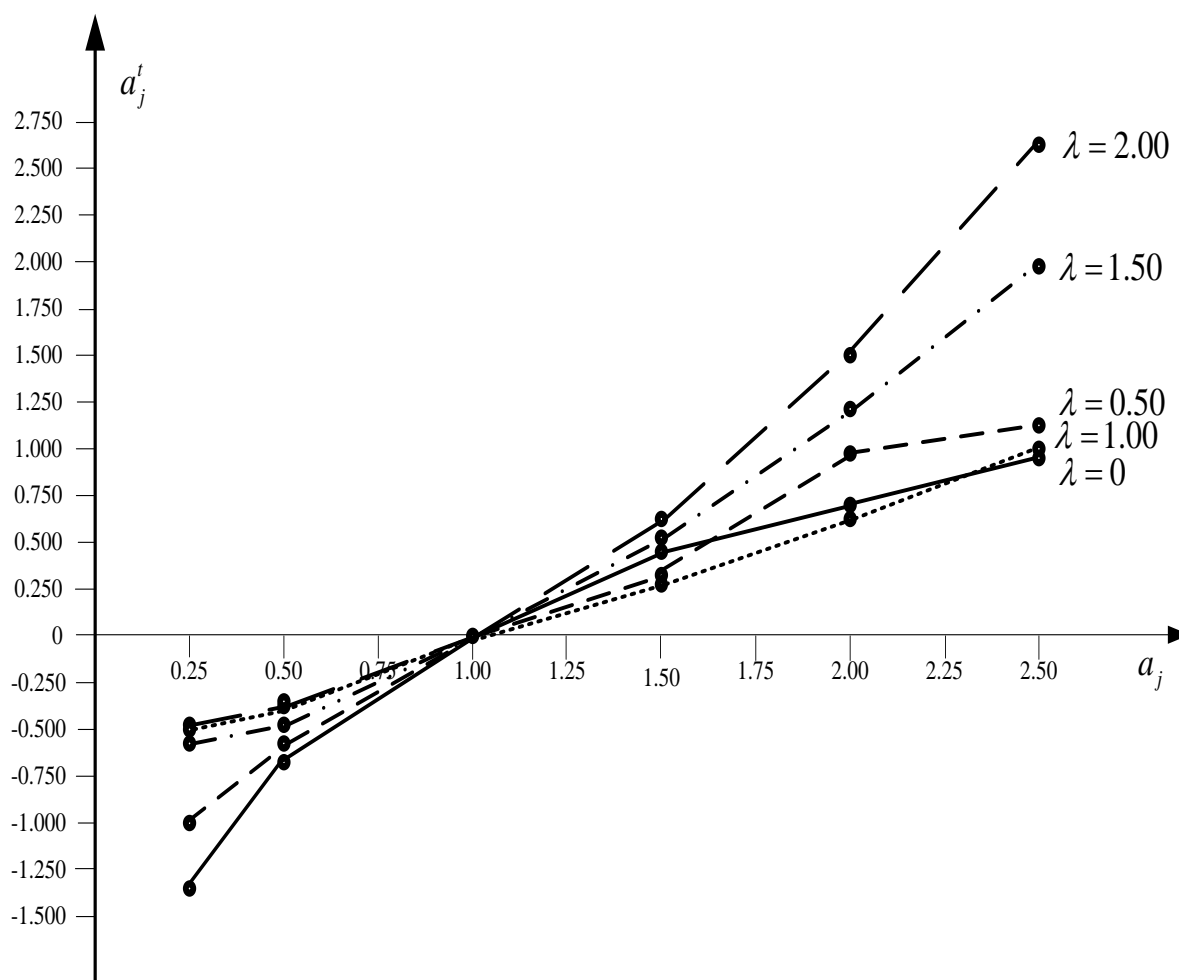
kur λ_2 -parametrs, ko izmanto sākotnējo atribūtu vērtību skalas izmaiņai, lai visas vērtības būtu lielākas par 1.

Tabulā 6.2.3 parādītas sākotnējās atribūtu a_j vērtības un to transformētās vērtības pie dažādām λ vērtībām. Transformētās atribūtu a_j vērtības ir grafiski attēlotas attēlā 6.2.4. No tabulas 6.2.3 un attēla 6.2.4 varam izdarīt pamatotu secinājumu, ka norādot dažādas λ vērtības, ir iespējams nodrošināt ļoti plašas sākotnējo atribūtu a_j vērtību transformēšanas iespējas, izmantojot Box-Cox transformācijas.

Tabula 6.2.3

Sākotnējo atribūtu a_j vērtību transformācijas rezultāti saskaņā ar vienādojumu (6.2.7) dažādām λ vērtībām

a_j	0.25	0500	1.00	1.50	2.00	2.50
$\lambda = 2.00$	-0.469	-0.375	0	0.625	1500	2.625
$\lambda = 1.50$	-0.583	-0.431	0	0.558	1.219	1.969
$\lambda = 1.00$	-0.500	-0.333	0	0.333	0.666	1000
$\lambda = 0.50$	-1000	-0.586	0	0.449	0.828	1.162
$\lambda = 0$	-1.386	-0.693	0	0.405	0.693	0.961



Attēls 6.2.4. Sākotnējo un transformēto atribūtu a_j vērtību grafisks attēlojums datiem no tabulas 6.2.3

6.3. Atribūtu vērtību transformācija normālā sadalījuma sasniegšanai

Ir piedāvātas dažādas pieejas šādu transformāciju veikšanai. Uz labo pusi vēršiem sadalījumiem bieži izmanto logaritmiskās transformācijas. Piemēra pēc mēs izmantosim logaritmiskās transformācijas atribūtu a_k vērtībām no tabulas P3.2.2 (pielikums P3.2). Šis sadalījums ir tipisks uz labo pusi vērsta sadalījuma pārstāvis. Tabulās P3.2.2 un 6.2.1 sniegto datu kopsavilkums ir parādīts tabulā 6.3.1.

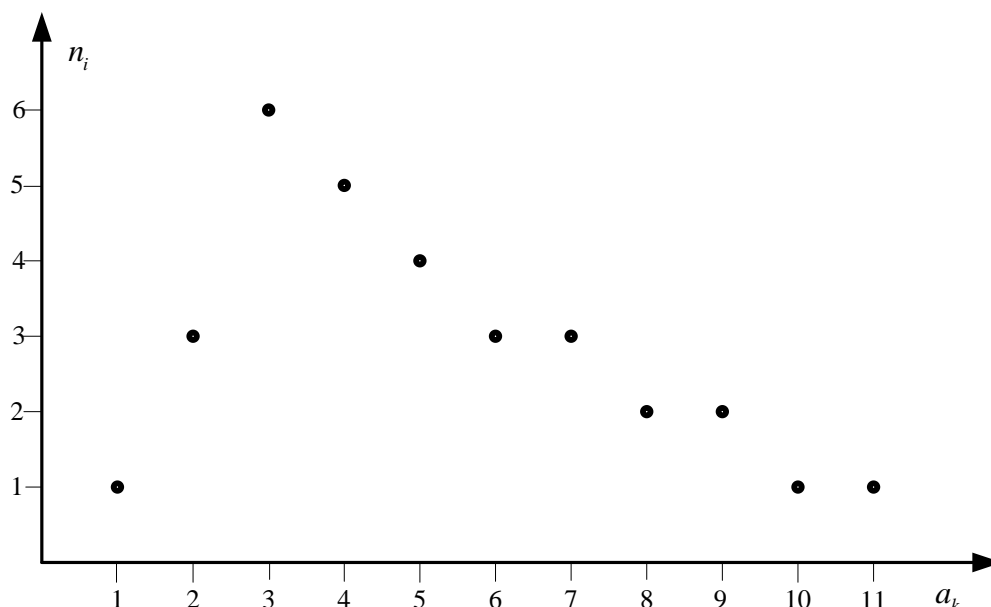
Tabula 6.3.1

Vispārināts datu kopsavilkums no tabulām P3.2.2 un 6.2.1

a_k	1	2	3	4	5	6	7	8	9	10	11
$\lg(a_k)$	0	0.301	0.477	0.602	0.699	0.778	0.845	0.903	0.954	1000	1.041
$\ln(a_k)$	0	0.693	1.097	1.386	1.609	1.792	1.946	2.079	2.197	2.302	2.398
$\log_2(a_k)$	0	1000	1.585	2000	2.322	2.585	2.807	3000	3.169	3.222	3.458
n_i	1	3	6	5	4	3	3	2	2	1	1

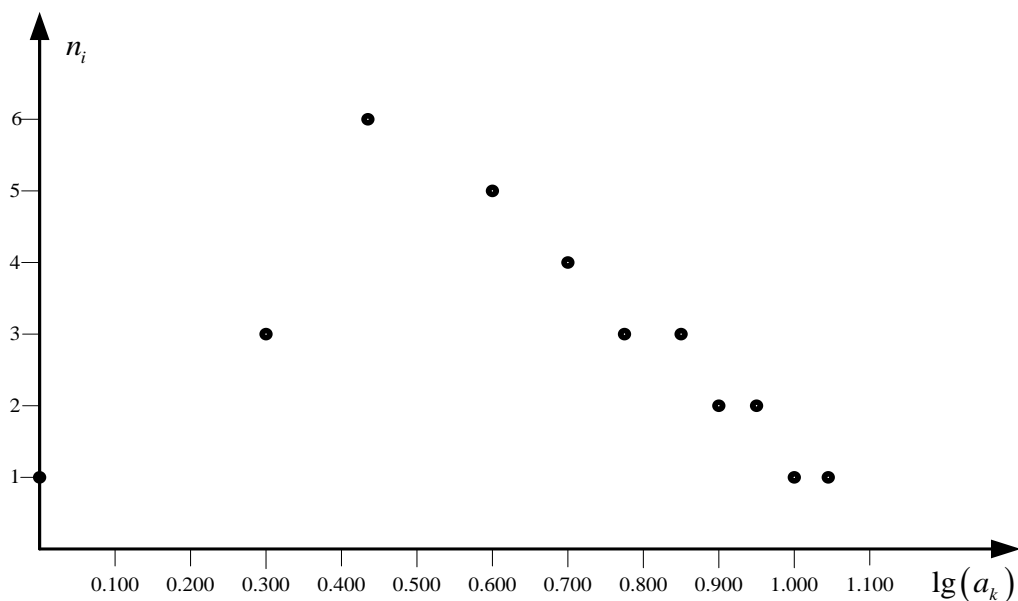
No tabulas P3.2.2 ir ņemtas atribūta sākotnējās vērtības a_k un tā vērtību atkārtotās gadījumu skaits n_i noteiktos intervālos.

Vērtības $\lg(a_k)$, $\ln(a_k)$, $\log_2(a_k)$ ir ņemtas no tabulas 6.2.1 un pievienotas vērtības priekš $a_k = 11$. Attēls P3.2.2 parāda atribūtu vērtību a_k histogrammu datiem no tabulas P3.2.2. Iegūto rezultātu tālākas salīdzināšanas nolūkos attēlā 6.2.1 parādīta atribūtu a_k vērtību sadalījuma histogramma atbilstoši tabulas 6.3.1 datiem.



Attēls 6.3.1. Atribūtu a_k vērtību sadalījuma histogramma datiem no tabulas 6.3.1

Attēlā 6.3.2 parādīta vērtību $\lg(a_k)$ sadalījuma histogramma.



Attēls 6.3.2. Vērtību $\lg(a_k)$ sadalījuma histogramma datiem no tabulas 6.3.1

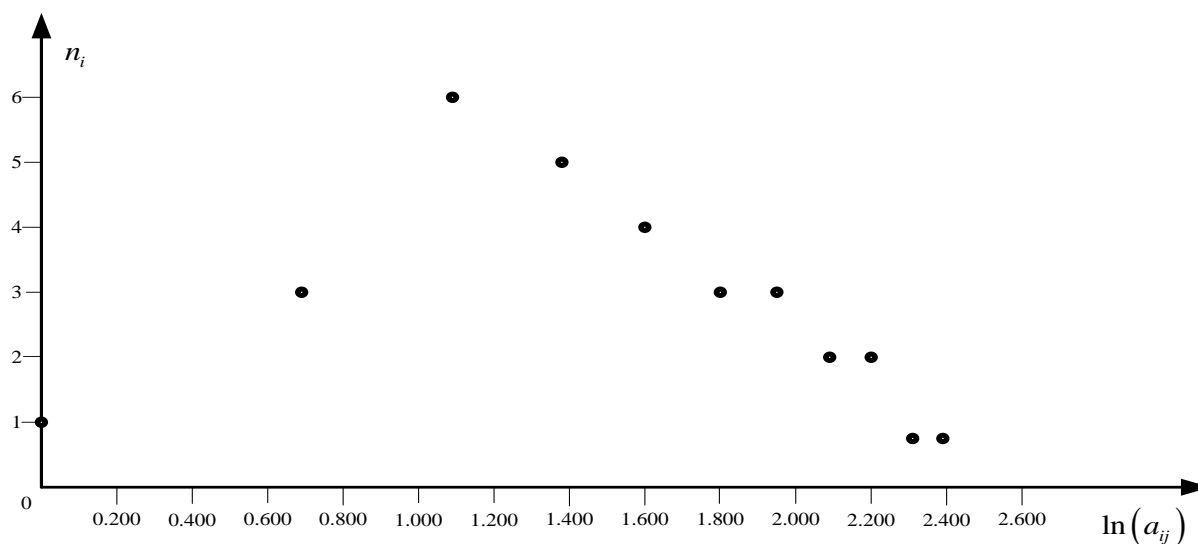
Tabulā 6.3.2 sniegti statistikas dati par vērtību $\lg(a_k)$ sadalījumiem datiem no tabulas 6.3.1 Turpmākajos aprēķinos tiks izmantoti pakotnes IBM SPSS Statistica aprēķini (detalizētāk skatīt pielikumā P3).

Tabula 6.3.2

Vērtību $\lg(a_k)$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.1

Descriptives		Statistic	Std. Error
Mean		,64887	,042686
95% Confidence Interval for	Lower Bound	,56169	
Mean	Upper Bound	,73605	
5% Trimmed Mean		,65801	
Median		,69900	
Variance		,056	
Std. Deviation		,237668	
Minimum		,000	
Maximum		1,041	
Range		1,041	
Interquartile Range		,368	
Skewness		-,603	,421
Kurtosis		,361	,821

Attēlā 6.3.3 parādīta histogramma vērtību $\ln(a_k)$ sadalījumam datiem no tabulas 6.3.1.



Attēls 6.3.3. Vērtību $\ln(a_k)$ sadalījuma histogramma datiem no tabulas 6.3.1

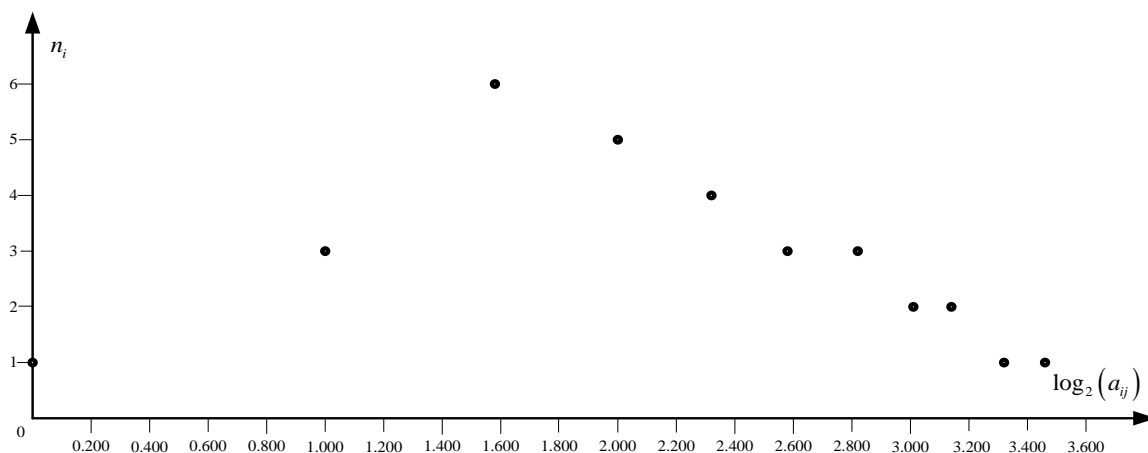
Tabulā 6.3.3 sniegti statistikas dati par vērtību $\ln(a_k)$ sadalījumu tabulas 6.3.1 datiem.

Tabula 6.3.3

Vērtību $\ln(a_k)$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.1

Descriptives		
	Statistic	Std. Error
Mean	1,49977	,099650
95% Confidence Interval for Mean	Lower Bound Upper Bound	1,29626 1,70329
5% Trimmed Mean	1,52142	
Median	1,60900	
Variance	,308	
Std. Deviation	,554826	
Minimum	,000	
Maximum	2,398	
Range	2,398	
Interquartile Range	,849	
Skewness	-,562	,421
Kurtosis	,282	,821

Attēlā 6.3.4 parādīta histogramma vērtību $\log_2(a_k)$ sadalījumam datiem no tabulas 6.3.1.



Attēls 6.3.4. Vērtību $\log_2(a_k)$ sadalījuma histogramma datiem no tabulas 6.3.1

Tabulā 6.3.4 sniegti statistikas dati par vērtību $\log_2(a_k)$ sadalījumu datiem no tabulas 6.3.1.

Tabula 6.3.4

Vērtību $\log_2(a_k)$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.1

Descriptives		Statistic	Std. Error
Mean		2,16103	,142847
95% Confidence Interval for Mean	Lower Bound	1,86930	
	Upper Bound	2,45276	
5% Trimmed Mean		2,19397	
Median		2,32200	
Variance		,633	
Std. Deviation		,795337	
Minimum		,000	
Maximum		3,458	
Range		3,458	
Interquartile Range		1,222	
Skewness		-,588	,421
Kurtosis		,308	,821

Analizēsim iegūtos rezultātus. Tabulā 6.3.5 ir parādītas sākotnējās un transformētās atribūtu a_k vērtību sadalījuma asimetrijas rādītāja (*skewness*) un ekscesa rādītāja (*kurtosis*) vērtības. Šo sadalījuma parametru definīcijas ir sniegtas pielikumā P3.3. Sākotnējo atribūtu vērtību asimetrijas un ekscesa vērtības ir ņemtas no tabulas pielikumā P3.4.

Tabula 6.3.5

Asimetrijas un ekscesa vērtības oriģinālajiem un transformētajiem a_k atribūtiem

	a_k	$\lg(a_k)$	$\ln(a_k)$	$\log_2(a_k)$
Asimetrija	0.581	0.361	0.282	0.308
Ekscess	-0.427	-0.603	-0.562	-0.588

Visiem transformēto atribūtu a_k vērtību sadalījumiem asimetrijas vērtības ir mazākas nekā sākotnējo atribūtu vērtību sadalījuma vērtības. Transformējot ar binārajiem logaritmiem, tika iegūta vismazākā asimetrijas vērtība. Tas nozīmē, ka asimetrijas korekcijas ziņā šī transformācija ir vislabākā.

Kas attiecas uz ekscesu, visi transformēto atribūtu a_k vērtību sadalījumi ir tālu no normālā sadalījuma nekā tā sākotnējo vērtību sadalījums.

Jāpiebilst, ka izmantojot logaritmiskās transformācijas, mēs sasniedzām savu galveno mērķi – samazinājām sākotnējā sadalījuma asimetriju.

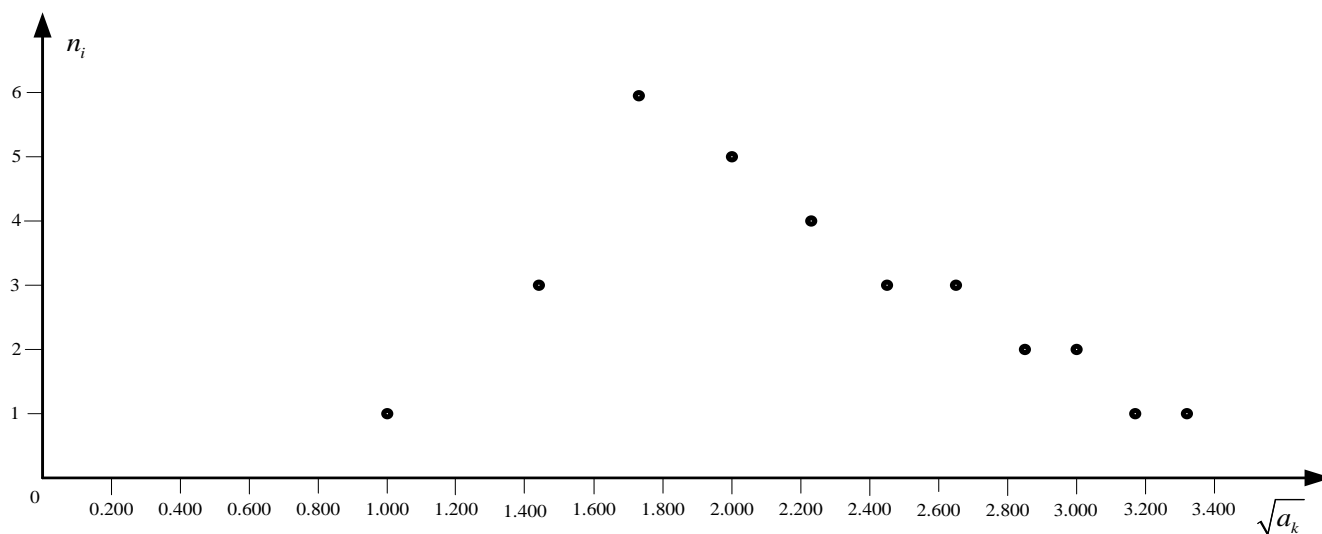
Uz labo pusi vērstajiem sadalījumiem ieteicams izmantot arī transformācijas ar pakāpes funkcijām ar pakāpēm, kas <1 . Piemēram, tabulā 6.3.6 ir parādīti sākotnējo atribūtu a_k vērtību pakāpju transformāciju rezultāti ar pakāpēm $\frac{1}{2}$ un $\frac{1}{3}$. Pretējā gadījumā transformētās atribūtu a_k vērtības tiek aprēķinātas, ņemot kvadrātsakni vai kubsakni no sākotnējām vērtībām.

Tabula 6.3.6

Sākotnējo atribūtu a_k vērtību transformācijas rezultāti, izmantojot kvadrātsakni un kubsakni.

a_k	1	2	3	4	5	6	7	8	9	10	11
$\sqrt{a_k}$	1	1.414	1.732	2.000	2.236	2.449	2.646	2.829	3.000	3.162	3.317
$\sqrt[3]{a_k}$	1	1.260	1.442	1.589	1.709	1.816	1.912	2.000	2.078	2.153	2.222
n_i	1	3	6	5	4	3	3	2	2	1	1

Transformēto atribūtu $\sqrt{a_k}$ vērtību histogramma parādīta attēlā 6.3.5.



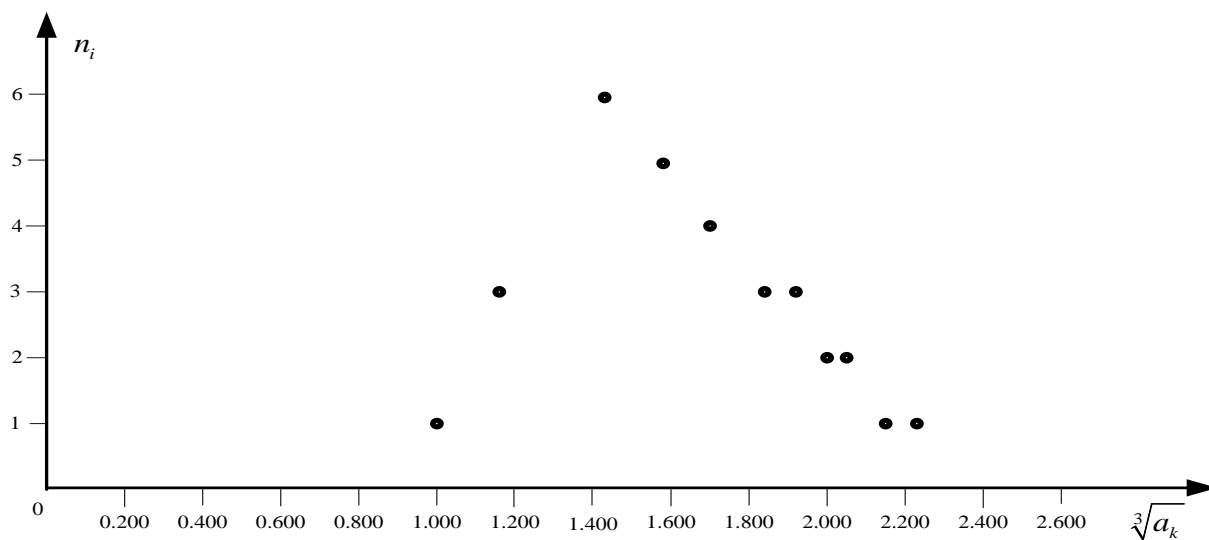
Attēls 6.3.5. Vērtību $\sqrt{a_k}$ sadalījuma histogramma datiem no tabulas 6.3.6

Tabulā 6.3.7 sniegti statistikas dati par vērtību $\sqrt{a_k}$ sadalījumu.

Vērtību $\sqrt{a_k}$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.6

Descriptives			Statistic	Std. Error
Mean			2,19355	,102855
95% Confidence Interval for Mean	Lower Bound		1,98349	
	Upper Bound		2,40361	
5% Trimmed Mean			2,19234	
Median			2,23600	
Variance			,328	
Std. Deviation			,572672	
Minimum			1,000	
Maximum			3,317	
Range			2,317	
Interquartile Range			,914	
Skewness			,089	,421
Kurtosis			-,599	,821

Transformēto atribūtu $\sqrt[3]{a_k}$ vērtību histogramma parādīta attēlā 6.3.6.



Attēls 6.3.6. Vērtību $\sqrt[3]{a_k}$ sadalījuma histogramma datiem no tabulas 6.3.6

Tabulā 6.3.8 sniegti statistikas dati par vērtību $\sqrt[3]{a_k}$ sadalījumu.

Vērtību $\sqrt[3]{a_k}$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.6

Descriptives			Statistic	Std. Error
Mean			1,67510	,053194
95% Confidence Interval for	Lower Bound		1,56646	
Mean	Upper Bound		1,78373	
5% Trimmed Mean			1,67845	
Median			1,70900	
Variance			,088	
Std. Deviation			,296173	
Minimum			1,000	
Maximum			2,222	
Range			1,222	
Interquartile Range			,470	
Skewness			-,111	,421
Kurtosis			-,457	,821

Analizēsim iegūtos rezultātus. Tabulā 6.3.9 ir parādītas asimetrijas un ekscesa vērtības atribūta sākotnējo vērtību a_k sadalījumam un transformētajām vērtībām $\sqrt{a_k}$, $\sqrt[3]{a_k}$.

Tabula 6.3.9.

Asimetrijas un ekscesa vērtības sākotnējam un transformētajam atribūtu a_k sadalījumiem

	a_k	$\sqrt{a_k}$	$\sqrt[3]{a_k}$
Ekscess	-0.427	-0.599	-0.457
Asimetrija	-0.581	0.086	-0.11

Transformēto atribūtu sadalījuma asimetrijas vērtība ievērojami samazinājās, salīdzinot ar sākotnējo sadalījuma asimetrijas vērtību. Transformācijai $\sqrt[3]{a_k}$ nozīmē to, ka asimetrijas vērtība pamainīja savu zīmi un kļuva pozitīva, t.i., vērtību sadalījums kļuva vērsts uz kreiso pusi.

Transformēto atribūtu sadalījumiem ir lielas ekscesa negatīvas vērtības.

Pamatojoties uz šiem un iepriekš iegūtajiem rezultātiem, mēs varam droši teikt, ka gan logaritmiskās, gan pakāpes transformācijas ievērojami samazina sadalījumu asimetriju, bet mazākā mērā ietekmē šo sadalījumu ekscesu. Tā kā šādu transformāciju primārais mērķis ir samazināt sākotnējo sadalījumu asimetriju, tad abu veidu transformācijas var sasniegt šo mērķi.

Uz kreiso pusi vērsto sadalījumu transformācija ir sarežģītāka. Šādos gadījumos ieteicams izmantot pakāpes transformācijas ar pakāpju rādītājiem 2 un 3.

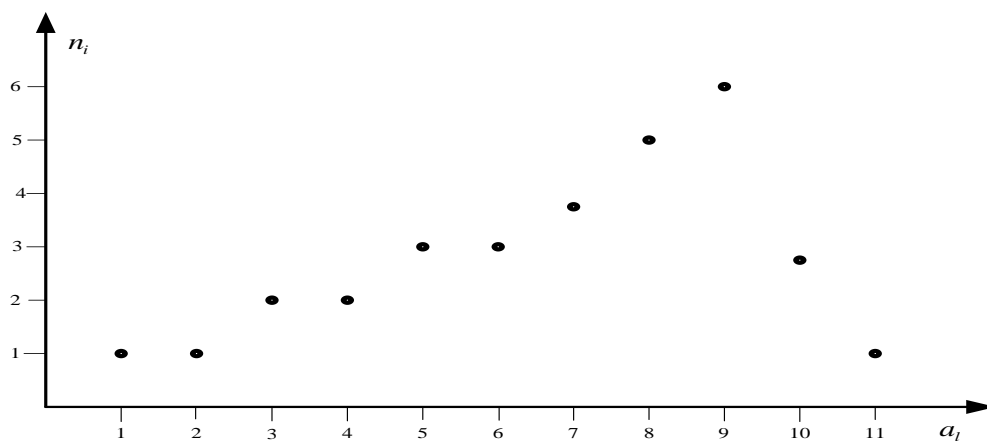
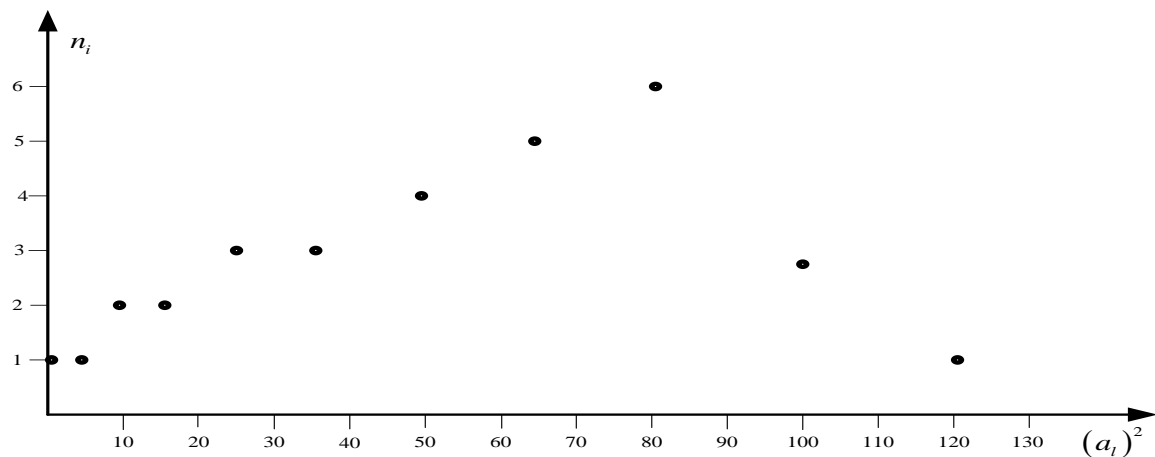
Tabulā 6.3.10 ir parādītas sākotnējās atribūtu a_i vērtības no tabulas P3.2.3 (pielikums P3.2), transformētās vērtības $(a_i)^2$ un $(a_i)^3$, un atbilstošo atribūtu vērtību biežums.

Attēlā 6.3.7 iegūto rezultātu salīdzināšanas un analīzes nolūkos parādīta sākotnējo atribūtu a_i vērtību sadalījuma histogramma.

Tabula 6.3.10

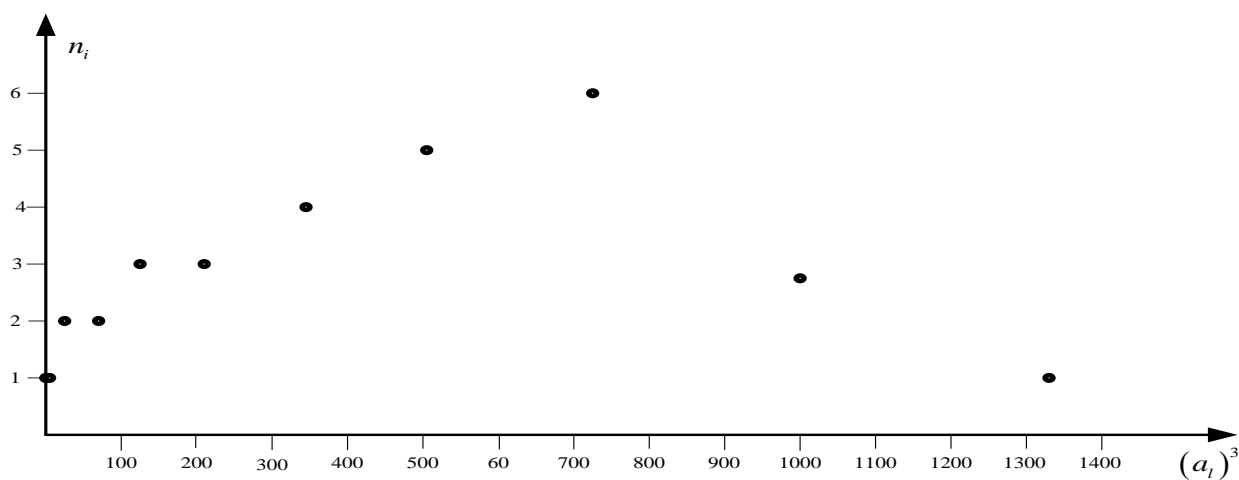
Sākotnējās atribūtu a_i vērtības, transformētās vērtības $(a_i)^2$, $(a_i)^3$ un biežums n_i

a_i	1	2	3	4	5	6	7	8	9	10	11
$(a_i)^2$	1	4	9	16	25	36	49	64	81	100	121
$(a_i)^3$	1	8	27	64	125	216	343	512	729	1000	1331
n_i	1	1	2	2	3	3	4	5	6	3	1

Attēls 6.3.7. Sākotnējā atribūtu a_i vērtību sadalījuma histogramma datiem no tabulas 6.3.10Attēlā 6.3.8 parādīta transformēto vērtību $(a_i)^2$ sadalījuma histogramma.Attēls 6.3.8. Transformēto vērtību $(a_i)^2$ sadalījuma histogramma datiem no tabulas 6.3.10Tabulā 6.3.11 sniegti statistikas dati par vērtību $(a_i)^2$ sadalījumu.Attēlā 6.3.9 parādīta transformēto vērtību $(a_i)^3$ sadalījuma histogramma.

Vērtību $(a_i)^2$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.10

Descriptives		Statistic	Std. Error
Mean		53,58	5,824
95% Confidence Interval for Mean	Lower Bound	41,69	
	Upper Bound	65,48	
5% Trimmed Mean		53,11	
Median		49,00	
Variance		1051,585	
Std. Deviation		32,428	
Minimum		1	
Maximum		121	
Range		120	
Interquartile Range		56	
Skewness		,103	,421
Kurtosis		-,898	,821



Attēls 6.3.9. Transformēto vērtību $(a_i)^3$ sadalījuma histogramma datiem no tabulas 6.3.10

Tabulā 6.3.12 sniegti statistikas dati par vērtību $(a_i)^3$ sadalījumu.

Analizēsim iegūtos rezultātus. Tabulā 6.3.13 ir parādītas asimetrijas un ekscesa vērtības sākotnējo atribūtu a_i vērtību sadalījumam un transformēto vērtību $(a_i)^2$, $(a_i)^3$ sadalījumiem. Parametru vērtības sākotnējo atribūtu vērtību sadalījumam ir ņemtas no tabulas pielikumā P3.4.

Vērtību $(a_i)^3$ sadalījuma statistikas kopsavilkums datiem no tabulas 6.3.10

Descriptives		Statistic	Std. Error
Mean		446,81	63,495
95% Confidence Interval for	Lower Bound	317,13	
Mean	Upper Bound	576,48	
5% Trimmed Mean		428,84	
Median		343,00	
Variance		124978,828	
Std. Deviation		353,523	
Minimum		1	
Maximum		1331	
Range		1330	
Interquartile Range		604	
Skewness		,614	,421
Kurtosis		-,319	,821

Tabula 6.3.13

Asimetrijas un ekscesa vērtības sākotnējo un transformēto vērtību a_i sadalījumiem

	a_i	$(a_i)^2$	$(a_i)^3$
Ekscess	-0.427	-0.898	-0.319
Asimetrija	-0.581	0.103	-0.111

Transformēto atribūtu sadalījuma asimetrijas vērtības a_i ir mazākas par asimetrijas vērtībām priekš sākotnējo vērtību sadalījuma. Transformēto vērtību $(a_i)^2$ sadalījums tagad ir ar nelielu asimetriju uz labo pusi. Tā kā praktiskās iespējas labot uz kreiso pusi vērsto sadalījumu ir ierobežotas, literatūrā tiek ieteikta šāda pieeja. Sākotnējo atribūtu a_{ii} vērtību vietā iesaka izmantot vērtības $c - a_{ii}$, kur c ir konstante. Šīs konstantes vērtība ir izvēlēta tā, lai varētu korekti aptvert visas sākotnējās atribūta a_i vērtības. Tādu transformāciju sauc par sākotnējā sadalījuma *attēlojumu*.

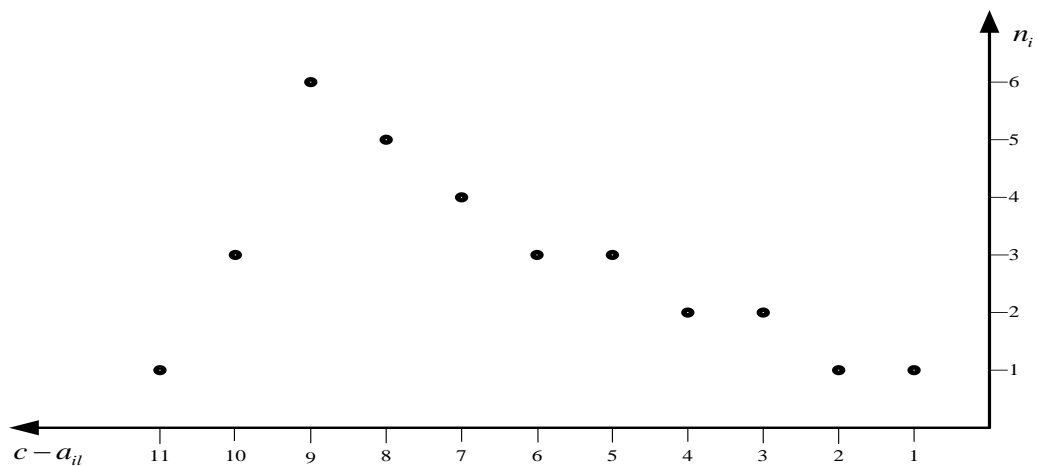
Demonstrēsim sadalījuma attēlošanas procesu ar piemēru. Tabulā 6.3.10 ir parādītas sākotnējās atribūtu a_i vērtības. Šīs vērtības ir uzdotas tabulas 6.3.14 pirmajā rindā. Izvēlēsimies konstantes vērtību $c = 12$ un no šī skaitļa atņemsim sākotnējās atribūtu vērtības. Transformētās vērtības $c - a_{ii}$ ir parādītas tabulas 6.3.14 otrajā rindā. Biezuma vērtības tabulas pēdējā rindā attiecas uz transformētajām atribūtu vērtībām.

Tabula 6.3.14

Sākotnējās un transformētās atribūtu a_i vērtības

a_i	1	2	3	4	5	6	7	8	9	10	11
$c - a_{ii}$	11	10	9	8	7	6	5	4	3	2	1
n_i	1	3	6	5	4	3	3	2	2	1	1

Transformēto atribūtu vērtību histogramma ir parādīta attēlā 6.3.10.



Attēls 6.3.10. Vērtību $c - a_{i_l}$ sadalījuma histogramma datiem no tabulas 6.3.14

Salīdzinot histogrammu attēlā 6.3.10 ar sākotnējo atribūtu a_i vērtību histogrammu attēlā 6.3.7, varam izdarīt acīmredzamu secinājumu, ka histogramma attēlā 6.3.10 ir attēla 6.3.7 histogrammas spoguļattēls. Bet vērtību $c - a_{i_l}$ sadalījums ir uz labo pusi vērsts sadalījums un šim sadalījumam var izmantot iepriekš minētās un citas piemērotas normalizācijas metodes.

Izmantojot sākotnējā sadalījuma attēlošanas pieeju, ir jāpārlicinās, ka visi rezultāti attiecas tieši uz pārveidoto sadalījumu. Tāpēc ir nepieciešama iegūto rezultātu korekcija, lai tos piemērotu sākotnējam sadalījumam.

Apkoposim rezultātus. Sākotnējo sadalījumu transformācijas, lai tos tuvinātu normāliem sadalījumiem, tiek veiktas dažādos veidos sadalījumiem, kas ir vērsti uz kreiso vai labo pusi.

Logaritmiskās transformācijas tiek plaši izmantotas, lai pārveidotu uz labo pusi vērstos sadalījumus. Šīs transformācijas ir īpaši piemērotas šādos gadījumos:

- iespējama eksponenciāla komponente datos;
- datu vērtības atšķiras pēc izkliedes pakāpes;
- galvenie efekti datos ir multiplikatīvi, bet ne additīvi.

Uz labo pusi vērstu sadalījumu transformācijas ar kvadrātsaknes palīdzību ir piemērotas šādos gadījumos:

- dati var būt aprēķinu rezultāti vai biežumu vērtības;
- datos ir daudz nulles vērtību un ārkārtīgi mazas vērtības.

Uz labo pusi vērstiem sadalījumiem Box-Cox transformācijas dod labus rezultātus, taču optimālās λ vērtības izvēle ir diezgan sarežģīta.

Uz kreiso pusi vērstiem sadalījumu transformācijām plaši izmanto kvadrātiskās un kubiskās transformācijas. Var izmantot arī sākotnējā sadalījuma attēlojumu, kas dod labus rezultātus. Bet tas prasa transformācijas rezultātu pielāgošanu.

Izmantojot pakāpes transformācijas un Box-Cox transformācijas, rodas problēma izvēlēties transformācijas parametru optimālās vērtības. Lai sasniegtu šādu optimizāciju, var ieteikt iestatīt dažādas attiecīgā parametra vērtības, veikt katrai izvēlētajai vērtībai transformāciju un salīdzināt rezultātus. Pēdējais variants ir tas, kas nodrošina vislabāko sākotnējā sadalījuma tuvinājumu normālajam sadalījumam. Kā variants ir salīdzināt asimetrijas un ekscesa aplēses visām transformācijas iespējām. Tā kā nepieciešamās transformācijas un rezultātu analīzi var veikt, izmantojot piemērotus programmatūras rīkus, daudzkārtīgu transformāciju veikšana un iegūto rezultātu analīze nav problēma.

Kā minēts iepriekš šajā sadaļā, asimetrisko sākotnējo sadalījumu transformāciju mērķis ir pārveidot šos sadalījumus gandrīz normālā formā.

Transformētajam atribūtu vērtību a_j sadalījumam var iegūt mūs interesējošos šī sadalījuma statistiskos raksturlielumus: sagaidāmā vērtība, standartnovirze, ticamības intervāls paredzamajai vērtībai un tamlīdzīgi. Tās ir absolūti korektas darbības. Iepriekš šajā sadaļā tika sniegti dažādu transformētu sadalījumu statistikas datu kopsavilkumi. Šie dati tika iegūti, izmantojot SPSS programmatūras pakotni. Iegūtos raksturlielumus var izmantot transformāciju rezultātu analīzei.

Problēma ar sākotnējo sadalījumu transformēšana ir saistīta ar citu šādu transformāciju aspektu. Lieta tāda, ka jebkura transformēta sadalījuma analīze nav tieši saistīta ar sākotnējo sadalījumu. Šāda analīze, kaut arī absolūti pareiza, neko nevar pateikt par sākotnējo sadalījumu, lai gan tas ir viens no analīzes galvenajiem mērķiem.

Lai atrisinātu šo problēmu, ir nepieciešams pārvērst transformētā sadalījuma statistisko analīzi atpakaļ uz sākotnējo sadalījumu. Kā to var panākt?

Transformētā sadalījuma sagaidāmo vērtību var pārvērst, izmantojot apgriezto transformāciju, kas tika piemērota sākotnējam sadalījumam. Tādējādi, ja tika izmantota transformācija ar decimālogaritmiem, tad transformēto atribūtu vērtību kopas sagaidāmā vērtība \bar{a}_j^t tiek pārvērsta sākotnējā datu telpā kā $\bar{a}_j = 10^{\bar{a}_j^t}$. Ja tika izmantota transformācija ar naturālajiem logaritmiem, tad $\bar{a}_j = e^{\bar{a}_j^t}$, ja tika izmantota transformācija ar binārajiem logaritmiem, tad $\bar{a}_j = 2^{\bar{a}_j^t}$.

Ja tika izmantota transformācija, kuras pamatā ir pakāpes funkcija, tad transformētā sadalījuma parametru vērtību apgrieztā transformēšana tiek veikta, pamatojoties uz, piemēram, apgriezto pakāpes funkciju.

- kvadrātiskajai transformācijai – pamatojoties uz kvadrātisko funkciju;
- kvadrātsaknes transformācijai – pamatojoties uz kvadrātsaknes funkciju;
- patvaļīgai x pakāpei - pamatojoties uz pakāpes funkcijas ar pakāpi $\frac{1}{x}$.

Jāņem vērā, ka transformēto sadalījumu standartnovirzes un standartklūdas netiek pārvērstas atpakaļ sākotnējo datu telpā, jo pēc transformēšanas šie raksturlielumi zaudē savu jēgu.

6.4. Atribūtu vērtību transformācija regresijas linearitātes sasniegšanai

Bez atribūtu vērtību transformācijas, lai normalizētu šo vērtību sākotnējos sadalījumus, datu pirmapstrādes jomā ir plaši pielietotas transformācijas lineāru atkarību starp atribūtu vērtību kopām sasniegšanai.

Standarta lineārās regresijas teorētiskais pamats ir dots pielikumā P1.3. Pieņemam, ka ir dotas divas atribūtu vērtību kopas a_j un a_k . Ja starp šīm vērtību kopām pastāv lineāra sakarība, tad šīs sakarības var aprakstīt ar lineārās regresijas vienādojumu (lineārā regresija jau tika apskatīta 3. nodaļas 4. sadaļā):

$$a_k = b_0 + b_1(a_j) + \varepsilon, \quad (6.4.1)$$

kur b_0 - atribūta a_k vērtība regresijas līknes krustpunktā ar vertikālo asi, kas attēlo atribūtu vērtības

a_k ;

b_1 - parametrs, kas nosaka regresijas līknes slīpumu pret horizontālo asi, kas attēlo atribūtu a_j vērtības;

ε - klūdas lielums, kas raksturo datu punktu novirzi no regresijas taisnes.

Ja atribūtu a_j un a_k sākotnējās vērtības ir tādas, ka attiecības starp tām var raksturot ar atbilstošu precizitātes pakāpi, izmantojot regresijas vienādojumu (6.4.1), tad šo vērtību

transformācijas nav nepieciešamas. Problēmas rodas, ja attiecības starp atribūtu vērtību kopām nav lineāras un tās nevar aprakstīt ar regresijas vienādojumu.

Šādā situācijā ir iespējamas šādas darbības:

- modelēt esošo atkarību ar piemērotu nelineārās regresijas vienādojumu un izmantot šo modeli turpmākajā datu apstrādes un analīzes procesā;
- transformēt atribūta a_k vai a_j vērtības (jeb abas), lai iegūtu vēlamo lineāro attiecību.

Šajā sadaļā mēs aplūkojam otro iespēju, proti, viena atribūta vai abu atribūtu vērtību transformāciju.

Visbiežāk, lai panāktu attiecīgu lineāru sakarību, tiek izmantotas logaritmiskās transformācijas. Pēc noklusējuma transformācijām tiek izmantoti naturālie logaritmi. Ir trīs šādu transformāciju veidi.

- lineāri-logaritmiskā transformācija

$$a_k = b_{01} + b_{11} \ln(a_j) + \varepsilon, \quad (6.4.2)$$

kur b_{01} un b_{11} ir attiecīgi krustošanās ar līkni un slīpuma parametri.

- logaritmiski-lineārā transformācija

$$\ln(a_k) = b_{02} + b_{12} a_j + \varepsilon. \quad (6.4.3)$$

- logaritmiski-logaritmiskā transformācija

$$\ln(a_k) = b_{03} + b_{13} \ln(a_j) + \varepsilon. \quad (6.4.4)$$

Nav speciālu ierobežojumu par to, kuras no iepriekš minētajām transformācijām piemērot konkrētajā datu priekšapstrādes uzdevumā. Izvēlei tiek pielietota izmēģinājumu un kļūdu metode. Sākotnējiem datiem tiek piemēroti visi trīs transformāciju veidi. Par galīgo tiek uzskatīts tas variants, kas dod vislabāko transformācijas rezultātu (mazākais kļūdas ε lielums).

Apskatīsim ilustratīvus piemērus atribūtu vērtību transformācijām, lai panāktu regresijas linearitāti.

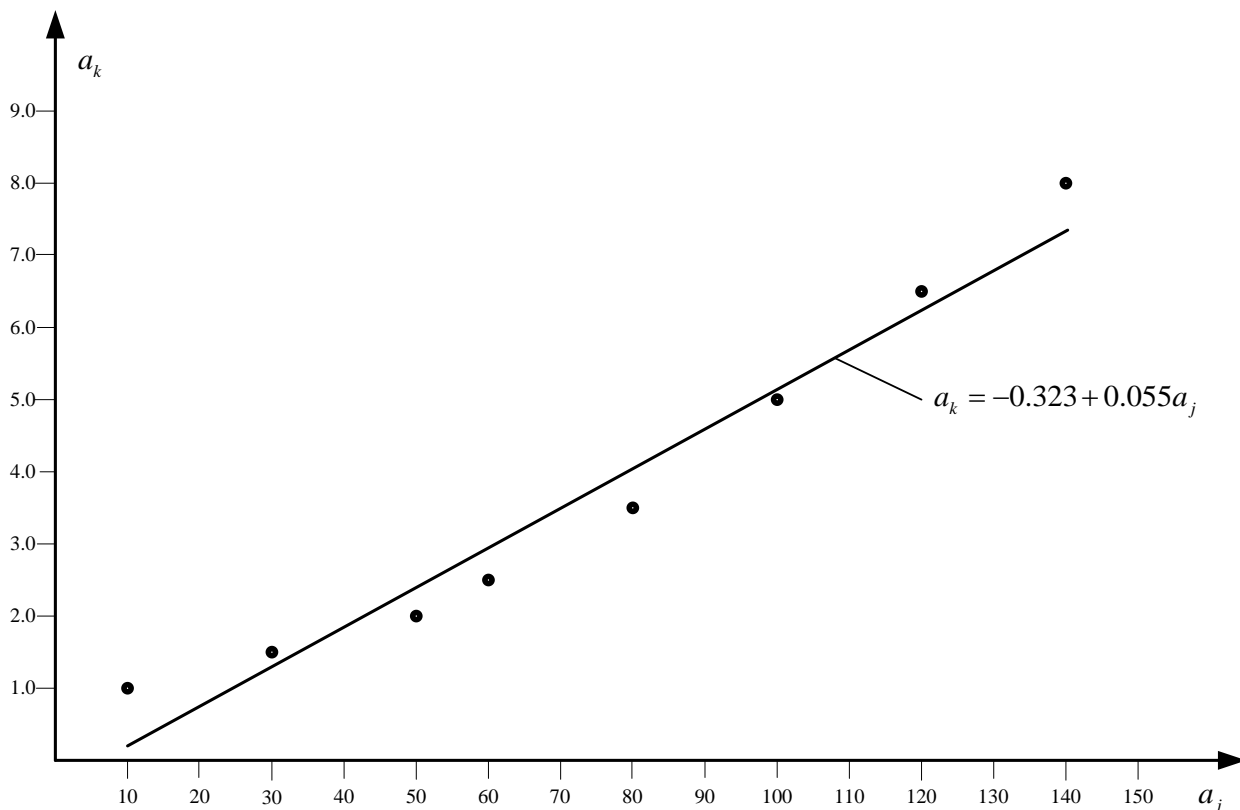
Piemērs 6.4.1. Tabulas 6.4.1 pirmajās divās rindās ir dotas atribūtu a_j , a_k sākotnējās vērtības, trešajā rindā ir parādītas transformētās vērtības $\ln(a_j)$, ceturtajā rindā ir parādītas transformētās vērtības $\ln(a_k)$.

Tabula 6.4.1

Sākotnējās un transformētās atribūtu a_j, a_k vērtības

a_j	10	30	50	60	80	100	120	140
a_k	1.0	1.5	2.0	2.5	3.5	5.0	6.5	8.0
$\ln(a_j)$	2.302	3.401	3.952	4.094	4.382	4.605	4.787	4.942
$\ln(a_k)$	0	0.405	0.693	0.916	1.253	1.609	1.872	2.079

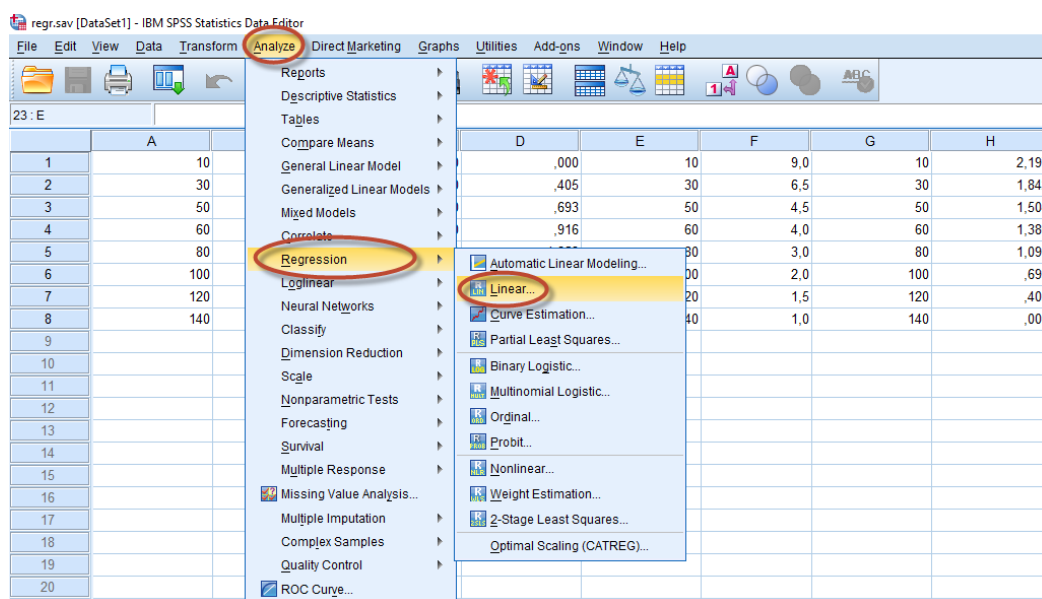
Sākotnējās atribūtu vērtības ir grafiski parādītas attēlā 6.4.1.



Attēls 6.4.1. Sākotnējo atribūtu a_j , a_k vērtību grafiskais attēlojums datiem no tabulas 6.4.1

Tā kā atribūta a_k vērtības palielinās, palielinoties atribūta a_j vērtībām, tad šeit ir pozitīva regresijas sakarība.

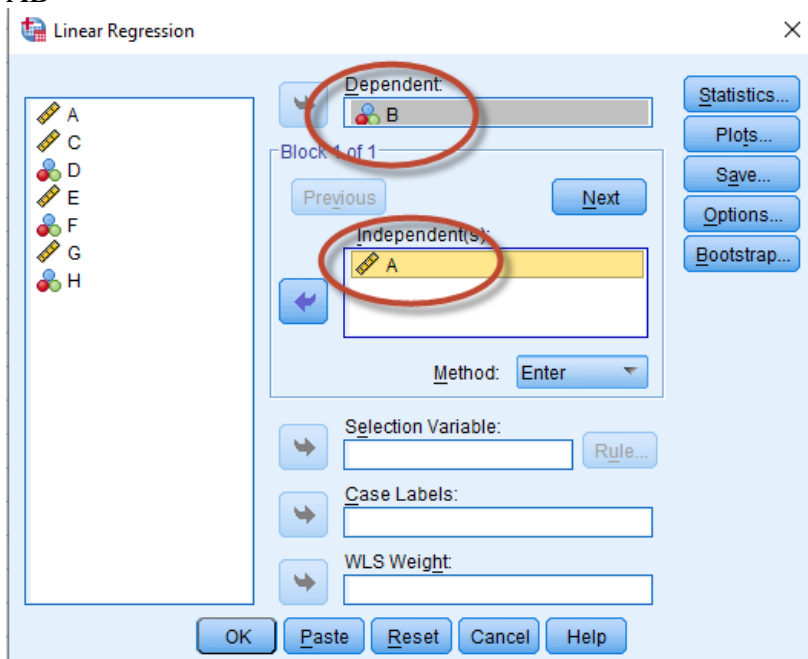
Attēlā 6.4.1 datu vizuālā analīze parāda, ka attiecības starp atribūtu a_j un a_k vērtībām nav lineāras. Mēģināsim aproksimēt šo sakarību, izmantojot lineārās regresijas vienādojumu. Veiksim aprēķinus, izmantojot SPSS programmatūras pakotni. Aprēķini sākas, izvēloties procedūru soļus **Analyze – Regression – Linear**.



Pēc tam jāievada dati. Mūsu gadījumā tiek ievadīti šādi dati: $A = a_j$, $B = a_k$.

Lai novērtētu regresiju starp vērtību kopām A un B, šie dati jāievada attiecīgajos laukos.

AB



Pēc aprēķinu veikšanas mums ir šādi rezultāti:

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,979 ^a	,959	,952	,5532

a. Predictors: (Constant), A

b. Dependent Variable: B

AB

Coefficients^a

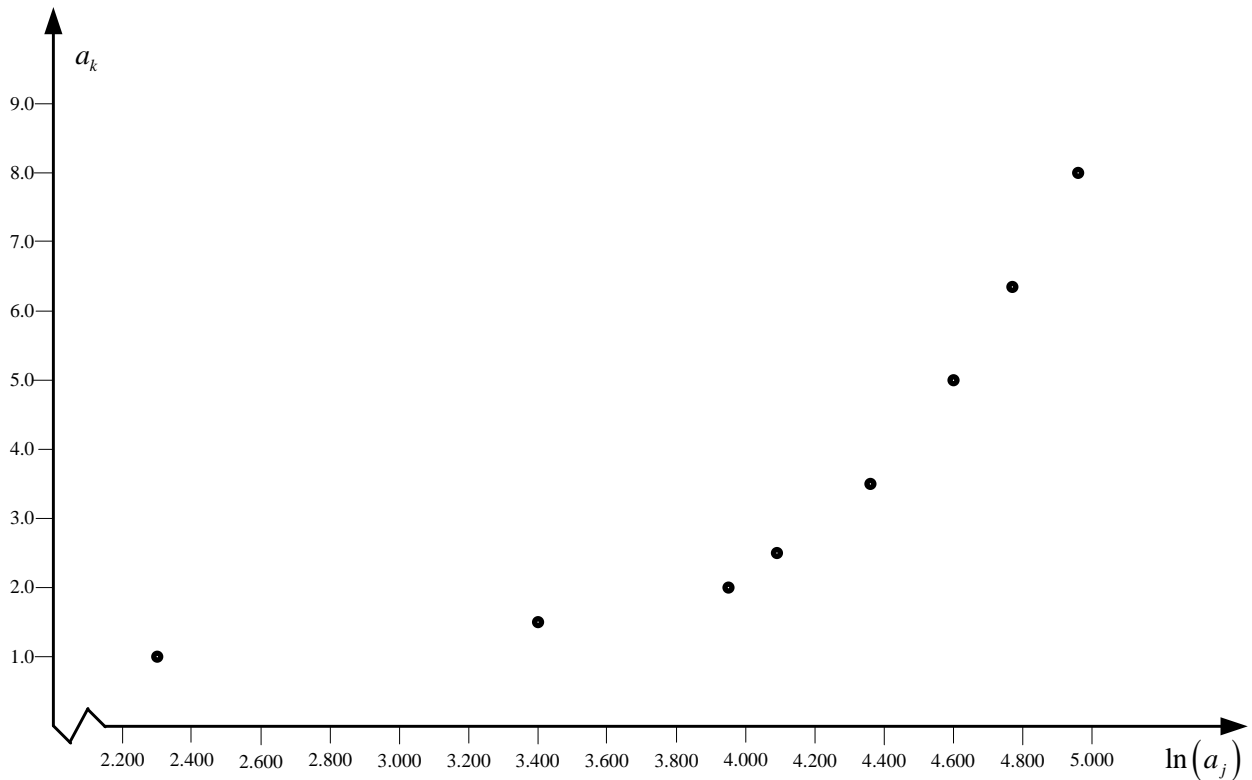
Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95,0% Confidence Interval for B	
		B	Std. Error	Beta			Lower Bound	Upper Bound
1	(Constant)	-,323	,397		-,815	,446	-1,293	,647
	A	,055	,005	,979	11,806	,000	,044	,067

a. Dependent Variable: B

Izmantojot pēdējās tabulas B ailē aprēķinātos datus, mēs iegūstam regresijas vienādojumu $a_k = -0.323 + 0.055a_j$.

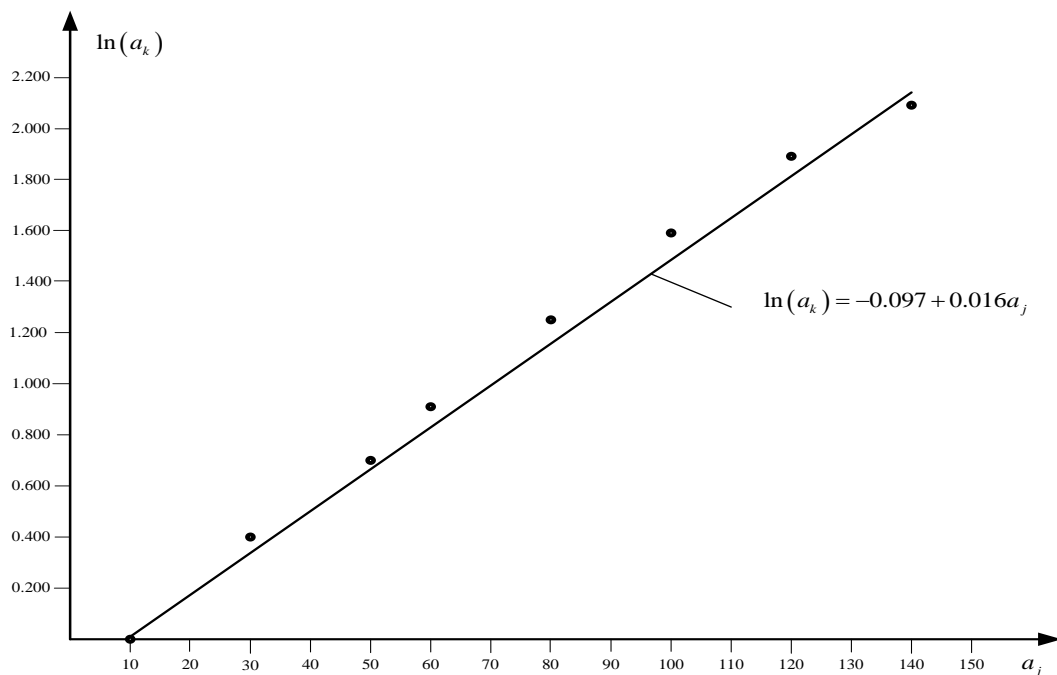
Izmantojot iegūto vienādojumu, mēs izveidojam regresijas līkni, kas parādīta attēlā. 6.4.1. Transformēsim atribūtu a_j vērtības logaritmiskā formā. Vērtību $\ln(a_j)$ un a_k grafiskais attēlojums parādīts attēlā 6.4.2.

Acīmredzams, ka atribūtu a_j vērtību pārveidošana logaritmiskā formā neizveidoja lineāru sakarību starp vērtībām $\ln(a_j)$ un a_k . Šis secinājums tieši izriet no attēla 6.4.2 vizuālās analīzes un neprasa īpašus komentārus.



Attēls 6.4.2. Vērtību $\ln(a_j)$, a_k grafisks attēlojums datiem no tabulas 6.4.1

Pārveidosim atribūtu a_k vērtības logaritmiskā formā. Vērtības a_j un $\ln(a_k)$ grafiski parādītas attēlā 6.4.3.



Attēls 6.4.3. Vērtību a_j , $\ln(a_k)$ grafisks attēlojums datiem no tabulas 6.4.1

Lai noteiktu attiecību izteiksmi starp vērtībām a_j , $\ln(a_k)$, mēs veiksīm nepieciešamos aprēķinus ar SPSS. Ieviesīsim šādus apzīmējumus: $C = a_j$, $D = \ln(a_j)$. Aprēķinu rezultāti ir parādīti divās nākamajās tabulās:

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,997 ^a	,993	,992	,065367

a. Predictors: (Constant), C

b. Dependent Variable: D

CD

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95,0% Confidence Interval for B	
		B	Std. Error	Beta			Lower Bound	Upper Bound
1	(Constant)	-,097	,047		-2,064	,085	-,211	,018
	C	,016	,001	,997	29,442	,000	,015	,018

a. Dependent Variable: D

Mums ir iegūts šāds regresijas vienādojums:

$$\ln(a_k) = -0.097 + 0.016a_j.$$

Izmantojot šo vienādojumu, mēs izveidojam regresijas likni, kas parādīta attēlā 6.4.3.

Analizēsim iegūtos rezultātus. Regresijas “kvalitātes” rādītāji ir koeficienta R^2 vērtība un standartkļūdas vērtība, kas atbilst kļūdas lielumam ε lineārās regresijas vienādojumos (6.4.1) – (6.1.4). Šo abu regresiju rādītāju kopsavilkums ir parādīts tabulā 6.4.2.

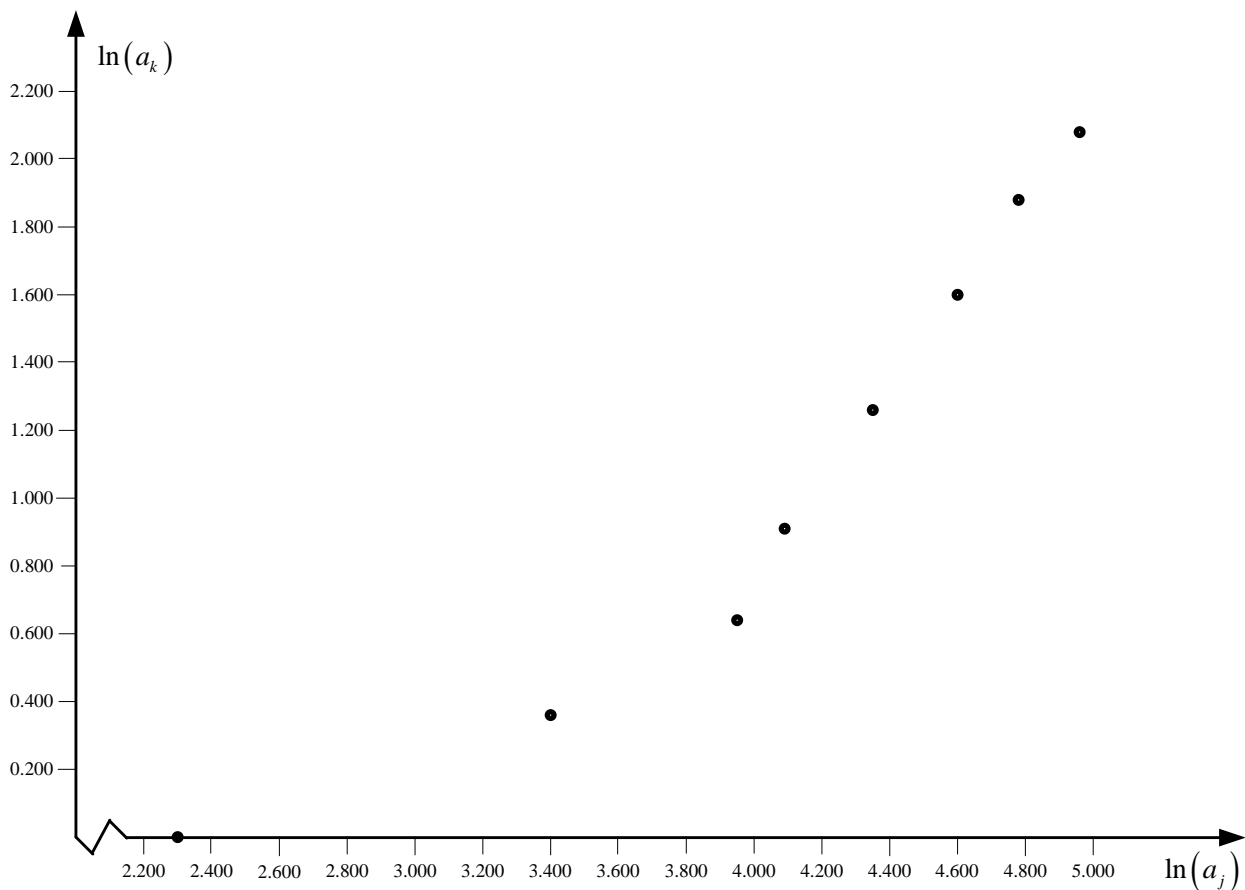
Tabula 6.4.2

Regresijas attiecību $a_j - a_k$, $a_j - \ln(a_k)$ rezultējošo rādītāju kopsavilkums no piemēra 6.4.1

Regresija	R^2	ε
$a_j - a_k$	0.959	0.5532
$a_j - \ln(a_k)$	0.993	0.0654

Regresijas attiecības $a_j - \ln(a_k)$ gadījumā vērtība $R^2 = 0.993$ ir lielāka par vērtību $R^2 = 0.959$ priekš attiecības $a_j - a_k$, un vērtība ir $\varepsilon = 0.0654$ ievērojami mazāka par vērtību $\varepsilon = 0.5532$. Tas nozīmē, ka atribūtu vērtību a_k logaritmiskā transformācija ļāva sasniegt ievērojami precīzāku regresijas $a_j - \ln(a_k)$ attiecību, nekā $a_j - a_k$.

Attēlā 6.4.4 vērtības $\ln(a_j)$, $\ln(a_k)$ tiek parādītas grafiskā veidā.



Attēls 6.4.4. Vērtību $\ln(a_j)$, $\ln(a_k)$ grafisks attēlojums datiem no tabulas 6.4.1

Attēlā 6.4.4 redzamo datu vizuālā analīze parāda, ka sakarība starp vērtībām $\ln(a_j)$, $\ln(a_k)$ nav lineāra.

Tādējādi sākotnējām atribūtu a_j , a_k vērtībām no tabulas 6.4.1 tikai atribūtu vērtību logaritmiskā transformācija a_k nodrošina lineāru sakarību starp vērtībām a_j un $\ln(a_k)$.

Piemērs 6.4.2. Tabulā 6.4.3 parādītas atribūtu a_l , a_m vērtības un transformētās vērtības $\ln(a_l)$, $\ln(a_m)$.

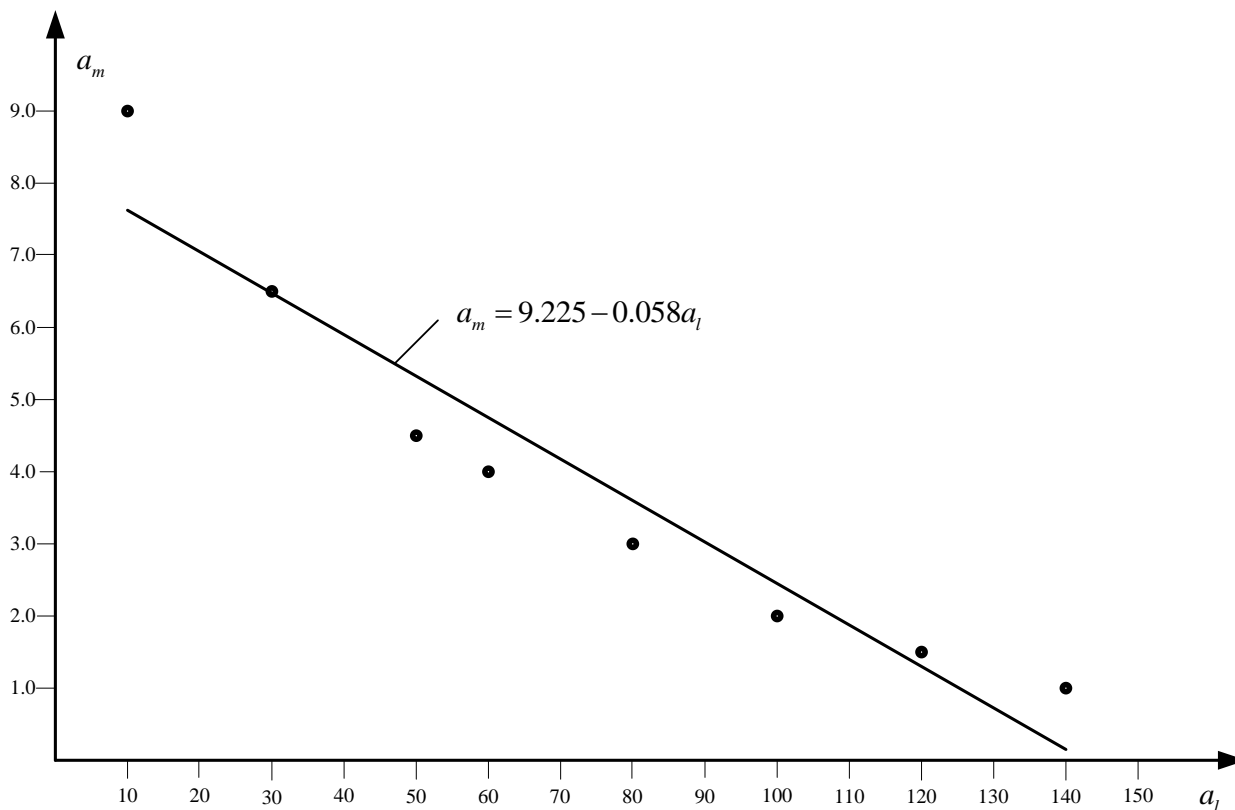
Tabula 6.4.3

Sākotnējās un transformētās atribūtu a_l, a_m vērtības

a_l	10	30	50	60	80	100	120	140
a_m	9.0	6.5	4.5	4.0	3.0	2.0	1.5	1.0
$\ln(a_l)$	2.302	3.401	3.952	4.094	4.382	4.605	4.787	4.942
$\ln(a_m)$	2.197	1.872	1.504	1.386	1.090	0.693	0.405	0

Tā kā atribūtu a_m vērtības samazinās, palielinoties atribūtu a_l vērtībām, tad šeit ir negatīva regresijas sakarība.

Attēls 6.4.5 grafiski attēlo atribūtu a_l , a_m vērtības datiem no tabulas 6.4.3.



Attēls 6.4.5. Atribūtu a_l, a_m vērtību grafisks attēlojums datiem no tabulas 6.4.3

Attēlā 6.4.5 redzamo datu vizuālā analīze parāda, ka sakarības starp atribūtu a_l un a_m vērtībām nav lineāras. Mēģināsim aproksimēt šo sakarību, izmantojot lineārās regresijas vienādojumu. Aprēķini tika veikti, izmantojot SPSS. Apzīmējam: $E = a_l$, $F = a_m$. Pēc attiecīgo aprēķinu veikšanas mums ir šādi rezultāti:

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,956 ^a	,914	,900	,8612

a. Predictors: (Constant), E

b. Dependent Variable: F

EF

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95,0% Confidence Interval for B	
		B	Std. Error	Beta			Lower Bound	Upper Bound
1	(Constant)	8,225	,617		13,323	,000	6,714	9,735
	E	-,058	,007	-,956	-7,983	,000	-,076	-,040

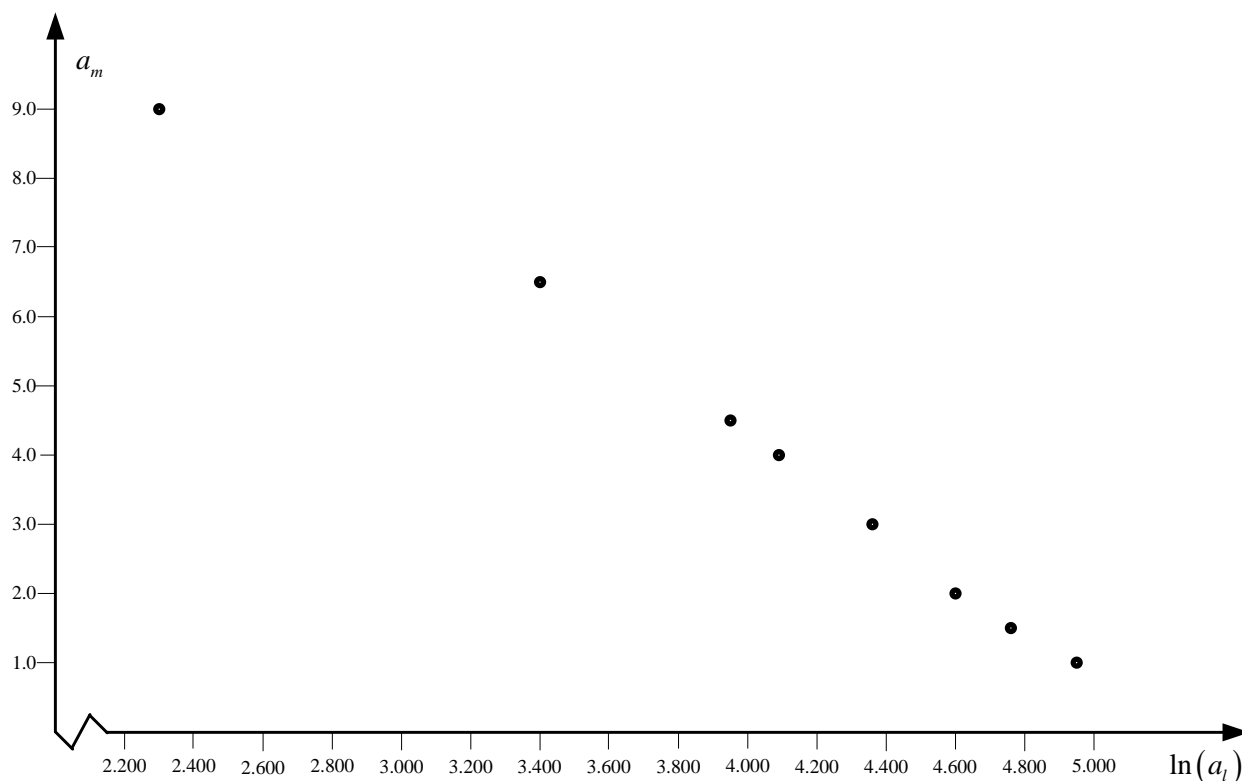
a. Dependent Variable: F

Mums ir iegūts šāds regresijas vienādojums:

$$a_m = 8.225 - 0.058a_l.$$

Izmantojot šo vienādojumu, mēs izveidojām regresijas līkni, kas parādīta attēlā 6.4.5.

Pārveidosim atribūtu a_l vērtības logaritmiskā formā. Vērtību $\ln(a_l)$ grafiskais attēlojums ir a_m parādīts attēlā 6.4.6.



Attēls 6.4.6. Vērtību $\ln(a_l)$, a_m grafisks attēlojums datiem no tabulas 6.4.2

Acīmredzams, ka sakarība starp vērtībām $\ln(a_l)$ un a_m ir “lineārāka” nekā sakarība starp atribūtu sākotnējām vērtībām a_l un a_m . Tomēr šī sakarība nav stingri lineāri atkarīga.

Mēs logaritmiski transformējam atribūtu a_m vērtības. a_l un $\ln(a_m)$ vērtības grafiski parādītas attēlā 6.4.7.

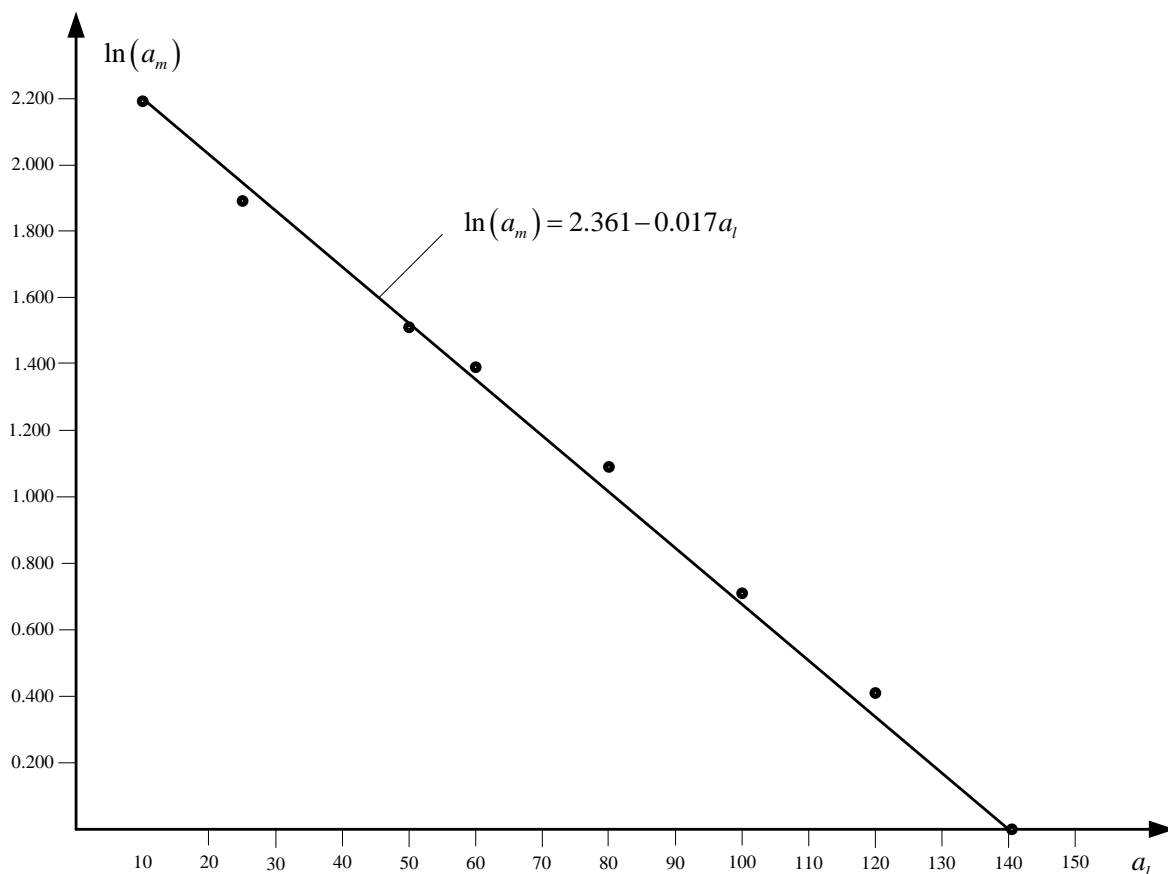
Attēlā 6.4.7 redzamo datu vizuālā analīze ļauj secināt, ka sakarība starp vērtībām a_l un $\ln(a_m)$ ir “aptuveni” lineāra Ieviesīsim apzīmējumu: $G = a_l$, $H = \ln(a_m)$. Ar SPSS palīdzību noteiksim lineārās regresijas vienādojumu starp šīm vērtībām un analizēsim sakarību statistiskos raksturlielumus.

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,999 ^a	,998	,998	,035376

a. Predictors: (Constant), G

b. Dependent Variable: H



Attēls 6.4.7. Vērtību $a_l, \ln(a_m)$ grafisks attēlojums datiem no tabulas 6.4.3

G.H.

Coefficients^a

Model	Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95,0% Confidence Interval for B	
	B	Std. Error	Beta			Lower Bound	Upper Bound
1 (Constant)	2,361	,025		93,113	,000	2,299	2,423
G	-,017	,000	-,999	-55,377	,000	-,017	-,016

a. Dependent Variable: H

Iegūts regresijas vienādojums

$$\ln(a_m) = 2.361 - 0.017a_l.$$

Izmantojot šo vienādojumu, mēs izveidojām regresijas līkni, kas parādīta attēlā 6.4.7.

Abu regresiju statistikas kopsavilkums ir parādīts tabulā 6.4.4.

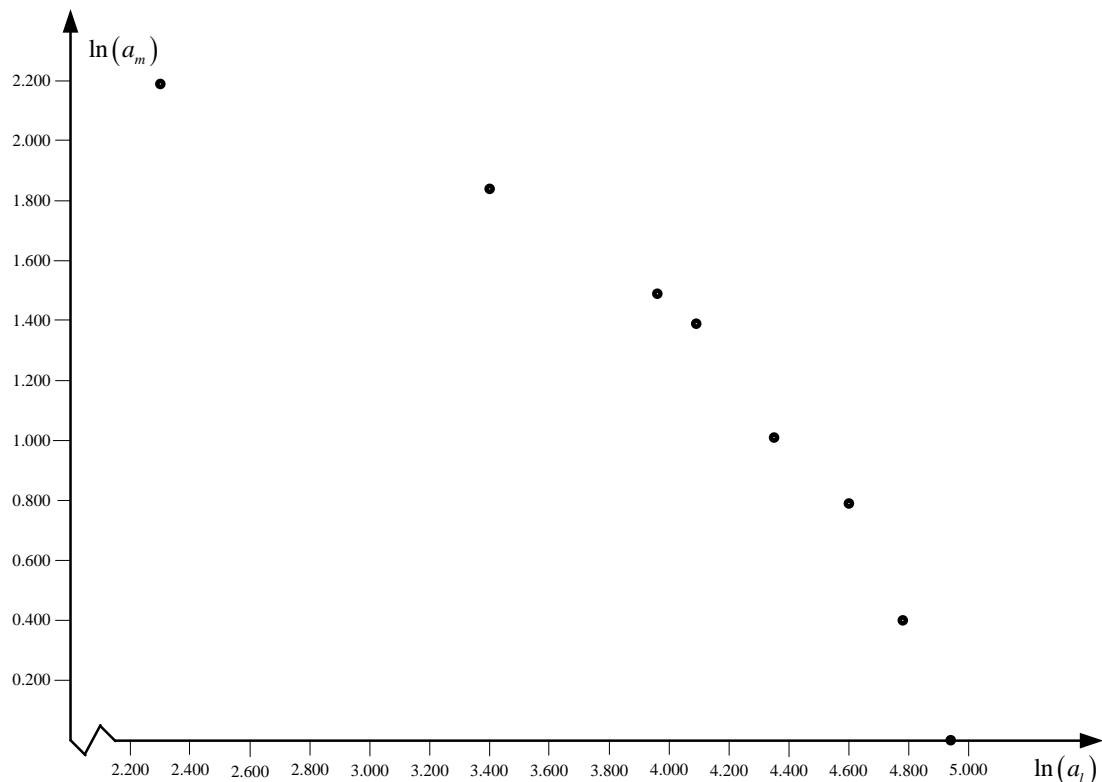
Tabula 6.4.4

Regresijas sakarību $a_l - a_m, a_l - \ln(a_m)$ rezultātu rādītāju kopsavilkums piemēram 6.4.2

Regresija	R^2	ε
$a_l - a_m$	0.914	0.8612
$a_l - \ln(a_m)$	0.998	0.0354

Ir pilnīgi skaidrs, ka atribūtu a_m vērtību logaritmiskā transformācija izraisīja stingru lineāru sakarību starp vērtībām a_l un $\ln(a_m)$.

Attēlā 6.4.8 grafiskā veidā tiek parādītas $\ln(a_l)$, $\ln(a_m)$ vērtības.



Attēls 6.4.8. Vērtību $\ln(a_l)$, $\ln(a_m)$ grafisks attēlojums datiem no tabulas 6.4.2

Attēlā 6.4.8 redzamo datu vizuālā analīze parāda, ka sakarība starp vērtībām $\ln(a_l)$ un $\ln(a_m)$ nav lineāra.

Abos iepriekšminētajos piemēros tika panākta regresijas linearitāte starp netransformētajām atribūtu a_j un logaritmiski transformētām vērtībām $\ln(a_k)$, $(\ln(a_m))$.

Kā minēts iepriekš, nav vispārīgu ieteikumu par to, kuras atribūtu vērtību transformācijas novedīs pie lineāras regresijas sakarības. Literatūrā bieži tiek izmantots abu atribūtu vērtību logaritmisko transformāciju “klasiskais” piemērs, lai panāktu lineāras regresijas sakarību starp to transformētajām vērtībām. Piemēram, viens no atribūtiem var būt ātri augošas priedes stumbra diametrs, ko mēra centimetros. Otrs atribūts ir priedes stumbra tilpums, ko mēra kubikcentimetros. Tikai ar abu atribūtu vērtību logaritmisku transformēšanu ir iespējams panākt lineāras regresijas sakarības starp šīm transformētajām vērtībām.

Apskatīsim lineārās regresijas koeficienta b_1 interpretāciju, izmantojot dažādas attiecīgo atribūtu vērtību transformācijas.

Ja izmantojam sākotnējo atribūtu a_j , a_k vērtību standarta lineāro regresiju (vienādojums 6.4.1), tad koeficienta b_1 vērtība parāda atribūta a_k vērtības izmaiņas, kad atribūta a_j vērtība mainās par 1 vienību. Acīmredzot ar pozitīvu korelāciju starp atribūtu a_j , a_k vērtībām, b_1 ir pozitīvs un palielinot b_1 vērtību par 1 vienību, atribūta a_k vērtība palielinās par b_1 vienībām. Ja starp atribūtu vērtībām a_j , a_k ir negatīva korelācija, tad atribūta a_j vērtībai palielinoties par 1 vienību, atribūta a_k vērtība samazinās par b_1 vienībām.

Koeficienta b_1 vērtības interpretācija standarta lineārā regresijā starp sākotnējām atribūtu vērtībām a_j, a_k ir saprotama un neprasa nekādus papildu paskaidrojumus.

Problēmas ar koeficienta b_1 interpretāciju rodas, izmantojot transformētās atribūtu a_j, a_k vērtības. Sniedzam šī koeficienta interpretācijas dažāda veida transformācijām.

1. *Lineāri-logaritmiskais modelis* (vienādojums 6.4.2).

Vērtības $\ln(a_j)$ palielināšana par 1 vienību izraisa atribūta a_k vērtības izmaiņas (palielinās vai samazinās atkarībā no koeficienta zīmes) par b_1 vienībām.

Kādā veidā var interpretēt vērtības $\ln(a_j)$ izmaiņas par 1 vienību saistībā ar atribūta a_j sākotnējām vērtībām? Lai to izdarītu, mēs izmantojam šādu sakarību:

$$\ln(a_j) + 1 = \ln(a_j) + \ln(1) = \ln(e * (a_j)).$$

No tās izriet, ka vērtības $\ln(a_j)$ palielināšana par 1 vienību ir līdzvērtīga naturālā logaritma izteiksmei $\ln(e * (a_j))$.

Alternatīvs veids, kā interpretēt operāciju vērtības a_j reizināšanai ar e , ir izteikt a_j vērtības procentuālo pieaugumu, izmantojot šādu sakarību:

$$100 * (2.72 - 1) = 172\%.$$

No tās izriet, ka koeficients b_1 atspoguļo atribūta a_j vērtības pieaugumu, kad vērtība $\ln(a_j)$ palielinās par 172%.

2. *Logaritmiski-lineārais modelis* (vienādojums 6.4.3).

Atribūta a_j vērtības palielināšana par 1 vienību izraisa $\ln(a_k)$ izmaiņas (palielinājumu vai samazinājumu) par b_1 vienībām. Sākotnējo atribūtu a_k vērtību izteiksmē tas nozīmē, ka atbilstošā atribūta a_k vērtība tiek reizināta ar e^{b_1} . Citiem vārdiem sakot, atribūta a_j vērtības palielināšana par 1 vienību izraisa a_k vērtības izmaiņas par e^{b_1} vienībām.

Cits atribūta a_j vērtības palielinājums par kādu citu skaitli, piemēram, c vienībām, izraisa sākotnējās atribūta a_k vērtības izmaiņas par e^{cb_1} vienībām.

3. *Logaritmiski-logaritmiskā transformācija* (vienādojums 6.4.4).

Šajā gadījumā atribūta a_j vērtības reizināšana ar e atbilst atribūta a_k vērtības reizināšanai ar e^{b_1} .

Izmantojot šo modeli, tiek izmantotas proporcionālas atribūta a_k vērtības izmaiņas (palielinājums vai samazinājums), kad atribūta a_j vērtība palielinās par noteiktu skaitu, piemēram, c procentiem. Lai to izdarītu, vispirms jāaprēķina

$$d = \ln(100 + c) / 100.$$

Šai d vērtībai atribūta a_k vērtība mainās par e^{db_1} procentiem.

Papildus logaritmiskajām transformācijām, lai panāktu lineāras regresijas sakarību, tiek izmantotas kvadrātsaknes transformācijas un Box-Cox transformācijas. Tomēr pēdējā gadījumā transformācijas parametra λ optimālās vērtības noteikšana ir sarežģīta matemātiska problēma.

IZMANTOJAMĀ LITERATŪRA

- Acuna E., Rodrigues C. (2004). The Treatment of Missing Values and its Effect on Classifier Accuracy. In book "Classification, Clustering and Data Mining Applications, pp. 639 – 647.
- Alasadi S. (2017). Review of Data Preprocessing Techniques in Data Mining. Journal of Engineering and Applied Sciences, 12 (16), pp. 4102 – 4107.
- Acuna E. (2011). Preprocessing in Data Mining. In Lovric M. (Ed.). International Encyclopedia of Ststistical Science, Springer, doi: 10.1007/978-3-642-04898-2_51.
- Aljuaid T., Sasi S. (2016). Intelligent Imputation Technique for Missing Values. Int. Conference on Advances in Computing, Communications and Informatics, September 21 – 24, 2016, Jaipur, India, pp. 2441 – 2445.
- Agresti A., Franklin Ch. (2013). Statistics. The Art and Science of Learning from data. (third edition). Pearson Education, Inc., 834 p.
- Anath C.V., Kleinbaum D.G. (1997). Regression Model for Ordinal Responses: A Review of Methods and Applications. Information Journal of Epidemiology, Vol. 26, (6), pp. 1323 – 1333.
- Agresti A. (2010). Modelling Ordinal Categorical data.
www.stat.ufl.edu/~aa/ordinal/ord.html.
- Agresti A. (2002). Categorical Data Analysis (second edition). John Wiley & Sons, Inc., Hoboken, New Jersey, 394 p. DOI:10.1002/0471249688.
- Agresti A. (2007). An Introduction to Categorical Data Analysis. John Wiley & Sons, Inc., Hoboken, New Jersey, 394 p., DOI:10.1002/0470114754.
- Alison P.D. (2005). Imputation of categorical variables with proc mi. Proceedings of the SUGI 30, 113 (30), pp. 1 -14.
- Allison P.D. (2000). Multiple imputation for missing data: A cautionary tale. Sociological Methods and Research, 28, pp. 301 – 309.
- Aggarwal Ch.C. (2016). Outlier Analysis (Second Edition).
<http://rd.springer.com/book/10.1007/978-3319-47578-3>.
- Aşikgil B., Erar A. (2009). Research into Multiple Outliers in Linear Regression Analysis. Hacetpere Journal of Mathematics and Statistics, 28 (2), pp. 185 – 198.
- Alma Ö.G. (2011). Comparison of Robust Regression Methods in Linear Regression. Int. J. Contemp. Math. Science, 6 (9), pp. 409 – 421.
- Anitha S., Metilda M., (2016). A Survey on Cluster Based Outlier Detection Techniques In Data Stream. Int. Journal of Data Mining Techniques and Applications, 5, (1), pp. 96 – 101.

- Batista G.E., Monard M.-C. (2003). An Analysis of Four Missing Data Treatment Methods for Supervised Learning. *Applied Artificial Intelligence*, Vol. 17, Issue 5 – 6, pp. 519 – 533.
- Batista G., Monard M.C. (2003). Experimental comparison of K-nearest neighbor and mean or mode imputation methods with the internal strategies used by C4.5 and CN2 to treat missing data. Technical Report, University of Sao Paulo.
- Batista G., Monard M.C. (2002). A study of K-nearest neighbor as an imputation method. In Abraham A. et al. (Eds.) *Hybrid Intell. Systems*, IOS Press, pp. 251 – 260.
- Barnard J., Rubin D.B. (1999). Small-sample degrees of freedom with multiple imputation. *Biometrika*, 86 (4), pp. 948 – 955, doi: 10.1093/biomet/86.4.948.
- Burgin M. (2004). Data, Information, and Knowledge. *Information*, 7 (1), pp. 47 – 57.
- Behrens J.T. (1997). Principles and Procedures of Exploratory Data Analysis. *Psychological Methods*, Vol. 2, No. 2, pp. 131 – 160.
- Batista G.E., Monard M.C. (2001). A Study of k-Nearest Neighbour as a Model-Based Method to Treat Missing Data. In *Proceedings of the Argentine Symposium on Artificial Intelligence*, Vol. 30, Buenos-Aires, Argentina, pp. 1 – 9.
- Bürkner P.-Ch., Vuorre M. (2019). Ordinal Regression Models in Psychology: A Tutorial. *Advances in Methods and Practices in Psychological Science*, Vol. 2 (1), pp. 77 – 101. DOI:10.1177/2515245918823199.
- Beretta L., Santaniello A. (2016). Nearest neighbor imputation algorithms: a critical evaluation. *Medical Informatics and Decision Making*, 16, pp. 197 – 208. DOI:10.1186/s12911-016-0318-z.
- Ben-Gal I. (2005). Outlier Detection. In Maimon O and Rockash (Eds.). *Data Mining and Knowledge Discovery Handbook: A Complete Guide for Practitioners and Researches*, Kluwer Academic Publishers, 16 p.
- Barnett C., Lewis T. (1994). *Outliers in Statistical Data*. (3rd edition). John Wiley & Sons, Chichester, 584 p.
- Bogashaw G.B., Yohannes Y.B. (2020). Review of Outlier Detection and Identifying Using Robust Regression Model. *Int Journal of Systems Science and Applied Mathematics*, 5 (1), pp. 4 – 11.
- Breunig M.M., Kriegel H.-P., Ng R.T., Sender J. (2000). LOF: identifying density-based local outliers. *ACM sigmod record*, 29 (2), pp. 93 – 104.
- Bishop Ch.M. (2006). *Pattern Recognition and Machine Learning*. Springer, 758 p.
- Boulle M. (2004). Khiops: A Statistical Discretization Method of Continuous Attributes. *Machine Learning*, 55, pp. 53 – 69.
- Box G.E.P., Cox D.R. (1964). An Analysis of Transformations. *Journal of the Royal Statistical Society, Series B (Methodological)*, 26 (2), pp. 211 – 252.

- Cios K.J., Pedrycz W., Swiniarski R.W., Kurgan L.A. (2007). *Data Mining. A Knowledge Discovery Approach*. Springer, 600 p.
- Cheng H.Y. (1994). *Abduction? Deduction? Induction? Is There a Logic of Exploratory Data Analysis?* Annual Meeting of the American Educational Research Association, New Orleans, LA, April 4 8, 28 p.
- Cheng H.Y. (2010). Exploratory Data Analysis in the Context of Data Mining and Resampling. *Int. Journal of Psychological Research*, 3, (1), pp. 9 – 22.
- Chen H.-Ch., Wang N.-Sh. (2014). The Assignment of Scores Procedure for Ordinal Categorical Data. *The Scientific World Journal*, Vol. 2014, 7 p, <http://dx.doi.org/101155/2014/304213>.
- Chen J., Shao J. (2000). Nearest neighbor imputation for survey data. *Journal of official Statistics*, 16,
- Cousineneau D. (2010). Outliers detection and treatment: a review. *Int. Journal of Psychological Research*, 5 (1), pp. 58 – 67.
- Çetin M., Toka O. (2011). The Comparing of S-estimators and M-estimators in Linear Regression. *Journal of Science, Gazi University*, 24 (4), pp. 747 – 752.
- Christy A.G., Gandhi M., Vaithyasubramanian (2015). Cluster Based Outlier Detection for Healthcare Data. *Procedia Computer Science*, 50, pp. 209 – 215.
- Chmielewski M.R., Grzymala-Busse J.W. (1996). Global Continuous Attributes as Preprocessing for Machine Learning. *Int. Journal of Approximate Reasoning*, 15, pp. 319 - 331.
- Ching J.Y., Wong A.K.C. (1995). Class-Dependent Discretization for Inductive Learning from Continuous and Mixed-Mode Data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17, (7), pp. 641 – 651.
- Dong Y., Peng Ch.-Y. J. (2013). *Principled missing data methods for researcher*. Springer Plus, 2, 17 p.
- Dempster A.P., Laird N.M., Rubin D.B. (1977). Maximum Likelihood from Incomplete Data via the EM algorithm. *Journal of the Royal Statistical Society, Series B (Methodological)*, 9 (1), pp. 1 – 38.
- De Mast J., Kemper B.P.H. (2009). Principles of Exploratory Data Analysis in Problem Solving: What Can We Learn from a Well-Known Case? *Quality Engineering*, 21, pp. 366 – 375.
- Dubani S.A. (1976). *Transactions. Systems Man Cybernetics*, 6 (4), pp. 325 – 327.
- Do Ch.B., Batzoglu S. (2008). What is the expectation maximization algorithm. *Primer. Nature Biotechnology*, 26 (8). pp. 897 – 899.

- Dempster A.P., Laird N.M., Rubin D.B. (1997). Maximum likelihood from incomplete data via the EM algorithm. *Journal Royal Statistical Society*, 39 (1), pp. 1 – 38.
- Dan E.D., Ijeona O.A. (2013). Statistical analysis/methods of detecting outliers in a univariate data in a regression analysis model. *Int. Journal of Education and Research*, 1 (5), 24 p.
- Davies L., Gather U. (1993). The identification of multiple outliers. *Journal of the American Statistical Association*, 88 (423), pp. 782 – 792.
- Dan E.D., Ijeoma O.A. (2013). Statistical Analysis/Methods Of Detecting Outliers in a Univariate Regression Analysis Model. *Int. Journal of Education and Research*, 1, 24 p.
- Dimitrova E.S., Licona M.P.V., McGee J., Laubenbacher R. (2010). Discretization of Time Series Data. *Journal of Computational Biology*, 17, (6), pp. 853 – 868.
DOI:10.1089/cmb.2008.0023.
- Dash R., Dash R. (2011). Comparative Analysis of Supervised and Unsupervised Discretization Techniques. *Int. Journal of Advances in Science and Technilogy*, 2 (3), pp. 29 – 37.
- Dougherty J., Kohavi R., Sahami M. (1995). Supervised and unsupervised discretization of continuous features. *Proceedings of the 12th Int. Conference on Machine Learning*. San Francisco, pp. 194 – 202.
- Dęlkowska K., Jarocka M. (2013). The impact of the methods of the data normalization on the result of linear ordering. *Acta Universitatis Lodzianensis, Polia Oeconomica*, 286, pp. 181 – 188.
- Forhangfar A., Kurgan L.A., Dy J. (2008). Impact of imputation of missing values on classification error for discrete data. *Pattern Recognition*, 41, pp. 3693 – 3705.
- Feng H., Guoshun C., Yang B., Chen Y. (2005). A SVM regression based approach to filling in missing values. In Khosla R., Howlett R.J., Jain L.C. (Eds.) "KES (3)", Vol. 3683 of lecture notes in computer science, Derlin, Springer, pp. 581 – 587.
- Farhangfar A., Kurgan L.A., Pedrycz W. (20007). A novel framework for imputation of missing values in databases. *IEEE Transaction on Systems Man Cybernetics, Part A*, 37 (5), pp. 692 – 709.
- Faisal Sh., Tutz G. (2016). Nearest Neighbor Imputation for Categorical Data by Weighting of Attributes. *Proceedings of the 31st International Workshop on Statistical Modeling*, Vol. 1, 28 p.
- Feng Ch., Wang H., Lu N., Chen T., He H., Lu Y., Tu X.M. (2014). Log-transformation and its implications for data analysis. *Shanghai Arch Psychiatry*, 26 (2), pp. 105 – 109.
- Fayyad U.M., Irani B. (1993). Multi-Interval Discretization of Continuous-Valued Attributes for Classification Learning. *Proceedings of the 13th Int. Joint Conference Artificial Intelligence*, pp. 1022 – 1027.
- Feng Ch., Wang H., Lu N., Chen T., He H., Lu Y., Tu X.M. (2014). Log-transformation and

its implications for data analysis. *Shanghai Arch Psychiatry*, 26 (2), pp. 105 – 109.

Grzymala-Buse J.W., Grzymala-Buse W.J. (2005). Handling missing attribute values. In Maimon O., Rokash L. (Eds.). *Data Mining and Knowledge Discovery Handbook*, Springer, Boston, MA, pp. 37 – 57.

Grzymala-Buse J.W., Ming Hu (2000). A comparison of Several Approaches to Missing Attribute Values in Data Mining. *Proceedings of Int. Conference on Rough Sets and Current Trends in Computing*, pp. 378 – 388.

Garía S., Luengo J., Herrera F. (2015). *Data Preprocessing in Data Mining*. Springer, 327 p.

Grzymala-Buse J., Goodwin L., Grzymala-Buse W., Zhang X. (2005). Handling missing attribute values in preterm birth data sets. *Proceedings of 10th Int. Conference of rough sets and fuzzy sets and data mining and granular computing*, pp. 342 – 352.

Geler Z., Kurbalija V., Radovanović M., Ivanović M. (2016). Comparison of different weighting schemes for the kNN classifier on time-series data. *Knowledge and Information Systems*, 48, pp. 331 – 378. <https://doi.org/10.1007/s10115-015-0881-0>.

Gökmen N., Erar A. (2018). Robust Estimators In Linear Regression: A Simulation Study And An Application On Health Care. *11 International Statistics Days Conference*, 3 – 7 October 2018, pp. 874 – 885.

Gan G., Ng K.-P. (2017). k-means clustering with outlier removal. *Pattern Recognition Letters*, 90, pp. 8 – 14.

Garsia S., Luengo J., Saez J.A., Lopez V. Herrera F. (2011). A Survey of Discretization Techniques: Taxonomy and Empirical Analysis in Supervised Learning. *IEEE Transaction On Knowledge and Data Engineering*. pp.

Grzymala-Buse J.W. (2004). Three strategies to rule induction from data with numerical attributes. *Lecture Notes in Computer Science* 3135, pp. 54 – 62.

Grünwald P.D. (2007). *The minimum description length principle*. The MIT Press, Cambridge, London, 703 p.

Howell D.C. (2007).. The analysis of missing data. In Outwaite W. & Turner S. *Handbook of Social Science Methodology*, London, Sage.

Hestie T., Tibshirani R., Friedman J. (2009). *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*. Second Edition, Springer, 764 p.

Haig B.B. (2005). An Abductive Theory of Scientific Method. *Psychological Methods*, Vol. 10, No. 4, pp. 371 – 388.

Han J., Kamber M., Pei J. (2012). *Data Mining Concepts and Techniques*, Third Edition, Morgan Kaufmann Publishers, 703 p.

Howell D. (2007). *The analysis of missing data*. SAGE Publication, London.

- Hechenbichler K., Schliep K. (2004). Weighted k-Nearest-Neighbor Techniques and Ordinal Classification. Ludwig-Maximilians-Universität München. Sonderforschungsbereich 386, Paper 399, 16 p.
- Hedge V.J. Augustin J. (2004). A Survey of Outlier Detection Methodologies. Sydney, Kluwer Academic Publishers.
- Hadi A.S. (1992). Identifying multiple outliers in multivariate data. Journal of the Royal Statistical Society, Series B, 54, pp. 761 – 771.
- Huber P.J. (1981). Robust Statistics. John Wiley & Sons, New York.
- Hawkings D.M. (1980). Identification of Outliers. Chapman and Hall-Science, 188 p.
- <https://towardsdatascience.com/11-dimensionality-reduction-techniques-you-should-know-in-2021-dcb9500d388b>.
- Hemada B., Lakshmi K.S.V. (2013). A Study On Discretization Techniques. Int Journal of Engineering Research & Technology, 2, (8), pp. 1887 – 1892.
- Ho K.M., Scott P.D. (1997). Zeta: A Global Method for Discretization of Continuous Variables. III Int. Conference on Knowledge Discovery and Data Mining, pp. 191 – 194.
- Holte R.C. (1993). Very simple classification rules perform well on most commonly used data sets. Machine Learning, 11, pp. 63 – 90.
- Iglewitz B., Hoaglin D.C. (1993). How to detect and handle outliers. ADQC Quality Press.
- Jović W A., Brkić K., Bogunović N. (2015). A review of feature selection methods with applications. Proceedings of the 38 Int. Convention on Information and Communication Technology, Electronics and Microelectronics, 6 p.
- Jain S., Yain K., Chodhary N (2016). A survey paper on missing data in data mining. Int. Journal of Innovations in Engineering Research and Technology, 3 (12), pp. 45 – 50.
- Jönsson P., Wohlin C. (2004). An Evaluation of k-nearest Neighbour Imputation Using Likert Data. Proceedings of the 10th International Symposium on Software Metrics, pp. 108 – 118.
- Jolliffe I.T. (2002). *Principal Component Analysis* (Second Edition). Springer-Verlag, New York, Berlin Heidelberg, 487 p.
- Jackson J.D. (1991). *A User's Guide to Principal Component Analysis*. John Willey & Sons, Inc.
- Jahan A., Edwards K.L. (2015). A state-of-the-art survey on the influence of normalization techniques in ranking: Improving the materials selection process in engineering design. Materials and Design, 65, pp. 335 – 342.
- Krzanowski W. (1986). Multiple discriminant analysis in the presence of mixed continuous and categorical data Computer & Mathematics Applications, 12 (2, part A), pp. 179 – 185.

- Kang H. (2013). The prevention and handling of the missing data. *Korean J. Anesthesiol*, 64 (5), pp. 402-406. doi: 10.4097/kjal.2013.64.5.402.
- Kuppusamy V., Paramasivam I. (2016). A Study of Impact on Missing Categorical Data - A Qualitative Review. *Indian Journal of Science and Technology*, (32), 6 p., doi:10.17485/ijst/2016/v9i32/83088.
- Kwak S.K., Kim J.H. (2017). Statistical data preparation management of missing values and outliers. *Korean Journal of Anesthesiology*, 79 (4), pp. 407 – 411.
- Knorr E.M., Ng R.T. (1999). Finding Intensional Knowledge of Distance-Based Outliers. *Proceedings of the 25th ULDB Conference*, Edinburg, Scotland, 12 p.
- Knorr E.M., Ng R.T. (1998). Algorithms for Mining Distance-Based Outlier in Large Datasets. *Proceedings of the 24th WLDB Conference*, New York, USA, pp. 392 – 403.
- Kumar V., Kumar S., Singh K. (2013). Outlier Detection: A Clustering-Based Approach. *Int Journal of Science and Modern Engineering*.
- Kaufman L., Rousseeuw P.J. (1987). Clustering by means of Medoids. In (Y.Dodge, ed.). *Statistical Data Analysis Based on the Norm and Related Methods*, North-Holland, pp. 405 – 416.
- Kotsiantis S., Kanellopou D. (2006). Discretization Techniques: A recent survey. *GESTS Int. Transactions on Computer Science and Engineering*, 32, (1), pp. 47-58.
- Kurgan L.A., Krzysztof J.C. (2004). CAIM Discretization Algorithm. *IEEE Transactions On Knowledge and Data Engineering*, 16, (2), .
- Kerber R. (1992). Chimerge: Discretization of numeric attributes. *Proceedings of the Ninth National Conference Artificial Intelligence*, pp. 123 – 128.
- Laencina P.J.G., Sancho-Gomez J.L., Figneras-Vidal A.R. (2010). Pattern classification with missing data: A review. *Neural Computing and Applications*, 19, (2), pp. 263 – 282.
- Little R.J.A., Rubin D.B. (2002). *Statistical analysis with missing data*. 2nd Edition, Wiley, New Jersey.
- Liu Ch.-H., Tsai Ch.-F., Sue K.-L., Huang M.-W. (2020). The Feature Selection Effect On Missing Value Imputation on Medical Datasets. *Applied Sciences*, 10, 12 p.
- Li D., Deogun J., Spaulding W., Shuart B. (2004). Towards missing data imputation: a study of fuzzy k-means clustering method. *Proceedings of 4th Conference of*
- Luengo J., Garcia S., Herrera F. (2011). On the choice the best imputation methods for missing values considering three groups of classification methods. *Knowledge and Information Systems*, 32, pp. 77 – 108.
- Liu B.Q. (2014). *The application of exploratory data analysis in auditing*. A dissertation submitted to the Graduate School-Newark, The State University of New Jersey for the degree Doctor of Philosophy.

- Little R.J.A., Schluchter M.D. (1985). Maximum likelihood estimation for mixed continuous and categorical data with missing values. *Biometrika*, 72, pp. 497 – 512.
- Liddell T., Krushke J.K. (2018). Analyzing ordinal data with metric models: What could possibly go wrong? *Journal of Experimental Social Psychology*, 79, pp. 328 – 348.
- Larose D.T. (2005). *Discovering Knowledge in data: An Introduction to Data Mining*. Wiley.
- Landerman L.R., Kenneth C.L., Pieper C.F. (1979). An Empirical Evaluation of the Predictive Mean Matching Method for Imputation Missing Values. *Sociological Methods & Research*, 26, pp. 3 – 23.
- LaLonde S.M. (2012). Transforming Variables for Normality and Linearity – When, How, Why and Why Not’s. *SAS Global Forum 2012. Statistics and Data Analysis*, 8 p.
- Lee D.K. (2020). Data transformation: a focus on the interpretation. *Korean Journal Anesthesiol*, 73 (6), pp. 503 – 508.
- LakshmiM.S. (2018). An Overview Study on Data Cleaning, its Types and its Methods for Data Mining. *Int. Journal of Pure and Applied Mathematics*, 119 (12), pp. 16837 - 16848.
- Loureiro A., Togo L. Soares C. (2004). Outlier detection using clustering methods: a data cleaning application. *Proceedings of KDNNet Symposium on Knowledge-based Systems for the Public Sector, Bonn, Germany*.
- Lin T.Y. (2002). Attribute Transformations for Data Mining I: Theoretical Explanations. *Int. Journal of Intelligent Systems*, 17, pp. 213 – 222.
- Liu H., Hussain F., Lim C., Dash M. (2002). Discretization. An Enabling Technique. *Data Mining and Knowledge Discovery*, 6 (4), pp.393 – 423.
- Lu D.K. (2020). Data transformation: a focus on the interpretation. *Korean Journal Anesthesiol*, 73 (6). pp. 503 – 508.
- LaLonde S.M. (2012). Transforming Variables for Normality and Linearity – When, How, Why, and Why Not’s. *SBS Global Forum 2012, Paper 430-2012*. 8 p.
- McGullgh P. (1980). Regression model for ordinal data. *Journal of the Royal Statistical Society, series B*, 42, pp. 109 – 142.
- Moon T.K. (1996). The expectation-maximization algorithm. *IEEE Signal Proceeding Magazine*, December 1996, pp. 47 – 60.
- Mielke A. (2016). *Robust Statistics*. Trufa Science Inside № 16 5p.
- Mofanavalli S.S., Poornima N.S. (2018). Outlier Detection using Clustering Techniques. *Int. Journal of Engineering & Technology*, 7, pp. 813 – 818.
- Moore Y.D., McCabe G. (1999). *The Practice of Statistics*, New York, Freeman.
- Mohanavalli S.S., Sriprija N., Poornima S. (2018). Outlier Detection using Clustering

Techniques. *Int. Journal of Engineering & Technology*, 7, pp. 813 – 818.

- Muhlenbach F., Rakotomalala R. (2005). Discretization of Continuous Attributes. *John Wang Encyclopedia of Data Warehousing and Mining*, Idea Group Reference, pp. 397 – 402.
- Nügroho W.H., Wardhani N.W.S., Fernandes A.R., Solimun (2020). Robust Regression Analysis Study for Data with Outliers at Some Significance Ltvels. *Mathematics and Statistics*, 8 (4), pp. 373 – 381.
- Nayak S.C., Misra B.B., Behera H.S. (2014). Impact of Data Normalization on Stock Index Forecasting. *Int. Journal of Computer Information Systems and Industrial Management Applications*, 6, pp. 257 – 269.
- Oba S., Sato M., Takamasa I., Monden M., Matsubara K. (2003). A bayesian missing value estimation method for gene expression profile data. *Biometrics*, 19 (16), pp. 2088 – 2096.
- Ozdemir O. (2016). A comparison Study of Data Transformation Methods to Achieve Normality. *Int. Journal of Mathematical and Computational Methods*, 1, pp. 382 – 383.
- Osborne J.W. (2002). Notes on the Use of Data Transformation. *Practical Assessment, Research & Evaluation*, 8 (6), 7 p.
- Pyle D. (1999). *Data preparation for data mining*. Morgan Kaufman, San Francisco, 560 p.
- Pitush K.A., Stevens J.P. (2016). *Applied Multivariate Statistics for the Social Sciences. Analyses with SAS and IBM's SPSS (sixth edition)*. Taylor & Francis, 814 p.
- Peng L., Stuart E.A., Allison D.B. (2015). Multiple Imputation: A Flexible Tool for Handling Missing Data. *JAMA*, 314 (18). pp. 1966 – 1967.
- Penny K.I., Jolliffe I.T. (2001). A comparison of multivariate outlier detection methods for clinical laboratory safety data. *The Statistician*, 50 (3), pp. 295 – 308.
- Porkodi F.A. (2014). A comparison of filter based feature selection algorithms: An overview. *Int. Journal of Innovative Research in Technology & Science*, 2, (2), pp. 1957 – 2000.
- Palczewski K., Salabun W. (2019). Influence of various normalization methods in PROMETHEE: an empirical study on the selection of the airport location. *Procedia Computer Science*, 159, pp. 2051 – 2060.
- Pandey A., Jain A. (2017). Analysis of KNN Algorithm using Various Normalization Techniques. *Int Journal Computer Network and Information Security*, 11, pp. 36 – 42.
- Quinlan J.R. (1986). Induction of Decision Trees. *Machine Learning*, 1 (1), pp. 81 – 106.
- Rubin D.B. (1976). Inference and missing data. *Biometrika*, 63 (3), pp. 581 – 592.
- Rubin D.B. (1996). Multiple Imputation after 18+ Years. *JASA*, 91, pp. 473 – 478, doi: 10.1080/01621459.1996.10476908.
- Rubin D.B. (1987). *Multiple imputation for nonresponse in surveys*. John Wiley & Sons,

Inc., New York.

- Rowley J. (2007). The wisdom hierarchy: representation of the DIKW hierarchy. *Journal of Information Science*, 33 (2), pp. 163 – 180, /doi: 10.1177/0165551506070706.
- Rubin D.B. (1996). Multiple Imputation After 18+ Years. *Journal of the American Statistical Association*, 91 (434), pp. 473 – 489.
- Rubin D.B. (1976). Inference and Missing Data. *Biometrika*, 63, pp. 581 – 592.
- Rubin D.B. (1977). Formalizing Subjective Notion about The Effect of Nonrespondents in Sample Surveys. *Journal of the American Statistical Association*, 72, pp. 538 – 543.
- Rubin D.b. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York, John Wiley & Sons.
- Rousseeuw P.J, Leory A. (1987). *Robust Regression and Outlier Detection*. Wiley Series In Probability and Statistics.
- Rousseeuw P.J., van Zomeren B.C. (1990). Unmasking Multivariate Outliers and Leverage Points. *Journal of the American Statistical Association*, 85, pp. 633 – 639.
- Rousseeuw P.J., Hubert M. (2011). Robust statistics for outlier detection. *WIREs Data Mining Knowledge Discovery*, 1, pp. 73 – 79.
- Rousseeuw P.J. (1987). Silhouettes: A Practical Aid to the Interpretation and Validation of Cluster Analysis. *Computational and Applied Mathematics*, 20, pp. 53 – 65.
- Rissanen J. (2007). *Information and complexity in statistical modeling*. New York, Springer, 142 p.
- Schafer J.L. (1997). *Analysis of incomplete multivariate data*. Chapman & Hall/CRC, London.
- Schafer J.L. (1999). Multiple imputation: a primer. *Stat. Methods in Medicine*. 8 (1), pp. 3 – 15, doi: 10.1177/096228029900800102.
- Schafer J.L., Olsen M.K. (1998). Multiple Imputation for Multivariate Missing-Data Problems. A Data Analyst's Perspective. *Multivar. Behav. Res.*, 33 (4), pp. 545 – 571, doi: 10.1207/s1532790mbr2204_5.
- Schafer J.L., Graham J.W. (2002). Missing data: Our view of the state of art. *Psychological Methods.*, 7 (2), pp. 147 – 177, doi: [org/10.1037/1082-989X.7.2.147](https://doi.org/10.1037/1082-989X.7.2.147).
- Schneider T. (2001). Analysis of incomplete climate data: estimation of mean values and covariance matrices and imputation of missing values. *Journal of Climate*, 14, pp. 853 – 871.
- Sanders J.D. (2016). Definition Terms: Data, information and Knowledge. *SAI Computing Conference*, 6 p. doi: 10.1109/SAI.2016.7555986.
- Stonier T. (1993). *The wealth of information*. London: Thames/Methuen.

- Stonier T. (1997). *Information and meaning – An Evolutionary perspective*. Berlin, Springer.
- Safarinejadan B., Menhaj M., Karrari M.B. (2010). A distributed EM algorithm to estimate the parameters of a finite mixture of components. *Knowl. Inf. Syst.*, pp. 267 – 292.
- Shylaja B., Kumar R.S. (2018). Traditional versus modern missing data handling techniques: an overview. *Int. Journal of Pure and Applied Mathematics*, 118 (14), pp. 77 – 83.
- Sim J., Lee J.S., Kwon O. (2015). Missing Values and Optimal Selection of an Imputation Method and Classification Algorithm to Improve the Accuracy of Ubiquitous Computing Applications. *Mathematical Problems in Engineering*, Vol. 2015, 14 p., doi: 10.1155/2015/538613.
- Schmitt P., Mandel J., Guedj M., (2015). A comparison of Six Methods for Missing Data Imputation. *Journal Biomet. Biostat.*, 6 (1), pp. 1 – 6.
- Schwenger H., Ickstadt K. (2008). Imputing missing genotypes with weighted k nearest neighbor. Technical Report< No. 2008,03, Technische Univesität Dortmund, 15 p.
- Sinharay S., Stern H.S., Russel D. (2001). The Use of Multiple Imputation for the Analysis of Missing Data. *Psychological Methods*, 6 (4), pp. 317 – 329.
- Schafer J.L. (1997). *Analysis of incomplete multivariate data*. Chapman & Hall, New York.
- Schafer J.L., Olsen M.K. (1998). Multiple imputation for multivariate missing-data problems: A data analysis perspective. *Multivariate Behavioral Research*, 33, pp. 545 – 571.
- Stephen R.S., Santhamarai K.K., (2017). Detection of Outliers in Regression Model for Medical Data. *Int. Journal of Medical Research & Health Sciences*, 6 (7), pp. 50 – 56.
- Salgato C., Arevedo C., Proença H.M., Vieira S. (2016). Noise Versus Outliers. In *Secondary Analysis of Electronic Health Records*. Springer, Cham, https://doi.org/10.1007/978-3-319-43742-2_2.
- Susanti Y., Prativi H., Sri Sulistijowati H., Liana T. (2014). M Estimation, S Estimation, and MM Estimation in Robust Regression. *Int. Journal of Pure and Applied Mathematics*, 91 (3), pp. 349 – 360.
- Soni K.G., Patel A. (2017). Comparative Analysis of k-means and k-medoids Algorithm on IRIS Data. *Int. Journal of Computational Intelligence Research*, 13, (5), pp. 899 – 906.
- Suarez-Alvarez M.N., Pham D.-T., Prostov M.Y., Prostov J.I. (2012). Statistical approach to normalization of feature vectors and clustering of mixed databases. *Proceedings of the Royal Society*, 468, pp. 2630 – 2651, <https://doi:10.1098/rspa.2011.0704>.
- Sakia R. (1992). The Box-Cox Transformation Technique: A Review. *Journal of the Royal Statistical Society. Series D (Statistician)*, pp. 169 – 178.
- Tabachnik B.G., Fidell L.s. (2012). *Using multivariate statistics*. 6th Edition, Allyn & Becon, Needhaum, Heights, MA.
- Tukay J.W. (1977). *Exploratory data analysis*. Reading, MA: Addison-Wesley.

- Tutz G., Ramzan Sh. (2014). Improved Methods for the Imputation of Missing Data by Nearest Neighbor Methods. Technical Report Number 172, Department of Statistics, University of Munich, 22 p.
- Trojanskaja O., Cantor M., Sherlock G., Brown P., Hastie T., Tibshirani R., Botstein D., Altman R.B. (2007). Missing value estimation method for DNA microarrays. *Bioinformatics*, 17, pp. 1 – 6.
- Tanner M.A. Wong W.H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association*, 82, pp. 528 – 550.
- Tukey J. (1997). *Exploratory data analysis*. Pearson.
- Tharwad A. (2017). *Principal component analysis – a tutorial*.
<https://www.academia.edu/36699625/Principal-Component-Analysis>.
- Thorndike R.L. (1953). Who Belongs in the Family. *Psychometrika*, 18, (4), pp. 267 – 271.
- Tukey J.W. (1957). The comparative anatomy of transformations. *Annals of Mathematical Statistics*, 28, pp. 602 – 632.
- Vafaei N., Ribeiro R.A., Camarinha-Matos L.M. (2018). Data normalization techniques in decision making: case study with TOPSIS method. *Int Journal Information and Decision Sciences*, 10 (1).
- Wong A.K.C., Chin D.K.Y. (1987). Synthesizing statistical knowledge from incomplete mixed-mode data. *IEEE Trans. Pattern Anal. Mach. Intell.*, 9 (6), pp. 796 – 805.
doi: 10.1109/tpami.1987.4767986.
- Winship Ch., Mare R.D. (1984). Regression Models with Ordinal Variables. *American Sociological Review*, 49 (4): 512, pp. 511 – 525, DOI:10.2307/2095465.
- Wu J. (2012). *Advances in k-means Clustering: A Data Mining Thinking*, Springer, 194 p.
- Xu G., Zong Yu, Yang Zh. (2013). *Applied Data Mining*. Taylor & Francis Group.
- Yamini C., Kowsalya (2015). Survey of classification of incomplete data handling technologies. *Int. Journal of Advance Research in Science and Engineering*, Vol 4, No. 4, Special Issue (01), pp. 111 – 118.
- Zins Ch. (2007). Conceptual Approaches for Defining Data, Information, and Knowledge. *Journal of the American Society for Information Science and Technology*, 58 (4), pp. 479 -493.
- Zaki M.J., Wagner M. Jr. (2014). *Data Mining and Analysis: Fundamental Concepts and Algorithms*. Cambridge University Press, 607 p.
- Ужга-Ревров О.И. (2021). Введение в анализ данных. Резекне.

PIELIKUMI

P1. REGRESIJAS ATRIBŪTU VĒRTĪBU SADALĪJUMA PARAMETRU APRĒĶINI

Atribūtu vērtību kopu var uzskatīt par šo vērtību sadalījumu. Datu priekšapstrādes un analīzes problēmās bieži vien ir nepieciešams raksturot atribūtu vērtību sadalījumu. Raksturošanu var saistīt ar sadalījuma formas noteikšanu, izvietojuma un diapazona parametru aprēķināšanu un attiecību novērtēšanu starp atsevišķām atribūtu vērtību kopām. Šajā pielikumā ir sniegtas pieejas izvietojuma parametru vērtību un atribūtu vērtību sadalījuma diapazona aprēķināšanai un regresijas noteikšanai starp atsevišķu atribūtu vērtību kopām.

Sadalījumu un regresijas parametru tipi ir atkarīgi no skalu veida, kurā tiek mērītas atribūtu vērtības. Turpmāk tiks parādītas aprēķinātās izteiksmes katram skalu tipam vispārpieņemtajā šo skalu hierarhijā.

P1.1. Nominālās skalas

Ja atribūtu vērtības mēra nominālā skalā, tas nozīmē tikai to, ka katrs no atbilstošajiem objektiem vai subjektiem pieder kādai no šīs skalas kategorijām. Nominālās skalas piemēri var ietvert subjektu dzimumu, politisko partiju piederību, subjektu matu krāsu un tamlīdzīgus elementus.

Nominālā skala ir pati vienkāršākā no visām mērījumu skalām. Tā sniedz ļoti ierobežotu informāciju par lietu būtību. Tomēr diezgan bieži tieši šajā skalā tiek mērītas atribūtu vērtības.

Kādi parametri var raksturot atribūtu vērtību sadalījumu, kas mērīts nominālā skalā? Vienīgais šāda sadalījuma parametrs ir *moda*. Moda ir statistisks raksturlielums, kas nosaka skaitli, kam ir vislielākais biežums kādā datu kopā.

Pieņemam, ka j -tais atribūts raksturo subjektu dzimumu. Kodēsim šī atribūta vērtības šādi: 1 – vīrietis, 2 – sieviete. Pieņemsim, ka ir dota šāda atribūtu a_j vērtību kopa:

$$a_j^T = 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1.$$

(Šajā piemērā mēs attēlojam atribūtu a_j vērtības formā a_j^T , jo datu tabulās atribūtu vērtības ir attēlotas šīs tabulas kolonās, un mūsu piemērā šīs vērtības ir attēlotas rindas veidā.)

No dotā piemēra ir skaidrs, ka dati ietver sevī 9 vīriešus un 6 sievietes, tātad sadalījuma moda ir vienāda ar $\text{mod}(a_j) = 1$.

Citiem vārdiem sakot, sadalījuma parametrs ir iestatīts uz vērtību 1 - vīrieši.

Definēsim regresijas jēdzienu starp diviem mainīgajiem (divām atribūtu vērtību kopām), kas mērīti nominālajā skalā. Tā kā standarta lineārās regresijas modelis šim gadījumam nav piemērots (skat. 11. nodaļu darbā [Pitush K.A., Stewens I.R., 2016] un 13.6 sadaļu darbā [Agresti A., Franklin Ch., 2013]), tiek izmantots *logistiskās regresijas modelis*. Šī modeļa būtība ir šāda. Dota faktoriālā pazīme (mainīgais) - \mathbf{X} un rezultējošā pazīme (mainīgais) - \mathbf{Y} . Mainīgās vērtības tiek izteiktas nominālajās skalās. Pieņemsim, ka mainīgie \mathbf{X} un \mathbf{Y} ietver divas nominālās kategorijas. (Vēlāk mēs paplašināsim šo ierobežojumu līdz patvaļīgam skaitam mainīgā lieluma \mathbf{X} nominālo kategoriju).

Logistiskās regresijas būtību skaidrosim ar konkrētu piemēru. Pirms gaidāmajām pilsētas vadības vēlēšanām tika veikta iedzīvotāju aptauja: kurš no diviem kandidātiem A vai B , ir piemērotākais pilsētas mēra amatam. Aptaujā kopumā piedalījās 100 respondenti, no tiem 70 bija vīrieši un 30 sievietes. Respondentu viedokļi dalījās šādi: vīrieši: par A - 50, par B - 20; sievietes: par A - 10, par B - 20.

Mūsu uzdevums ir noteikt attiecības starp katra kandidāta uzvaras varbūtību un respondentu viedokļiem katrā nominālajā kategorijā.

Sākotnējos datus kodēsim šādi:

- $y = 1$ - kandidāta A uzvara vēlēšanās;
 $y = 2$ - kandidāta B uzvara vēlēšanās;
 $x = 1$ - aptaujas dalībnieku vīriešu kategorija;
 $x = 2$ - aptaujas dalībnieku sievietes kategorija.

Uzdomam sākotnējos datus tabulas P1.1.1 veidā.

Tabula P1.1.1

Iedzīvotāju aptaujas rezultāti par diviem kandidātiem

Respondentu grupas	Kandidāti		
	A ($y = 1$)	B ($y = 2$)	Summa rindās
Vīrieši ($x = 1$)	50	10	60
Sievietes ($x = 2$)	20	20	40

Pārveidosim aptaujas rezultātus no tabulas P1.1.1 ar biežuma (varbūtību) aplēsēm (tabula P1.1.2).

Tabula P1.1.2

Kandidāta uzvaras varbūtību aplēses, pamatojoties uz iedzīvotāju aptaujas rezultātiem

Respondentu grupas	Kandidāti		
	A ($y = 1$)	B ($y = 2$)	Summa rindās
Vīrieši ($x = 1$)	0.833	0.167	1.000
Sievietes ($x = 2$)	0.500	0.500	1.000

Šajā tabulā varbūtību vērtības ir aprēķinātas kā vīriešu īpatsvars, kuri attiecīgi deva priekšroku kandidātiem A un B , kā arī sievietes īpatsvars, kas deva priekšroku kandidātiem A un B .

Aprēķināsim kandidātu izredzes uzvarēt, pamatojoties uz tabulā P1.1.2 dotajiem, šādi.

$$\text{Ch}(y = j)/(x = i) = \frac{p(y = j)/(x = i)}{1 - p(y = j)/(x = i)}, \quad i, j = 1, 2 \quad (\text{P.1.1.1})$$

kur $p(y = j)/(x = i)$ - kandidāta j uzvaras iespējamība pamatojoties uz respondentu grupas i aptaujas rezultātiem.

Mums ir iegūti šādi rezultāti:

$$\text{Ch}(y = 1)/(x = 1) = \frac{0.873}{1 - 0.873} = \frac{0.873}{0.167} = 4.988;$$

$$\text{Ch}(y = 1)/(x = 2) = \frac{0.500}{1 - 0.500} = \frac{0.500}{0.500} = 1.000;$$

$$\text{Ch}(y = 2)/(x = 1) = \frac{0.167}{1 - 0.167} = \frac{0.167}{0.833} = 0.200;$$

$$\text{Ch}(y=2)/(x=2) = \frac{0.500}{1-0.500} = \frac{0.500}{0.500} = 1.000;$$

Iegūtie rezultāti ir apkopoti tabulā P.1.1.3.

Tabula P1.1.3

Kandidātu izredžu aplēses, pamatojoties uz aptaujas rezultātiem

Respondentu grupas	Kandidāti	
	A (y = 1)	B (y = 2)
Vīrieši (x = 1)	4.988	1.000
Sievietes (x = 2)	0.200	1.000

Ja ir zināmas kandidātu iespējas uzvarēt noteiktos apstākļos, tad uz tā pamata var aprēķināt viņu uzvaras varbūtību. Piemēram, priekš $x = 1$ un $y = 1$

$$p(y=1)/(x=1) = \frac{p(y=1)/(x=1)}{1+p(y=1)/(x=1)}. \quad (\text{P1.1.2})$$

Ievietojot skaitliskās vērtības, iegūstam

$$p(y=1)/(x=1) = \frac{4.988}{1+4.988} = \frac{4.988}{5.988} = 0.714.$$

Aprēķinātā varbūtības vērtība sakrīt ar tās vērtību tabulā P1.1.2. Pārējās varbūtības tiek aprēķinātas pēc analogijas.

Jāņem vērā, ka varbūtības aprēķiniem, izmantojot vienādojumu (P1.1.2) mūsu piemērā nav jēgas, jo pašas varbūtības vērtības tiek aprēķinātas, pamatojoties uz aplēstajām varbūtības vērtībām. Tomēr šo procedūru var izmantot praktiskos loģistiskās regresijas pielietojumos.

Ieviesīsim *izredžu attiecības* jēdzienu.

$$\xi(y=1) = \frac{\text{Ch}(y=1)/(x=1)}{\text{Ch}(y=1)/(x=2)}; \quad (\text{P1.1.3.a})$$

$$\xi(y=2) = \frac{\text{Ch}(y=2)/(x=1)}{\text{Ch}(y=2)/(x=2)}. \quad (\text{P1.1.3.b})$$

Aizstājot skaitliskās vērtības, mums ir

$$\xi(y=1) = \frac{4.988}{1.000} = 4.988; \quad \xi(y=2) = \frac{0.200}{1.000} = 0.200$$

Būtībā izredžu attiecība parāda, cik reizi izredzes uzvarēt kandidātam *A* vai *B* pēc vīriešu vērtējuma, ir lielākas vai mazākas nekā šo kandidātu izredzes uzvarēt, pamatojoties uz sieviešu vērtējumiem.

Kā minēts iepriekš, loģistiskās regresijas mērķis ir noteikt rezultatīvā mainīgā lieluma kategoriju varbūtības (mūsu piemērā: kandidāts *A* vai kandidāts *B* uzvarēs gaidāmajās vēlēšanās) atkarībā no tiem dotajām priekšrocībām (mūsu piemērā: vīrieši vai sievietes). Aprēķinātās varbūtības tabulā P1.1.2 un izredžu aplēšu tabulā P1.1.3 nav adekvātas aplēses, lai noteiktu attiecīgo regresiju (sīkāku informāciju skatīt darba 11. nodaļā [Pitush K.A., Stewens I.R., 2016] un 13.6 sadaļā [Agresti A., Franklin Ch., 2013]). Lineārās regresijas modelis var būt piemērots gadījumos, kad rezultatīvā

mainīgā vērtības ir sadalītas pār normālā sadalījuma likuma. Pretējā gadījumā vērtībām jāatrodas simetriski attiecībā pret centrālo vērtību un jābūt neierobežotām (t.i., jāatrodas diapazonā $-\infty, +\infty$).

Ja mēs varam piemērotā veidā transformēt izredžu aplēses, lai iegūtu šo transformēto aprēķinu sadalījumu, kas tuvinās normālajam sadalījumam, tad var izveidot lineārās regresijas modeli.

Transformācija, kas var atrisināt šo problēmu, ir izredžu aprēķinu transformācija, izmantojot naturālos logaritmus

$$e^{\ln Ch} = Ch. \quad (P1.1.4)$$

Izmantojot šādu transformāciju, mēs varam izveidot lineāru sakarību starp faktoriālā mainīgā kategorijām un rezultatīvā mainīgā lieluma \ln kategoriju izredzēm. Piemēram, iepriekš minētajam piemēram un $y = 1$

$$\ln(\text{Ch}(y = 1)) = \beta_0 + \beta_1 x, \quad (P1.1.5)$$

kur β_0, β_1 - regresijas koeficienti.

Aprēķinot vērtību $\ln(\text{Ch}(y = j))$, izmantojot vienādojumu (P1.1.4), var aprēķināt iespēju vērtību $\text{Ch}(y = 1)$ un ar izteiksmi (P1.1.2) aprēķināt vēlamās varbūtības vērtību $p(y = j)$, $j = 1, 2$.

Iepriekš dotajā piemērā mēs darbojamies ar nominālām aptaujāto iedzīvotāju kategorijām: vīriešiem un sievietēm. Bet patiesībā skaitļi tabulā P1.1.1 ir norādīti nepārtrauktā reālo skaitļu x skalā. Mēs varam aprēķināt varbūtību, ka kandidāts A uzvarēs, dodot priekšroku 50, 100, 200 vai jebkuram citam vīriešu skaitam, vai patvaļīgam sieviešu skaitam, vai jebkuram patvaļīgam vīriešu un sieviešu skaitam. Pretējā gadījumā mums ir vajadzīgas vispārīgas attiecības starp kandidātam A piešķirto priekšrocību skaitu katrā respondentu kategorijā vai abās kategorijās.

Tas pats pamatojums attiecas uz aplēsēm par iespējamību, ka kandidāts B uzvarēs vēlēšanās.

No tā izriet, ka jāizmanto mainīgais x , kura vērtības atspoguļo mūs interesējošo iedzīvotāju skaitu gan katrai nominālajai kategorijai "vīrieši", "sievietes", gan to kopējo vērtību. Šāds mainīgais šajā kontekstā bieži tiek saukts par fiktīvu mainīgo (*dummy variable*).

Vispārīgā gadījumā loģistiskās regresijas grafikam ir viena no S-veida formām, kas parādīta attēlā R1.2.1. Attēls P1.2.1.a atbilst gadījumam $\beta_1 > 0$, attēls P1.2.1.b atbilst gadījumam $\beta_1 < 0$.

Iepriekš minētajā piemērā faktoriālais mainīgais un rezultatīvais mainīgais bija bināri nominālie mainīgie. Bet loģistiskās regresijas jēdzienu var attiecināt arī uz gadījumu, kad rezultatīvais mainīgais ir binārs nominālais mainīgais un faktoriālajam mainīgajam ir m kategorijas ($m > 2$). Šajā gadījumā loģistiskās regresijas vienādojums iegūst formu

$$\ln(\text{Ch}(y = 1)) = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m, \quad (P1.1.6)$$

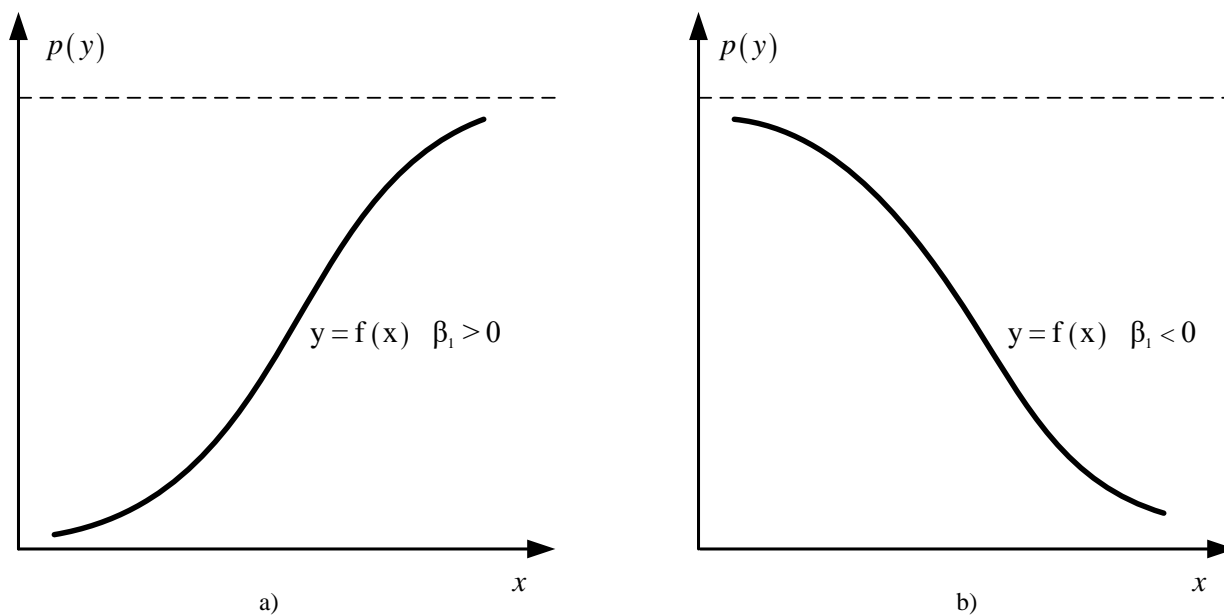
kur β_i , $i = 0, 1, \dots, m$ ir regresijas koeficienti.

Ja ir zināmas konstantes β_0 un regresijas koeficienta vērtības β_1 regresijas vienādojumā (P1.2.5), tad noteiktai fiktīvā mainīgā vērtībai x atbilstošās nominālās kategorijas varbūtības vērtību y_j var tieši aprēķināt, izmantojot šādu vienādojumu

$$p(y_j) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}. \quad (P1.2.7)$$

Ja faktoriālajam mainīgajam X ir m vērtības, tad atbilstošās nominālās kategorijas varbūtības vērtību y_j var tieši aprēķināt no vienādojuma

$$p(y_j) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m}}. \quad (P1.2.8)$$



Attēls P2.1.1. Loģistiskās regresijas diagrammas a) $\beta_1 > 0$; b) $\beta_1 < 0$

Kādā veidā var tikt aprēķināti konstantu un regresijas koeficientu vērtības vienādojumos (P1.2.7), (P1.2.8)? Šī ir sarežģīta skaitļošanas problēma, kuras pamatā ir maksimālās varbūtības princips. Praksē šī problēma tiek atrisināta, izmantojot atbilstošu programmatūru.

P1.2. Ordinālās skalas

Ordinālās skalas kopā ar nominālajām skalām veido kategorisko (kategorijas) skalu klasi. Atšķirībā no nominālajām skalām, ordinālās skalas ļauj sakārtot kategorijas datus atbilstoši noteiktam atribūtam.

Tipisks ordinālās piemērs ir Likerta skala. Šo skalu ļoti bieži izmanto dažāda veida aptauju organizēšanai, kuru mērķis ir noskaidrot respondentu viedokļus par kādu konkrētu priekšlikumu, plānu, reformu u.c. Piemēram, varētu organizēt aptauju, lai noskaidrotu iedzīvotāju attieksmi pret pašvaldības ierosinājumu mainīt pilsētas sabiedriskā transporta kustības sarakstus. Šīs skalas kategorijas var būt:

- 1 – kategoriski nepiekrītu;
- 2 – nepiekrītu;
- 3 – neitrāls;
- 4 – piekrītu;
- 5 – pilnībā piekrītu.

Pilsētas iedzīvotājs, kas piedalās aptaujā, vienkārši atzīmē kategoriju, kas atbilst viņa attieksmei pret vērtējamo piedāvājumu. Ja aptauja tiks veikta noteiktu iedzīvotāju kategoriju vidū (piemēram, dažādu pilsētas rajonu iedzīvotāji), tad mums būs daudz atbilžu kategoriju daudzām iedzīvotāju kategorijām. Visus var kompakti attēlot tabulas veidā, kurā katra rinda atbilst kādai no aptaujāto iedzīvotāju kategorijām, bet katra kolonna atbilst kādai no atbilžu kategorijām. Skaitlis rindas i un kolonas j krustpunktā apzīmē to iedzīvotāju skaitu no i -tās kategorijas, kuri kā atbildi izvēlējušies kategoriju j no visu vērtēšanas kategoriju saraksta.

Vēl viens mērīšanas piemērs ordinālajā skalā ir respondentu apmierinātības vērtējums par viņu dalību kādā pasākumā, izstādē, teātra vai kino apmeklējumos utml. Šādas skalas kategorijas var būt:

- 1 - pilnīgi neapmierināts;

- 2 - neapmierināts;
- 3 - vienaldzīgs;
- 4 - apmierināts;
- 5 - pilnībā apmierināts.

Vēl viens mērīšanas piemērs ordinālajā skalā. Atkarībā no personas vecuma viņu var iedalīt vienā no šādām kategorijām:

- 1 – bērns;
- 2 – pusaudzis;
- 3 – jauniešs;
- 4 – pusmūža cilvēks;
- 5 – vecāka gadagājuma cilvēks.

Visiem iepriekš sniegtajiem kārtas skalu piemēriem raksturīga iezīme ir tā, ka visas kategorijas var sakārtot: pēc sakritības ar pašvaldības priekšlikumu, pēc apmierinātības pakāpes un pēc vecuma. Bet par šādām skalām nav iespējams pateikt, *cik ļoti* vai *cik daudz* kategorijas atšķiras. Tādējādi nevar konstatēt, kādā mērā atšķirība starp bērnu un pusaudzi vai pusaudzi un jauniešs ir lielāka vai mazāka.

Ja mēs nosakām skaitliskas robežas starp vecuma kategorijām un sadalām noteiktas izlases indivīdus (piemēram, novada iedzīvotājus) iegūtājās kategorijās, mēs atkal iegūsim mērījumu pēc ordinālās skalas. Lieta tāda, ka personas tiek iedalītas vecuma kategorijās, pamatojoties uz viņu faktisko vecumu. Tādējādi tiek izveidoti ordināli dati, un šie dati ir jāapstrādā saskaņā ar ordinālo skalu noteikumiem.

(Nedrīkst jaukt indivīdu vecuma salīdzināšanu ar viņu vecuma kategoriju salīdzināšanu. Pieņemam, ka mums ir divi indivīdi *A* un *B*. Indivīda *A* vecums ir 13 gadi un 6 mēneši, indivīda *B* vecums ir 18 gadi un 3 mēneši. Varam droši apgalvot, ka indivīds *B* ir 4 gadus un 9 mēnešus vecāks par indivīdu *A*. Bet, ja mēs zinām tikai to, ka indivīds *A* pieder kategorijai “Pusaudži” un indivīds *B* pieder kategorijai “Jauniešs”, tad vienīgais, ko mēs varam teikt ir tas, ka indivīds *B* ir vecāks par indivīdu *A*).

Lai atvieglotu datu analīzi, to ordinālās kategorijas var atbilstoši kodēt. Iepriekš minētajos piemēros kategoriju numurus var izmantot kā šo kategoriju kodus.

Tipisks ordinālās skalas kategoriju kodēšanas piemērs ir studentu snieguma rādītāji:

- 2 – ļoti vājš;
- 3 – vājš;
- 4 – gandrīz apmierinošs;
- 5 – apmierinošs;
- 6 – gandrīz labi;
- 7 – labi;
- 8 – ļoti labi;
- 9 – teicami;
- 10 – lieliski.

Skaitlisko snieguma novērtējumus var interpretēt kā atbilstošo verbālo un kategorisko kategoriju kodus, tāpēc mēģinājums novērtēt studenta vidējo sniegumu, aprēķinot viņa punktu vidējo vērtību, ir nekorekti.

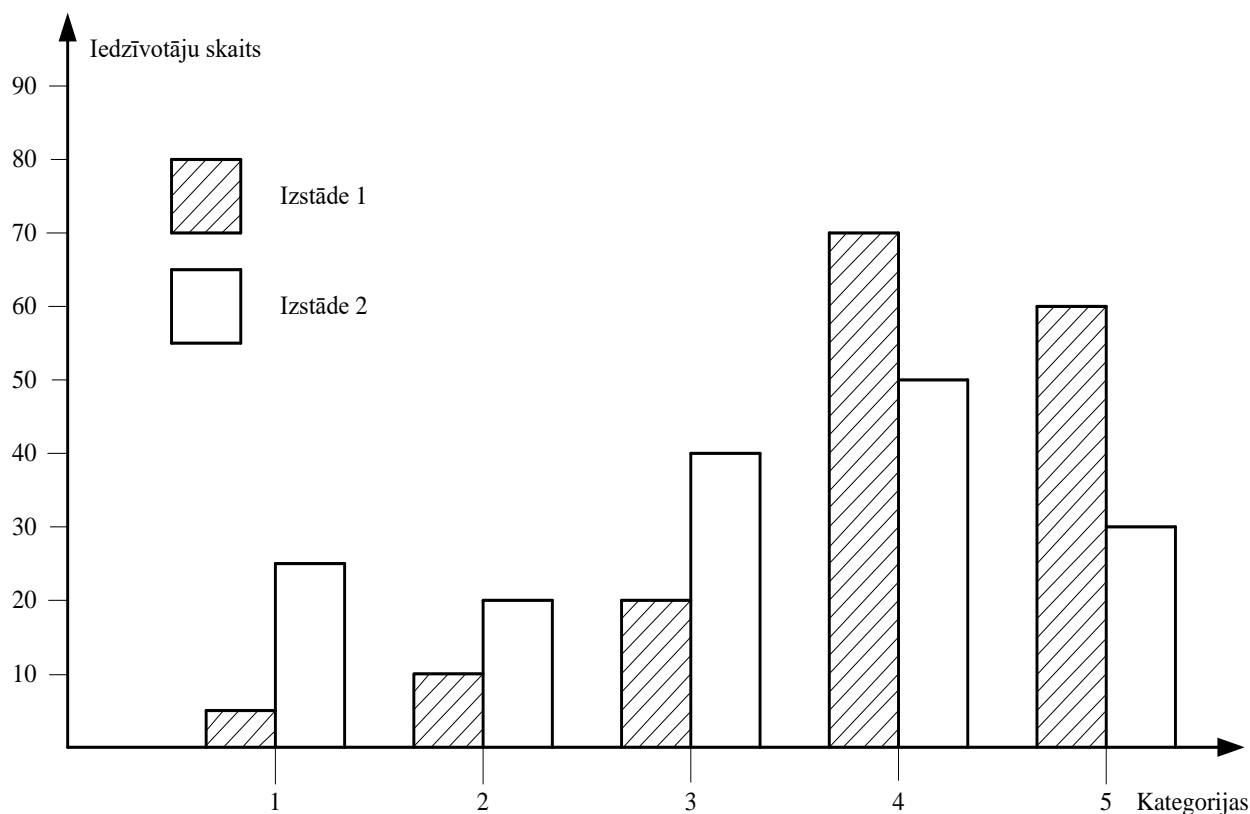
Kā var parādīt un vizualizēt datus, kas mērīti ordinālajā skalā? Mēs piedāvāsim šīs problēmas risinājumu uz piemēra pamata.

Pieņemam, ka pilsētā tiek rīkotas divas izstādes vienlaikus. Izstādes apmeklētāji tika lūgti novērtēt savu apmierinātības līmeni ar izstādēm, atzīmējot vienu no kategorijām: “1. Pilnīgi neapmierināts..., 5 - Pilnīgi apmierināts”, kas tika parādītas iepriekš šajā sadaļā. Apkopotie rezultāti ir parādīti tabulā P1.2.1.

Apmierinātības rezultāti izstādes apmeklētājiem

Izstādes	Apmierinātības kategorijas					Summa pa rindām
	1	2	3	4	5	
1	5	10	20	70	60	165
2	25	20	40	50	30	165

Šie dati ir grafiski vizualizēti attēlā P1.2.1.



Attēls P1.2.1. Datu grafiskā vizualizācija datiem no tabulas P1.2.1

Vienkārša datu analīze tabulā P1.2.1 un attēlā P1.2.1 ļauj izdarīt acīmredzamu secinājumu, ka 1. izstādes apmeklētāji ir apmierinātāki nekā 2. izstādes apmeklētāji.

Ordinālo datu analīzē var būt lietderīgi aprēķināt un vizualizēt rezultātu proporcijas (biežumu) katrā kategorijā. Iepriekš minētajā piemērā vērtējumu daļas biežumi ir parādīti tabulā P1.2.2.

Tabula P1.2.2

Respondentu vērtējumu proporcijas (biežums) pa kategorijām saskaņā ar datiem no tabulas P1.2.1

Izstādes	Apmierinātības kategorijas					Summa pa rindām
	1	2	3	4	5	
1	0.030	0.061	0.121	0.424	0.364	1000
2	0.152	0.121	0.242	0.303	0.182	1000

Lai vizuāli attēlotu novērtējumu biežumu pa kategorijām, varat aprēķināt uzkrātās biežumu vērtības. Šīs vērtības tabulas P1.2.2 datiem ir parādītas tabulā P1.2.3.

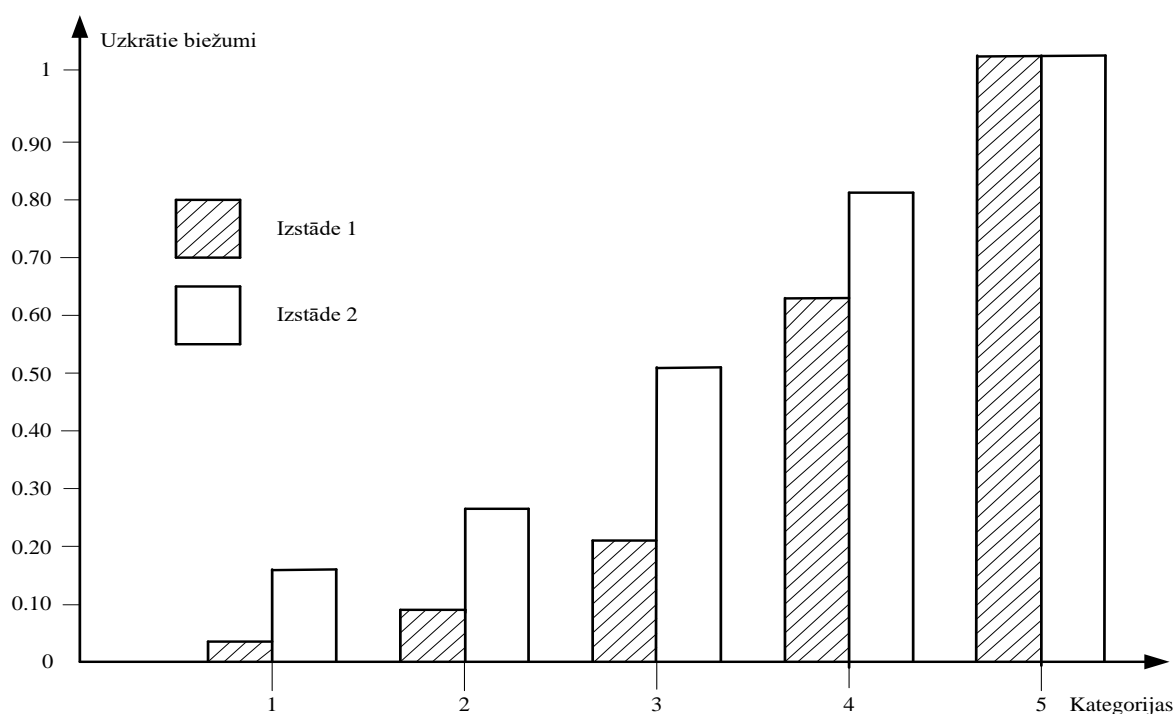
Uzkrātās biežumu vērtības datiem no tabulas P1.2.2

Izstādes	Apmierinātības kategorijas				
	1	2	3	4	5
1	0.030	0.091	0.212	0.636	1000
2	0.152	0.273	0.512	0.818	1000

Dati no tabulas P1.2.3 ir grafiski vizualizēti attēlā P1.2.2.

Uzkrāto biežumu grafiki šajā attēlā skaidri parāda, ka 1. izstādes apmeklētājiem uzkrātie biežumi lēnām pieaug mazo kategoriju numuru zonā un strauji palielinās 4. un 5. kategorijā. Tas ir mazā neapmierināto apmeklētāju skaita sekas un lielais apmierināto apmeklētāju skaits izstādē 1. Apmeklētājiem 2. izstāde mums ir cita aina. Uzkrātie biežumi strauji pieaug mazo kategoriju numuru zonā un lēnām palielinās 4. un 5. kategorijā. Tas ir lielā neapmierināto apmeklētāju skaita sekas.

Kā tiek noteikts izvietojuma parametrs ordinālās skalas kategoriju sadalījumam? Pirmā iespēja ir izmantot modu kā kategoriju ar lielāko vērtību skaitu. Apmierinātības kategoriju sadalījumam starp izstādes apmeklētājiem moda ir 4. kategorija abām izstādēm.



Attēls P1.2.2. Uzkrāto biežumu grafiskā vizualizācija datiem no tabulas P1.2.3

Otra iespēja ir izmantot sadalījuma mediānu. Mediānas vērtību nosaka pēc vienādojumiem (n ir kategoriju skaits):

$$med = \frac{n}{2}, \quad n - \text{pāra}; \quad (\text{P1.2.1.a})$$

$$med = \frac{n+2}{2}, \quad n - \text{nepāra}. \quad (\text{P 1.2.1.b})$$

Dotajā piemērā ar izstādēm sadalījuma mediāna ir 3. kategorija, jo $n=5$.

Ordinālo kategoriju sadalījumam principā var noteikt diapazona parametrus: rangu un procentuālo diapazonu. Bet šie parametri pārāk informatīvi, tāpēc praksē tos izmanto reti.

Ordinālo kategoriju sadalījumu analīze, pamatojoties uz šo sadalījumu parametriem, šķiet problemātiska, jo šie parametri sniedz maz būtiskas informācijas. Tāpēc šādu sadalījumu analīzei tiek plaši izmantotas tiešās metodes, tostarp metodes, kuru pamatā ir sadalījumu vizualizācija (sk. attēlus P1.2.1, P1.2.2).

Jāatzīmē, ka praksē ordinālo skalu dati bieži tiek analizēti, izmantojot metodes, kas paredzētas kvantitatīviem datiem, piemēram, aprēķinot šādu datu vidējo vērtību un variācijas. Tomēr tas noved pie nelabvēlīgiem rezultātiem.

Šādu metožu trūkumi ir detalizēti analizēti darbā [Liddell T., Krushke J. K., 2018]. Pirmkārt, analizējot ordinālos datus statistiskajiem modeļiem, kas paredzēti metriskajiem mainīgajiem, piemēram, T-testiem un dispersijas analīzei, var rasties zems pareizas atpazīšanas līmenis vai ietekmes lieluma aplēšu kļūdas. Otrkārt, ordinālo mainīgo sadalījums var būt neparasts, jo īpaši, ja dominē ļoti mazas vai ļoti lielas vērtības. Treškārt, variācijas mainīgajos var atšķirties dažādās grupās, apstākļos, laika diapazonos utt.

Tas noved pie nepārprotama secinājuma: ordinālie mainīgie jāanalizē tikai ar metodēm, kas paredzētas tieši šādiem mainīgajiem.

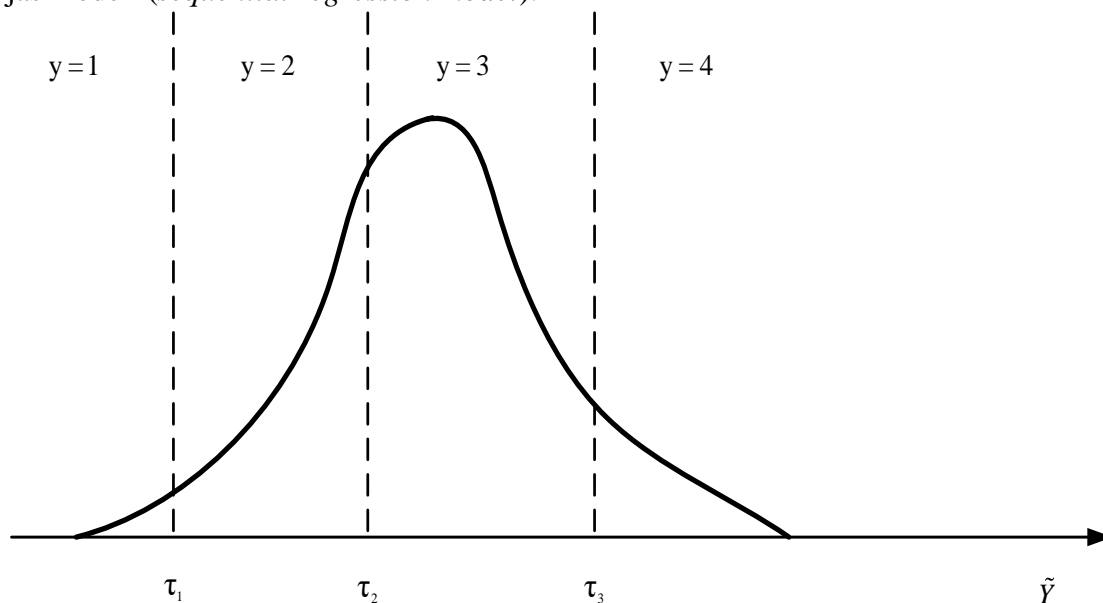
Kā var definēt regresiju ordinālajiem mainīgajiem? P1.1 sadaļā mēs definējam loģistisko regresiju nominālajiem mainīgajiem. Tā kā ordinālā skala kopējā skalu hierarhijā ieņem augstāku vietu nekā nominālās skala, tad ordinālajiem mainīgajiem ir iespējams definēt plašāku regresijas modeļu klasi.

Kopumā var definēt 6 regresijas modeļus ordinālajiem mainīgajiem [Ananth G.V., Kleinbaum D.G., 1997]. Šajā sadaļā prezentēsim trīs vispārinātas šādu modeļu klases [Bürkner P.-Ch., Vuorre M., 2019].

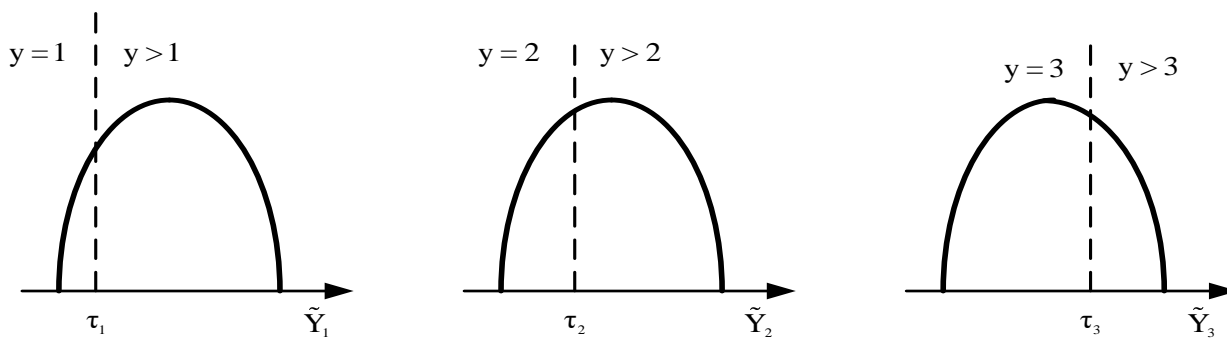
Dota faktoriālā pazīme (mainīgais) - \mathbf{X} un rezultējošā pazīme (mainīgais) - \mathbf{Y} . Ieviešam slēpto (latento) mainīgo \tilde{Y} , kura vērtības tiek iegūtas, atbilstoši transformējot rezultējošo mainīgo \mathbf{Y} , un šīs transformētās vērtības tiek sadalītas saskaņā ar normālā sadalījuma likumu. Tad, ja mainīgā lieluma \tilde{Y} sadalījuma blīvuma funkcijai ir tāda forma, kā parādīts attēlā P1.2.3, mums ir darīšana ar kumulatīvo regresijas modeli (*cumulative logistic regression models*).

Šajā attēlā y_j , $j=1,2,3,4$ ir rezultējošā mainīgā \mathbf{Y} ordinālā mainīgā kategorijas, τ_j , $j=1,2,3$ ir robežas starp kategorijām. Šajā gadījumā kategorijas y_j tiek uzdotas tā, lai tām atbilstošā latentā mainīgā \tilde{Y} vērtības tiktu sadalītas saskaņā ar normālo sadalījuma likumu.

Apsvērsim vēl vienu modeli. Pieņemam, ka mums ir tāds rezultējošais mainīgais \mathbf{Y} , lai katrai ordinālajai kategorijai varētu izveidot latentu mainīgo \tilde{Y} , kura vērtības tiek sadalītas saskaņā ar normālā sadalījuma likumu. Šī situācija ir grafiski parādīta attēlā P1.2.4. Šo modeli sauc par secīgas regresijas modeli (*sequential regression model*).

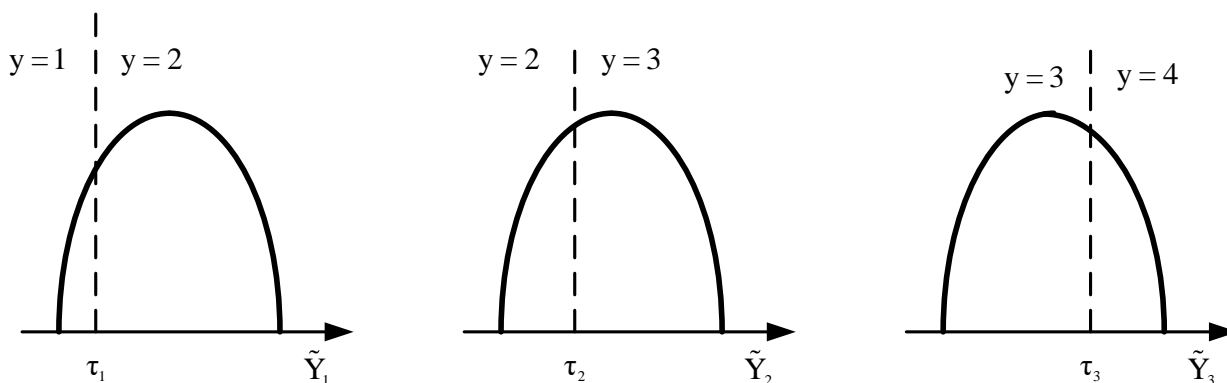


Attēls P1.2.3. Latentā mainīgā \tilde{Y} vērtību blīvuma sadalījuma grafiks kumulatīvajā regresijas modelī



Attēls P1.2.4. Latentā mainīgā \tilde{Y} vērtību blīvuma sadalījuma grafiks secīgas regresijas modelī

Trešo regresijas modeli sauc par konjugētās kategorijas modeli (*adjacent-category model*). Tā atšķirība no iepriekšējiem modeļiem ir tāda, ka latentie mainīgie \tilde{Y}_j tiek veidoti dažādām rezultējošā mainīgā Y ordinālo mainīgo kategorijām. Latento mainīgo sadalījuma funkciju grafiki \tilde{Y}_j šādam regresijas modelim ir parādīti attēlā P1.2.5.



Attēls P1.2.5. Latentā mainīgā \tilde{Y} vērtību blīvuma sadalījuma grafiks konjugētās kategorijas modelī

Interesentiem, kuri vēlas iegūt sīkāku informāciju par apspriestajiem jautājumiem, ieteicams izmantot darbu [Bükner P.-Ch., Vuorre M., 2019].

Savukārt kumulatīvās regresijas modeli var konstruēt dažādos veidos (sīkāk darbā [Ananth C.V., Kleinbaum D.G., 1997]). Šajā sadaļā mēs iepazīstināsim ar loģistiskā modeļa paplašinājumu, proti, kumulatīvo logit modeli (*Cumulative Logit Model*), jeb precīzāk, vienu no tā variantiem, ko sauc par proporcionālo izredžu modeli. Tālāk sniegtais materiāls ir balstīts uz datiem, kas sniegti darbā [Agresti A., 2002; Agresti A., 2007; Agresti A., 2010]. Iepazīstinām ar šī regresijas modeļa pamatprocedūrām, izmantojot konkrētu piemēru.

Pieņemsim, ka pirms pašvaldību vēlēšanām tika aptaujātas trīs pilsētas iedzīvotāju grupas, lai noskaidrotu, kā viņi vērtē domes darbu iepriekšējā periodā.

Kodēsim pilsētas iedzīvotājus atspoguļojošā faktoriālā mainīgā X ordinālās skalas kategorijas šādi:

- x_1 - jaunieši;
- x_2 - pusmūža cilvēki;
- x_3 - gados veci cilvēki.

Mainīgais Y ir saistīts ar domes darba kvalitātes novērtējumu ordinālām kategorijām. Kodēsim šīs kategorijas šādi:

- y_1 - labi;
- y_2 - apmierinoši;
- y_3 - slikti.

Apkopotie iedzīvotāju aptaujas rezultāti atspoguļoti tabulā P1.2.4.

Tabula P1.2.4

Pilsētas iedzīvotāju aptaujas rezultāti par iepriekšējās domes darba kvalitāti

Iedzīvotāju kategorijas	Domes darba kvalitātes novērtējuma kategorijas			Summa pa rindām
	y_1	y_2	y_3	
x_1	50	40	10	100
x_2	40	40	20	100
x_3	30	40	30	100

Kā interpretēt skaitļus šajā tabulā? Iedzīvotāju kategoriju x_i skaita apzīmēšanai izmantosim indeksu i , $i = 1, 2, 3$ un ar indeksu j , $j = 1, 2, 3$ apzīmēsīm novērtēšanas kategoriju y_j numurus. Tad cipars $n_{11} = 50$ pirmās rindas un pirmās kolonas krustpunktā parāda to iedzīvotāju skaitu no kategorijas x_1 “Jaunieši”, kuri pilsētas domes darbu novērtējuši ar kategoriju y_1 “Labi”. Skaitlis $n_{33} = 30$ trešās rindas un trešās kolonas krustpunktā parāda to iedzīvotāju skaitu no kategorijas x_3 “Gados veci cilvēki”, kuri pilsētas domes darbu novērtējuši kā y_3 “Vāji”.

Kopumā skaitlis n_{ij} atspoguļo iedzīvotāju skaitu no kategorijas x_i , $i = 1, 2, 3$, kuri pilsētas domes darba kvalitāti novērtēja ar kategoriju y_j , $j = 1, 2, 3$.

(Turpmākajā izklāstā izmantosim pilno kategoriju kodu vietā to indeksus i un j . Tas nesagādās grūtības, apgūstot materiālu).

Tabulas P1.2.4 pēdējā kolonā ir norādītas skaitļu summas tabulas rindās - n_i , $i = 1, 2, 3$. Katrs skaitlis apzīmē kopējo attiecīgās kategorijas iedzīvotāju skaitu, kas piedalījās aptaujā. Mūsu gadījumā visi šie skaitļi ir vienādi ar 100.

Aprēķināsim katras vērtējuma kategorijas nosacītās varbūtības katrai respondentu grupai

$$p_{j/i} = \frac{n_{ij}}{n_i}, \quad i, j = 1, 2, 3. \quad (\text{P1.2.2})$$

Aprēķinu rezultāti tabulu P1.2.4 datiem ir parādīti tabulā P2.1.5.

Tabula P1.2.5

Domes darba kvalitātes novērtējumu kategoriju nosacītās varbūtības pēc respondentu kategorijas

Iedzīvotāju kategorijas	Domes darba kvalitātes novērtējuma kategorijas			Summa pa rindām
	y_1	y_2	y_3	
x_1	0.500	0.400	0.100	1.000
x_2	0.400	0.400	0.200	1.000
x_3	0.300	0.400	0.300	1.000

Aprēķināsim nosacīto varbūtību $F_{j/i}$ kumulatīvās vērtības katrai iedzīvotāju kategorijai. Lai to izdarītu, mēs izmantosim datus no tabulas P1.2.5.

$$F_{1/1} = p_{1/1} = 0.500; F_{2/1} = p_{1/1} + p_{2/1} = 0.500 + 0.400 = 0.900;$$

$$F_{3/1} = p_{1/1} + p_{2/1} + p_{3/1} = 0.500 + 0.400 + 0.100 = 1.000.$$

$$F_{1/2} = p_{1/2} = 0.400; F_{2/2} = p_{1/2} + p_{2/2} = 0.400 + 0.400 = 0.800;$$

$$F_{3/2} = p_{1/2} + p_{2/2} + p_{3/2} = 0.400 + 0.400 + 0.200 = 1.000.$$

$$F_{1/3} = p_{1/3} = 0.300; F_{2/3} = p_{1/3} + p_{2/3} = 0.300 + 0.400 = 0.700;$$

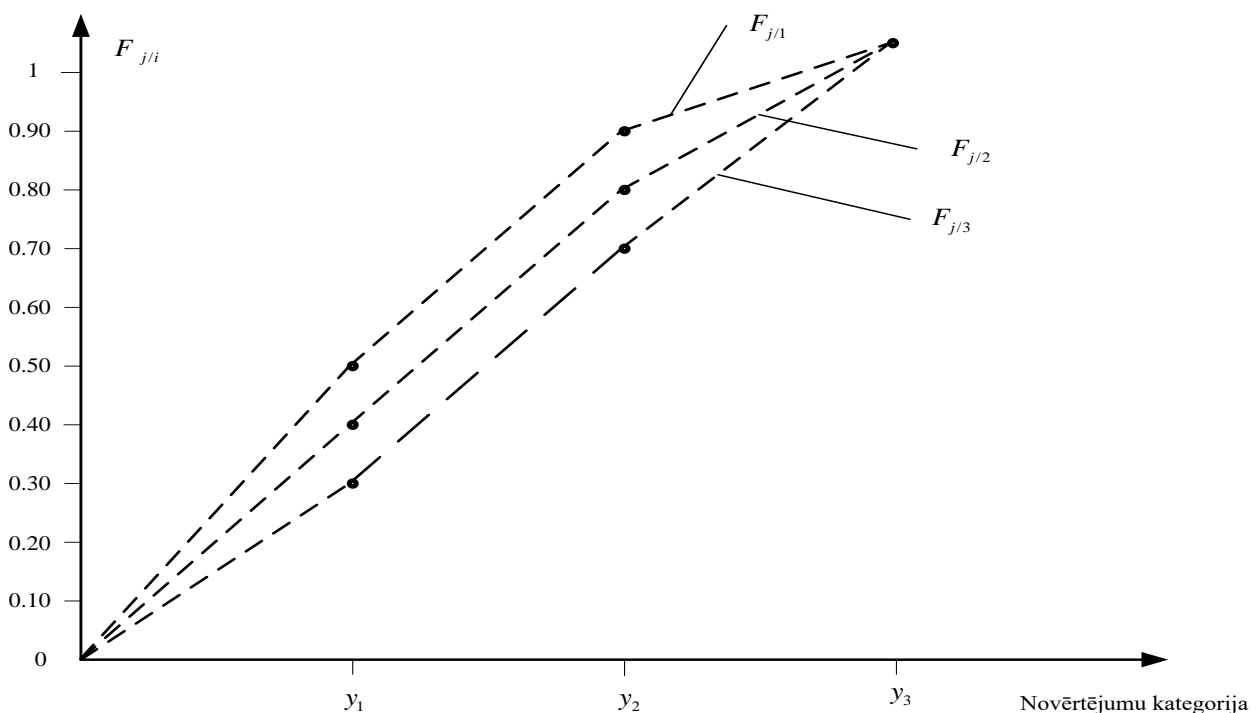
$$F_{3/3} = p_{1/3} + p_{2/3} + p_{3/3} = 0.300 + 0.400 + 0.300 = 1.000.$$

Nosacīto varbūtību kumulatīvās vērtības ir apkopotas tabulā P1.2.6.

Tabula P1.2.6

Iedzīvotāju kategorijas	Domes darba kvalitātes novērtējuma kategorijas		
	y_1	y_2	y_3
x_1	0.500	0.900	1000
x_2	0.400	0.800	1000
x_3	0.300	0.700	1000

Skaidrības labad iegūtie rezultāti ir grafiski parādīti attēlā P1.2.6.



Attēls P1.2.6. Nosacīto varbūtību kumulatīvo vērtību $F_{j/i}$ grafiki datiem no tabulas P.1.2.6

Mēs varam aprēķināt beznosacījumu varbūtību vērtības p_{ij} , izmantojot vienādojumu

$$p_{ij} = \frac{n_{ij}}{n}, \quad i, j = 1, 2, 3, \quad (\text{P1.2.3})$$

kur $n = \sum_i n_i = 300$ ir kopējais iedzīvotāju skaits, kas uzrādīts tabulā P1.2.4. Aprēķinu rezultāti ir parādīti tabulā P1.2.7.

Tabula P1.2.7

Beznosacījumu varbūtības p_{ij} vērtības datiem no tabulas P1.2.4

Iedzīvotāju kategorijas	Domes darba kvalitātes novērtējuma kategorijas			Summa pa rindām
	y_1	y_2	y_3	
x_1	0.167	0.133	0.033	0.333
x_2	0.133	0.133	0.067	0.333
x_3	0.100	0.133	0.133	0.333

Ieviesīsim izredžu attiecības jēdzienu. Atšķirībā no loģistiskās regresijas, kas aplūkota sadaļā P1.1, logit-regresijas modelī, izmantojot tabulā P1.2.4 esošos datus, varam uzreiz aprēķināt izredžu attiecības [Agresti A., 2010]. Datiem tabulai ar izmēriem $r \times c$ var aprēķināt $(r-1)(c-1)$ izredžu attiecības. Tabulai P1.2.4 var aprēķināt $2 \times 2 = 4$ izredžu attiecības.

1. Lokālās izredžu attiecības.

$$\xi_{ij}^L = \frac{n_{ij}n_{i+1,j+1}}{n_{i,j+1}n_{i+1,j}}. \quad (\text{P1.2.3})$$

Aprēķināsim lokālās izredžu attiecības tabulas P1.2.4 datiem, izmantojot vienādojumu (P1.2.3).

$$\xi_{11}^L = \frac{n_{11}n_{22}}{n_{12}n_{21}} = \frac{50 \cdot 40}{40 \cdot 40} = 1.250;$$

$$\xi_{12}^L = \frac{n_{12}n_{23}}{n_{13}n_{22}} = \frac{40 \cdot 20}{10 \cdot 40} = 2.000;$$

$$\xi_{21}^L = \frac{n_{21}n_{32}}{n_{22}n_{31}} = \frac{40 \cdot 40}{40 \cdot 30} = 1.333;$$

$$\xi_{22}^L = \frac{n_{22}n_{33}}{n_{23}n_{32}} = \frac{40 \cdot 30}{20 \cdot 40} = 1.500.$$

2. Globālās izredžu attiecības.

$$\xi_{ij}^G = \frac{\left(\sum_{a<i} \sum_{b \leq j} n_{ab}\right) \left(\sum_{a>i} \sum_{b>j} n_{ab}\right)}{\left(\sum_{a \leq i} \sum_{b>j} n_{ab}\right) \left(\sum_{a>i} \sum_{b \leq j} n_{ab}\right)}. \quad (\text{P1.2.4})$$

Aprēķināsim globālās izredžu attiecības tabulas P1.2.4 datiem, izmantojot vienādojumu (P1.2.4).

$$\xi_{11}^G = \frac{n_{11}(n_{22} + n_{23} + n_{32} + n_{33})}{(n_{12} + n_{13})(n_{21} + n_{31})} = \frac{50 * (40 + 20 + 40 + 30)}{(40 + 10) * (40 + 30)} = \frac{50 * 130}{50 * 70} = \frac{6500}{3500} = 1.857;$$

$$\xi_{12}^G = \frac{(n_{11} + n_{12})(n_{23} + n_{33})}{n_{13}(n_{21} + n_{22} + n_{31} + n_{32})} = \frac{(50 + 40) * (20 + 30)}{10 * (40 + 40 + 30 + 10)} = \frac{90 * 50}{10 * 120} = \frac{3750}{1200} = 3.750;$$

$$\xi_{21}^G = \frac{(n_{11} + n_{21})(n_{32} + n_{33})}{(n_{12} + n_{13} + n_{22} + n_{23})(n_{31} + n_{32})} = \frac{(50 + 40) * (40 + 30)}{(40 + 10 + 40 + 20) * (30 + 40)} = \frac{90 * 70}{110 * 70} = \frac{6300}{7700} = 0.818;$$

$$\xi_{22}^G = \frac{(n_{11} + n_{12} + n_{21} + n_{22})n_{33}}{(n_{13} + n_{23})(n_{31} + n_{32})} = \frac{(50 + 40 + 40 + 40) * 30}{(10 + 20) * (30 + 40)} = \frac{170 * 30}{30 * 70} = \frac{5100}{2100} = 2.429.$$

3. Kumulatīvās izredžu attiecības.

$$\xi_{ij}^C = \frac{\left(\sum_{b \leq j} n_{ib}\right) \left(\sum_{b > j} n_{i+1,b}\right)}{\left(\sum_{b > j} n_{ib}\right) \left(\sum_{b \leq j} n_{i+1,b}\right)}. \quad (\text{P1.2.5})$$

Aprēķināsim kumulatīvās izredžu attiecības tabulas P1.2.4 datiem, izmantojot vienādojumu (P1.2.5).

$$\xi_{11}^C = \frac{n_{11}(n_{22} + n_{23})}{(n_{12} + n_{13})n_{21}} = \frac{50 * (40 + 20)}{(40 + 10) * 40} = \frac{50 * 60}{50 * 40} = \frac{3000}{2000} = 1.500;$$

$$\xi_{12}^C = \frac{(n_{11} + n_{12})n_{23}}{n_{13}(n_{21} + n_{22})} = \frac{(50 + 40) * 20}{10 * (40 + 40)} = \frac{90 * 20}{10 * 80} = \frac{1800}{800} = 2.250;$$

$$\xi_{21}^C = \frac{n_{21}(n_{31} + n_{32})}{(n_{22} + n_{23})n_{31}} = \frac{40 * (30 + 40)}{(40 + 20) * 30} = \frac{40 * 70}{60 * 30} = \frac{2400}{1400} = 1.715.$$

Kumulatīvās izredžu attiecības var arī aprēķināt, izmantojot kumulatīvās nosacīto varbūtību vērtības $F_{j/i}$.

$$\xi_{ij}^C = \frac{F_{j/i} / (1 - F_{j/i})}{F_{j/i+1} / (1 - F_{j/i+1})}. \quad (\text{P1.2.6})$$

Aprēķināsim kumulatīvās izredžu attiecības tabulas P1.2.6 datiem, izmantojot vienādojumu (P1.2.6).

$$\xi_{11}^C = \frac{F_{1/1} / (1 - F_{1/1})}{F_{1/2} / (1 - F_{1/2})} = \frac{0.500 / (1 - 0.500)}{0.400 / (1 - 0.400)} = \frac{1.000}{0.667} = 1.500;$$

$$\xi_{12}^C = \frac{F_{2/1} / (1 - F_{2/1})}{F_{2/2} / (1 - F_{2/2})} = \frac{0.900 / (1 - 0.900)}{0.800 / (1 - 0.800)} = \frac{9.000}{4.000} = 2.250;$$

$$\xi_{21}^C = \frac{F_{1/2} / (1 - F_{1/2})}{F_{1/3} / (1 - f_{1/3})} = \frac{0.400 / (1 - 0.400)}{0.300 / (1 - 0.300)} = \frac{0.667}{0.429} = 1.556;$$

$$\xi_{22}^C = \frac{F_{2/2} / (1 - F_{2/2})}{F_{2/3} / (1 - f_{2/3})} = \frac{0.800 / (1 - 0.800)}{0.700 / (1 - 0.700)} = \frac{4.000}{2.233} = 1.715.$$

Kumulatīvo izredžu attiecību vērtības, kas aprēķinātas, izmantojot vienādojumus (P1.2.5), (P1.2.6), sakrīt.

Ieviesīsim kumulatīvo regresijas logit-modeli ar proporcionālajām izredzēm.

Pieņemam, ka mums ir rezultējošais mainīgais \mathbf{Y} , kas satur c ordinālas kategorijas y_1, y_2, \dots, y_c . Turpmāk apzīmējumā y_j simbolu y izlaidīsim un aprēķinu izteiksmēs vienkārši izmantosim indeksus j . Tātad, tā vietā, lai apzīmētu nosacīto varbūtību $p(Y \leq y_j)$ mēs vienkārši izmantosim vienkāršoto apzīmējumu $p(Y \leq j)$. Tas vienkāršos aprēķinu izteiksmes un neradīs nekādus pārpratumus.

Sekojo [Agresti A., 2010], katras ordinālās kategorijas nosacītās varbūtības $p_{j/i}$ apzīmējam ar π_j . Kategorijas j kumulatīvā nosacītā varbūtība tiek aprēķināta kā visu iepriekšējo kategoriju nosacīto varbūtību summa:

$$p(Y \leq j) = \pi_1 + \dots + \pi_j, \quad j = 1, \dots, c.$$

Šīs varbūtību vērtības veido sakārtotu secību:

$$\pi_1 \leq \pi_2 \leq \dots \leq \pi_c.$$

Ieviesīsim logitu (*logits*) jēdzienu uzkrātajām nosacītajām varbūtībām.

$$\log it(p(Y \leq j)) = \ln \left(\frac{p(Y \leq j)}{1 - p(Y \leq j)} \right) = \ln \frac{\pi_1 + \pi_2 + \dots + \pi_j}{\pi_{j+1} + \pi_{j+2} + \dots + \pi_c}, \quad j = 1, \dots, c-1. \quad (\text{P1.2.7})$$

Vērtības $\log it(p(Y \leq j))$ sauc par *kumulatīviem logitiem*. Piemēram, priekš $c = 3$

$$\log it(p(Y \leq 1)) = \ln \frac{\pi_1}{\pi_2 + \pi_3}; \quad \log it(p(Y \leq 2)) = \ln \frac{\pi_1 + \pi_2}{\pi_3}.$$

Apskatāmā kumulatīvā logita modeļa vispārējā ideja ar proporcionālajām izredzēm ir šāda. Katrai mainīgā j kategorijai var izmantot šādu regresijas vienādojumu

$$\log it(p(Y \leq j)) = \alpha_j + \beta x, \quad j = 1, \dots, c-1. \quad (\text{P1.2.8})$$

Vienādojumā (P1.2.8) regresijas koeficientu β vērtības ir vienādas visām kategorijām y_j . Tā ir apskatāmā regresijas logit-modeļa raksturīga iezīme. Tomēr konstantu α_j nozīme visām kategorijām y_j ir atšķirīga. Tas nozīmē, ka veidojot regresijas līknes visām kategorijām y_j , šīm līknēm būs vienāda forma, bet tās tiks nobīdītas viena pret otru proporcionāli α_j vērtībām.

Mēs varam iegūt šādu izredžu attiecību j mainīgā kategorijai diviem mainīgajiem x_1 un x_2 :

$$\frac{p(Y \leq j / x_2) / p(Y > j / x_2)}{p(Y \leq j / x_1) / p(Y > j / x_1)}.$$

Šīs attiecības naturālais logaritms ir starpība starp kumulatīvo logitu vērtībām x_1 un x_2 . Šī starpība ir vienāda ar $\beta(x_2 - x_1)$ un proporcionāla attālumam starp vērtībām x_1 un x_2 .

Šim izredžu attiecību logaritmam tas pats proporcionalitātes koeficients β attiecas uz visām uzkrātajām varbūtībām. No šejienes arī ir cēlies nosaukums *uzkrāto izredžu modelis*.

Pieņemsim, ka izmantojot iepriekš uzdoto modeli, mēs esam aprēķinājuši visu uzkrāto varbūtību $p(Y \leq j)$ vērtības. Tad visu kategoriju y_j *neuzkrātās varbūtības* vērtību aprēķina kā

$$p(y_j) = p(Y \leq j) - p(Y \leq j-1). \quad (\text{P1.2.9})$$

Izmantojot vienādojumu (P1.2.9), var aprēķināt visu mūs interesējošo kategoriju varbūtības.

Izmantojot konstantes α_j un koeficienta β vērtību vienādojumā (P1.2.8), uzkrāto varbūtību vērtības var aprēķināt ar vienādojumu

$$p(Y \leq j) = \frac{e^{\alpha_j + \beta x}}{1 + e^{\alpha_j + \beta x}}. \quad (\text{P1.2.10})$$

Faktoriālais mainīgais kumulatīvajā logit-modelī var būt nepārtraukts vai kategorisks. (Iepriekš sniegtajā piemērā par domes darba kvalitātes novērtēšanu mainīgais ir ordināls mainīgais, kura kategorijas ir iedzīvotāju vecuma gradācijas).

Nepārtrauktu faktoriālā mainīgā vērtību gadījumā formulā (P1.2.10) tiek iekļauta x vērtība un veikti attiecīgie aprēķini. Ja faktoriālais mainīgais ir kategoriskais mainīgais, tad x vērtības vienkārši tiek aizstātas ar kategoriju numuriem un tiek veikti attiecīgie aprēķini.

Ja viena faktoriālā mainīgā \mathbf{X} vietā tiek izmantota mainīgo \mathbf{X}_r kopa $r, r = 1, \dots, c$, tad vienādojums (P1.2.10) iegūst šādu formu

$$p(Y \leq j) = \frac{e^{\alpha_j + \beta_1 x_1 + \dots + \beta_r x_r}}{1 + e^{\alpha_j + \beta_1 x_1 + \dots + \beta_r x_r}}. \quad (\text{P1.2.11})$$

No skaitļošanas viedokļa visgrūtākais ir konstantes α_j un regresijas koeficienta β (vai regresijas koeficientu β_r izteiksmē (P1.2.11) aprēķināšanas uzdevums. To var atrisināt tikai ar piemērotu programmatūras rīku palīdzību.

P1.3. Intervālu skalas

Intervālu skalas ir kvantitatīvās vai skaitliskās skalas. Šo skalu raksturīga iezīme ir spēja noteikt, par cik viena vērtība ir lielāka vai mazāka par citu vērtību. Piemēram, veicot divus temperatūras mērījumus, 10^0 C un 15^0 C , varam droši apgalvot, ka otrā vērtība ir par 5^0 C lielāka nekā pirmā vai pirmā vērtība ir par 5^0 C mazāka nekā otrā vērtība. Vēl viena intervālu skalu iezīme ir tāda, ka vienādiem intervāliem starp vērtībām ir tāda pati nozīme. Piemēram, temperatūras diapazonam $10^0 \text{ C} - 20^0 \text{ C}$ ir tāda pati nozīme kā intervālam ir $35^0 \text{ C} - 45^0 \text{ C}$ jeb intervālam $96^0 \text{ C} - 106^0 \text{ C}$.

Vēl viens mērījumu piemērs intervālu skalā ir indivīdu vecuma mērīšana. Lai gan šos mērījumus var interpretēt kā attiecību skalas mērījumus, indivīdu vecuma salīdzināšana ir tipisks intervālu skalas mērījumu piemērs.

Dažreiz Likerta skala tiek interpretēta kā intervālu skala. Bet tā ir nekorekta interpretācija, jo mēs nevaram skaitliski izteikt atšķirības starp skalas kategorijām. Tāpēc Likerta skala ir tipiska ordināla skala.

Tas pats attiecas uz IQ mērīšanas skalu. Tā ir tipiska ordināla skala, nevis intervālu skala.

Intervālu skalām nav noteikta nulles punkta. Piemēram, Celsija skalā ūdens sasalšanas punkts tiek uzskatīts par nulli, bet Fārenheita skalā ūdens sasalšanas punkts atbilst 32^0 F .

Vērtības intervālu skalā var saskaitīt un atņemt, bet ne reizināt un dalīt.
 Par parametriem intervālu skalā var izmantot vidējo vērtību, mediānu un modu.
 Apskatīsim šādu piemēru. Tika izmērītas šādas temperatūru vērtības:

$$X = \{12, 16, 14, 17, 14, 18, 14, 10, 19, 16\}.$$

Vidējo vērtību aprēķina šādi:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}. \quad (\text{P1.3.1})$$

Mūsu datu kopai

$$\bar{x} = \frac{12+16+14+17+14+18+14+10+19+16}{10} = \frac{150}{10} = 15.$$

Lai aprēķinātu mediānas vērtību, sākotnējās vērtības sakārtojam nedilstošā secībā:

$$X = \{10, 12, 14, 14, 14, 16, 16, 17, 18, 19\}.$$

Tā kā vērtību skaits ir pāra ($n=10$), tad par mediānas vērtību pieņem skaitli, kas atrodas pozīcijā $n/2 = 10/2 = 5$. Šī vērtība ir 14, tāpēc $med(X) = 14$.

Par modu pieņem visbiežāk izmantoto vērtību šajā kopā. Mūsu piemērā visizplatītākā vērtība ir 14, tāpēc $mod(X) = 14$.

Visbiežāk izmantotie diapazona parametri datiem, ko mēra intervālu skalā, ir dispersija un standarta novirze.

Vērtību kopas dispersija tiek aprēķināta, izmantojot vienādojumu

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}, \quad (\text{P1.3.2})$$

kur \bar{x} ir vidējā vērtība, kas aprēķināta, izmantojot vienādojumu (P1.3.1).

Aprēķināsim temperatūras vērtību kopas dispersiju. Sākotnējie dati un aprēķinu rezultāti ir parādīti tabulā P1.3.1.

Tabula P1.3.1

Temperatūras vērtību kopas dispersijas aprēķināšanai

x_i	12	16	14	17	14	18	14	10	19	16	
$x_i - \bar{x}$	-3	1	-1	2	-1	3	-1	-5	4	1	
$(x_i - \bar{x})^2$	9	1	1	4	1	9	1	25	16	1	$\sum = 67$

No vienādojuma (P1.3.1) iegūstam:

$$S^2 = \frac{67}{9} = 7.444.$$

Standartnovirzes vērtību aprēķina, izmantojot vienādojumu

$$s = \sqrt{S^2}. \quad (\text{P1.3.3})$$

Mūsu piemērā:

$$s = \sqrt{7.444} = 2.728.$$

Apskatīsim, kā definēt lineāru regresiju starp divām vērtību kopām, no kurām katra tiek mērīta intervālu skalā. Lineārās regresijas modeli var izteikt šādā lineārā formā

$$E(Y/x) = \beta_0 + \beta_1 x. \quad (\text{P1.3.4})$$

Vienādojumā (P1.3.4) β_0 ir konstante, jeb brīvais loceklis, β_1 ir slīpuma koeficients. Skaitliski β_1 nosaka, par cik vienībām izmainīsies Y , ja X tiks palielināts par vienu vienību.

Līknei, kas atbilst vienādojumam (P1.3.4), jābūt tādai, lai sākotnējo vērtību novirzes no šīs līknes būtu minimālas. To var panākt, pielāgojot koeficientus β_0 un β_1 .

Pastāv divas konceptuālas pieejas lineārās regresijas risināšanai: maksimālās līdžības metode un mazāko kvadrātu metode. Sekojošās aprēķinu izteiksmes ir balstītas uz mazāko kvadrātu metodes principa.

Uzdodam vienādojumu (P1.3.4) uzskatāmākā formā. Regresijas koeficienti b_0 un b_1 tiek aprēķināti, pamatojoties uz sākotnējiem datiem.

$$y = b_0 + b_1 x. \quad (\text{P1.3.5})$$

Koeficienta b_1 vērtības tiek aprēķinātas, izmantojot vienādojumu

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}. \quad (\text{P1.3.6})$$

Šajā vienādojumā \bar{x} , \bar{y} ir vērtību kopu X un Y vidējās vērtības.

Koeficienta b_0 vērtību aprēķina, izmantojot vienādojumu

$$b_0 = \bar{y} - b_1 \bar{x}. \quad (\text{P1.3.7})$$

Apskatīsim vienkāršu ilustratīvu piemēru. Karstā ūdens temperatūra mājās ar centrālo apkuri ir atkarīga no ārējā gaisa temperatūras. Temperatūras mērījumu dati ir parādīti tabulā P1.3.2.

Tabula P1.3.2

Āra temperatūras un ūdens temperatūras mērījumu dati

X	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
	-10	-15	-7	-20	-18	-12	-10	-9	-16	-10
Y	y_1	y_2	y_3	y_4	y_5	y_6	y_7	y_8	y_9	y_{10}
	72	74	70	80	78	74	71	68	77	74

Ir nepieciešams izveidot regresijas modeli, kas parāda sakarību starp vērtību kopām X un Y .

Aprēķināsim vidējās temperatūras.

$$\bar{x} = \frac{(-10) + (-15) + (-7) + (-20) + (-18) + (-12) + (-10) + (-9) + (-16) + (-10)}{10} = \frac{-127}{10} \approx -13.$$

$$\bar{y} = \frac{72+74+70+80+78+74+71+68+77+74}{10} = \frac{738}{10} \approx 74.$$

(Mēs esam noapaļojuši aprēķinu rezultātus, lai vienkāršotu turpmākos aprēķinus).

Veiksim starpaprēķinus, lai noteiktu regresijas koeficientus b_0 un b_1 . Aprēķinu rezultāti ir apkopoti tabulā P1.3.3.

Tabula P1.3.3

Starpaprēķinu rezultāti, lai noteiktu nepieciešamo regresijas atkarība saskaņā ar tabulu P1.3.2

Mērījumi	1	2	3	4	5	6	7	8	9	10	Σ
$x_i - \bar{x}$	3	-2	6	-7	-5	1	3	5	-3	3	
$y_i - \bar{y}$	-1	0	-4	6	4	0	-3	-6	3	0	
$(x_i - \bar{x}) * (y_i - \bar{y})$	-3	0	-24	-42	-20	0	-9	-30	-9	0	-137
$(x_i - \bar{x})^2$	9	4	36	49	25	1	9	25	9	9	176

Izmantojot starpaprēķinu rezultātus un vienādojumu (P1.3.6), aprēķinām koeficienta b_1 vērtību.

$$b_1 = \frac{-137}{176} = -0.778.$$

Izmantojot vienādojumu (P1.3.7), aprēķinām koeficienta b_0 vērtību.

$$b_0 = 74 - (-0.778) * (-13) = 74 - 10.150 = 63.850.$$

Galīgais regresijas vienādojums izskatās šādi:

$$y = 63.850 - 0.778x.$$

Lai izveidotu regresijas līkni, mēs aprēķinām y vērtības divām uzdotām x vērtībām.

(Faktiski, vienai no vērtībām vajadzētu būt vērtībai $x = 0$, pie kuras $y = b_0 = 63.850$. Taču šī vērtība ir problemātiska grafiskajam attēlojumam (sk. attēlu P1.3.1). Tāpēc mēs izmantojam tādas x vērtības, kuras var attēlot grafiski).

$$x = -10$$

$$y = 63.850 - 0.778 * (-10) = 63.850 + 7.780 = 71.630.$$

$$x = -20$$

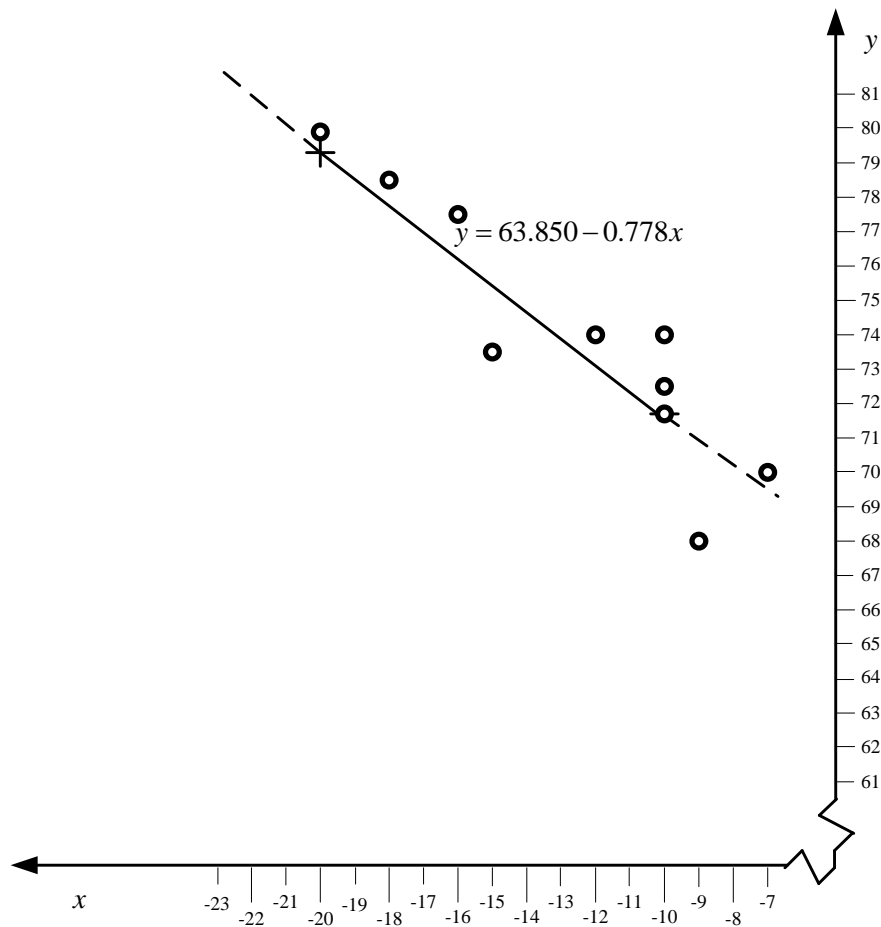
$$y = 63.850 - 0.778 * (-20) = 63.850 + 15.560 = 79.410.$$

Attēlā P1.3.1 parādīta konstruētā regresijas līkne. Punkti, kas atbilst vērtībām $x = -10$, $y = 71.630$ un $x = -20$, $y = 79.410$, šajā attēlā ir apzīmēti ar krustiņiem.

Šis grafiks parāda izteikti negatīvu korelāciju starp temperatūras vērtībām. Palielinoties ārējā gaisa temperatūras negatīvajām vērtībām, paaugstinās ūdens temperatūras vērtība.

Reālās centrālās apkures sistēmās, mainoties āra gaisa temperatūrai, istabas ūdens temperatūra svārstās daudz plašākā diapazonā nekā iepriekš dotajā piemērā. Bet šis piemērs ir sniegts tikai ar

mērķi ilustrēt lineārās regresijas attiecības starp diviem mainīgajiem konstruēšanas principus. Tāpēc piemērā tiek izmantotas mākslīgas temperatūras vērtības.



Attēls P1.3.1. Regresijas grafiks starp ūdens temperatūru un āra temperatūru

P1.4. Attiecību skalas

Šīs skalas ieņem augstāko vietu vispārējā mērījumu skalu hierarhijā. Attiecību skalas ir kvantitatīvās (skaitliskās) skalas ar definētu nulles punktu. Ar vērtībām attiecību skalā var veikt saskaitīšanu, atņemšanu, reizināšanu un dalīšanu. Ņemot divas vērtības šajā skalā, ir iespējams noteikt *par cik* viena vērtība ir pārāka par otru.

Ar attiecību skalas vērtībām var noteikt vidējo vērtību, mediānu un modu. Tāpat var izskaitļot dispersiju un standarta novirzi, var izveidot lineāru regresiju, kas attiecas uz abām vērtību kopām.

Mēs nesniedzam aprēķinu izteiksmes un ilustratīvus piemērus, jo šeit ir spēkā viss, kas tika parādīts iepriekšējā sadaļā par vērtībām, kas mērītas intervālu skalās.

P2. METRIKAS DATU TELPĀ

P2.1. Kas ir metrika?

Šajā pielikumā esošais materiāls ir balstīts uz materiāliem, kas sniegti darbā [Užga-Rebrovs O.I., 2021].

Datu analīzē bieži vien ir jānovērtē, cik tuvu vai tālu atrodas datu elementi viens no otra. Kā var novērtēt interesējošos attālumus starp punktiem vai vektoriem sākotnējā datu kopā? Apzīmēsim ar a un b divus datu elementus (objektus, subjektus vai citas entītijas). Mēs apzīmējam attālumu starp šiem elementiem ar $d(a,b)$.

Formulēsīm prasības, kurām jāatbilst jebkuram pareizam attāluma novērtējumam datu telpā.

1. *Negativitāte*: jebkurš rezultāts $d(a,b)$ nav negatīvs:

$$d(a,b) \geq 0.$$

2. *Identitāte*: $d(a,b) = 0$ tad un tikai tad, ja a un b ir vienādi.

3. *Simetrija*: attālumam starp a un b jābūt vienādam ar attālumu starp b un a

$$d(a,b) = d(b,a).$$

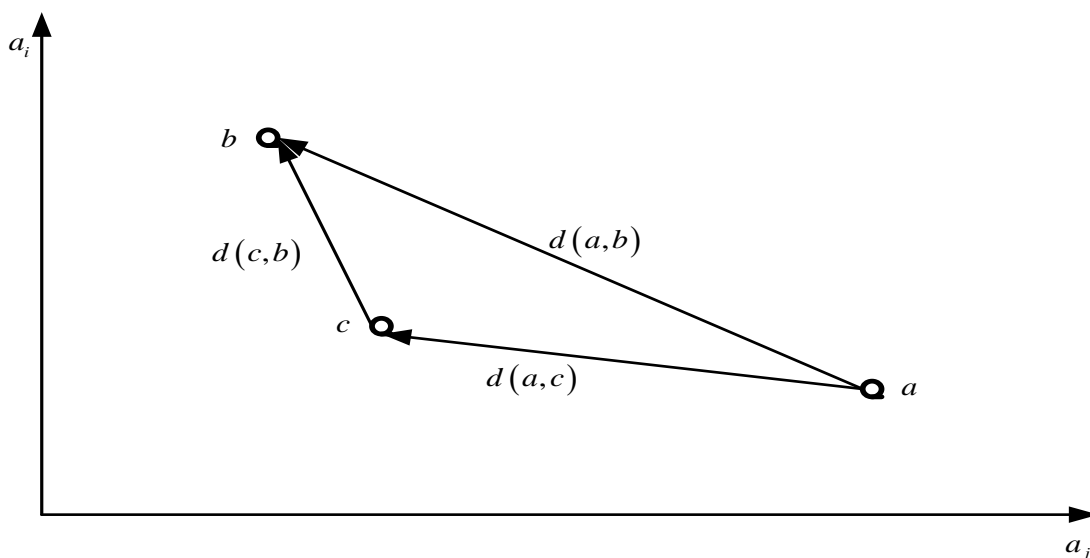
4. *Trijstūra nevienādība*:

$$d(a,b) \leq d(a,c) + d(c,b).$$

Pēdējā prasība ir grafiski interpretēta attēlā P2.1.1 divu atribūtu vērtību telpā a_j, a_k .

Jebkuru attālumu mērīšanas metodi datu telpā, kas atbilst (1) – (4) prasībām, sauc par metriku jeb distanci.

Attālumu mērīšanas metodi, kas neatbilst dažām no šīm prasībām, sauc par semimetriku. Metodi, kas neatbilst (2) prasībai, sauc par pseidometriku, un metodi, kas neatbilst (3) prasībai, sauc par kvazimetriku.



Attēls P2.1.1. Trijstūra noteikuma grafiskā interpretācija divu atribūtu vērtību telpā

Nākamajās sadaļās tiks aplūkotas izplatītās metodes attālumu (metrikas) mērīšanai datu telpās.

P2.2. Eiklīda attālums

Pieņemsim, sākotnējie dati ir uzdoti tabulas veidā. Katrā tabulas rindā tiek parādīts viens datu elements un tam atbilstošais atribūtu vērtību vektors. Katrā kolonā tiek parādīta atbilstošo atribūtu vērtību kopa visiem datu vienumiem.

Apskatām divus datu elementus o_i, o_k . Šos elementus raksturo atribūtu vērtību vektori:

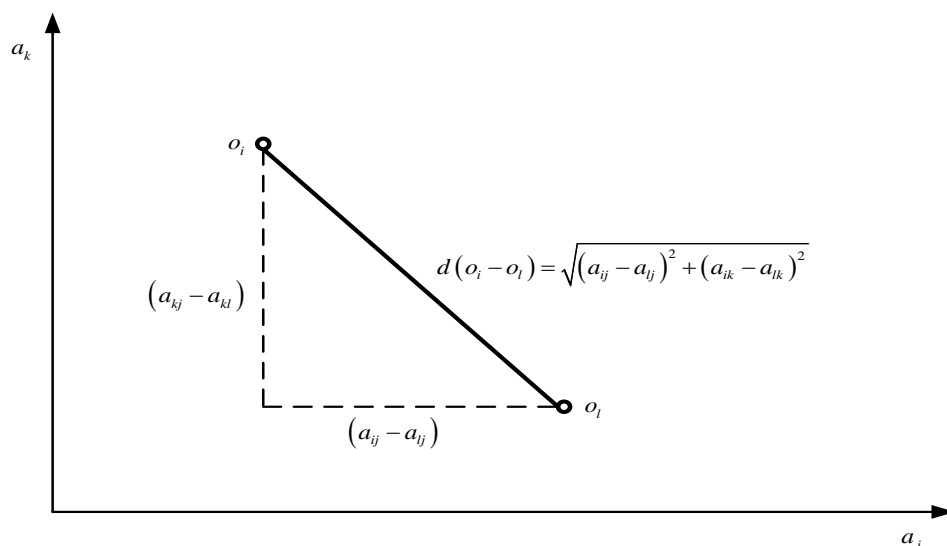
$$o_i: (a_{i1}, a_{i2}, \dots, a_{ij}, \dots, a_{in-1}, a_{in});$$

$$o_k: (a_{k1}, a_{k2}, \dots, a_{kj}, \dots, a_{kn-1}, a_{kn}).$$

Eiklīda attālumu starp šiem elementiem aprēķina pēc vienādojuma

$$d(o_i, o_k) = \sqrt{\sum_{j=1}^n (a_{ij} - a_{kj})^2}. \quad (\text{P2.2.1})$$

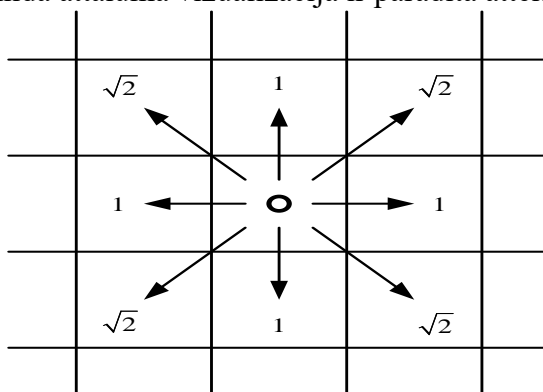
Attēlā P2.2.1 parādīta grafiska interpretācija Eiklīda attālumam starp diviem datu elementiem divu atribūtu vērtību telpā.



Attēls P2.2.1. Eiklīda attāluma grafiskā interpretācija divu atribūtu telpā

Acīmredzams, ka Eiklīda attālums ir mazākais attālums starp elementiem o_i, o_k .

Lai vizualizētu aprēķinātās izteiksmes attālumiem datu telpā, dažkārt tiek izmantota šaha galdiņa interpretācija. Šī Eiklīda attāluma vizualizācija ir parādīta attēlā P2.2.2.



Attēls P2.2.2. Eiklīda attāluma vizualizācija ar šaha galdiņu

Tiek pieņemts, ka attālums starp blakus esošajām šūnām horizontāli un vertikāli ir vienāds ar 1. Šie attālumi ir parādīti attēlā P2.2.2 ar skaitli 1. Lai pārietu no centrālās šūnas uz jebkuru no diagonālajām šūnām, jāveic viens solis horizontāli un viens solis vertikāli. Saskaņā ar vienādojumu (P2.2.1) attālums starp centrālo šūnu un jebkuru no diagonālajām šūnām ir vienāds ar $\sqrt{2}$.

Apskatīsim ilustratīvu piemēru Eiklīda attālumu aprēķināšanai starp elementu (objektu) pāriem.

Piemērs P2.2.1. Sākotnējā datu kopa ir parādīta tabulā P2.2.1. Ir nepieciešams aprēķināt Eiklīda attālumus starp visiem objektu pāriem.

Tabula P2.2.1

Sākotnējie dati piemērā P2.2.1

	a_{i1}	a_{i2}	a_{i3}	a_{i4}
o_1	7	3	2	10
o_2	9	2	5	6
o_3	12	5	8	8
o_4	5	4	6	9

Risinājums. Izmantojot vienādojumu (P2.2.1), mēs aprēķinām attālumus no objekta o_1 līdz citiem objektiem.

$$\begin{aligned} d(o_1 - o_2) &= \sqrt{(7-9)^2 + (3-2)^2 + (2-5)^2 + (10-6)^2} = \sqrt{(-2)^2 + 1^2 + (-3)^2 + 4^2} = \\ &= \sqrt{4+1+9+16} = \sqrt{30} = 5.48; \end{aligned}$$

$$\begin{aligned} d(o_1 - o_3) &= \sqrt{(7-12)^2 + (3-5)^2 + (2-8)^2 + (10-8)^2} = \sqrt{(-5)^2 + (-2)^2 + (-6)^2 + 2^2} = \\ &= \sqrt{25+4+36+4} = \sqrt{69} = 8.31; \end{aligned}$$

$$\begin{aligned} d(o_1 - o_4) &= \sqrt{(7-5)^2 + (3-4)^2 + (2-6)^2 + (10-9)^2} = \sqrt{2^2 + (-1)^2 + (-4)^2 + 1^2} = \\ &= \sqrt{4+1+16+1} = \sqrt{22} = 4.69. \end{aligned}$$

Nav nepieciešams aprēķināt attālumus $d(o_2 - o_1)$, $d(o_3 - o_1)$, $d(o_4 - o_1)$, jo pamatojoties uz simetrijas prasību $d(o_2 - o_1) = d(o_1 - o_2)$, $d(o_3 - o_1) = d(o_1 - o_3)$, $d(o_4 - o_1) = d(o_1 - o_4)$

Aprēķinām attālumus no objekta o_2 līdz objektiem o_3, o_4

$$\begin{aligned} d(o_2 - o_3) &= \sqrt{(9-12)^2 + (2-5)^2 + (5-8)^2 + (6-8)^2} = \sqrt{(-3)^2 + (-3)^2 + (-3)^2 + (-2)^2} = \\ &= \sqrt{9+9+9+4} = \sqrt{31} = 5.57; \end{aligned}$$

$$\begin{aligned} d(o_2 - o_4) &= \sqrt{(9-5)^2 + (2-4)^2 + (5-6)^2 + (6-9)^2} = \sqrt{4^2 + (-2)^2 + (-1)^2 + (-3)^2} = \\ &= \sqrt{16+4+1+9} = \sqrt{30} = 5.48. \end{aligned}$$

Aprēķinām attālumu no objekta o_3 līdz objektam o_4

$$\begin{aligned} d(o_3 - o_4) &= \sqrt{(12-5)^2 + (5-4)^2 + (8-6)^2 + (8-9)^2} = \sqrt{7^2 + 1^2 + 2^2 + (-1)^2} = \\ &= \sqrt{49+1+4+1} = \sqrt{55} = 7.42. \end{aligned}$$

Visi nepieciešamie aprēķini ir veikti. Citas vērtības tiek noteiktas, pamatojoties uz simetrijas principu. Aprēķinu rezultāti ir parādīti tabulā P2.2.2.

Tabula P2.2.2

Attālumu starp objektiem aprēķinu rezultāti piemērā P2.2.1

	o_1	o_2	o_3	o_4
o_1	0	5.48	8.31	4.69
o_2	5.48	0	5.57	5.48
o_3	8.31	5.57	0	7.42
o_4	4.69	5.48	7.42	0

Šūnās gar tabulas diagonāli ir nulles, jo šīs šūnas parāda objektu attālumus pašas uz sevi. Pamatojoties uz simetrijas principu, vērtības tabulas šūnās ir simetriskas attiecībā pret diagonāli.

Piemērā P2.2.1 ir aprēķināti tabulā P2.2.1 uzrādīto objektu attālumi telpā. Tomēr attālumus atribūtu vērtību telpā var aprēķināt, izmantojot to pašu principu. Mēs to parādīsim ar šādu piemēru.

Piemērs P2.2.2. Izmantojot sākotnējos datus no tabulas P2.2.1, ir nepieciešams aprēķināt attālumus starp atribūtiem, izmantojot vienādojumu (P2.2.1).

Risinājums.

$$\begin{aligned} d(a_1 - a_2) &= \sqrt{(7-3)^2 + (9-2)^2 + (12-5)^2 + (5-4)^2} = \sqrt{4^2 + 7^2 + 7^2 + 1^2} = \\ &= \sqrt{16 + 49 + 49 + 1} = \sqrt{115} = 10.72; \end{aligned}$$

$$\begin{aligned} d(a_1 - a_3) &= \sqrt{(7-2)^2 + (9-5)^2 + (12-8)^2 + (5-6)^2} = \sqrt{5^2 + 4^2 + 4^2 + (-1)^2} = \\ &= \sqrt{25 + 16 + 16 + 1} = \sqrt{58} = 7.62; \end{aligned}$$

$$\begin{aligned} d(a_1 - a_4) &= \sqrt{(7-10)^2 + (9-6)^2 + (12-8)^2 + (5-9)^2} = \sqrt{(-3)^2 + 3^2 + 4^2 + (-4)^2} = \\ &= \sqrt{9 + 9 + 16 + 16} = \sqrt{50} = 7.07; \end{aligned}$$

$$\begin{aligned} d(a_2 - a_3) &= \sqrt{(3-2)^2 + (2-5)^2 + (5-8)^2 + (4-6)^2} = \sqrt{1^2 + (-3)^2 + (-3)^2 + (-2)^2} = \\ &= \sqrt{1 + 9 + 9 + 4} = \sqrt{23} = 4.80; \end{aligned}$$

$$\begin{aligned} d(a_2 - a_4) &= \sqrt{(3-10)^2 + (2-6)^2 + (5-8)^2 + (4-9)^2} = \sqrt{(-7)^2 + (-4)^2 + (-3)^2 + (-5)^2} = \\ &= \sqrt{49 + 16 + 9 + 25} = \sqrt{99} = 9.95; \end{aligned}$$

$$\begin{aligned} d(a_3 - a_4) &= \sqrt{(2-10)^2 + (5-6)^2 + (8-8)^2 + (6-9)^2} = \sqrt{(-8)^2 + (-1)^2 + 0 + (-3)^2} = \\ &= \sqrt{64 + 1 + 0 + 9} = \sqrt{74} = 8.60. \end{aligned}$$

Aprēķinu rezultāti ir parādīti tabulā P2.2.3.

Attālumu starp atribūtiem aprēķinu rezultāti piemērā P2.2.2

	a_1	a_2	a_3	a_4
a_1	0	10.72	7.62	7.07
a_2	10.72	0	4.80	9.95
a_3	7.62	4.80	0	8.60
a_4	7.07	9.95	8.60	0

Eiklīda attāluma divkāršā mērogošana.

Tās būtība ir šāda. Datu tabulas j -tajai kolonai tiek aprēķināta starpība

$$\chi_j = (a_j^{\max} - a_j^{\min}), \quad j = 1, \dots, n,$$

kur a_j^{\max} , a_j^{\min} ir atribūta a_j maksimālās un minimālās vērtības.

Tiek aprēķināta mērogotā Eiklīda attāluma vērtība pēc vienādojuma

$$d^*(o_i - o_l) = \sqrt{\sum_{j=1}^n \left(\frac{(a_{ij} - a_{lj})^2}{\chi_j} \right)}, \quad i, l = 1, \dots, m, \quad (\text{P2.2.2})$$

Visbeidzot, divkāršā mēroga Eiklīda attāluma vērtība tiek aprēķināta pēc vienādojuma

$$d^{**}(o_i - o_l) = \frac{\sqrt{\sum_{j=1}^n \left(\frac{(a_{ij} - a_{lj})^2}{\chi_j} \right)}}{\sqrt{n}}, \quad i, l = 1, \dots, m, \quad i \neq l, \quad (\text{P2.2.3})$$

kur n ir atribūtu (kolonu) skaits sākotnējo datu tabulā.

Piemērs P2.2.3. Nepieciešams aprēķināt mērogotos un divkārši mērogotos Eiklīda attālumus starp objektiem no tabulas P2.2.1.

Risinājums. Izmantojot sākotnējos datus no tabulas P2.2.1, mēs aprēķinām vērtības χ_j , $j = 1, 2, 3, 4$.

$$\chi_1 = (12 - 5) = 7, \quad \chi_2 = (5 - 2) = 3, \quad \chi_3 = (8 - 2) = 6; \quad \chi_4 = (10 - 6) = 4.$$

Izmantojot vienādojumu (P2.2.2), mēs aprēķinām mērogotos Eiklīda attālumus.

$$\begin{aligned} d^*(o_1 - o_2) &= \sqrt{\frac{(7-9)^2}{7} + \frac{(3-2)^2}{3} + \frac{(2-5)^2}{6} + \frac{(10-6)^2}{4}} = \sqrt{\frac{4}{7} + \frac{1}{3} + \frac{9}{6} + \frac{16}{4}} = \\ &= \sqrt{0.571 + 0.333 + 1.500 + 4.000} = \sqrt{6.404} = 2.53; \end{aligned}$$

$$\begin{aligned} d^*(o_1 - o_3) &= \sqrt{\frac{(7-12)^2}{7} + \frac{(3-5)^2}{3} + \frac{(2-8)^2}{6} + \frac{(10-8)^2}{4}} = \sqrt{\frac{25}{7} + \frac{4}{3} + \frac{36}{6} + \frac{4}{4}} = \\ &= \sqrt{3.571 + 1.333 + 6.000 + 1.000} = \sqrt{11.904} = 3.45; \end{aligned}$$

$$d^*(o_1 - o_4) = \sqrt{\frac{(7-5)^2}{7} + \frac{(3-4)^2}{3} + \frac{(2-6)^2}{6} + \frac{(10-9)^2}{4}} = \sqrt{\frac{4}{7} + \frac{1}{3} + \frac{16}{6} + \frac{1}{4}} =$$

$$= \sqrt{0.571 + 0.333 + 2.667 + 0.250} = \sqrt{3.821} = 1.95;$$

$$d^*(o_2 - o_3) = \sqrt{\frac{(9-12)^2}{7} + \frac{(2-5)^2}{3} + \frac{(5-8)^2}{6} + \frac{(6-8)^2}{4}} = \sqrt{\frac{9}{7} + \frac{9}{3} + \frac{9}{6} + \frac{4}{4}} =$$

$$= \sqrt{1.286 + 3.000 + 1.500 + 1.000} = \sqrt{6.786} = 2.60;$$

$$d^*(o_2 - o_4) = \sqrt{\frac{(9-5)^2}{7} + \frac{(2-4)^2}{3} + \frac{(5-6)^2}{6} + \frac{(6-9)^2}{4}} = \sqrt{\frac{16}{7} + \frac{4}{3} + \frac{1}{6} + \frac{9}{4}} =$$

$$= \sqrt{2.286 + 1.333 + 0.167 + 2.250} = \sqrt{6.036} = 2.46;$$

$$d^*(o_3 - o_4) = \sqrt{\frac{(12-5)^2}{7} + \frac{(5-4)^2}{3} + \frac{(8-6)^2}{6} + \frac{(8-9)^2}{4}} = \sqrt{\frac{49}{7} + \frac{1}{3} + \frac{4}{6} + \frac{1}{4}} =$$

$$= \sqrt{7.000 + 0.333 + 0.667 + 0.250} = \sqrt{8.250} = 2.87.$$

Iegūtie rezultāti ir apkopoti tabulā P2.2.4.

Tabula P4.2.4

Mērogoto Eiklīda attālumu vērtības piemērā P2.2.3

	o_1	o_2	o_3	o_4
o_1	0	2.53	3.45	1.95
o_2	2.53	0	2.60	2.45
o_3	3.45	2.60	0	2.87
o_4	1.95	2.46	2.87	0

Izmantojot vienādojumu (P2.2.3), mēs aprēķinām divkārsi mērogotos Eiklīda attālumus. Tā kā mūsu piemērā atribūtu skaits $n = 4$, $\sqrt{n} = \sqrt{4} = 2$, lai aprēķinātu nepieciešamos attālumus, pietiek dalīt ar 2 attāluma vērtības no tabulas P2.2.4. Galīgie rezultāti ir parādīti tabulā P2.2.5.

Tabula P2.2.5

Divkārsi mērogoto Eiklīda attālumu vērtības piemērā P2.2.3

	o_1	o_2	o_3	o_4
o_1	0	1.27	1.72	0.98
o_2	1.27	0	1.30	1.23
o_3	1.72	1.30	0	1.43
o_4	0.98	1.23	1.43	0

Lai noteiktu attālumus starp atribūtiem, var izmantot relatīvos, mērogotos un divkārsi mērogotos Eiklīda attālumus.

Noslēdzot šo sadaļu, mēs atzīmējam, ka Eiklīda attālums ir Minkovska attāluma īpašs gadījums

$$d(o_i - o_l) = \sum_{j=1}^n \left((a_{ij} - a_{lj})^p \right)^{\frac{1}{p}}. \quad (\text{P2.2.5})$$

Pie $p = 2$ Minkovska attālumā tiek iegūta standarta Eiklīda distance.

P2.3. Manhetenas attālums

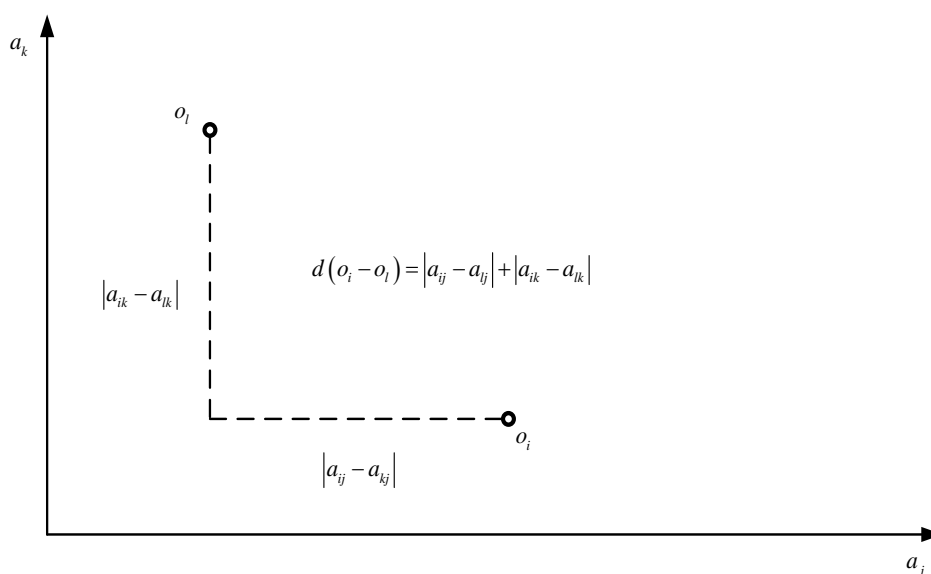
Pieņem, ka dati attēloti tabulas veidā, kurā katra rinda attēlo objektu un tam atbilstošu atribūtu vērtību vektoru (a_{ij}) . Manhetenas attālums starp objektiem a_i , a_k tiek aprēķināts kā

$$d(o_i - o_l) = \sum_{j=1}^n |a_{ij} - a_{lj}|, \quad (\text{P2.3.1})$$

kur n ir atribūtu skaits.

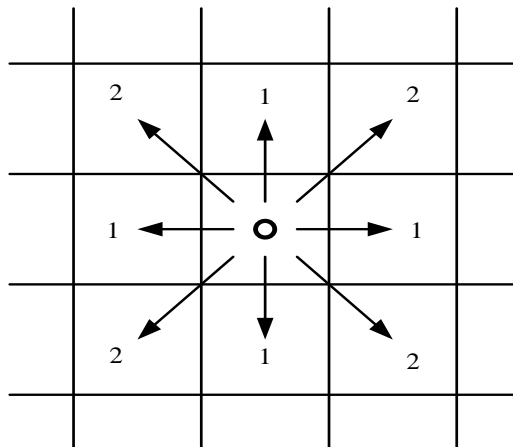
Šo attālumu var interpretēt kā ceļa garumu, kas jāiet pa pilsētas ielām (jo īpaši Ņujorkas Manhetenas rajonā), lai no punkta a_i sasniegtu punktu a_k . Līdz ar to radies šī attāluma nosaukums.

Attēls P2.3.1 attēlo Manhetenas attāluma grafisku interpretāciju divu atribūtu telpā.



Attēls P2.3.1. Manhetenas attāluma grafiskā interpretācija divu atribūtu telpā

Manhetenas attāluma vizualizācija, izmantojot šaha galdiņu, ir parādīta attēlā P2.3.2. Lai nokļūtu no centrālās šūnas uz jebkuru blakus esošo šūnu horizontāli vai vertikāli, ir jāveic viens solis. Lai nokļūtu no centrālās šūnas uz jebkuru pa diagonāli blakus esošo šūnu, ir jāveic divas darbības: viena horizontāli un viena vertikāli.



Attēls P2.3.2. Manhetenas attāluma vizualizācija ar šaha galdiņu

Tā kā pēc definīcijas soļa garums ir vienāds ar 1, summējot divu soļu garumus saskaņā ar vienādojumu (P2.3.1), iegūstam, ka visi šie attālumi ir vienādi ar 2. Tādējādi Manhetenas attālums ir metrika.

Apskatīsim piemēru, kas ilustrē Manhetenas attālumu aprēķinus objektu telpā.

Piemērs P2.3.1. Nepieciešams aprēķināt Manhetenas attālumus starp tabulā P2.2.1 dotajiem objektiem.

Risinājums. Ērtības labad sākotnējie dati no tabulas P2.2.1 ir reproducēti tabulā P2.3.1.

Tabula P2.3.1

Sākotnējie dati piemērā P2.3.1

	a_{i1}	a_{i2}	a_{i3}	a_{i4}
o_1	7	3	2	10
o_2	9	2	5	6
o_3	12	5	8	8
o_4	5	4	6	9

Aprēķināsim nepieciešamos attālumus, izmantojot vienādojumu (P2.3.1).

$$d(o_1 - o_2) = |7 - 9| + |3 - 2| + |2 - 5| + |10 - 6| = 2 + 1 + 3 + 4 = 10;$$

$$d(o_1 - o_3) = |7 - 12| + |3 - 5| + |2 - 8| + |10 - 8| = 5 + 2 + 6 + 2 = 15;$$

$$d(o_1 - o_4) = |7 - 5| + |3 - 4| + |2 - 6| + |10 - 9| = 2 + 1 + 4 + 1 = 8;$$

$$d(o_2 - o_3) = |9 - 12| + |2 - 5| + |5 - 8| + |6 - 8| = 3 + 3 + 3 + 2 = 11;$$

$$d(o_2 - o_4) = |9 - 5| + |2 - 4| + |5 - 6| + |6 - 9| = 4 + 2 + 1 + 3 = 10;$$

$$d(o_3 - o_4) = |12 - 5| + |5 - 4| + |8 - 6| + |8 - 9| = 7 + 1 + 2 + 1 = 11.$$

Aprēķinu rezultāti ir parādīti tabulā P2.3.2.

Tabula P2.3.2

Attālumu starp objektiem aprēķinu rezultāti piemērā P2.3.1

	o_1	o_2	o_3	o_4
o_1	0	10	15	8
o_2	10	0	11	10
o_3	15	11	0	11
o_4	8	10	11	0

Manhetenas attālumus var arī aprēķināt, izmantojot normalizētas atribūtu vērtības.

P2.4. Alternatīvi attāluma mēri starp objektiem

Eiklīda attālums un Manhetenas attālums ir populāri attāluma mēri, kurus plaši izmanto datu analīzē. Tomēr papildus tiek izmantoti arī citi attāluma mēri. Šajā sadaļā ir aplūkoti daži no šiem alternatīvajiem mēriem.

1. Breja-Kurtisa koeficients (Bray-Curtis coefficient).

$$d(o_i - o_l) = 1 - \frac{2 \sum_{j=1}^n \min(a_{ij}, a_{lj})}{\sum_{j=1}^n a_{ij} + \sum_{j=1}^n a_{lj}}, \quad i, l = 1, \dots, m, \quad i \neq l, \quad (\text{P2.4.1})$$

kur \min ir minimuma operators.

Piemērs P2.4.1. Ir nepieciešams aprēķināt Breja-Kurtisa koeficientu vērtības pēc vienādojuma (P2.4.1) objektu kopai piemērā P2.2.1 (Tabula P2.2.1).

Risinājums. Tabulā P2.4.1 tiek atjaunoti dati no tabulas P2.2.1. Tabulas pēdējā kolona parāda atribūtu vērtību summas visās tabulas rindās.

Tabula P2.4.1

Sākotnējie dati piemērā P2.4.1

	a_{i1}	a_{i2}	a_{i3}	a_{i4}	$\sum_{j=1}^4 a_{ij}$
o_1	7	3	2	10	22
o_2	9	2	5	6	22
o_3	12	5	8	8	33
o_4	5	4	6	9	24

Veiksim nepieciešamos aprēķinus, izmantojot vienādojumu (P2.4.1).

$$d(o_1 - o_2) = 1 - \frac{2 * (7 + 2 + 2 + 6)}{22 + 22} = 1 - \frac{34}{44} = 1 - 0.77 = 0.23;$$

$$d(o_1 - o_3) = 1 - \frac{2 * (7 + 3 + 2 + 8)}{22 + 33} = 1 - \frac{40}{55} = 1 - 0.73 = 0.27;$$

$$d(o_1 - o_4) = 1 - \frac{2 * (5 + 3 + 2 + 9)}{22 + 24} = 1 - \frac{38}{46} = 1 - 0.83 = 0.17;$$

$$d(o_2 - o_3) = 1 - \frac{2 * (9 + 2 + 5 + 6)}{22 + 33} = 1 - \frac{44}{55} = 1 - 0.80 = 0.20;$$

$$d(o_2 - o_4) = 1 - \frac{2 * (5 + 2 + 5 + 6)}{22 + 24} = 1 - \frac{36}{46} = 1 - 0.78 = 0.22;$$

$$d(o_3 - o_4) = 1 - \frac{2 * (5 + 4 + 6 + 8)}{33 + 24} = 1 - \frac{46}{57} = 1 - 0.81 = 0.19.$$

Aprēķinu rezultāti ir apkopoti tabulā P2.4.2.

Tabula P2.4.2

Attālumu starp objektiem aprēķinu rezultāti piemēram P2.4.1

	o_1	o_2	o_3	o_4
o_1	0	0.23	0.27	0.17
o_2	0.23	0	0.20	0.22
o_3	0.27	0.20	0	0.19
o_4	0.17	0.22	0.19	0

2. Relatīvais Sorensena attālums (Sorensen Distance).

$$d(o_i - o_l) = 1 - \sum_{j=1}^n \min \left(\frac{a_{ij}}{\sum_{j=1}^n a_{ij}}, \frac{a_{lj}}{\sum_{j=1}^n a_{lj}} \right), \quad i, l = 1, \dots, m, \quad i \neq l, \quad (\text{P2.4.2})$$

Piemērs P2.4.2. Izmantojot sākotnējos datus no tabulas P2.4.1, ir nepieciešams aprēķināt relatīvos Sorensena attālumus starp objektiem, izmantojot vienādojumu (P2.4.2).

Risinājums. Veiksim starpposma aprēķinus, pamatojoties uz tabulas P2.4.1 datiem. Šo aprēķinu rezultāti ir parādīti tabulā P2.4.3.

Tabula P2.4.3

Starpposma aprēķinu rezultāti piemēram P2.4.2

	$\frac{a_{i1}}{\sum_{j=1}^4 a_{ji}}$	$\frac{a_{i2}}{\sum_{j=1}^4 a_{ji}}$	$\frac{a_{i3}}{\sum_{j=1}^4 a_{ji}}$	$\frac{a_{i4}}{\sum_{j=1}^4 a_{ji}}$
o_1	0.32	0.14	0.09	0.45
o_2	0.41	0.09	0.23	0.27
o_3	0.36	0.15	0.24	0.24
o_4	0.21	0.17	0.25	0.37

Veiksim aprēķinus, izmantojot vienādojumu (P2.4.2).

$$d(o_1 - o_2) = 1 - (0.32 + 0.09 + 0.09 + 0.27) = 1 - 0.77 = 0.23;$$

$$d(o_1 - o_3) = 1 - (0.32 + 0.14 + 0.09 + 0.24) = 1 - 0.79 = 0.21;$$

$$d(o_1 - o_4) = 1 - (0.21 + 0.14 + 0.09 + 0.37) = 1 - 0.81 = 0.19;$$

$$d(o_2 - o_3) = 1 - (0.36 + 0.09 + 0.23 + 0.24) = 1 - 0.92 = 0.08;$$

$$d(o_2 - o_4) = 1 - (0.21 + 0.09 + 0.23 + 0.27) = 1 - 0.80 = 0.20;$$

$$d(o_3 - o_4) = 1 - (0.21 + 0.15 + 0.24 + 0.24) = 1 - 0.84 = 0.16.$$

Aprēķinu rezultāti ir parādīti tabulā P2.4.4.

Tabula P2.4.4

Attālumu starp objektiem aprēķinu rezultāti piemēram P2.4.2

	o_1	o_2	o_3	o_4
o_1	0	0.23	0.21	0.19
o_2	0.23	0	0.08	0.20
o_3	0.21	0.08	0	0.16
o_4	0.19	0.20	0.16	0

Relatīvais Sorensena attālums ir semimetrisks.

Vienādojuma (P2.4.2) variants, izmantojot starpību absolūtās vērtības operatora min vietā, ir vienādojums

$$d(o_i - o_l) = \frac{1}{2} \sum_{j=1}^n \left| \frac{a_{ij}}{\sum_{j=1}^n a_{ij}} - \frac{a_{lj}}{\sum_{j=1}^n a_{lj}} \right| \dots i, l = 1, \dots, m, \quad i \neq l \quad (\text{P2.4.3})$$

Piemērs P2.4.3. Izmantojot starposma aprēķinu rezultātus tabulā P2.4.3 ir jāaprēķina attālumu vērtības starp objektiem, izmantojot vienādojumu (P2.4.3).

Risinājums. Veiksim aprēķinus, izmantojot vienādojumu (P2.4.3).

$$\begin{aligned} d(o_1 - o_2) &= \frac{1}{2} * (|0.32 - 0.41| + |0.14 - 0.09| + |0.09 - 0.23| + |0.45 - 0.27|) = \\ &= \frac{1}{2} * (0.09 + 0.05 + 0.14 + 0.18) = \frac{1}{2} * 0.46 = 0.23; \end{aligned}$$

$$\begin{aligned} d(o_1 - o_3) &= \frac{1}{2} * (|0.32 - 0.36| + |0.14 - 0.15| + |0.09 - 0.24| + |0.45 - 0.24|) = \\ &= \frac{1}{2} * (0.04 + 0.01 + 0.15 + 0.21) = \frac{1}{2} * 0.41 = 0.21. \end{aligned}$$

Pārējie aprēķini tiek veikti līdzīgi. Aprēķinu rezultāti ir parādīti tabulā P2.4.5.

Tabula P2.4.5

Attālumu starp objektiem aprēķinu rezultāti piemēram P2.4.3

	o_1	o_2	o_3	o_4
o_1	0	0.23	0.21	0.19
o_2	0.23	0	0.08	0.20
o_3	0.21	0.08	0	0.16
o_4	0.19	0.20	0.16	0

P2.5. χ^2 attālums

Atgriezīsimies pie Eiklīda attāluma, kas bija aplūkots sadaļā P2.2. Apzīmēsim ar \mathbf{a}_i objekta o_i atribūtu vērtību vektoru un ar \mathbf{a}_l objekta o_l atribūtu vērtību vektoru. Vispārīgā veidā kvadrātiskā Eiklīda attāluma aprēķināšanai starp atribūtu vērtību vektoriem (\mathbf{a}_i) , (\mathbf{a}_l) (objektiem o_i, o_l) var uzdot vektoru reizinājuma formā.

$$d^2 d^2(o_i, o_l) = d^2(\mathbf{a}_i, \mathbf{a}_l) = (\mathbf{a}_{ij} - \mathbf{a}_{lj})^T (\mathbf{a}_{ij} - \mathbf{a}_{lj}) = \sum_{j=1}^n (a_{ij} - a_{lj})^2,$$

kur T ir vektora transponēšanas simbols.

(Jāņem vērā, ka šeit ir runa par kvadrātisko Eiklīda attālumu. Šis attālums tiek izmantots arī dažādos praktiskos pielietojumos. Lai iegūtu standarta Eiklīda attāluma aprēķinu, vienkārši jāņem kvadrātsakne no Eiklīda kvadrātiskā attāluma aprēķina).

Aprēķinot kvadrātisko (un standarta) Eiklīda attālumu, visām atribūtu vērtībām ir vienāds svars. Tomēr var ieviest paplašināto Eiklīda attāluma versiju, kurā atribūtu vērtību vektoriem datu tabulas kolonās tiek piešķirti dažādi svari. Tādā veidā var iegūt kvadrātiskā attāluma χ_q^2 aprēķinus. Vispārīgā veidā kvadrātisko attālumu χ_q^2 starp atribūtu $\mathbf{a}_i, \mathbf{a}_l$ vērtību vektoriem (objektiem o_i, o_l) var definēt kā

$$\chi_q^2(o_i - o_l) \chi_q^2(\mathbf{a}_i, \mathbf{a}_l) = (\mathbf{a}_{ij} - \mathbf{a}_{lj})^T \mathbf{W} (\mathbf{a}_{ij} - \mathbf{a}_{lj}). \quad (\text{P2.5.1})$$

Vienādojumā (P2.5.1) simbols \mathbf{W} apzīmē atsevišķiem atribūtiem piešķirto svaru matricu.

Šajā sadaļā mēs apskatīsim divas svāra matricas \mathbf{W} formēšanas iespējas.

Būtiska attāluma χ_q^2 aprēķinu iezīme ir tā, ka tie balstās nevis uz faktiskajām atribūtu vērtībām, bet gan uz šo vērtību profiliem. Lai pārietu no atribūtu vērtībām uz to profilu vērtībām, katrai datu tabulas rindai tiek aprēķināta atribūtu vērtību summa

$$a_i^* = \sum_{j=1}^n a_{ij}, \quad i = 1, \dots, m.$$

Profila vērtības tiek aprēķinātas kā faktisko atribūtu vērtību attiecība pret to summām a_i^*

$$p_{ij} = \frac{a_{ij}}{a_i^*} \dots \quad i = 1, \dots, m \quad j = 1, \dots, n.$$

Pirmajā variantā svaru matricas \mathbf{W} komponentu vērtības tiek aprēķinātas šādi. Katrā datu tabulas kolonā tiek aprēķināta atribūtu vērtību summa:

$$a_j^* = \sum_{i=1}^m a_{ij}, \quad j = 1, \dots, n.$$

Par svaru vērtībām tiek pieņemtas vērtības, kas ir a_{ij}^* apgrieztās vērtības:

$$w_j = \frac{1}{a_j^*}.$$

Svāra matrica \mathbf{W} ir diagonāla matrica:

$$\mathbf{W} = \begin{bmatrix} \frac{1}{a_1^*} & \dots & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & \frac{1}{a_j^*} & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & 0 & \dots & 0 & \dots & \frac{1}{a_n^*} \end{bmatrix}.$$

Formāli priekš pirmā varianta var formulēt šādu aprēķina vienādojumu

$$\chi_1^2(o_i - o_l) = \sum_{j=1}^n \frac{1}{a_j^*} \left(\frac{a_{ij}}{a_i^*} - \frac{a_{lj}}{a_l^*} \right)^2 = \sum_{j=1}^n \frac{1}{a_j^*} (p_{ij} - p_{lj})^2, \quad i, l = 1, \dots, m, \quad i \neq l, \quad (\text{P2.5.2})$$

kur p_{ij} , p_{lj} ir atribūtu a_{ij} , a_{lj} vērtību profili.

Apskatīsim piemēru, kas ilustrē χ_1^2 attālumu aprēķināšanas metodi.

Piemērs P2.5.1. Ņemsim par pamatu sākotnējos datus no tabulas P2.2.1. Ir nepieciešams aprēķināt kvadrātisko un standarta attālumu χ_1^2 , χ_1 vērtības starp objektiem (to atribūtu vērtību vektoriem), izmantojot pirmo variantu svaru matricas veidošanai.

Risinājums. Tabulas P2.5.1 kreisajā daļā ir parādīti sākotnējie dati no tabulas P2.2.1.

Tabula P2.5.1

Sākotnējie dati un provizorisko aprēķinu rezultāti piemēram P2.5.1

	Sākotnējās atribūtu vērtības				a_i^*	Atribūtu vērtību profili			
	a_{i1}	a_{i2}	a_{i3}	a_{i4}		p_{i1}	p_{i2}	p_{i3}	p_{i4}
o_1	7	3	2	10	22	0.318	0.136	0.091	0.455
o_2	9	2	5	6	22	0.409	0.091	0.228	0.273
o_3	12	5	8	8	33	0.364	0.152	0.242	0.242
o_4	5	4	6	9	24	0.208	0.167	0.250	0.375
a_j^*	33	14	21	33					
$w_1 = \frac{1}{a_j^*}$	0.030	0.071	0.048	0.030					

Tabulas labajā daļā ir redzamas atribūtu p_{ij} profila vērtības, kas tiek aprēķinātas, dalot atribūtu a_{ij} vērtības ar atbilstošajām a_i^* vērtībām.

Tabulas priekšpēdējā rindā dotas a_j^* vērtības, kas ir atribūtu vērtību summas attiecīgajās tabulas kolonās.

Tabulas pēdējā rindā ir norādīti svāri $w_j = \frac{1}{a_j^*}$.

Veicam aprēķinus, izmantojot vienādojumu (P2.5.2).

$$\begin{aligned}\chi_1^2(o_1 - o_2) &= 0.030 * (0.318 - 0.409)^2 + 0.071 * (0.136 - 0.091)^2 + 0.048 * (0.091 - 0.228)^2 + \\ &+ 0.030 * (0.455 - 0.273) = \\ &= 0.030 * 0.00828 + 0.071 * 0.00202 + 0.048 * 0.01877 + 0.030 * 0.03312 = \\ &= 0.00025 + 0.00014 + 0.00090 + 0.00010 \approx 0.0014;\end{aligned}$$

$$\begin{aligned}\chi_1^2(o_1 - o_3) &= 0.030 * (0.318 - 0.364)^2 + 0.071 * (0.136 - 0.152)^2 + 0.048 * (0.091 - 0.242)^2 + \\ &+ 0.030 * (0.455 - 0.242) = \\ &= 0.030 * 0.00212 + 0.071 * 0.00026 + 0.048 * 0.02280 + 0.030 * 0.04537 = \\ &= 0.00006 + 0.00002 + 0.00109 + 0.00136 \approx 0.0025.\end{aligned}$$

Pārējie aprēķini tiek veikti tādā pašā veidā. Aprēķinu rezultāti ir parādīti tabulā P2.5.2.

Tabula P2.5.2

Kvadrātisko attālumu χ_1^2 aprēķinu rezultāti piemēram P2.5.1

	o_1	o_2	o_3	o_4
o_1	0	0.0014	0.0025	0.0018
o_2	0.0014	0	0.0004	0.0019
o_3	0.0025	0.0004	0	0.0012
o_4	0.0018	0.0019	0.0012	0

Lai aprēķinātu standarta attālumus χ_1 , vienādojums (P2.5.2) jāpārvērš formā

$$\chi_1(o_i - o_l) = \sqrt{\sum_{j=1}^n \frac{1}{a_j^*} \left(\frac{a_{ij}^*}{a_i^*} - \frac{a_{lj}^*}{a_l^*} \right)^2} = \sqrt{\sum_{j=1}^n \frac{1}{a_j^*} (p_{ij} - p_{lj})^2} \dots i, l = 1, \dots, m, i \neq l. \quad (\text{P2.5.3})$$

Šajā piemērā nav nepieciešams veikt aprēķinus, izmantojot vienādojumu (P2.5.3). Lai iegūtu nepieciešamās vērtības, pietiek ar kvadrātsakni no katras vērtības no tabulas P2.5.2. Gala rezultāti ir parādīti tabulā P2.5.3.

Tabula P2.5.3

Standarta attālumu χ_1 aprēķinu rezultāti piemēram P2.5.1

	o_1	o_2	o_3	o_4
o_1	0	0.0374	0.0500	0.0424
o_2	0.0374	0	0.0200	0.0436
o_3	0.0500	0.0200	0	0.0346
o_4	0.0424	0.0436	0.0346	0

Apskatīsim otro variantu kvadrātisko attālumu χ^2 aprēķināšanai. Tāpat kā pirmajā variantā, tiek aprēķinātas atribūtu vērtību summas datu tabulas rindās un kolonās: $a_i^* = \sum_{j=1}^n a_{ij}$, $i = 1, \dots, m$, $a_j^* = \sum_{i=1}^m a_{ij}$, $j = 1, \dots, n$. Pamatojoties uz aprēķinātajām a_i^* vērtībām, tiek veidoti atribūtu vērtību profili.

Atšķirība starp variantiem slēpjas svaru matricas \mathbf{W} veidošanas principā. Otrajā variantā tiek aprēķināta visu atribūtu vērtību summa tabulā: $a_{ij}^* = \sum_{i=1}^m \sum_{j=1}^n a_{ij}$. Pēc tam atribūtu smaguma centru vērtības aprēķina kā attiecību

$$c_j = \frac{a_j^*}{a_{ij}^*}, \quad j = 1, \dots, n. \quad (\text{P2.5.4})$$

Diagonālās svaru matricas \mathbf{W} komponentu vērtības tiek aprēķinātas kā

$$w_j = \frac{1}{c_j}. \quad (\text{P2.5.5})$$

Tādējādi svaru matrica ir veidā

$$\mathbf{W} = \begin{bmatrix} \frac{1}{c_1} & \dots & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & \frac{1}{c_j} & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & 0 & \dots & 0 & \dots & \frac{1}{c_n} \end{bmatrix}.$$

Kvadrātisko attālumu χ_2^2 noteikšanai tiek izmantots vienādojums

$$\chi_2^2(o_i - o_l) = \sum_{j=1}^n w_j (p_{ij} - p_{lj})^2, \quad i, l = 1, \dots, m, \quad l \neq k, \quad (\text{P2.5.6})$$

kur p_{ij} , p_{lj} ir atribūtu a_{ij} , a_{lj} vērtību profili;

w_j tiek aprēķināts, izmantojot vienādojumu (P2.5.5).

Piemērs P2.5.2. Ņemot P2.5.1 piemēra sākotnējos datus, ir jāaprēķina kvadrātisko un standarta attālumu χ_2^2 , χ_2 vērtības starp objektiem (atribūtu vērtību vektoriem), izmantojot otro variantu atribūtu svaru matricas formēšanai.

Risinājums. Tabulas P2.5.4 kreisajā daļā ir sākotnējie dati, bet labajā pusē ir atribūtu vērtību profili.

Tabula P2.5.4

Sākotnējie dati un provizorisko aprēķinu rezultāti piemēram P2.5.2

	Sākotnējās atribūtu vērtības				a_i^*	Atribūtu vērtību profili			
	a_{i1}	a_{i2}	a_{i3}	a_{i4}		P_{i1}	P_{i2}	P_{i3}	P_{i4}
o_1	7	3	2	10	22	0.318	0.136	0.091	0.455
o_2	9	2	5	6	22	0.409	0.091	0.228	0.273
o_3	12	5	8	8	33	0.364	0.152	0.242	0.242
o_4	5	4	6	9	24	0.208	0.167	0.250	0.375
a_j^*	33	14	21	33	101				
c_j	0.327	0.138	0.208	0.327					
$w_j = \frac{1}{c_j}$	3.058	7.246	4.808	3.058					

Skaitlis $a_{ij}^* = 101$ rindas a_i^* un kolonas a_j^* krustpunktā ir visu tabulas atribūtu vērtību summa. Vērtības c_j , $j=1, \dots, 4$, aprēķina, dalot atbilstošo vērtību a_j^* ar vērtību $a_{ij}^* = 101$. Vērtības w_j ir vērtību c_j apgrieztās vērtības.

Veiksim aprēķinus, izmantojot vienādojumu (P2.5.6). Nav nepieciešams veikt aprēķinus pilnā mērā. Izmantosim aprēķinu rezultātus no piemēra P2.5.1. Mums ir

$$\chi_2^2(o_1 - o_2) = 3.058 * 0.00828 + 7.246 * 0.00202 + 4.808 * 0.01877 + 3.058 * 0.03312 = 0.02532 + 0.01464 + 0.09025 + 0.10128 \approx 0.2315;$$

$$\chi_2^2(o_1 - o_3) = 3.058 * 0.00212 + 7.246 * 0.00026 + 4.808 * 0.02280 + 3.058 * 0.04537 = 0.00648 + 0.00188 + 0.10962 + 0.13874 \approx 0.2567.$$

Pārējie aprēķini tiek veikti līdzīgi. Aprēķinu rezultāti ir parādīti tabulā P2.5.5.

Lai aprēķinātu standarta attālumus χ_2 saskaņā ar otro variantu, vienādojums (P2.5.6) ir jāpārvērš šādā formā

$$\chi_2(o_i - o_l) = \sqrt{\sum_{j=1}^n w_j (p_{ij} - p_{lj})^2} \dots i, l = 1, \dots, m, i \neq l. \quad (\text{P2.5.7})$$

Nav nepieciešamības veikt aprēķinus, izmantojot vienādojumu (P2.5.7). Lai iegūtu standarta attālumus χ_2 šajā piemērā, pietiek ņemt kvadrātsaknes no tabulā P2.5.5 uzrādītajām $\chi_2^2(o_i - o_l)$ vērtībām. Galīgie rezultāti ir parādīti tabulā P2.5.6.

Tabula P2.5.5

Kvadrātisko attālumu χ_2^2 aprēķinu rezultāti piemēram P2.5.2

	o_1	o_2	o_3	o_4
o_1	0	0.2315	0.2567	0.1877
o_2	0.2315	0	0.0371	0.1995
o_3	0.2567	0.0371	0	0.1167
o_4	0.1877	0.1995	0.1167	0

Tabula P2.5.6

Standarta attālumu χ_2 aprēķinu rezultāti piemēram P2.5.2

	o_1	o_2	o_3	o_4
o_1	0	0.4811	0.5067	0.4532
o_2	0.4811	0	0.1926	0.4466
o_3	0.5067	0.1926	0	0.3416
o_4	0.4332	0.4466	0.3416	0

Kurš no diviem iepriekš parādītajiem variantiem ir labāks? Atbilde uz šo jautājumu ir nepārprotama: abi varianti ir līdzvērtīgi. Aprēķinu apjoms abos variantos ir vienāds. Otrā varianta neliela priekšrocība ir tā, ka iegūtajiem attāluma aprēķiniem ir ievērojami lielākas absolūtās vērtības, tāpēc aprēķinu rezultātu noapaļošana mazāk ietekmē šīs aplēses.

Noslēdzot šo sadaļu, mēs atzīmēsim, ka vērtības c_j tiek sadalītas saskaņā ar χ^2 likumu, tāpēc arī ir tāds attāluma aprēķinu nosaukums.

P2.6. Mahalanobisa distance

Iepriekšējā sadaļā tika apskatītas divas Eiklīda attāluma versijas, kurās tika izmantoti atribūtu svāri. Šo svāru vērtības tika aprēķinātas īpašā veidā, pamatojoties uz atribūtu vērtību kombinācijām.

Šajā sadaļā ir apskatīts cits paplašinātā Eiklīda attāluma variants. Šo versiju ierosināja indiešu zinātnieks P. Mahalanobis savā darbā [Mahalanobis P.S., 1936]. Autors pētīja līdzības un atšķirības starp dažādām iedzīvotāju grupām Indijā un nonāca pie secinājuma, lai aprēķinātu attālumus starp attiecīgajām grupām daudzdimensiju datus, ir jāņem vērā atribūtu vērtību izkliedes pakāpe (variācija) katrā no grupām un kovariācija starp atsevišķiem atribūtu vērtību vektoriem katram grupu pārim. Šajā sadaļā mēs noteiksim Mahalanobisa attālumu, pamatojoties tikai uz atribūtu vērtību izmaiņām.

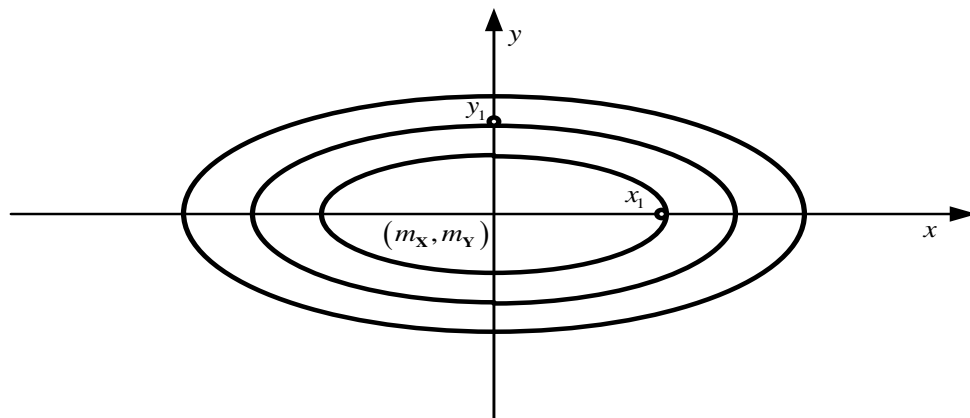
Lai iepazīstinātu ar Mahalanobisa attāluma definīcijas vispārējo ideju, apskatīsim šādu piemēru. Pieņemam, ka mums ir divdimensiju mainīgais $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$. Pamatojoties uz sākotnējo statistisko materiālu, tiek noteikta šī mainīgā lieluma vērtību sadalījuma funkcija $f(\mathbf{X}, \mathbf{Y})$. Attēlā P2.6.1 parādīta šādas nosacītas sadalījuma funkcijas projekcijas mainīgo sākotnējo vērtību telpā \mathbf{X} un \mathbf{Y} . Koordinātu sākuma punkts ir (m_x, m_y) , kur m_x, m_y ir mainīgo \mathbf{X} un \mathbf{Y} vidējās vērtības.

Katra ovāla līnija šajā attēlā apzīmē divdimensiju sadalījuma blīvuma funkcijas vērtību kopu $f(\mathbf{X}, \mathbf{Y})$.

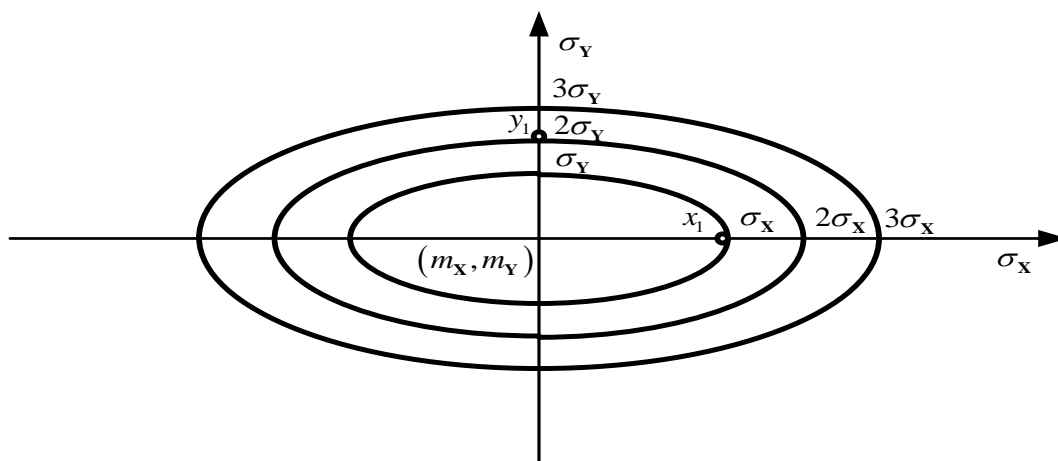
Attēlā P2.6.2 parādītas tās pašas projekcijas, bet koordinātu sistēmā, kuras sastāvdaļas ir standartnovirzes σ_x, σ_y .

Apskatīsim divus punktus x_1, y_1 , kas parādīti attēlā P2.6.1. Acīmredzams, ka koordinātu sistēmā x, y punkts y_1 ir tuvāk koordinātu centram nekā punkts x_1 . Pārnēsīsim šos punktus uz koordinātu sistēmu σ_x, σ_y attēlā P2.6.2. Jaunajā koordinātu sistēmā punktam x_1 ir koordinātes $(\sigma_x, 0)$, punktam y_1 ir koordinātes $(0, 2\sigma_y)$.

Ja par attālumu mērīšanas mērvienībām ņemam vērtības σ_x, σ_y , redzams, ka punkts x_1 ir tuvāk koordinātu sākumam nekā punkts y_1 . No tā izriet, ka standarta novirzes σ (un attiecīgi variācijas σ^2) var kalpot kā specifiski attāluma mēri, kuru pamatā ir gadījuma lielumu vērtību izkliedes pakāpe attiecībā pret noteiktiem centrālajiem punktiem. Šī ideja arī ir Mahalanobisa attāluma definīcijas pamatā.



Attēls P2.6.1. Divdimensiju sadalījuma blīvuma funkcijas $f(\mathbf{X}, \mathbf{Y})$ projekciju grafiskais attēlojums vērtību \mathbf{X} un \mathbf{Y} telpā



Attēls P2.6.2. Funkciju $f(\mathbf{X}, \mathbf{Y})$ projekciju grafiskais attēlojums koordinātu sistēmā σ_x, σ_y

Katra atribūta vērtības var korekti interpretēt kā gadījuma mainīgo realizāciju objektu kopā. Izmantojot atribūtu vērtību variācijas attiecībā pret to vidējām vērtībām kā atribūtu svarus, mēs varam formulēt šādu vispārīgu vienādojumu kvadrātiskajam attālumam starp objektiem

$$d^2(o_i - o_l) = (\mathbf{a}_{ij} - \mathbf{a}_{lj})^T \mathbf{W} (\mathbf{a}_{ij} - \mathbf{a}_{lj}), \quad (\text{P2.6.1})$$

kur $\mathbf{a}_{ij}, \mathbf{a}_{lj}$ ir atribūtu vērtību vektori objektiem o_i, o_k (šo vektoru vērtības parādītas sākotnējo datu tabulas rindās o_k, o_i ;

T - matricas transponēšanas simbols;

\mathbf{W} - svaru diagonālā matrica, kuras elementi ir atbilstošo atribūtu variāciju apgrieztās vērtības.

Datu tabulai, kas sastāv no n kolonām, matricai \mathbf{W} ir forma

$$\mathbf{W} = \begin{pmatrix} \frac{1}{s_1^2} & \dots & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & \frac{1}{s_j^2} & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \dots & 0 & \dots & 0 & \dots & \frac{1}{s_n^2} \end{pmatrix} \cdot \cdot$$

Praktisku aprēķinu veikšanai vispārīgo vienādojumu (P2.6.1) var uzrādīt formā

$$d^2(o_i - o_l) \sum_{j=1}^n \frac{1}{s_j^2} (a_{ij} - a_{lj})^2, \quad i, l = 1, \dots, m, \quad i \neq l, \quad (\text{P2.6.2})$$

kur s_j^2 ir atribūta j vērtību variācijas novērtējums;

a_{ij}, a_{lj} - j -tā vērtības objektiem o_i, o_l .

(Jāņem vērā, ka vienādojums (P2.6.2) tiek izmantots, lai aprēķinātu kvadrātiskos Mahalanobisa attālumus, neņemot vērā kovariācijas starp atsevišķu atribūtu vērtībām).

Apskatīsim piemēru, kas ilustrē Mahalanobisa attālumu starp objektiem aprēķinu.

Piemērs P2.6.1. Ņemsim par pamatu sākotnējos datus no tabulas 4.2.1. Šie dati ir reproducēti tabulā P2.6.1. Nepieciešams aprēķināt kvadrātiskos Mahalanobisa attālumus starp visiem objektiem, izmantojot vienādojumu (P2.6.2).

Risinājums. Aprēķināsim atribūtu vidējās vērtības tabulas P2.6.1 kolonās pēc vienādojuma

$$m_j = \frac{\sum_{i=1}^4 a_{ij}}{4}, \quad j = 1, \dots, 4.$$

Tabula P2.6.1

Sākotnējie dati piemēram P2.6.1

	a_{i1}	a_{i2}	a_{i3}	a_{i4}
o_1	7	3	2	10
o_2	9	2	5	6
o_3	12	5	8	8
o_4	5	4	6	9

$$m_1 = \frac{7+9+12+5}{4} = \frac{33}{4} = 8.25; \quad m_2 = \frac{3+2+5+4}{4} = \frac{14}{4} = 3.50;$$

$$m_3 = \frac{2+5+8+6}{4} = \frac{21}{4} = 5.25; \quad m_4 = \frac{10+6+8+9}{4} = \frac{33}{4} = 8.25.$$

Variāciju vērtības tabulas kolonās tiek aprēķinātas, izmantojot vienādojumu

$$s_j^2 = \frac{\sum_{i=1}^4 (a_{ij} - m_j)^2}{3}.$$

$$\begin{aligned} s_1^2 &= \frac{1}{3} * \left((7-8.25)^2 + (9-8.25)^2 + (12-8.25)^2 + (5-8.25)^2 \right) = \\ &= \frac{1}{3} (1.5625 + 0.5625 + 14.0625 + 10.5625) = \frac{1}{3} * 26.2675 = 8.9167; \end{aligned}$$

$$\begin{aligned} s_2^2 &= \frac{1}{3} * \left((3-3.50)^2 + (2-3.50)^2 + (5-3.50)^2 + (4-3.50)^2 \right) = \\ &= \frac{1}{3} (0.2500 + 2.2500 + 2.2500 + 0.2500) = \frac{1}{3} * 5.0000 = 1.6667; \end{aligned}$$

$$\begin{aligned} s_3^2 &= \frac{1}{3} * \left((2-5.25)^2 + (5-5.25)^2 + (8-5.25)^2 + (6-5.25)^2 \right) = \\ &= \frac{1}{3} * (10.5625 + 0.0625 + 7.5625 + 0.5625) = \frac{1}{3} * 18.750 = 6.2500; \end{aligned}$$

$$\begin{aligned} s_4^2 &= \frac{1}{3} * \left((10-8.25)^2 + (6-8.25)^2 + (8-8.25)^2 + (9-8.25)^2 \right) = \\ &= \frac{1}{3} * (3.0625 + 5.0625 + 0.0625 + 0.5625) = \frac{1}{3} * 8.7500 = 2.9167. \end{aligned}$$

Aprēķināsim s_j^2 apgrieztās vērtības, $j = 1, 2, 3, 4$.

$$\frac{1}{s_1^2} = \frac{1}{8.9167} = 0.1121; \quad \frac{1}{s_2^2} = \frac{1}{1.6667} = 0.6000;$$

$$\frac{1}{s_3^2} = \frac{1}{6.2500} = 0.1600; \quad \frac{1}{s_4^2} = \frac{1}{2.9167} = 0.3429.$$

Veiksim nepieciešamos aprēķinus, izmantojot vienādojumu (P2.6.2).

$$d^2(o_1 - o_2) = 0.1121 * (7-9)^2 + 0.6000 * (3-2)^2 + 0.1600 * (2-5)^2 + 0.3429 * (10-6)^2 =$$

$$= 0.1121 * 4 + 0.6000 * 1 + 0.1600 * 9 + 0.3429 * 16 =$$

$$= 0.4484 + 0.6000 + 1.4400 + 5.4864 \approx 7.975;$$

$$d^2(o_1 - o_3) = 0.1121 * (7-12)^2 + 0.6000 * (3-5)^2 + 0.1600 * (2-8)^2 + 0.3429 * (10-8)^2 =$$

$$= 0.1121 * 25 + 0.6000 * 4 + 0.1600 * 36 + 0.3429 * 4 =$$

$$= 2.8025 + 2.4000 + 5.7600 + 1.3716 \approx 12.234;$$

$$d^2(o_1 - o_4) = 0.1121 * (7-5)^2 + 0.6000 * (3-4)^2 + 0.1600 * (2-6)^2 + 0.3429 * (10-9)^2 =$$

$$= 0.1121 * 4 + 0.6000 * 1 + 0.1600 * 16 + 0.3429 * 1 =$$

$$= 0.4484 + 0.6000 + 2.5600 + 0.3429 \approx 3.951;$$

$$d^2(o_2 - o_3) = 0.1121 * (9-12)^2 + 0.6000 * (2-5)^2 + 0.1600 * (5-8)^2 + 0.3429 * (6-8)^2 =$$

$$= 0.1121 * 9 + 0.6000 * 9 + 0.1600 * 9 + 0.3429 * 4 =$$

$$= 1.0089 + 5.4000 + 1.4400 + 1.3716 \approx 9.220;$$

$$d^2(o_2 - o_4) = 0.1121 * (9-5)^2 + 0.6000 * (2-4)^2 + 0.1600 * (5-6)^2 + 0.3429 * (6-9)^2 =$$

$$= 0.1121 * 16 + 0.6000 * 4 + 0.1600 * 1 + 0.3429 * 9 =$$

$$= 1.7936 + 2.4000 + 0.1600 + 3.0861 \approx 7.440;$$

$$d^2(o_3 - o_4) = 0.1121 * (12-5)^2 + 0.6000 * (5-4)^2 + 0.1600 * (8-6)^2 + 0.3429 * (8-9)^2 =$$

$$= 0.1121 * 49 + 0.6000 * 1 + 0.1600 * 4 + 0.3429 * 1 =$$

$$= 5.4929 + 0.6000 + 0.6400 + 0.3429 \approx 7.076.$$

Iegūtie rezultāti ir apkopoti tabulā P2.6.2.

Lai aprēķinātu standarta Mahalanobisa attālumu, vienādojums (P2.6.2) ir jāpārvērš formā

$$d(o_i - o_l) = \sqrt{\sum_{j=1}^n \frac{1}{s_j^2} (a_{ij} - a_{lj})^2} \dots i, l = 1, \dots, m, i \neq l. \quad (\text{P2.6.3})$$

Tabula P2.6.2

Kvadrātisko Mahalanobisa attālumu aprēķinu rezultāti piemēram P2.6.1

	o_1	o_2	o_3	o_4
o_1	0	7.975	12.334	3.951
o_2	7.975	0	9.220	7.440
o_3	12.334	9.220	0	7.076
o_4	3.951	7.440	7.076	0

Piemērs P2.6.2. Izmantojot P2.6.1 piemērā iegūtos rezultātus, ir jāaprēķina standarta Mahalanobisa attālumi starp visiem objektu pāriem.

Risinājums. Šajā piemērā nav jāveic pilni aprēķini. No vienādojuma (P2.6.3) analīzes izriet, ka, lai aprēķinātu standarta Mahalanobisa attālumu starp objektu pāriem, pietiek aprēķināt kvadrātsaknes vērtību no atbilstošās kvadrātiskās Mahalanobisa attāluma vērtības tabulā P2.6.2. Aprēķinu rezultāti ir parādīti tabulā P2.6.3.

Tabula P2.6.3

Standarta Mahalanobisa attālumu aprēķinu rezultāti piemēram P2.6.2

	o_1	o_2	o_3	o_4
o_1	0	2.824	3.512	1.988
o_2	2.824	0	3.036	2.728
o_3	3.512	3.036	0	2.660
o_4	1.988	2.728	2.660	0

P3. NORMĀLĀ SADALĪJUMA LIKUMA PĀRBAUDE

6. nodaļā tika apskatītas datu transformācijas. Viens no šādu transformāciju mērķiem ir pārveidot uz labo vai kreiso pusi vērstus sadalījumus normālos vai gandrīz normālos sadalījumos.

Pirms atribūtu vērtību transformācijas ir jāpārlicinās, kā šīs vērtības tiek sadalītas. Šeit ir divas iespējas: (1) attiecīgais sadalījums ir normāls vai tuvu normālam un nav nepieciešama atribūtu vērtību transformācija; (2) attiecīgais sadalījums nav normāls un ir nepieciešama tā transformācija.

Var atzīmēt, ka daudzas cilvēku populācijas raksturojošas īpašības pēc savas būtības ir sadalītas pēc normālā sadalījuma likuma (svars, augums, vecums, IQ koeficients).

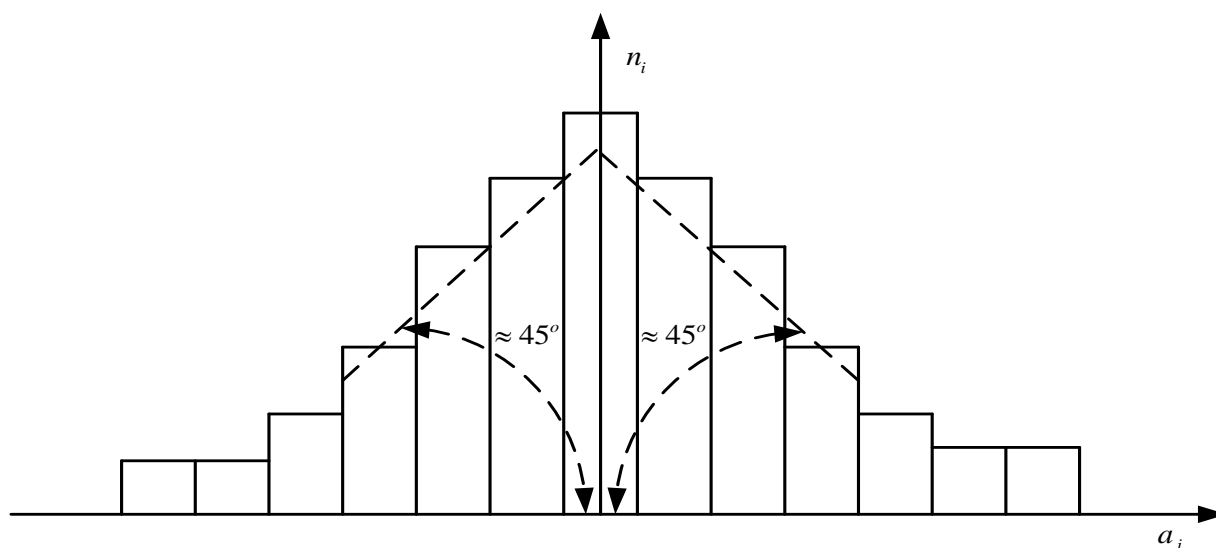
Šeit rodas vēl viens jautājums. Pieņemam, ka reģiona iedzīvotāju vidū šo atribūtu vērtības patiešām ir sadalītas saskaņā ar normālo likumu. Bet socioloģiskajos un citos pētījumos tiek izmantoti attiecīgās populācijas izlases dati. Šādos gadījumos attiecīgo atribūtu sadalījuma likumi var lielākā vai mazākā mērā atšķirties no normālā sadalījuma. Turklāt sadalījuma parametri dažādiem paraugiem var atšķirties. Šī problēma ir saistīta ar izlases datu reprezentativitāti un ir jāatrisina pētījuma organizēšanas procesā.

Šajā pielikumā ir sniegtas dažādas pieejas normālā sadalījuma pārbaudei un dažādas aplēses, kas raksturo konkrētā sadalījumu novirzes pakāpi no normālā sadalījuma prasībām. Ērtības labad turpmāk bieži lietosim terminu *normalitāte*, kas arī raksturo sadalījuma atbilstību normālajam sadalījuma likumam.

P3.1. Vienkāršākās metodes sadalījumu pārbaudei uz normalitāti

Vienkāršākais veids, kā novērtēt sadalījuma iespējamo normalitāti, ir grafiski attēlot atribūtu vērtību kopu histogrammas veidā. Viens šāds nosacīts attēlojums ir parādīts attēlā P3.1.1.

Sadalījuma histogramma tiek veidota šādi. Atribūtu a_j vērtību diapazons ir sadalīts noteiktā intervālu skaitā. Šie intervāli ir parādīti attēlā uz horizontālās ass. Tālāk tiek uzskaitīts atribūtu vērtību skaits, kas ietilpst katrā intervālā. Šis skaitlis tiek uzdots kā kolonas augstums, kas attēlots katrā intervālā.

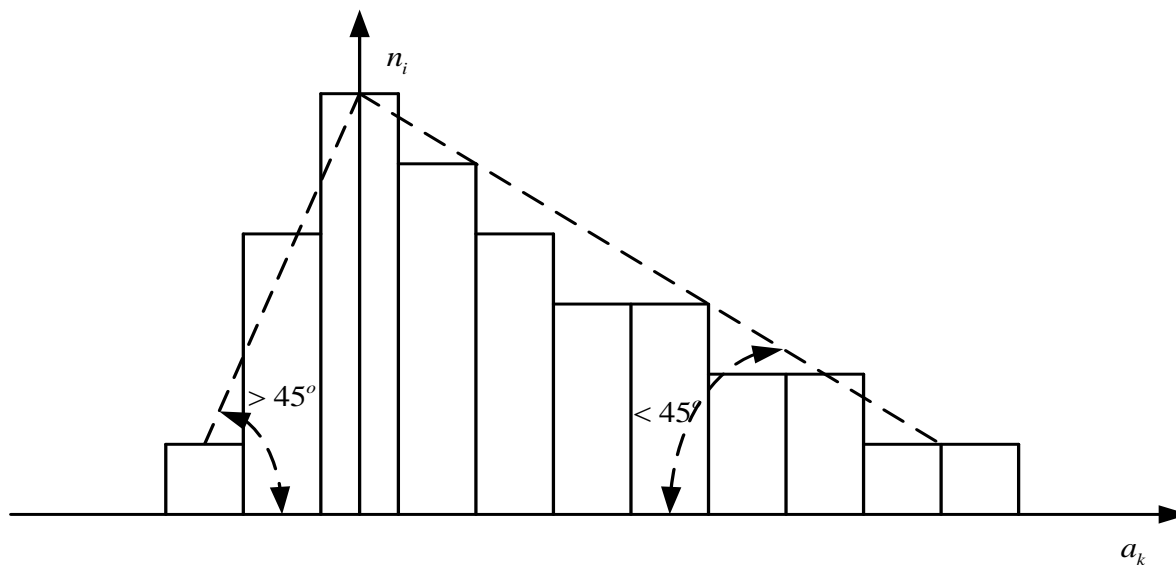


Attēls P3.1.1. Atribūtu a_j vērtību nosacītā sadalījuma histogramma

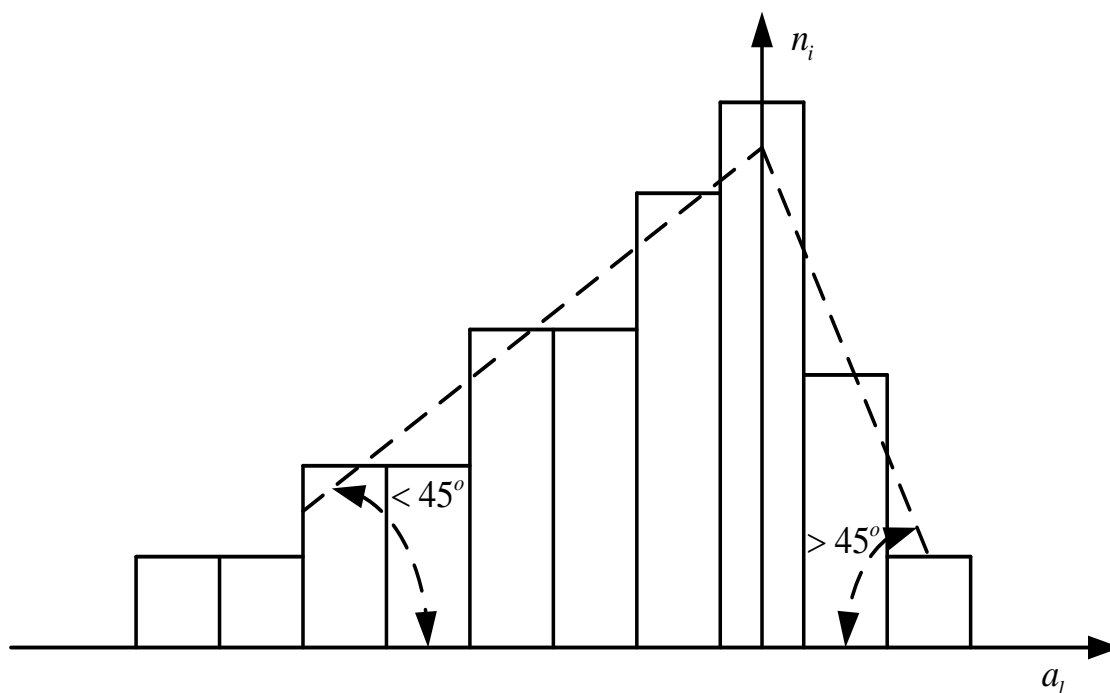
Kādus secinājumus var izdarīt no histogrammas analīzes attēlā P3.1.1? Pirmkārt, atribūtu a_j vērtības ir simetriski sadalītas abās pusēs no centrālās vērtības. Otrkārt, taisnajām līnijām, kas atspoguļo vispārējās tendences atribūtu vērtību izmaiņās to vidējo vērtību zonās, ir aptuveni 45°

slīpums attiecībā pret horizontālo asi. Šīs abas pazīmes norāda uz to, ka atribūtu vērtības a_j var tikt sadalītas pēc normālā sadalījuma likuma.

Attēlos P3.1.2 un P3.1.3 parādītas atribūtu a_k , a_l vērtību sadalījuma histogrammas.



Attēls P3.1.2. Atribūtu a_k vērtību nosacītā sadalījuma histogramma



Attēls P3.1.3. Atribūtu a_l vērtību nosacītā sadalījuma histogramma

Attēlā P3.1.2 redzams, ka atribūtu a_k vērtību sadalījums nav simetrisks. Datu punktu skaits sadalījuma labajā pusē ievērojami pārsniedz datu punktu skaitu tā kreisajā pusē. Līniju slīpumi, kas raksturo vispārējās atribūtu vērtību izmaiņu tendences, ir ļoti atšķirīgi. Līnijas slīpums sadalījuma kreisajā pusē ir ievērojami lielāks par 45^0 , bet līnijas slīpums sadalījuma labajā pusē ir ievērojami mazāks par 45^0 .

Attēlā P3.1.3 redzams, ka atribūtu a_l vērtību sadalījums arī nav simetrisks. Datu punktu skaits sadalījuma kreisajā pusē ievērojami pārsniedz datu punktu skaitu labajā pusē. Līnijas slīpums

sadalījuma kreisajā pusē ir ievērojami mazāks par 45° , savukārt līnijas slīpums sadalījuma labajā pusē ir ievērojami lielāks par 45° .

Histogramma attēlā P3.1.2 parāda, ka atribūtu a_k vērtību sadalījums ir tipisks uz labo pusi vērsts sadalījums, histogramma attēlā P3.1.3 parāda, ka atribūtu a_l vērtību sadalījums ir tipisks uz kreiso pusi vērsts sadalījums.

Mēs sakām, ka atribūtu a_j vērtību sadalījums attēlā P3.1.1 var būt normāls tādā nozīmē, ka tas potenciāli var būt normāls, taču ir jāveic turpmāki novērtējumi, lai apstiprinātu šo faktu. Atribūtu vērtību a_k, a_l sadalījumi attēlos P3.1.2, P3.1.3 nevar būt normāli un nav nepieciešams cits šī fakta apstiprinājums.

Tādējādi vienkāršāko metožu mērķi sadalījuma pārbaudei uz normalitāti ir šādi: (1) noteikt, ka attiecīgais sadalījums var būt normāls un turpināt to pārbaudīt; (2) pārliecināties, ka attiecīgais sadalījums nav normāls un izstrādāt stratēģiju tā transformācijai.

P3.2. Sadalījuma normalitātes pārbaude pamatojoties uz sadalījuma parametriem

Pieņem, ka visiem objektiem var norādīt sākotnējo atribūtu a_j vērtību kopu. Tad mēs vienmēr varam aprēķināt tā vērtību sadalījuma parametrus.

- *sadalījuma izvietojuma parametrs* – sagaidāmā vērtība

$$\bar{a}_j = \sum_{i=1}^n a_{ij} p_i, \quad (\text{P3.2.1})$$

kur n - intervālu skaits, kas veidojas atribūtu a_j vērtību izmaiņu diapazonā;

a_{ij} - atribūta a_j vidējā vērtība i -tajā intervālā $i = 1, \dots, n$;

p_i - varbūtība, ka atribūta a_j vērtība nonāks i -tajā intervālā

$$p_i = \frac{n_i}{n}, \quad (\text{P3.2.2})$$

kur n_i - atribūtu a_j vērtību skaits i -tajā intervālā;

n - kopējais atribūtu a_j vērtību skaits.

- *sadalījuma skalas parametrs* – atribūtu a_j vērtību standartnovirze no sagaidāmās vērtības

$$s_j = \sqrt{\sum_{i=1}^n (a_{ij} - \bar{a}_j)^2 p_i}. \quad (\text{P3.2.3})$$

Apskatīsim dažus ilustratīvus piemērus.

Piemērs P3.2.1. Tabulā P3.2.1 ir parādīta sākotnējā atribūtu a_j vērtību kopa un dati par to vērtību sadalījumu. Ir nepieciešams novērtēt šī sadalījuma izvietojuma parametrus un novērtēt iespēju, ka sadalījums var būt normāls.

Tabulas pirmajā rindā ir norādīti atribūtu a_j vērtību intervālu numuri. Šajā piemērā tiek pieņemts, ka katra intervāla garums ir vienāds ar vienu atribūtu a_j vērtību vienību. Otrajā tabulas rindā ir norādītas katra intervāla robežas.

Katrā intervālā tiek pieņemts, ka visas atribūtu a_j vērtības, kas ietilpst šajā intervālā, tiek aizstātas ar atribūta vidējo vērtību šajā intervālā. Tabulas trešajā rindā ir parādītas vidējās atribūtu a_j vērtības katrā intervālā.

Tabula P3.2.1

Sākotnējo atribūtu a_j vērtību kopa un vērtību sadalījums

Intervāla numurs	1	2	3	4	5	6	7	8	9	10	11	12	13
Intervālu robežas	0.5-1.5	1.5-2.5	2.5-3.5	3.5-4.5	4.5-5.5	5.5-6.5	6.5-7.5	7.5-8.5	8.5-9.5	9.5-10.5	10.5-11.5	11.5-12.5	12.5-13.5
Vidējā vērtība intervālā	1	2	3	4	5	6	7	8	9	10	11	12	13
Vērtību skaits intervālā	1	1	2	3	4	6	8	6	4	3	2	1	1
Varbūtība, ka atribūta vērtība nonāks intervālā	0.024	0.024	0.048	0.071	0.095	0.143	0.190	0.143	0.095	0.071	0.048	0.024	0.024

Tabulas ceturtajā rindā ir parādīts atribūtu vērtību skaits katrā no intervāliem.

Tā kā šajā piemērā mēs runājam par atribūtu vērtību sadalījumu pa intervāliem, mums ir nepieciešamas atribūtu vērtību a_j varbūtības (biežumi) katrā intervālā. Aprēķināsim tabulas ceturtais rindas skaitļu summu. Tas ir vienāds ar 42. Lai aprēķinātu mūs interesējošās varbūtības, katram intervālam ir jādala tabulas 4. rindas atbilstošajā šūnā esošais skaitlis ar 42. Iegūtie varbūtības aprēķini ir parādīti pēdējā tabulas rindā.

Atribūtu vērtību a_j sadalījuma histogramma ir parādīta attēlā P3.2.1.

Aprēķināsim sadalījuma parametru vērtības:

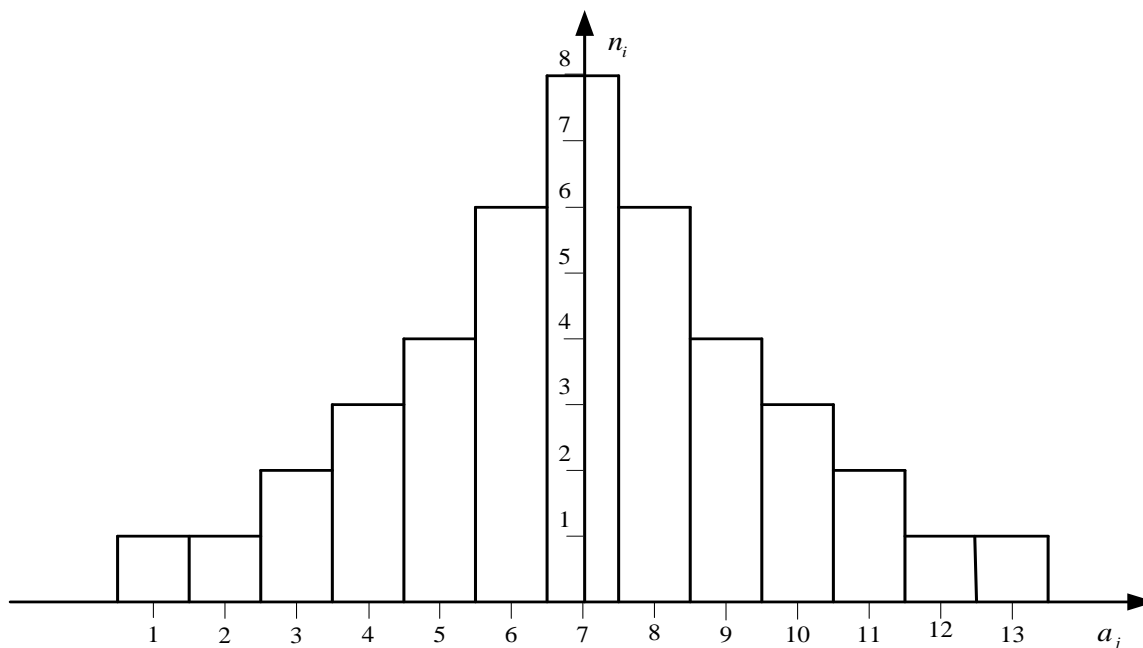
$$\begin{aligned} \bar{a}_j &= 1*0.024 + 2*0.024 + 3*0.048 + 4*0.071 + 5*0.095 + 6*0.143 + 7*0.190 + \\ &+ 8*0.143 + 9*0.095 + 10*0.071 + 11*0.048 + 12*0.024 + 13*0.024 = \\ &= 0.024 + 0.048 + 0.144 + 0.284 + 0.476 + 0.858 + 1.330 + 1.144 + \\ &+ 0.855 + 0.710 + 0.528 + 0.288 + 0.312 = 6.926 \approx 7. \end{aligned}$$

Lai aprēķinātu atribūtu a_j vērtību standartnovirzi, vispirms mēs aprēķinām šī sadalījuma izkliedi (variāciju), izmantojot vienādojumu

$$D_j = \sum_{i=1}^n (a_{ij} - \bar{a}_j)^2 p_i. \quad (\text{P3.2.4})$$

$$\begin{aligned} D_j &= (1-7)^2 * 0.024 + (2-7)^2 * 0.024 + (3-7)^2 * 0.048 + (4-7)^2 * 0.071 + (5-7)^2 * 0.095 + \\ &+ (6-7)^2 * 0.143 + (7-7)^2 * 0.190 + (8-7)^2 * 0.143 + (9-7)^2 * 0.095 + (10-7)^2 * 0.071 + \\ &+ (11-7)^2 * 0.048 + (12-7)^2 * 0.024 + (13-7)^2 * 0.024 = \\ &= 0.864 + 0.600 + 0.768 + 0.639 + 0.380 + 0.143 + 0 + 0.143 + 0.380 + \\ &+ 0.639 + 0.768 + 0.600 + 0.864 = 6.788. \end{aligned}$$

$$s_j = \sqrt{D_j} = \sqrt{6.788} = 2.60.$$



Attēls P3.2.1 Atribūtu a_j vērtību sadalījuma histogramma datiem no tabulas P3.2.1

Kādā veidā, pamatojoties uz sadalījuma parametriem, var izdarīt secinājumus par sadalījuma normalitāti?

Jebkura simetriska sadalījuma raksturīga iezīme ir tā, ka tā sagaidāmā vērtība, moda un mediāna sakrīt. Sagaidāmā vērtība tika aprēķināta iepriekš: $\bar{a}_j = 7$.

Aprēķināsim šī sadalījuma mediānu. Lai to izdarītu, mēs sakārtojam visas atribūtu a_j vērtības nedilstošā secībā:

1, 2, 3, 3, 3, 4, 4, 4, 5, 5, 5, 5, 6, 6, 6, 6, 6, 6, 7, 7, 7, 7, 7, 7, 7, 7, 7, 7, 8, 8, 8, 8, 8, 8, 9, 9, 9, 9, 10, 10, 10, 11, 11, 12, 13

Kopumā mums ir 42 atribūtu a_j vērtības. Tā kā tas ir pārskaitlis, sadalījuma mediānu aprēķina kā vidējo vērtību no skaitļiem 21. un 22. pozīcijā iepriekš norādītajā vērtību secībā. Šajās pozīcijās mums ir atribūtu a_j vērtības, kas vienādas ar 7. Tāpēc

$$\text{med}(a_j) = \frac{7+7}{2} = 7.$$

Moda ir atribūta a_j vērtība, kurai datu punktu skaits pa kreisi no šīs vērtības ir vienāds ar datu punktu skaitu pa labi no šīs vērtības. Acīmredzot mūsu piemērā $\text{mod}(a_j) = 7$.

Tā kā dažādo sadalījuma parametru vērtības sakrīt, var pamatoti secināt, ka atribūtu a_j vērtību sadalījums ir simetrisks sadalījums. Tas nozīmē, ka šis sadalījums var būt (bet ne vienmēr ir) normāls sadalījums.

Ja novērtējamais sadalījums patiešām ir normāls sadalījums, ir jāievēro šādas prasības:

- 68% atribūtu vērtību atrodas standarta novirzes robežās no sagaidāmās vērtības;
- 95% atribūtu vērtību atrodas divu standartnoviržu robežās no sagaidāmās vērtības;
- 99,7% atribūtu vērtību atrodas trīs standartnoviržu robežās no sagaidāmās vērtības.

Mūsu piemērā $\bar{a}_j = 7$, $s_j = 2.60$. Tas nozīmē, ka intervālā $\pm s_j$ no sagaidāmās vērtības ir intervāli 5, 6, 7, 8, 9. Šajā intervālā ir $4 + 6 + 8 + 6 + 4 = 28$ atribūtu a_j vērtības jeb 67% no visām tā vērtībām.

Intervālā $\pm 2s_j$ ir 38 atribūtu a_j vērtības jeb 91% no visām tā vērtībām.

Intervāls $\pm 3s_j$ satur 100% atribūtu a_j vērtību.

Pamatojoties uz šiem rezultātiem, var apgalvot, ka atribūtu a_j vērtību sadalījums no tabulas P3.2.1 tuvināti ir normāls sadalījums.

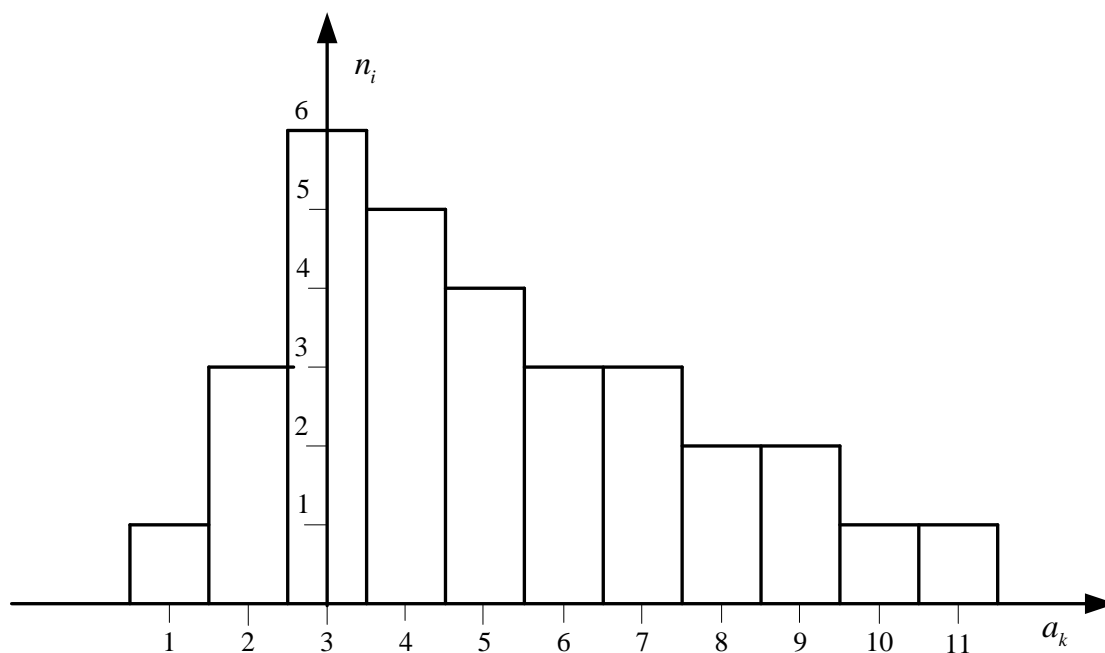
Piemērs P3.2.2. Tabulā P3.2.2 ir parādītas dažādas sākotnējās atribūta a_k vērtības un citi dati par šī atribūta vērtību sadalījumu. Šīs tabulas datu interpretācija ir līdzīga tabulas P3.2.1 datu interpretācijai.

Tabula P3.2.2

Sākotnējā atribūta a_k vērtību kopa un vērtību sadalījums

Intervāla numurs	1	2	3	4	5	6	7	8	9	10	11
Intervālu robežas	0.5-1.5	1.5-2.5	2.5-3.5	3.5-4.5	4.5-5.5	5.5-6.5	6.5-7.5	7.5-8.5	8.5-9.5	9.5-10.5	10.5-11.5
Vidējā vērtība intervālā	1	2	3	4	5	6	7	8	9	10	11
Vērtību skaits intervālā	1	3	6	5	4	3	3	2	2	1	1
Varbūtība, ka atribūta vērtība nonāks intervālā	0.032	0.098	0.194	0.161	0.128	0.098	0.098	0.064	0.064	0.32	0.032

Atribūtu a_k vērtību sadalījuma histogramma ir parādīta attēlā P3.2.2.



Attēls P3.2.2. Atribūtu a_k vērtību sadalījuma histogramma datiem no tabulas P3.2.2

Aprēķināsim atribūtu a_k vērtību sadalījuma parametrus.

$$\begin{aligned}\bar{a}_k &= 1*0.032 + 2*0.098 + 3*0.194 + 4*0.161 + 5*0.128 + 6*0.098 + \\ &+ 7*0.098 + 8*0.064 + 9*0.064 + 10*0.032 + 11*0.032 = \\ &= 0.032 + 0.196 + 0.582 + 0.644 + 0.640 + 0.588 + 0.686 + 0.512 + \\ &+ 0.576 + 0.320 + 0.352 = 5.127 \approx 5.\end{aligned}$$

$$\begin{aligned}D_k &= (1-5)^2 * 0.032 + (2-5)^2 * 0.098 + (3-5)^2 * 0.194 + (4-5)^2 * 0.161 + (5-5)^2 * 0.129 + \\ &+ (6-5)^2 * 0.098 + (7-5)^2 * 0.098 + (8-5)^2 * 0.064 + (9-5)^2 * 0.064 + \\ &+ (10-5)^2 * 0.032 + (11-5)^2 * 0.032 = 0.512 + 0.576 + 0.776 + 0.111 + 0 + \\ &+ 0.098 + 0.392 + 0.576 + 0.800 + 1.152 = 5.349.\end{aligned}$$

$$s_k = \sqrt{5.349} = 2.31.$$

Noteiksim sadalījuma mediānu. Sakārtosim atribūtu a_k vērtības nedilstošā secībā:

12,2,2,3,3,3,3,3, 4,4,4,4,4,5,5,5,5,6,6,6,7,7,7, 8, 8,9,9,10,11.

Mums ir 31 atribūtu vērtība, tāpēc mediāna būs vērtība, kas atrodas virknes 16. pozīcijā.

$$med(a_k) = 5.$$

Šajā piemērā $mod(a_k) = 4$. Tā kā visi parametru aprēķini nav vienādi, atribūtu a_k vērtību sadalījums nav simetrisks un nevar būt normāls sadalījums.

Intervālā $\pm s_j$ no sagaidāmās vērtības atrodas intervāli 3, 4, 5, 6, 7. Šajā intervālā atrodas 21 atribūta a_k vērtība jeb 67% no visām tā vērtībām.

Intervālā $\pm 2s_j$ ir 29 atribūtu a_j vērtības jeb 93% no visām tā vērtībām.

Intervāls $\pm 3s_j$ satur 100% atribūtu a_j vērtību.

Šie rezultāti arī norāda, ka atribūtu a_k vērtību sadalījums nav normāls sadalījums.

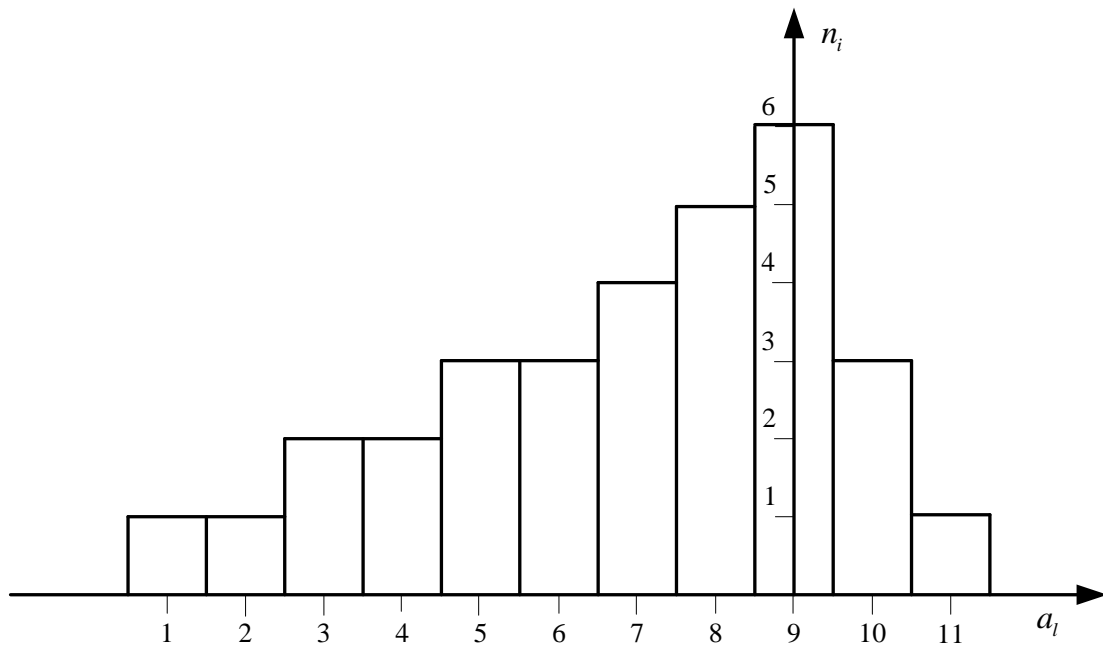
Piemērs P3.2.3. Tabulā P3.2.3 ir norādītas atribūta a_l sākotnējās vērtības un citi dati par šī atribūta vērtību sadalījumu. Aprēķināsim šī atribūta parametru vērtības.

Tabula P3.2.3

Sākotnējā atribūta a_l vērtību kopa un vērtību sadalījums

Intervāla numurs	1	2	3	4	5	6	7	8	9	10	11
Intervālu robežas	0.5-1.5	1.5-2.5	2.5-3.5	3.5-4.5	4.5-5.5	5.5-6.5	6.5-7.5	7.5-8.5	8.5-9.5	9.5-10.5	10.5-11.5
Vidējā vērtība intervālā	1	2	3	4	5	6	7	8	9	10	11
Vērtību skaits intervālā	1	1	2	2	3	3	4	5	9	3	1
Varbūtība, ka atribūta vērtība nonāks intervālā	0.032	0.032	0.064	0.064	0.097	0.097	0.129	0.161	0.194	0.098	0.032

Atribūtu a_i vērtību sadalījuma histogramma ir parādīta attēlā P3.2.3.



Attēls P3.2.3. Atribūtu a_i vērtību sadalījuma histogramma datiem no tabulas P3.2.3

Aprēķināsim atribūtu a_i vērtību sadalījuma parametrus.

$$\begin{aligned} \bar{a}_k &= 1 * 0.032 + 2 * 0.032 + 3 * 0.064 + 4 * 0.064 + 5 * 0.097 + 6 * 0.097 + \\ &+ 7 * 0.129 + 8 * 0.161 + 9 * 0.194 + 10 * 0.098 + 11 * 0.032 = \\ &= 0.032 + 0.064 + 0.192 + 0.256 + 0.485 + 0.582 + 0.903 + 0.1.288 + \\ &+ 1.746 + 0.980 + 0.352 = 6.880 \approx 7. \end{aligned}$$

$$\begin{aligned} D_i &= (1-7)^2 * 0.032 + (2-7)^2 * 0.032 + (3-7)^2 * 0.064 + (4-7)^2 * 0.064 + (5-7)^2 * 0.097 + \\ &+ (6-7)^2 * 0.097 + (7-7)^2 * 0.129 + (8-7)^2 * 0.161 + (9-7)^2 * 0.194 + \\ &+ (10-7)^2 * 0.098 + (11-7)^2 * 0.032 = 1.152 + 0.800 + 1.024 + 0.570 + 0.388 + \\ &+ 0.097 + 0 + 0.161 + 0.776 + 0.882 + 0.512 = 6.368. \end{aligned}$$

$$s_i = \sqrt{6.368} = 2.52.$$

$$\text{med}(a_i) = 7.$$

$$\text{mod}(a_i) = 8.$$

No tā izriet, ka atribūtu a_i vērtību sadalījums nevar būt normāls sadalījums.

67% atribūtu vērtību ietilpst vienas standartnovirzes intervālā. 93% atribūtu vērtību ietilpst divu standartnoviržu intervālā. 100% atribūtu vērtību ietilpst trīs standarta noviržu intervālā. Tas arī norāda, ka atribūtu a_i vērtību sadalījums nav normāls sadalījums.

P3.3. Sadalījuma normalitātes pārbaude pamatojoties uz augstākas kārtas parametriem

Atribūtu a_j vērtību sadalījuma asimetrijas rādītāju (*skewness*) S aprēķina pēc vienādojuma

$$S(a_j) = \frac{n \sum_{i=1}^n (a_{ij} - \bar{a}_j)^3}{s_j^3 (n-1)(n-2)}, \quad (\text{P3.3.1})$$

kur a_{ij} - atribūta a_j vērtība i -tajā intervālā ;

\bar{a}_j - atribūta a_j sagaidāmā vērtība;

s_j - atribūtu a_j vērtību standarta novirze no sagaidāmās vērtības;

n - kopējais atribūtu a_j vērtību skaits.

Stingri simetriskam sadalījumam S vērtība būs vienāda ar nulli. Uz kreiso pusi vērstam sadalījumam S vērtība ir mazāka par nulli, uz labo pusi vērstam sadalījumam S vērtība ir lielāka par nulli.

Simetrijas prasība ir nepieciešama, bet nav pietiekama, lai sadalījumu novērtētu kā normālu. Otrs faktors, no kura ir atkarīgi sadalījuma pārbaudes rezultāti attiecībā uz normalitāti, ir tā forma. Normālajam sadalījumam ir īpaša zvanveida forma (Gausa sadalījums). Bet ir iespējami simetriski sadalījumi, kuriem ir “plakanāka” vai “asāka” forma nekā parastajam sadalījumam.

Lai novērtētu sadalījuma formas atbilstību normalitātes prasībām, tiek izmantots īpašs parametrs – ekscesa rādītājs (*kurtosis*). Ekscesa vērtība tiek aprēķināta, izmantojot vienādojumu

$$K(a_j) = \frac{\sum_{i=1}^n (a_{ij} - \bar{a}_j)^4 / n}{s_j^4} - 3, \quad (\text{P3.3.2})$$

kur a_{ij} - atribūtu a_j vērtības i -tajam objektam;

\bar{a}_j - atribūtu a_j sagaidāmā vērtība;

s_j - atribūtu a_j vērtību standarta novirze no sagaidāmās vērtības;

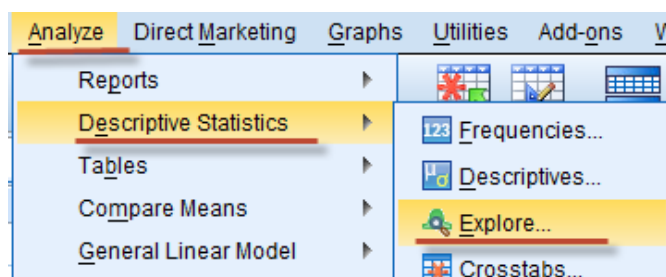
n - kopējais atribūtu a_j vērtību skaits.

Ja sadalījums ir “plakanāks” par parasto sadalījumu, ekscesa vērtība ir mazāka par nulli. Stingri normālam sadalījumam ekscesa vērtība ir nulle. Ja sadalījums ir “asāks” par parasto sadalījumu, ekscesa vērtība ir lielāka par nulli.

Sniegsim ilustratīvus piemērus asimetrijas pakāpes novērtēšanai, pamatojoties uz parametru S .

Piemērs P3.3.1. Pēc tabulas P3.2.1 datiem nepieciešams novērtēt atribūtu a_j vērtību sadalījuma asimetrijas un ekscesa rādītājus.

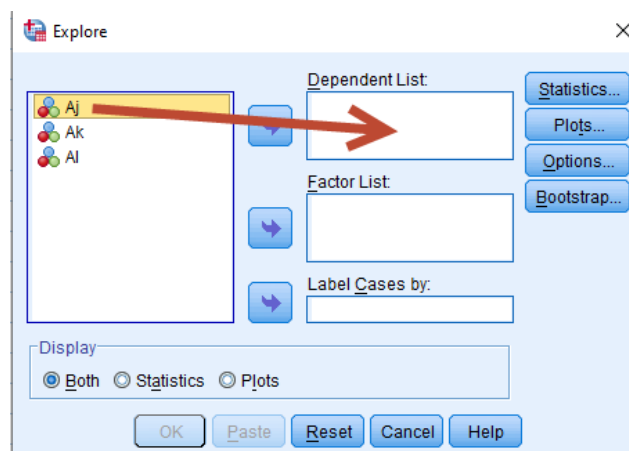
Mēs izmantojam SPSS programmatūras pakotni. Asimetrijas un ekscesa aprēķināšanas procedūras tiek veiktas šādi. Pēc SPSS palaišanas secīgi ir jāizvēlas opcijas “Analyze”, “Descriptive Statistics” un “Explore”.



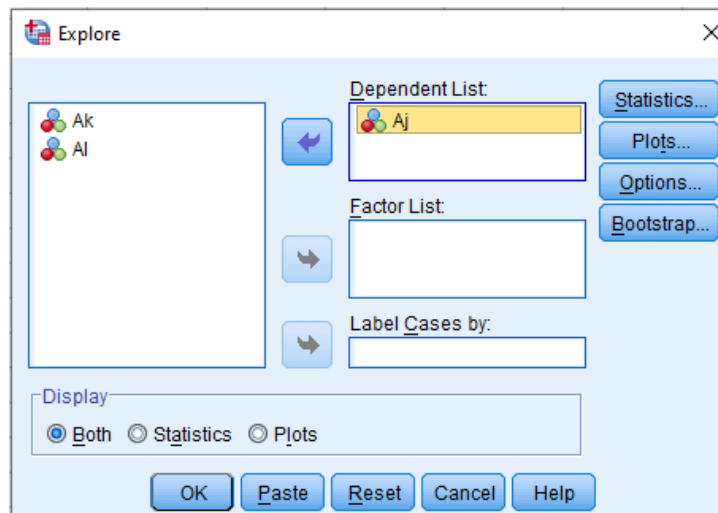
Logā, kas parādās, jāievada sākotnējie dati (*atribūts vai atribūtu vērtības, bet ne biežuma dati*).

	Aj	Ak	Al
1	1	1	1
2	2	2	2
3	3	2	3
4	3	2	3
5	4	3	4
6	4	3	4
7	4	3	5
8	5	3	5
9	5	3	5
10	5	3	6
11	5	4	6
12	6	4	6
13	6	4	7
14	6	4	7
15	6	4	7
16	6	5	7
17	6	5	8
18	7	5	8
19	7	5	8
20	7	6	8
21	7	6	8
22	7	6	9
23	7	7	9
24	7	7	9
25	7	7	9
26	8	8	9
27	8	8	9
28	8	9	10
29	8	9	10
30	8	10	10
31	8	11	11
32	9	.	.
33	9	.	.
34	9	.	.
35	9	.	.
36	10	.	.
37	10	.	.

(Tehnisku iemeslu dēļ visas atribūtu a_j vērtības nav parādītas).
Pēc noklikšķināšanas uz pogas “OK” parādās logs



Kreisajā panelī redzami testējamo atribūtu nosaukumi (mūsu gadījumā A_j , A_k , A_l). Izvēlamies mūs interesējošo atribūtu A_j un ar bultiņas palīdzību pārsūtām to uz labo paneli.



Pēc pogas “OK” noklikšķināšanas tiek iegūti šādi rezultāti:

Descriptives			Statistic	Std. Error
Mean			7,00	,406
95% Confidence Interval for Mean	Lower Bound		6,18	
	Upper Bound		7,82	
5% Trimmed Mean			7,00	
Median			7,00	
Variance			6,927	
Std. Deviation			2,632	
Minimum			1	
Maximum			13	
Range			12	
Interquartile Range			4	
Skewness			,000	,365
Kurtosis			-,001	,717

Tabulas apakšējās divās rindās ir norādītas asimetrijas un ekscesa vērtības: $S(a_j) = 0.000$, $K = -0.001$. Tas norāda uz to, ka atribūtu a_j vērtības ir sadalītas pēc normālā sadalījuma likuma.

Piemērs P3.3.2. Atribūtu a_k vērtību sadalījuma asimetrijas vērtību nepieciešams novērtēt datiem no tabulas P3.2.2. No šīs tabulas mums ir $S(a_k) = 0.581$, $K(a_k) = -0.427$. Atribūtu a_k vērtību sadalījums ir uz labo pusi vērsts sadalījums.

Izmantojot SPSS programmatūras pakotni, mēs iegūstam šādus rezultātus.

Descriptives			Statistic	Std. Error
Mean			5,13	,461
95%	LowerBound		4,19	
ConfidenceIntervalforMean	UpperBound		6,07	
5% TrimmedMean			5,03	
Median			5,00	
Variance			6,583	
Std. Deviation			2,566	
Minimum			1	
Maximum			11	
Range			10	
InterquartileRange			4	
Skewness			,581	,421
Kurtosis			-,427	,821

Piemērs P3.3.3. Atribūtu a_i vērtību sadalījuma asimetrijas vērtību nepieciešams novērtēt datiem no tabulas P3.2.3.

Izmantojot SPSS programmatūras pakotni, mēs iegūstam šādus rezultātus.

Descriptives			Statistic	Std. Error
Mean			6,87	,461
95%	LowerBound		5,93	
ConfidenceIntervalforMean	UpperBound		7,81	
5% TrimmedMean			6,97	
Median			7,00	
Variance			6,583	
Std. Deviation			2,566	
Minimum			1	
Maximum			11	
Range			10	
InterquartileRange			4	
Skewness			-,581	,421
Kurtosis			-,427	,821

No šīs tabulas mums ir $S(a_i) = -0.581$, $K(a_i) = -0.427$. Atribūtu a_i vērtību sadalījums ir uz kreiso pusi vērsts sadalījums.

P3.4. Sadalījuma normalitātes neparametriskā pārbaude

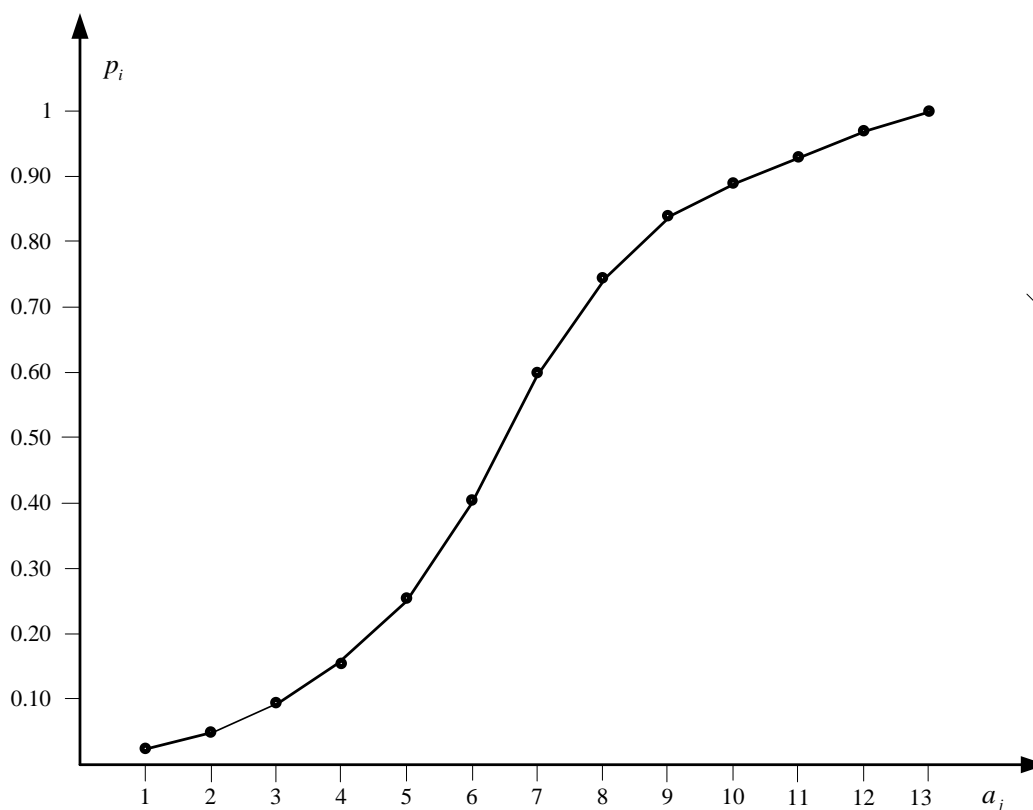
Sadaļās P3.2, P3.3 tika izmantoti parametriskie testi sadalījumu normalitātes pārbaudei. Sadaļā P3.2 šim nolūkam tika izmantoti pirmās kārtas sadalījumu parametri, P3.3 sadaļā – augstākās kārtas sadalījumu parametri.

Šajā sadaļā ir parādīti divi plaši izmantoti neparametriskie normalitātes testi: Kolmogorova-Smirnova tests (*Kolmogorov-Smirnov test*) un Šapiro-Vilksa tests (*Shapiro-Wilk test*).

1. Kolmogorova-Smirnova tests.

Ja ir norādīts atribūtu a_j vērtību sadalījums, tad šim sadalījumam vienmēr var izveidot sadalījuma funkciju. Attēlā P3.4.1 parādīts atribūtu a_j vērtību sadalījuma funkcijas grafiks datiem no tabulas P3.2.1 (piemērs P3.2.1).

Grafiks tiek veidots šādi. Uz horizontālās ass attēlotas atribūtu a_j vērtības (mūsu piemērā tās ir vidējās atribūtu vērtības katrā intervālā). Vertikālā ass parāda to atribūtu vērtību varbūtības (biežumi), kas ietilpst atbilstošajos intervālos. Šeit iet runa par varbūtību uzkrātajām vērtībām.



Attēls P3.4.1. Atribūtu a_j vērtību sadalījuma funkcijas grafiks datiem no tabulas P3.2.1

Pirmajam intervālam grafikā tiek attēlota tā varbūtības vērtība $p_1 = 0.024$. Šī vērtība ir atzīmēta ar punktu. Otrajam intervālam punkts parāda uzkrāto varbūtību $p_1 + p_2 = 0.024 + 0.024 = 0.048$ vērtības 1. un 2. intervālā. Trešais punkts parāda uzkrāto varbūtību vērtību $p_1 + p_2 + p_3 = 0.096$.

Pārējie punkti tiek noteikti līdzīgā veidā. Acīmredzams, ka pēdējā intervāla kumulatīvā varbūtības vērtība ir 1. Tas nozīmē, ka jebkura atribūta vērtība noteikti ietilpst tās vērtību diapazonā no 1 līdz 13.

Kolmogorova-Smirnova testa ideja ir šāda. Atribūtu a_j vērtību diapazonā tiek konstruēta “ideālā” normālā sadalījuma funkcija. Šim normālajam sadalījumam ir tāda pati sagaidāmā vērtība un standarta novirze kā pārbaudāmajam empīriskajam sadalījumam. “Ideālā” sadalījuma funkcijas vērtības tiek salīdzinātas ar empīriskās sadalījuma funkcijas vērtībām dotajos punktos. Pēc tam testa statistiku aprēķina, izmantojot vienādojumu

$$D = \sup_{a_j} \left| \left(F_0(a_{ij}) - F_e(a_{ij}) \right) \right|, \quad i = 1, \dots, n, \quad (\text{P3.4.1})$$

kur $F_0(a_{ij})$ - “ideālā” sadalījuma funkcijas vērtība i -tajā punktā;

$F_e(a_{ij})$ - empīriskā sadalījuma funkcijas vērtība i -tajā punktā;

n - datu punktu skaits (mūsu gadījumā vērtību intervālu skaits).

Citiem vārdiem sakot, vērtība D ir maksimālā atšķirība starp “ideālo” un empīrisko sadalījuma funkciju vērtību kādā i -tajā punktā.

Aprēķināto vērtību D salīdzina ar kritisko vērtību c . Ja $D < c$, nulles hipotēze tiek pieņemta: atribūtu a_j vērtības tiek sadalītas saskaņā ar normālo sadalījuma likumu. Ja $D > c$, nulles hipotēze tiek noraidīta, kas nozīmē, ka atribūtu a_j vērtības nav normāli sadalītas.

Kritiskās vērtības c ir tabulētas dažādiem nozīmīguma līmeņiem α . Šīs kritiskās vērtības ir norādītas tabulā P3.4.1.

Tabula P3.4.1

Kolmogorova-Smirnova normalitātes testa kritiskās vērtības c

n	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.10$	$\alpha = 0.15$	$\alpha = 0.20$
1	0.995	0.975	0.950	0.925	0.900
2	0.929	0.842	0.776	0.726	0.684
3	0.828	0.708	0.642	0.597	0.565
4	0.733	0.624	0.564	0.525	0.494
5	0.669	0.565	0.510	0.474	0.446
6	0.618	0.521	0.470	0.436	0.410
7	0.577	0.486	0.438	0.405	0.381
8	0.543	0.457	0.411	0.381	0.358
9	0.514	0.432	0.388	0.360	0.339
10	0.490	0.410	0.368	0.342	0.322
11	0.468	0.391	0.352	0.326	0.307
12	0.450	0.375	0.338	0.313	0.295
13	0.433	0.361	0.325	0.302	0.284
14	0.418	0.349	0.314	0.292	0.274
15	0.404	0.338	0.304	0.283	0.266
16	0.382	0.328	0.295	0.274	0.258
17	0.381	0.318	0.286	0.266	0.250
18	0.371	0.309	0.278	0.259	0.244
19	0.363	0.301	0.272	0.252	0.237
20	0.356	0.294	0.264	0.246	0.231
25	0.320	0.270	0.240	0.220	0.210
30	0.290	0.240	0.220	0.200	0.190
35	0.270	0.230	0.210	0.190	0.180
40	0.250	0.210	0.190	0.180	0.170
45	0.240	0.200	0.180	0.170	0.160
50	0.230	0.190	0.170	0.160	0.150
>50	$\frac{1.63}{\sqrt{n}}$	$\frac{1.36}{\sqrt{n}}$	$\frac{1.22}{\sqrt{n}}$	$\frac{1.14}{\sqrt{n}}$	$\frac{1.07}{\sqrt{n}}$

Kolmogorova-Smirnova testa priekšrocības ir:

1. Neatkarība no esošās empīriskās funkcijas sadalījuma.
2. Rezultātu precizitāte.
3. Izlases apjomam nav ierobežojumu.

Testa trūkumi ir:

1. Tas attiecas tikai uz nepārtrauktiem sadalījumiem.
2. Tests ir jūtīgāks pret vērtībām empīriskā sadalījuma centrā nekā pret vērtībām tā galos.

2. Šapiro - Vilksa tests.

Šis tests ir balstīts uz empīriskās blīvuma funkcijas $f_e(a_{ij})$ vērtību salīdzināšanu ar "ideālas" normālā blīvuma funkcijas vērtībām $f_0(a_{ij})$, kas konstruētas noteiktā atribūtu a_j vērtību diapazonā. "Ideālajam" sadalījumam ir tāda pati sagaidāmā vērtība un standarta novirze kā empīriskajam sadalījumam.

Šapiro-Vilksa testa nulles hipotēze ir tāda, ka atribūta a_j sākotnējās vērtības ir izlase no populācijas, kurā šī atribūta vērtības ir normāli sadalītas. Pārbaudes rezultātā tiek aprēķināta varbūtības p vērtība. Kritiskā vērtība ir $p = 0.5$. Ja aprēķinātā varbūtības vērtība ir $p < 0.05$, nulles hipotēze tiek noraidīta. Ja $p > 0.05$, tiek pieņemta nulles hipotēze, t.i., sākotnējās atribūtu a_j vērtības nāk no normālas populācijas.

Nosacījumu, ka atribūtu a_j vērtības nāk no normālas populācijas, var viegli apiet. Nelielas izmaiņas aprēķinu izteiksmēs ļauj tieši pārbaudīt sākotnējo sadalījumu bez atsauces uz pamatpopulāciju.

Tā kā abu testu praktiskai izpildei nepieciešams liels skaits aprēķinu, praksē šie testi tiek veikti, izmantojot piemērotus skaitļošanas rīkus, piemēram, statistikas pakotni SPSS. Lietojot P3.4. sadaļā aprakstītās procedūras, programma kopā ar izvades statistiku, balstoties uz Kolmogorova-Smirnova un Šapiro-Vilksa testiem, rāda arī normalitātes sadalījuma pārbaudes rezultātus. Tomēr šie testi ir paredzēti sadalījumiem, kas satur lielu skaitu atbilstošu atribūtu vērtību. Ja sadalījums satur nelielu skaitu atribūtu vērtību, testa rezultāti būs ļoti neprecīzi un neuzticami, tāpēc nav jēgas izmantot šos testus iepriekš sniegtajiem ilustratīvajiem piemēriem.