

ÕPIOBJEKT

„Sissejuhatus üldiste lineaarsete mudelite teoriasse“

Tanel Kaart

http://www.eau.ee/~ktanel/lineaarne_mudel/


www.emu.ee

Euroopa Liit
Euroopa Sotsiaalfond

Eesti tuleviku heaks

Õpiobjektid -> Sissejuhatus üldiste lineaarsete mudelite teoriasse

SISSEJUHATUS ÜLDISTE LINEAARSETE MUDELITE TEORIASSE

Õpiobjekti kirjeldus

Õpijuhatis

Sissejuhatus

1. Põhimõisted

2. Üldiste lineaarsete mudelite ja faktorite liigitus

- * Diskreetsed ja pidevad faktorid
- * Faktorite vahetõrke mudelid
- * Juhuslikud ja fikseeritud faktorid
- * Tasakaalus ja mittetasakaalus andmed (mudelid)

3. Üldise lineaarse mudeli esitused

- * Mudeli esitus objektiviisi
- * Üldine lineaarne mudel maatrikskuul
- * Üldine lineaarne segamudel maatrikskuul

4. Keskväärtused ja dispersioonid

- * Vaatluste keskväärtused ja kovariatsioonistruktuur fikseeritud mudeli korral
- * Vaatluste keskväärtused ja kovariatsioonistruktuur segamudeli korral

5. Eeldused ja kitsendused

- * Jaotused
- * Kovariatsioonistruktuur
- * Täisastakuga ja mitte täisastakuga mudelid, reparametriseerimis-tingimused

6. Fikseeritud efektide hindamine (BLUE)

- * Vähimruutude hinnangud
- * Hinnatavad funktsioonid
- * Vähimruutkeskmised

7. Juhuslike faktorite realiseerunud väärtuste prognoosimine (BLUP)

Õpiobjekti kirjeldus

Õppekava: Loomakasvatuse (449)

Õppeaine: VL.0192 Loomade aretusväärtuse hindamine ja aretusprogrammid

Maht: 10 tundi

Sihtrühm: Loomakasvatuse õppekava magistrandid

Eesmärk: Õpiobjekti eesmärk on toetada õppeaine omandamist

Õpiobjekti läbinu:

- mõistab üldiste lineaarsete mudelite esitust nii objektiviisi kui ka maatrikskuul;
- oskab hinnata mudelite parameetrite väärtuseid lihtsamate mudelite korral ja mõistab nende tähendust;
- omab võimalust enesekontrolliks.

Sisu ja tehniline teostus: Tanel Kaart

Eesti Maaülikool
kevadsemester 2012

[Järgmine >](#)

ÜLDINE LINEAARNE MUDEL

Eelmises materjalides kirjeldatud selektsiooniindeksite rakendamine suurtes reaalsetes populatsioonides on sageli pärsitud mitmel põhjusel.

- Esiteks on keerulisemate põlvnemisskeemide puhul tülikas ja vägagi töömahukas kõikvõimalikke sugulussidemeid arvestavate indeksite välja töötamine.
- Teiseks eeldab selektsiooniindeksite teooria kõigi mitte aditiivgeneetiliste faktorite mõjude täpset teadmist ja mõõtmistulemuste nende suhtes korrigeerimist, mida on aga suuremamahuliste, mitmeid aastaid ja loomakasvatusevõtteid hõlmavate uuringute korral pea võimatu teostada. Seetõttu võivad selektsiooniindeksite abil saadud aretusväärtuste hinnangud küll rahuldada parima lineaarse prognoosi (BLP) tingimusi, aga olla kokkuvõttes ikkagi positiivselt kallutatud paremais ettevõtteis või aastatel mõõdetud loomade suhtes (tulemuseks on nihkega hinnangud).
- Kolmandaks ei saa siiski alati lugeda selektsiooniindeksites kasutatavaid populatsiooni iseloomustavaid geneetilisi parameetreid (päritavuskoeffitsient, korduvus, geneetilised korrelatsioonikordajad jne) teadaolevaiks, nende väärtuste hindamiseks selektsiooniindeksid aga enamasti ei sobi.

Kõigi kolme probleemi lahendus on kasutada **üldisi lineaarseid segamudeleid** (inglise keeles *general linear mixed models*, GLMM). Nende mudelite olemuses on nii tundmatute keskkonnamõjude kui ka aretusväärtuste samaaegne üksteise mõju arvesse võttev hindamine, mistõttu on tulemuseks nihketa hinnangud – nii keskkonnamõjude **parimad lineaarsed nihketa hinnangud** (*best linear unbiased estimate*, BLUE) kui ka aretusväärtuste **parimad lineaarsed nihketa prognoosid** (*best linear unbiased predictor*, BLUP). Samuti võimaldavad need mudelid sobival valitud kovariatsioonistruktuuri abil arvesse võtta kõikvõimalikke sugulussidemeid ning hinnata polügeensete tunnuste geneetilist determineeritust ja omavahelist seotust populatsioonis iseloomustavaid parameetreid.

Mitmetel populatsioonigeneetikas ja aretusteoorias kasutatavatel üldiste lineaarsete segamudelite erijuhtudel peatutakse pikemalt järgmistes materjalides. Käesolev materjal esitab üldiste lineaarsete mudelite ja üldiste lineaarsete segamudelite korrektseks rakendamiseks ja tulemuste tõlgendamiseks vajalikud põhiteadmised.

1 PÕHIMÕISTED

Üldise lineaarse mudeli rakendamiseks jagatakse registreeritud tunnused kahte ossa – uuritavateks e sõltuvateks tunnusteks (need, mille käitumine huvi pakub) ja argument- e sõltumatuteks tunnusteks e faktoriteks (need, mille mõju uuritavatele tunnustele soovitakse selgitada).

Definitsiooni kohaselt on **üldine lineaarne mudel** (*general linear model*, GLM) uurija poolt eeldatav faktorite ja uuritavate tunnuste vahekorra üldskeem (mudel), mille detailid tuleb hinnata algandmetest (valimist).

Üldise lineaarse mudeli alusel tehtud järeldused on õiged kui

- a) uuritav mudel vastab tegelikkusele – on õigesti paika pandud faktorite olemus ja vahekord;
- b) kehtivad vajalikud matemaatilised eeldused (normaaljaotus, vaatluste sõltumatus jm).

Tänu oma lihtsusele ja tõsiasjale, et faktorite mõjude lineaarne kombinatsioon vastab ligilähedaselt tegelikkusele paljudes reaalsetes analüüsides ning sobib lähendama lineaarsest märksa keerulisemaid funktsioone, on üldine lineaarne mudel tänapäeva eluteadustes enim rakendatav matemaatilise statistika meetod.

Iga üldise lineaarse mudeli püstitamine eeldab kolme asja:

- a) mudeli esitust faktorite mõjude summamana,
- b) uuritava tunnuse ning faktorite keskvärtuste ja dispersioonimaatriksite struktuuri fikseerimist vastavalt andmestiku ülesehitusele ja faktorite olemusele, ning
- c) tõenäosusjaotuslike eelduste tegemist mudeli liikmete kohta võimaldamaks määrata hinnangute täpsust ja kontrollida hüpoteese.

2 ÜLDISTE LINEARSETE MUDELITE JA FAKTORITE LIIGITUS

2.1 Diskreetsed ja pidevad faktorid

Faktortunnuse erinevaid väärtusi nimetatakse **tasemeteks** e **nivoodeks** (inglise keeles *levels*). Iga faktor jaguneb vastavalt oma tasemete iseloomule **diskreetseks** või **pidevaks**, arvuliseks (kvantitatiivne) või klassifitseerivaks (kvalitatiivne).

Näiteks on lehma sünniaasta, laktatsioon jne diskreetsed arvulised faktorid; farm tasemetega (väärtustega) 'Vorbuse', 'Ülenurme' jne on diskreetne klassifitseeriv faktor; laktatsiooni pikkus, piimatoodang, pekipaksus jne (möödetud tunnused) on aga pidevad faktorid.

Üldine lineaarne mudel on **dispersioonanalüüsi**, **regressioonanalüüsi** või **kovariatsioonanalüüsi** mudel vastavalt sellele, kas faktorid on diskreetsed, pidevad või esineb mõlemaid.

2.2 Faktorite vahekord mudelis

Faktorid on lihtsad ja tuletatud. Lihtsate faktorite väärtused on vahetult möödetud või registreeritud, tuletatud faktorid moodustatakse lihtsatest. Tüüpilised tuletatud faktorid on **interaktsioonid** e. **koosmõjud** ning arvuliste faktorite korrutised.

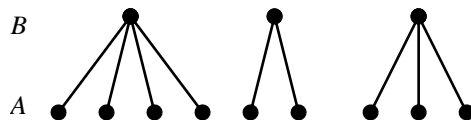
Näiteks on farm, isa ja laktatsiooni pikkus lihtsad faktorid; farm*isa (koosmõju) ja laktatsioon*laktatsioon (arvulise faktori kõrgem järk) aga tuletatud faktorid.

Praktikas on faktorite vahel sageli ka alluvusseosed. Faktor A **allub** faktorile B – tähistatakse enamasti $A(B)$ –, kui A iga nivoo (tase) esineb koos vaid ühe B nivooaga (joonis 1).

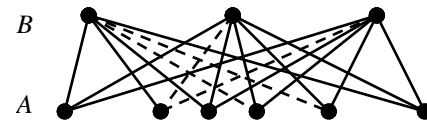
Näiteks võime me tavaliselt lugeda farmi allutatuks maakonnale; kui iga ema on ristatud kindla isaga, on ema allutatud isale.

Faktorid A ja B on **ristseoses**, kui A iga nivoo kombineerub (saab põhimõtteliselt kombineeruda) B kõigi nivooodega (joonis 2).

Näiteks kui viiel aastal on uuritud pullide tütarde jõudlusandmeid ja igal pullil on igal aastal tütreid, on pull ja aasta ristseoses; kui aga igal aastal on valitud uued pullid, allub pull aastale.



Joonis 1. Alluvusseoses faktorid



Joonis 2. Ristseoses faktorid

2.3 Juhuslikud ja fikseeritud faktorid

Sõltuvalt uurija eesmärkidest ja analüüsitavaatest andmetest omistatakse igale faktorile kindel tüüp – **juhuslik** või **fikseeritud**.

Fikseeritud faktoril on:

- vähe nivooideid,
- iga nivoo pakub iseseisvat huvi ja on valitud mittejuhuslikult,
- andmetes on või saavad põhimõtteliselt olla esindatud kõik nivood.

Seetõttu käsitletakse fikseeritud faktorite tasemetele vastavaid efekte kui konstante.

Juhuslikul faktoril on

- potentsiaalselt väga palju (lõpmatu hulk) nivooideid,
- andmetes on neist esindatud juhuslik valim,
- huvi pakub kõigi (ka andmetes esindamata) nivooide keskmine mõju e see, kui suur osa uuritava tunnuse koguvarieeruvusest on kirjeldatud antud faktori poolt.

Juhuslike faktorite tasemetele vastavaid efekte käsitletakse kui mingi teoreetilise jaotusega juhuslike suuruste realiseerunud väärtusi, kus selle teoreetilise jaotuse näol mõistetakse üldjuhul normaaljaotust.

Näiteks, kui andmestik sisaldab kümne spetsiaalselt valitud isaslooma järglaste andmeid, loetakse isa mõju (faktor 'isa') fikseerituks, kuna me soovime võrrelda just antud isasid ega soovi teha järeldusi teiste isade kohta. Kui aga huvi pakub isa kui geneetilise faktori üldine mõju (isade poolt järglastele pärandatavate juhuslike geenidekomplektide mõju) ja järglased on valitud juhuslikult, siis loetakse isa mõju juhuslikuks.

See, kas faktor on juhuslik või fikseeritud, ei kajastu mudeli võrrandis, küll aga sõltuvad sellest mudeli liikmete keskväärtused ja dispersioonid ning seeläbi ka mõjude hinnangud ja nende statistiline olulisus.

Vastavalt sellele, kas üldine lineaarne mudel sisaldab üksnes fikseeritud faktoreid, üksnes juhuslikke faktoreid või siis mõlemaid, nimetatakse kasutatavat mudelit kas **fikseeritud mudeliks** (*fixed model*), **juhuslikuks mudeliks** (*random model*) või **segamudeliks** (*mixed model*). Aretusväärtuste ja teiste geneetiliste parameetrite hindamisel rakendatakse just viimast.

2.4 Tasakaalus ja mittetasakaalus andmed (mudelid)

Tasakaalus (*balanced*) andmed on sellised, kus kõigil faktorite ja nende kombinatsioonide tasemetel on sooritatud võrdne arv mõõtmisi/vaatlusi/vmt.

Mittetasakaalus (*unbalanced*) andmed on sellised, kus vastupidiselt tasakaalulisele juhule ei vasta kõigile faktorite tasemetele võrdne arv objekte. Andmete mittetasakaalulisus eeldab vähegi keerulise- ma struktuuriga mudelite puhul ligikaudsete arvutusmeetodite kasutamist, mistõttu tulemused ei pruugi olla enam ühesed.

Näide 1. Järgnevalt toodud andmetabeleist esimene kujutab tasakaalus ja teine mittetasakaalus andmestikku. Esimeses on mõlemast majandist ühepalju eesti holsteini ja eesti punast tõugu lehma, lisaks on ühepalju kõikvõimalikke majandite, tõugude ja laktatsioonide kombinatsioone. Teise andmestiku korral pole enamusi nimetatud tingimustest täidetud.

Tasakaalus andmed			Mittetasakaalus andmed		
Farm	Tõug	Laktatsioon	Farm	Tõug	Laktatsioon
Põlula	EHF	1	Põlula	EHF	1
Põlula	EHF	2	Põlula	EHF	2
Põlula	EHF	3	Põlula	EPK	2
Põlula	EPK	1	Põlula	EPK	1
Põlula	EPK	2	Põlula	EPK	2
Põlula	EPK	3	Põlva	EHF	1
Põlva	EHF	1	Põlva	EHF	2
Põlva	EHF	2	Põlva	EHF	2
Põlva	EHF	3	Põlva	EHF	3
Põlva	EPK	1	Põlva	EPK	1
Põlva	EPK	2	Põlva	EPK	3
Põlva	EPK	3			

3 ÜLDISE LINEAARSE MUDELI ESITUSED

3.1 Mudeli esitus objektiviisi

Mudeli objektiviisi esitamisel avaldatakse igal objektil mõõdetud väärtus üldkeskmise (nn vabaliikme, inglise keeles *intercept*), objektile vastavate mõjude ja konkreetse objekti omapära e faktoritega kirjeldamata jääva osa (mudeli vea) summana.

Erinevate objektide ja faktorite tasemete üheseks identifitseerimiseks kasutatakse indeksite abi: igale diskreetsele faktorile omistatakse põhiindeks (harilikult i, j, k, l, \dots), mis näitab iga objekti kohta just temale omase faktori taseme mõju, peale selle tähistab üks lisaindeks objekti (mõõtmist) ennast. Faktorite koosmõjude korral kirjutatakse indeksiks mõlema faktori indeksid, allutatud faktorite puhul tuuakse tavaliselt faktori enese põhiindeksi järel sulgudes ära kõigi nende faktorite indeksid, millele antud faktor allub. Pideva faktori tähis omab kogu indeksite komplekti, sest märgib konkreetset objektile ja diskreetsete faktorite kombinatsioonil sooritatud mõõtmise tulemust, mille mõju samal objektile ja diskreetsete faktorite kombinatsioonil mõõdetud uuritava tunnuse väärtusega väljendab vastav regressioonikordaja.

Näide 2. Uuritavaks tunnuseks on tallede võõrutusmass ja analüüsi eesmärgiks on võrrelda ühes lambafarmis kasutatud kahte isa. Objektiviisemate tulemuste huvides soovitakse arvestada ka talle soo ja pesakonna suuruse kui diskreetsete faktorite ning võõrutusvanuse kui pideva faktori mõjudega. Andmed on kirjas järgnevas tabelis.

Sugu	Pesakonna suurus	Võõrutusmass	Võõrutusvanus	Isa kood
1	3	22	91	1025
1	1	39	90	1025
1	2	24	86	1025
2	2	27	101	1025
2	2	24	86	1025
1	2	31	117	1027
1	2	30	152	1027
2	1	31	110	1027
2	2	26	107	1027

Üldise lineaarse mudeli võib toodud ülesande püstituse korral esitada järgmiselt:

$$y_{ijkl} = \mu + I_i + S_j + P_k + b \times v_{ijkl} + e_{ijkl}, \quad (1)$$

kus parameetrit μ võib mõista kui nõ keskmist tallede 100 päeva kehamassi (mis siiski ei pruugi võrduda aritmeetilise keskmisega, sest on tegelikult faktorite keskmine mõju);

y_{ijkl} tähistab objektile l sooritatud mõõtmist e antud juhul l talle 100 päeva kehamassi (ülejäänud indeksid fikseerivad täpselt ära vaatlusaluse talle isa, sünniaasta jne), $l = 1, \dots, 9$;

I_i tähistab isa i mõju, $i = 1025, 1027$;

S_j märgib soo j mõju, $j = 1, 2$;

P_k on pesakonna suuruse k mõju, $k = 1, 2, 3$;

v_{ijkl} tähistab konkreetse talle võõrutusvanust ja b on võõrutusmassi lineaarset sõltuvust võõrutusvanusest iseloomustav regressioonikordaja;

e_{ijkl} on talle l omapära (e juhuslik viga).

3.2 Üldine lineaarne mudel maatrikskujul

Üldiste lineaarsete mudelite maatrikskujul esitamiseks kasutatakse eriliste maatriksite abi, mis seostavad iga objekti just temal sooritatud mõõtmistega või temale vastava faktori tasemega. Selliseid maatrikseid, mille ridade arv võrdub uuritud objektide arvuga ja veergude arv mudelist hinnatavate efektide arvuga, nimetatakse **plaani-** e **disainimaatriksiteks**.

Regressioonivõrrand, kirja panduna objekti i kohta,

$$y_i = \mu + bx_i + e_i,$$

näeb maatrikskujul esitatuna välja järgmiselt:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad (2)$$

kus $\mathbf{y} = (y_1 \ y_2 \ \dots \ y_n)^T$ ja $\mathbf{e} = (e_1 \ e_2 \ \dots \ e_n)^T$ on vastavalt funktsioontunnuse väärtuste ja prognoosivigade vektorid, $\boldsymbol{\beta} = (\mu \ b)^T$ on hinnatavate parameetrite vektor ja plaanimaatriks $\mathbf{X} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix}^T$

sisaldab ühte veergu kummagi hinnatava parameetri tarvis – esimene, ühtedest koosnev veerg tähendab, et kõigil objektidel mõõdetud uuritava tunnuse väärtusi piüütakse prognoosida sama keskmise μ abil, teine, argumenttunnuse väärtusi sisaldav veerg seob igal objektil mõõdetud uuritava tunnuse väärtuse läbi kordaja b just temale vastava argumenttunnuse väärtusega.

Dispersioonanalüüsi korral koosneb plaanimaatriks vaid nullidest ja ühtedest – mingile objektile vastavas reas ja faktori tasemele vastavas veerus on 1, kui mõõtmine antud objektil on sooritatud just sellel tasemel, ja 0 vastupidisel juhul. Näiteks ühefaktorilise dispersioonanalüüsi mudel

$$y_{ij} = \mu + \alpha_i + e_{ij}$$

on maatriksite abil esitatav kujul

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad \text{ehk} \quad \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{1}_{n_1} & \mathbf{0} & \dots & \mathbf{0} \\ 1 & \mathbf{0} & \mathbf{1}_{n_2} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \mathbf{0} & \mathbf{0} & \dots & \mathbf{1}_{n_q} \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_q \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}.$$

Siin n_i tähistab faktori i -ndale tasemele vastavate objektide arvu, q on faktori tasemete arv ($i = 1, \dots, q$), $\mathbf{1}_{n_i}$ on ühtedest koosnev $n_i \times 1$ -vektor: $\mathbf{1}_{n_i} = \underbrace{(1 \ 1 \ \dots \ 1)}_{n_i \text{ tk}}^T$, ja $\mathbf{0}$ on sobiva dimensiooniga nullidest

koosnev vektor.

Kovariatsioonanalüüsi mudel, milles on nii mitteamulisi kui ka arvilisi faktortunnuseid, sisaldab plaanimaatriksis nii nulle ja ühtesid sisaldavaid veerge kui ka pidevate argumenttunnuste väärtuseid.

Kokkuvõttes kujutab üldise lineaarse mudeli maatriksesitus enesest tavalise lineaarvõrranditesüsteemi maatriksesitust, kus kordajate maatriksi rollis on plaanimaatriks, tundmatud parameetrid on koondatud vektorisse $\boldsymbol{\beta}$ ja vabaliikmete veerg moodustub vahest $\mathbf{y} - \mathbf{e}$.

Näide 3. Paneme näite 2 mudeli (1) kirja ka maatrikskujul (2).

Vektor \mathbf{y} on tallede võõrutusmasside vektor kujul

$$\mathbf{y} = (22 \ 39 \ 24 \ 27 \ 24 \ 31 \ 30 \ 31 \ 26)^T.$$

Plaanimaatriks \mathbf{X} sisaldab ühte rida iga vaatluse ja ühte veergu iga vektoris $\boldsymbol{\beta}$ kirjas oleva mudeli parameetri kohta.

Hinnatavate parameetrite vektor $\boldsymbol{\beta}$ sisaldab 9 elementi (parameetrite tähistused on üle võetud sama mudeli objektiviisi esitusest):

$$\boldsymbol{\beta} = (\mu \ I_{1025} \ I_{1027} \ S_1 \ S_2 \ P_1 \ P_2 \ P_3 \ b)^T$$

Plaanimaatriks \mathbf{X} , mis seob iga vaatluse just temale vastava faktori tasemega või väärtusega, näeb praeguse tallede järjestuse korral välja järgmine:

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & 1 & 0 & 0 & 0 & 1 & 91 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 90 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 86 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 101 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 86 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 117 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 152 \\ 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 110 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 107 \end{pmatrix} \quad (3)$$

(parema jälgitavuse huvides on erinevatele faktoritele vastavad blokid eraldatud punktiirjoontega).

Juhuslike vigade vektor on kujul

$$\mathbf{e} = (e_1 \ e_2 \ e_3 \ e_4 \ e_5 \ e_6 \ e_7 \ e_8 \ e_9)^T.$$

Näitamaks, et eelnevalt defineeritud plaanimaatriks tõepoolest iga talle võõrutusmassi just temale vastavate faktorite tasemetega seob, asendame kõik väljakirjutatud vektorid ja maatriksid mudelisse (2) ning sooritame korrutustehte $\mathbf{X}\boldsymbol{\beta}$. Tulemuseks saame maatriksvõrduse

$$\begin{pmatrix} 22 \\ 39 \\ 24 \\ 27 \\ 24 \\ 31 \\ 30 \\ 31 \\ 26 \end{pmatrix} = \begin{pmatrix} \mu + I_{1025} + S_1 + P_3 + b \times 91 \\ \mu + I_{1025} + S_1 + P_1 + b \times 90 \\ \mu + I_{1025} + S_1 + P_2 + b \times 86 \\ \mu + I_{1025} + S_2 + P_2 + b \times 101 \\ \mu + I_{1025} + S_2 + P_2 + b \times 86 \\ \mu + I_{1027} + S_1 + P_2 + b \times 117 \\ \mu + I_{1027} + S_1 + P_2 + b \times 152 \\ \mu + I_{1027} + S_2 + P_1 + b \times 110 \\ \mu + I_{1027} + S_2 + P_2 + b \times 107 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ e_7 \\ e_8 \\ e_9 \end{pmatrix}.$$

3.3 Üldine lineaarne segamudel maatrikskujul

Juhul, kui üldine lineaarne mudel sisaldab nii fikseeritud kui ka juhuslikke faktoreid (tegu on segamudeliga), esitatakse need eraldi liidetavatena, jagades kaheks nii hinnatavate parameetrite vektori kui ka plaanimaatriksi:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}, \quad (4)$$

siin $\boldsymbol{\beta}$ ja \mathbf{u} tähistavad vastavalt fikseeritud ja juhuslike efektide vektoreid ning \mathbf{X} ja \mathbf{Z} vastavaid plaanimaatrikseid.

Näide 4. Oletame nüüd, et näites 2 vaadeldud jäärade alusel soovitakse iseloomustada isa kui geneetilise faktori üldist mõju ning jäärad on valimisse sattunud juhuslikult. Sellisel juhul tuleb isa mõjusid käsitleda juhuslike efektidena ning mudeli võib välja kirjutada valemi (4) alusel:

$$\begin{pmatrix} 22 \\ 39 \\ 24 \\ 27 \\ 24 \\ 31 \\ 30 \\ 31 \\ 26 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 1 & 91 \\ 1 & 1 & 0 & 1 & 0 & 0 & 90 \\ 1 & 1 & 0 & 0 & 1 & 0 & 86 \\ 1 & 0 & 1 & 0 & 1 & 0 & 101 \\ 1 & 0 & 1 & 0 & 1 & 0 & 86 \\ 1 & 1 & 0 & 0 & 1 & 0 & 117 \\ 1 & 1 & 0 & 0 & 1 & 0 & 152 \\ 1 & 0 & 1 & 1 & 0 & 0 & 110 \\ 1 & 0 & 1 & 0 & 1 & 0 & 107 \end{pmatrix} \begin{pmatrix} \mu \\ S_1 \\ S_2 \\ P_1 \\ P_2 \\ P_3 \\ b \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} \underbrace{\begin{pmatrix} I_{1025} \\ I_{1027} \end{pmatrix}}_{\mathbf{u}} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ e_7 \\ e_8 \\ e_9 \end{pmatrix}. \quad (5)$$

4 KESKVÄÄRTUSED JA DISPERSIOONID

Üldise lineaarse mudeli esituse teine osa seisneb mudeli juhuslike liikmete keskväärtuste ja dispersioonimaatriksite defineerimises (et fikseeritud mõjusid käsitletakse konstantsetena, ei ole nende keskväärtuste ja dispersioonide leidmisel mõtet).

4.1 Vaatluste keskväärtused ja kovariatsioonistruktuur fikseeritud mudeli korral

Igas lineaarses mudelis on vähemalt kaks juhuslikku liiget (mis eeldatakse traditsiooniliselt käituvat vastavalt normaaljaotuse seaduspäradele) – uuritav tunnus, kui populatsioonist juhuslikult valitud indiviididel sooritatud mõõtmiste tulemusi koondav suurus, ja juhuslik viga. Juhul, kui mudel sisaldab faktoritena üksnes fikseeritud mõjusid, on need kaks ka mudeli ainsad mittekonstantsed ja seega mittenullilist dispersiooni omavad suurused.

On loomulik eeldada, et juhuslike vigade keskväärtus võrdub nulliga (mudeli keskmine viga on null): $E(\mathbf{e}) = \mathbf{0}$. Uuritava tunnuse keskväärtus avaldub fikseeritud mudeli (2) korral kujul

$$E(\mathbf{y}) = E(\mathbf{X}\boldsymbol{\beta} + \mathbf{e}) = E(\mathbf{X}\boldsymbol{\beta}) + E(\mathbf{e}) = \mathbf{X}\boldsymbol{\beta}$$

ja dispersioonimaatriks võrdusena

$$\text{var}(\mathbf{y}) = \text{var}(\mathbf{X}\boldsymbol{\beta} + \mathbf{e}) = \underbrace{\text{var}(\mathbf{X}\boldsymbol{\beta})}_0 + \text{var}(\mathbf{e}) = \text{var}(\mathbf{e}) \quad (6)$$

4.2 Vaatluste keskväärtused ja kovariatsioonistruktuur segamudeli korral

Segamudeli (4) korral tuleb arvesse võtta ka juhuslike efektide keskväärtust ja varieeruvust. Analoogselt juhuslike vigadega eeldatakse ka juhuslike faktorite mõjude keskväärtuste nulliga võrdumist: $E(\mathbf{u}) = \mathbf{0}$. Seetõttu jääb uuritava tunnuse keskväärtus endiselt määratuks üksnes fikseeritud mõjude poolt:

$$E(\mathbf{y}) = E(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}) = E(\mathbf{X}\boldsymbol{\beta}) + E(\mathbf{Z}\mathbf{u}) + E(\mathbf{e}) = \mathbf{X}\boldsymbol{\beta}.$$

Arvestades dispersioonide defineerimisel ka juhuslike mõjude varieeruvusega, saame, et

$$\text{var}(\mathbf{y}) = \text{var}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}) = 0 + \text{var}(\mathbf{Z}\mathbf{u}) + \text{var}(\mathbf{e}) - \underbrace{2\text{cov}(\mathbf{Z}\mathbf{u}, \mathbf{e}^T)}_0 = \mathbf{Z}\text{var}(\mathbf{u})\mathbf{Z}^T + \text{var}(\mathbf{e}).$$

Eelnevas võrduste reas on eeldatud, nagu üldiste lineaarsete mudelite teoorias tavaks, et faktorite mõjud ja juhuslikud vead on sõltumatud: $\text{cov}(\mathbf{u}, \mathbf{e}^T) = \mathbf{0}$.

Tuues sisse tähistused $\mathbf{V} = \text{var}(\mathbf{y})$, $\mathbf{G} = \text{var}(\mathbf{u})$ ja $\mathbf{R} = \text{var}(\mathbf{e})$, saame segamudeli (4) keskväärtused ja kovariatsioonistruktuuri määrata maatriksitena

$$E \begin{pmatrix} \mathbf{y} \\ \mathbf{u} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \quad (7)$$

ja

$$\text{var} \begin{pmatrix} \mathbf{y} \\ \mathbf{u} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{Z}\mathbf{G}\mathbf{Z}^T + \mathbf{R} & \mathbf{G}\mathbf{Z}^T & \mathbf{R} \\ \mathbf{Z}\mathbf{G} & \mathbf{G} & \mathbf{0} \\ \mathbf{R} & \mathbf{0} & \mathbf{R} \end{pmatrix}. \quad (8)$$

Traditsioonilises segamudelis eeldatakse, et nii kõik uuritavad objektid kui ka kõik juhuslikud efektid on omavahel sõltumatud, mistõttu on nii jääkdispersioonimaatriks \mathbf{R} kui ka juhuslike efektide dispersioonimaatriks \mathbf{G} diagonaalmaatriksid:

$$\mathbf{R} = \mathbf{I}_n \sigma_e^2, \quad (9)$$

$$\mathbf{G} = \mathbf{I}_a \sigma_u^2, \quad (10)$$

siin n tähistab uuritavate objektide arvu ja a juhusliku faktori tasemete arvu (enam kui ühe juhusliku faktori korral on maatriks \mathbf{G} blokkdiagonaalne – igale faktorile vastab tema tasemete arvule vastava dimensiooniga diagonaalmaatriks).

Dispersioone σ_e^2 ja σ_u^2 valemeis (9) ja (10) nimetatakse **dispersioonikomponentideks** (*variance components*), sest nad kujutavad enesest vaatluste dispersiooni komponente: $\text{var}(y_i) = \sigma_u^2 + \sigma_e^2$.

Näide 5. Paneme kirja mudeli (5) kovariatsioonimaatriksid (isade mõjusid käsitleme juhuslike efektidena).

Vastavalt valemitele (9) ja (10) saame, et

$$\text{var}(\mathbf{e}) = \mathbf{R} = \begin{pmatrix} \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 \end{pmatrix}, \quad \text{var}(\mathbf{u}) = \begin{pmatrix} \sigma_u^2 & 0 \\ 0 & \sigma_u^2 \end{pmatrix}. \quad (11)$$

Valemist (8) järeldub, et $\text{var}(\mathbf{y}) = \mathbf{V} = \mathbf{ZGZ}^T + \mathbf{R}$, seega

$$\mathbf{V} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \sigma_u^2 & 0 \\ 0 & \sigma_u^2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} + \begin{pmatrix} \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 \end{pmatrix} \\ = \begin{pmatrix} \sigma_e^2 + \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & 0 & 0 & 0 & 0 & 0 \\ \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & 0 & 0 & 0 & 0 & 0 \\ \sigma_u^2 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & \sigma_u^2 & 0 & 0 & 0 & 0 & 0 \\ \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & 0 & 0 & 0 & 0 & 0 \\ \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 \\ 0 & 0 & 0 & 0 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 \\ 0 & 0 & 0 & 0 & \sigma_u^2 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & \sigma_u^2 & \sigma_u^2 \\ 0 & 0 & 0 & 0 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 & \sigma_u^2 \\ 0 & 0 & 0 & 0 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \sigma_e^2 + \sigma_u^2 \end{pmatrix}.$$

Viimasest maatriksist ilmneb ka, et vaatluste dispersioonid esituvad summana $\text{var}(y_i) = \sigma_u^2 + \sigma_e^2$ ja kovariatsioonid kujul

$$\text{cov}(y_i, y_j) = \begin{cases} \sigma_u^2, & \text{kui } i \neq j \text{ ja loomad } i \text{ ja } j \text{ on sama isa järglased,} \\ 0, & \text{kui } i \neq j \text{ ja loomad } i \text{ ja } j \text{ on erinevate isade järglased.} \end{cases}$$

Märkus: käsitledes isade mõjusid fikseerituna, avaldub vaatluste dispersioonimaatriks seosena (6):

$$\text{var}(\mathbf{y}) = \mathbf{V} = \begin{pmatrix} \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \sigma_e^2 \end{pmatrix} = \mathbf{R} = \text{var}(\mathbf{e}),$$

millest omakorda järeldub, et $\text{var}(y_i) = \sigma_e^2$ ja $\text{cov}(y_i, y_j) = 0$, kui $i \neq j$.

5 EELDUSED JA KITSENDUSED

5.1 Jaotused

Üldiste lineaarsete mudelite analüüsil eeldatakse uuritava(te) tunnus(t)e jaotumist vastavalt normaaljaotuse seaduspäradele (pt 2.3.1). Matemaatiliselt väljendudes tähendab see, et vektor \mathbf{y} on normaaljaotusega keskvärtusega $\mathbf{X}\boldsymbol{\beta}$ ja dispersioonimaatriksiga \mathbf{V} :

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}).$$

Taolise eelduse puhul on garanteeritud, et vähimruutude printsiibil (pt 2.4.2) leitud efektide hinnangud on tõepoolest väikseima varieeruvusega, osutub võimalikuks faktorite mõjude täpsuse ja statistilise olulisuse kontrollimine (F -test dispersioon- ja regressioonanalüüsil) ning suurima tõepära meetodi rakendamine dispersioonikomponentide ja seeläbi ka päritavuse ja geneetiliste korrelatsioonide hindamiseks.

5.2 Kovariatsioonistruktuur

Traditsioonilisel üldiste lineaarsete mudelite analüüsil (nagu näiteks ka tavalise regressioon- ja dispersioonanalüüsi korral) eeldatakse dispersioonimaatriksite diagonaalset kuju. Segamudelite tarvis on selline lihtne kovariatsioonistruktuur defineeritud seostega (9) ja (10). Et aga tegelikkus (uuritava(te) tunnus(t)e varieeruvus) sageli nii triviaalselt kirjeldatav pole, leiab reaalse andmete analüüsil kasutust suur hulk rohkem või vähem keerukaid dispersioonistruktuure, mis on faktorite mõjude ja juhuslike vigade sõltumatus eeldusel kokkuvõtvalt esitatavad dispersioonimaatriksina (8).

Fikseeritud mudelite puhul on levinuimad kõiksugu korduvate mõõtmiste analüüsi mudelid, kus püütakse arvesse võtta ühel ja samal objektil sooritatud mõõtmistele vastavate mudeli jääkide korreleeritust.

Näide 6. Olgu meil uuritavaiks objektideks lehmad, kellel kõigil on fikseeritud esimese kolme laktatsiooni piimatoodangud (väljavõte andmetabelist on toodud teksti paremas servas).

Ükskõik, milliste faktortunnuste mõju me ka uurida tahame, on loomulik eeldada, et sama lehma piimatoodangud ei ole sõltumatud. Mudelisse saab viimase kirja panna jääkdispersioonimaatriksi \mathbf{R} struktuuri abil, kus endiselt vastab kõigile vaatlustele juhuslike vigade dispersioon σ_e^2 , aga lisaks võib defineerida (ja lasta arvutil ka andmetest hinnata) samale loomale vastavate mudeli jääkliikmete vahelist varieeruvust peegeldava dispersioonikomponendi σ^2 , misjärel jääkdispersioonimaatriks esitub kompaundsümmeetrilisel kujul

$$\mathbf{R} = \begin{pmatrix} \sigma_e^2 + \sigma^2 & \sigma^2 & \sigma^2 & 0 & 0 & 0 & \dots \\ \sigma^2 & \sigma_e^2 + \sigma^2 & \sigma^2 & 0 & 0 & 0 & \dots \\ \sigma^2 & \sigma^2 & \sigma_e^2 + \sigma^2 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & \sigma_e^2 + \sigma^2 & \sigma^2 & \sigma^2 & \dots \\ 0 & 0 & 0 & \sigma^2 & \sigma_e^2 + \sigma^2 & \sigma^2 & \dots \\ 0 & 0 & 0 & \sigma^2 & \sigma^2 & \sigma_e^2 + \sigma^2 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}.$$

ID	Lakt.	Piim, kg
3396	1	4119
3396	2	5857
3396	3	6260
3990	1	3106
3990	2	3934
3990	3	5171
4390	1	2473
4390	2	3301
4390	3	2958
...		

Sammu võrra keerulisema dispersioonistruktuuri saame, kui eeldame, et mida suurem on mõõtmiste ajaline vahe, seda nõrgem on nende vaheline seos, ehk antud näite puhul – 1. ja 2. laktatsiooni toodangud on omavahel tugevamini seotud, kui 1. ja 3. laktatsiooni toodangud. Üks võimalus seda mudelisse kirja panna on kasutada jääkdispersioonimaatriksi defineerimisel esimest järku autoregressiivset struktuuri kujul

$$\mathbf{R} = \sigma_e^2 \begin{pmatrix} 1 & \rho & \rho^2 & 0 & 0 & 0 & \dots \\ \rho & 1 & \rho & 0 & 0 & 0 & \dots \\ \rho^2 & \rho & 1 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 1 & \rho & \rho^2 & \dots \\ 0 & 0 & 0 & \rho & 1 & \rho & \dots \\ 0 & 0 & 0 & \rho^2 & \rho & 1 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix},$$

kus parameetrit ρ võib mõista kui ühe ja sama looma järjestikuste vaatluste vahelist autokorrelatsioonikordajat, $|\rho| < 1$.

Segamudelite puhul tuleb arvestada ka juhuslike faktorite mõjude võimaliku korreleeritusega, mis tähendab dispersioonimaatriksi \mathbf{G} struktuuri erinevust traditsioonilisest diagonaalsest kujust (10). Loomade aretuses väljendub see enamasti kõikvõimalike sugulussidemete matemaatilisel kujul väljendamises ja mudelisse kaasamises – taolisel kovariatsioonistruktuuri defineerimisel põhineb näiteks aretusväärtuste prognoosimine looma mudelist.

5.3 Täisastakuga ja mitte täisastakuga mudelid, reparametriseerimistingimused

Üldine lineaarne mudel on **täisastakuga**, kui täisastakuga (kõik read-veerud lineaarselt sõltumatud, pööratav) on tema plaanimaatriksi transponeeritud maatriksi ja plaanimaatriksi enese korrutismaatriks, näiteks $\mathbf{X}^T \mathbf{X}$. Nagu järgmisest peatükist nähtub, kujutab just taoline korrutismaatriks enesest mudeli parameetrite hindamiseks konstrueeritud võrrandi kordajatemaatriksit.

Mudeli parameetrid on **üheselt hinnatavad** üksnes täisastakuga mudeli korral, sest vaid siis leidub kasutataval kordajatemaatriksil ühene pöördmaatriks; mitte täisastakuga mudeli korral on erinevaid parameetrite hinnanguid palju (teoreetiliselt lõpmatu arv).

Tavalise fikseeritud mudeli (2) korral piisab mudeli parameetrite üheseks hinnanguks ka plaanimaatriksi enese veergude lineaarsest sõltumatusest. Täisastakuga on näiteks regressioonanalüüsi mudel, sest plaanimaatriksi ühtedest koosnev vabaliikmele vastav veerg pole üldjuhul mingi lineaarkombinatsiooni tulemusena teisendatav argumenttunnuse väärtuste veeruks.

Näide 7. Püstitades näite 2 alusel lihtsa lineaarse regressioonimudeli prognoosimaks tallede võõrutusmassi võõrutusvanuse abil, saame regressioonivõrrandi maatrikskujul

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

kus

$$\mathbf{X} = \begin{pmatrix} 1 & 91 \\ 1 & 90 \\ 1 & 86 \\ 1 & 101 \\ 1 & 86 \\ 1 & 117 \\ 1 & 152 \\ 1 & 110 \\ 1 & 107 \end{pmatrix} \text{ ja } \boldsymbol{\beta} = \begin{pmatrix} \mu \\ b \end{pmatrix}.$$

Plaanimaatriksi \mathbf{X} veeruastak on 2, sest vabaliikmele ja regressioonikordajale vastavad veerud on lineaarselt sõltumatud, ja seega on parameetervektor $\boldsymbol{\beta}$ üheselt hinnatav.

Täisastakuga pole aga dispersioonanalüüsi mudel, sest näiteks mistahes faktori kõigile tasemetele vastavate veergude summeerimine annab tulemuseks ühtedest koosneva vabaliikmele vastava veeru. Seega pole ka dispersioonanalüüsi korral mudeli parameetrite ühene hindamine võimalik. Lahenduseks on hinnata mudeli parameetrite funktsioone, mis on üheselt hinnatavad (nn **hinnatavad funktsioonid** – *estimable functions*, vt pt 6.6.2) või kasutada mudeli parameetrite **reparametriseerimist**. Viimane tähendab parameetritele mingi lisakitsenduse rakendamist, mille all võib mõista näiteks iga faktortunnuse viimase taseme mõju nulliga võrdsustamist ja teiste tasemete mõjude hindamist selle

suhtes (kasutab näiteks statistikapakett SAS ja kasutatakse ka antud õppevahendi järgnevais näidetes) või traditsioonilisema reparametriseerimistingimusena iga faktori mõjude summa nulliga võrdsustamist. Seejuures peab selliseid lisakitsendusi olema sama palju, kui on plaanimaatriksis lineaarselt sõltuvaid veerge.

Näide 8. Mudel (1) näites 2 ei ole täisastakuga, sest plaanimaatriksi

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & 1 & 0 & 0 & 0 & 1 & 91 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 90 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 86 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 101 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 86 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 117 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 152 \\ 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 110 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 107 \end{pmatrix}$$

veeruastak on 6, mis on väiksem maatriksi \mathbf{X} veergude arvust (9) – veergude lineaarne sõltuvus ilmneb selles, et kõigi diskreetsete faktorite korral annavad neile vastavad plaanimaatriksi veerud summeeritult esimese, vabaliikmele vastava ühtedest koosneva veeru.

Faktorite mõjude üheseks hindamiseks on vajalik teatud kitsenduste püstitamine. Näiteks võime nõuda kõigi faktorite korral nende mõjude summa nulliga võrdumist ($I_{1025} + I_{1027} = 0$, $S_1 + S_2 = 0$, $R + P_2 + P_3 = 0$) või siis iga faktori viimase efekti nulliga võrdumist ($I_{1027} = 0$, $S_2 = 0$, $P_3 = 0$).

6 FIKSEERITUD EFEKTIDE HINDAMINE (BLUE)

6.1 Vähimruutude hinnangud

Faktorite mõjud mudelist (2) hinnatakse vähimruutude meetodil, st et parameetervectori β hinnang $\hat{\beta}$ valitakse selliselt, et mudeli vigade ruudud oleks minimaalsed. Matrikskujul on vähimruutude tingimus mudeli (2) tarvis väljendatav seosena

$$\min_{\beta}(\mathbf{e}^T\mathbf{e}) = \min_{\beta}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) = (\mathbf{y} - \mathbf{X}\hat{\beta})^T(\mathbf{y} - \mathbf{X}\hat{\beta}).$$

Hinnangu $\hat{\beta}$ avaldamiseks tuleb mudeli (2) vigade ruutude summast

$$\mathbf{e}^T\mathbf{e} = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) = \mathbf{y}^T\mathbf{y} - \mathbf{y}^T\mathbf{X}\beta - \beta^T\mathbf{X}^T\mathbf{y} + \beta^T\mathbf{X}^T\mathbf{X}\beta = \mathbf{y}^T\mathbf{y} - 2\beta^T\mathbf{X}^T\mathbf{y} + \beta^T\mathbf{X}^T\mathbf{X}\beta$$

võtta tuletis β järgi ja võrdsustada tulemus nulliga. Diferentseerimise tagajärjel saame

$$\partial\mathbf{e}^T\mathbf{e}/\partial\beta = -2\mathbf{X}^T\mathbf{y} + 2\mathbf{X}^T\mathbf{X}\beta$$

ning viimase avaldise nulliga võrdsustamisest järeldub, et

$$(\mathbf{X}^T\mathbf{X})\hat{\beta} = \mathbf{X}^T\mathbf{y}, \quad (12)$$

millest

$$\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}. \quad (13)$$

Võrrand (12) on tuntud kui **normaalvõrrand** (*normal equation*) ja hinnang (13) kujutab enesest parameetervectori β **vähimruutude hinnangut** (OLS, *ordinary least squares*), mis diagonaalse võrdsete dispersioonikomponentidega jääkdispersioonimaatriksi (9) korral on **parim lineaarne nihketa hinnang** (BLUE, *best linear unbiased estimator*). Sõna “parim” hinnangu nimetuses tähendab, et tegu on täpseima e väikseima dispersiooniga $\text{var}(\beta - \hat{\beta})$ hinnanguga (efektiivse hinnanguga, pt 2.4.3) ning sõna “nihketa” märgib hinnangu nn keskmist õigsust (süsteemaatilise vea puudumist).

Kui meil on aga tegu keerulisema kovariatsioonistruktuuriga, kus $\mathbf{V} = \mathbf{R} = \text{var}(\mathbf{e})$ ei ole võrdsete dispersioonikomponentidega diagonaalmaatriks, siis pole ka valemiga (13) defineeritud hinnang enam BLUE parameetervectorile β . Hädast aitab välja kaval teisendus. Nimelt, kui teame dispersioonimaatriksit \mathbf{V} , mis tänu dispersiooni definitsioonile on alati positiivselt määratud ja mistõttu leiduvad ka \mathbf{V}^{-1} ja selline maatriks $\mathbf{V}^{1/2}$, et $\mathbf{V} = \mathbf{V}^{1/2}\mathbf{V}^{1/2}$, võime võrrandi (2) mõlemaid pooli vasakult korrutada maatriksiga $\mathbf{V}^{-1/2}$, saades tulemuseks mudeli

$$\mathbf{V}^{-1/2}\mathbf{y} = \mathbf{V}^{-1/2}\mathbf{X}\beta + \mathbf{V}^{-1/2}\mathbf{e}.$$

Tähistades $\mathbf{y}^* = \mathbf{V}^{-1/2}\mathbf{y}$, $\mathbf{X}^* = \mathbf{V}^{-1/2}\mathbf{X}$ ja $\mathbf{e}^* = \mathbf{V}^{-1/2}\mathbf{e}$ saame viimase mudeli esitada võrrandina

$$\mathbf{y}^* = \mathbf{X}^*\beta + \mathbf{e}^*, \quad (14)$$

kusjuures

$$\mathbf{V}^* = \text{var}(\mathbf{e}^*) = \text{var}(\mathbf{V}^{-1/2}\mathbf{e}) = \mathbf{V}^{-1/2}\text{var}(\mathbf{e})\mathbf{V}^{-1/2} = \underbrace{(\mathbf{V}^{-1/2}\mathbf{V}^{1/2})}_{\mathbf{I}}\underbrace{(\mathbf{V}^{-1/2}\mathbf{V}^{1/2})}_{\mathbf{I}} = \mathbf{I}.$$

Et nüüd on meil dispersioonimaatriks \mathbf{V}^* diagonaalsel kujul, nagu on vaja BLUE esitamiseks valemi (13) abil, rakendame viimast mudelile (14) ja saame parameetervectori hinnangu kujul

$$\hat{\beta} = (\mathbf{X}^{*T}\mathbf{X}^*)^{-1}\mathbf{X}^{*T}\mathbf{y}^* = (\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}. \quad (15)$$

Viimast hinnangut nimetatakse **üldistatud vähimruutude hinnanguks** (GLS, *generalized least squares*).

Juhul, kui analüüsitav mudel ei ole täisastakuga, ei ole normaalvõrrandi kordajate maatriksid $\mathbf{X}^T\mathbf{X}$ ja $\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X}$ (vastavalt tavalise ja üldistatud vähimruutude hinnangu puhul) pööratavad ja kasutada tuleb üldistatud pöördmaatrikseid. Vähimruutude hinnangud (13) ja (15) saavad siis üldisemad kujud

$$\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \quad (16)$$

ja

$$\hat{\beta} = (\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}. \quad (17)$$

Näide 9. Vaatleme edasi tallede võõrutusmassi andmestikku ja hindame talle isa, soo, pesakonna suuruse ning võõrutusvanuse kui fikseeritud faktorite mõjud. Garanteerimaks

hinnangute ühesust, võrdsustame nulliga iga faktori viimase efekti ($I_{1027} = 0$, $S_2 = 0$, $P_3 = 0$). Ülejäänud efektide hindamiseks paneme kirja täisastakuga (ilma nulliga võrdsustatud efektidele vastavate veergudeta) plaanimaatriksi \mathbf{X}^* :

$$\mathbf{X}^* = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 91 \\ 1 & 1 & 1 & 1 & 0 & 90 \\ 1 & 1 & 1 & 0 & 1 & 86 \\ 1 & 1 & 0 & 0 & 1 & 101 \\ 1 & 1 & 0 & 0 & 1 & 86 \\ 1 & 0 & 1 & 0 & 1 & 117 \\ 1 & 0 & 1 & 0 & 1 & 152 \\ 1 & 0 & 0 & 1 & 0 & 110 \\ 1 & 0 & 0 & 0 & 1 & 107 \end{pmatrix}.$$

Et meil vaatluste dispersioonimaatriks on diagonaalsel kujul, $\mathbf{V} = \mathbf{I}_9 \sigma_e^2$, ja plaanimaatriks \mathbf{X}^* on täisveeruastakuga, saame faktorite mõjude vektori

$$\boldsymbol{\beta}^* = (\mu \ I_{1025} \ S_1 \ P_1 \ P_2 \ b)^T$$

hindamiseks kasutada valemit (13):

$$\underbrace{\begin{pmatrix} \hat{\mu} \\ \hat{I}_{1025} \\ \hat{S}_1 \\ \hat{P}_1 \\ \hat{P}_2 \\ \hat{b} \end{pmatrix}}_{\hat{\boldsymbol{\beta}}^*} = \underbrace{\begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ 91 & 90 & 86 & 101 & 86 & 117 & 152 & 110 & 107 \end{pmatrix}}_{\mathbf{X}^T} \underbrace{\begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 91 \\ 1 & 1 & 1 & 1 & 0 & 90 \\ 1 & 1 & 1 & 0 & 1 & 86 \\ 1 & 1 & 0 & 0 & 1 & 101 \\ 1 & 1 & 0 & 0 & 1 & 86 \\ 1 & 0 & 1 & 0 & 1 & 117 \\ 1 & 0 & 1 & 0 & 1 & 152 \\ 1 & 0 & 0 & 1 & 0 & 110 \\ 1 & 0 & 0 & 0 & 1 & 107 \end{pmatrix}}_{\mathbf{X}}^{-1}$$

$$\times \underbrace{\begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ 91 & 90 & 86 & 101 & 86 & 117 & 152 & 110 & 107 \end{pmatrix}}_{\mathbf{X}^T} \underbrace{\begin{pmatrix} 22 \\ 39 \\ 24 \\ 27 \\ 24 \\ 31 \\ 30 \\ 31 \\ 26 \end{pmatrix}}_{\mathbf{y}} = \begin{pmatrix} 11,38 \\ 1,10 \\ 3,30 \\ 14,59 \\ 6,03 \\ 0,07 \end{pmatrix}.$$

Peale kitsenduste lisamist saame mudeli (1) parameetrite hinnangud kujul

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\mu} \\ \hat{I}_{1025} \\ \hat{I}_{1027} \\ \hat{S}_1 \\ \hat{S}_2 \\ \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_3 \\ \hat{b} \end{pmatrix} = \begin{pmatrix} 11,4 \\ 1,1 \\ 0 \\ 3,3 \\ 0 \\ 14,6 \\ 6,0 \\ 0 \\ 0,07 \end{pmatrix}. \quad (18)$$

Võrreldes faktorite tasemete kaupa leitud aritmeetiliste keskmistega (vt tabelit näite lõpus) hakkab silma jäärade paremusjärjestuse muutus – kui järglaste keskmiste väärtuste alusel võinuks eelistada jäära numbriga 1027, siis peale teiste faktorite mõjude arvesse võtmist on ilmselge hoopis jäära nr 1025 paremus. Andmetabelist (pt 6.3.1) põhjust otsides hakkab eelkõige silma, et jäära nr 1027 järglased on võõrutatud hiljem ja saanud seega kauem aega kasvada (positiivset seost võõrutusvanuse ja võõrutusmassi vahel näitab ka regressioonikordaja positiivne väärtus, $\hat{b}=0,07$). Jäärade mõjude hindamine mudeli (1) abil võimaldab nende seostega arvestada.

Faktorite tasemed ja neile vastavad keskmised võõrutusmassid:

Isa		Sugu		Pesakonna suurus	
1025	27,2	1	29,2	1	35,0
1027	29,5	2	27,0	2	27,0
				3	22,0

6.2 Hinnatavad funktsioonid

Olgugi, et üldjuhul parameetervektor β ise ei ole üheselt hinnatav, on seda hulk tema elementide lineaarkombinatsioone $k^T\beta$. Viimaseid nimetatakse **hinnatavateks funktsioonideks** (*estimable functions*). Hinnatavuse olulisemad omadused lineaarsete mudelite korral on, et

1. vabaliige μ ei ole üheselt hinnatav (va regressioonanalüüsi korral);
2. erinevused ühe faktori tasemete vahel on üheselt hinnatavad tingimusel, et mudelis puuduvad selle ja teiste faktorite vahelised interaktsioonid ning uuritava faktori tasemete mõjud ei kattu mõne teise faktori tasemete mõjudega;
3. sisulist tähendust omavad üksnes hinnatavate funktsioonide väärtused.

Näide 10. Et tallede võõrutusmassi mõjutavate faktorite tasemete hinnangute ei ole üheselt hinnatavad, saame hinnangutele $\hat{\beta}_1$ (18) lisaks leida suure hulga erinevaid hinnanguid. Näiteks järgmised hinnangute vektorid on saadud kasutades erinevaid reparametriseerimistingimusi või rakendades erinevaid üldistatud pöördmaatrikseid valemis (16):

$$\hat{\beta}_2 = \begin{pmatrix} \hat{\mu} \\ \hat{I}_{1025} \\ \hat{I}_{1027} \\ \hat{S}_1 \\ \hat{S}_2 \\ \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_3 \\ \hat{b} \end{pmatrix} = \begin{pmatrix} 30,4 \\ 0 \\ -1,1 \\ 0 \\ -3,3 \\ 0 \\ -8,6 \\ -14,6 \\ 0,07 \end{pmatrix}, \quad \hat{\beta}_3 = \begin{pmatrix} \hat{\mu} \\ \hat{I}_{1025} \\ \hat{I}_{1027} \\ \hat{S}_1 \\ \hat{S}_2 \\ \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_3 \\ \hat{b} \end{pmatrix} = \begin{pmatrix} 20,45 \\ 0,55 \\ -0,55 \\ 1,65 \\ -1,65 \\ 7,71 \\ -0,84 \\ -6,87 \\ 0,07 \end{pmatrix}, \quad \hat{\beta}_4 = \begin{pmatrix} \hat{\mu} \\ \hat{I}_{1025} \\ \hat{I}_{1027} \\ \hat{S}_1 \\ \hat{S}_2 \\ \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_3 \\ \hat{b} \end{pmatrix} = \begin{pmatrix} 20,7 \\ 1,1 \\ 0 \\ 0 \\ -3,3 \\ 8,5 \\ 0 \\ -6,0 \\ 0,07 \end{pmatrix} \text{ jne.}$$

Hinnatavad funktsioonid on näiteks

$$\underbrace{(0 \mid 0 \mid 0 \mid 1 \mid -1 \mid 0 \mid 0 \mid 0 \mid 0)}_{k_1^T} \hat{\beta} = \hat{S}_1 - \hat{S}_2 = 0,$$

mis hindab sugude vahelist erinevust ja on sama väärtusega mistahes parameetrite hinnangute vektori $\hat{\beta}$ korral, või

$$\underbrace{(0 \mid 0 \mid 0 \mid 0 \mid 0 \mid 0 \mid 1 \mid -1 \mid 0)}_{k_2^T} \hat{\beta} = \hat{P}_2 - \hat{P}_3 = -2,$$

mis hindab kahe- ja kolmetallelisest pesakonnast pärit tallede võõrutusmasside erinevust ja on jällegi ühene mistahes parameetrite hinnangute vektori $\hat{\beta}$ korral.

6.3 Vähimruutkeskmised

Vähimruutkeskmised (VRK; inglise keeles *least square means*, LSM), mis esitatakse sageli teadusartiklites ja mille hindamisvõimalust pakuvad paljud statistikaprogrammid, on oma olemuselt samuti hinnatavad parameeterfunktsioonid. Nende laialdane kasutamine on tingitud sellest, et ühelt poolt on tegu tõepoolest lihtsalt mõistetavate hinnangutega uuritava tunnuse keskmistele väärtustele faktorite tasemete alusel moodustatud gruppides, teiselt poolt on aga vähimruutkeskmiste hindamise meetodika märksa täiuslikum võrreldes tavaliste aritmeetiliste keskmiste arvutamisega – VRK-d leitakse baseeruvana mudelil ja on korrigeeritud mudelis sisalduvate segavate faktorite mõjude suhtes ega ole tundlikud vaatluste arvu erinevusele võrreldavais gruppides (seetõttu erinevad vähimruutkeskmised üldjuhul tavalistest keskmistest).

Näide 11. Illustreerime vähimruutkeskmiste leidmist tallede võõrutusmassi näite varal. Üks variant mudeli (1) parameetrite hinnanguist on esitatud vektorina (18).

Kõikvõimalikud vähimruutkeskmised koos arvutusvalemitega (ja võrdluseks ka tavalised aritmeetilised keskmised) on esitatud järgnevas tabelis.

Efekt	\bar{x}	VRK
I_{1025}	27,2	$LSM(\hat{I}_{1025}) = \hat{\mu} + \hat{I}_{1025} + \frac{1}{2}(\hat{S}_1 + \hat{S}_2) + \frac{1}{3}(\hat{P}_1 + \hat{P}_2 + \hat{P}_3) + \hat{b}_1 \times \bar{v}$ $= 11,4 + 1,1 + \frac{1}{2}(3,3 + 0) + \frac{1}{3}(14,6 + 6,0 + 0) + 0,07 \times 104,4 = 28,14$
I_{1027}	29,5	$LSM(\hat{I}_{1027}) = \hat{\mu} + \hat{I}_{1027} + \frac{1}{2}(\hat{S}_1 + \hat{S}_2) + \frac{1}{3}(\hat{P}_1 + \hat{P}_2 + \hat{P}_3) + \hat{b}_1 \times \bar{v} + \hat{b}_2 \times \bar{v}$ $= 11,4 + 0 + \frac{1}{2}(3,3 + 0) + \frac{1}{3}(14,6 + 6,0 + 0) + 0,07 \times 104,4 = 27,04$
S_1	29,2	$LSM(\hat{S}_1) = \hat{\mu} + \frac{1}{2}(\hat{I}_{1025} + \hat{I}_{1027}) + \hat{S}_1 + \frac{1}{3}(\hat{P}_1 + \hat{P}_2 + \hat{P}_3) + \hat{b}_1 \times \bar{v} = 29,24$
S_2	27,0	$LSM(\hat{S}_2) = \hat{\mu} + \frac{1}{2}(\hat{I}_{1025} + \hat{I}_{1027}) + \hat{S}_2 + \frac{1}{3}(\hat{P}_1 + \hat{P}_2 + \hat{P}_3) + \hat{b}_1 \times \bar{v} = 25,94$
P_1	35,0	$LSM(\hat{P}_1) = \hat{\mu} + \frac{1}{2}(\hat{I}_{1025} + \hat{I}_{1027}) + \frac{1}{2}(\hat{S}_1 + \hat{S}_2) + \hat{P}_1 + \hat{b}_1 \times \bar{v} = 35,30$
P_2	27,0	$LSM(\hat{P}_2) = \hat{\mu} + \frac{1}{2}(\hat{I}_{1025} + \hat{I}_{1027}) + \frac{1}{2}(\hat{S}_1 + \hat{S}_2) + \hat{P}_2 + \hat{b}_1 \times \bar{v} = 26,75$
P_3	22,0	$LSM(\hat{P}_3) = \hat{\mu} + \frac{1}{2}(\hat{I}_{1025} + \hat{I}_{1027}) + \frac{1}{2}(\hat{S}_1 + \hat{S}_2) + \hat{P}_3 + \hat{b}_1 \times \bar{v} = 20,72$

7 JUHUSLIKE FAKTORITE REALISEERUNUD VÄÄRTUSTE PROGNOOSIMINE (BLUP)

Et juhuslikud efektid pole konstandid, nagu fikseeritud efektid, vaid kujutavad enesest mingist jaotusest pärinevaid juhusliku suuruse (juhusliku faktori) realiseerunud väärtusi, tuleb nende hindamisel arvestada ka selle jaotuse dispersiooniga (ehk teisiti väljendatuna, juhuslike efektide poolt määratud osaga uuritava tunnuse koguvarieeruvusest). Mudeli (4) korral tähendab see dispersioonimaatriksite (8) teadmist.

Juhul, kui dispersioonimaatriks (8) on teada, avaldub juhuslike faktorite realiseerunud väärtuste \mathbf{u} parim lineaarne nihketa prognoos (BLUP, *best linear unbiased prediction*) kujul

$$\hat{\mathbf{u}} = \mathbf{GZ}^T\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}), \quad (19)$$

kus $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}$ on fikseeritud efektide üldine vähimruutude hinnang. Seejuures tehakse siin ingliskeelses kirjanduses selget vahet sõnadel hinnang (*estimate*) ja prognoos (*prediction*) – hinnatakse midagi fikseeritud, konstantset, prognoositakse aga midagi, mille ilmumine sõltub juhusest.

Fikseeritud efektide hindamine ning juhuslike efektide prognoosimine valemite (17) ja (19) alusel nõuab vaatluste dispersioonimaatriksi \mathbf{V} (8) pöördmaatriksi leidmist, milline operatsioon võib vähegi suurema andmestiku korral problemaatiliseks kujuneda. Üks 20. sajandi suurimaid segamudelite ja aretusteooria arendajaid C. Henderson pakkus 1950. aastal välja normaalvõrrandile (12) sarnase maatriksvõrduse, mis võimaldab korraga leida nii BLUE($\boldsymbol{\beta}$) kui ka BLUP(\mathbf{u}) ja seda ilma vaatluste dispersioonimaatriksit pööramata. Tänapäeval tuntaksegi järgnevat maatriksvõrdust **Hendersoni segamudeli võrrandina** (või lihtsalt segamudeli võrrandina, inglise keeles *mixed model equation*, MME):

$$\begin{pmatrix} \mathbf{X}^T\mathbf{R}^{-1}\mathbf{X} & \mathbf{X}^T\mathbf{R}^{-1}\mathbf{Z} \\ \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{X} & \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{y} \end{pmatrix}. \quad (20)$$

Traditsioonilise segamudeli korral, kus vealiikmele ja juhuslikele efektidele vastavad dispersioonimaatriksid \mathbf{R} ja \mathbf{G} esituvad kujul (9) ja (10), on võrrand (20) kirjutatav ka kujul

$$\begin{pmatrix} \mathbf{X}^T\mathbf{X} & \mathbf{X}^T\mathbf{Z} \\ \mathbf{Z}^T\mathbf{X} & \mathbf{Z}^T\mathbf{Z} + \frac{\sigma_e^2}{\sigma_u^2}\mathbf{I} \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T\mathbf{y} \\ \mathbf{Z}^T\mathbf{y} \end{pmatrix}, \quad (21)$$

millest

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \text{BLUE}(\boldsymbol{\beta}) \\ \text{BLUP}(\mathbf{u}) \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T\mathbf{X} & \mathbf{X}^T\mathbf{Z} \\ \mathbf{Z}^T\mathbf{X} & \mathbf{Z}^T\mathbf{Z} + \frac{\sigma_e^2}{\sigma_u^2}\mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}^T\mathbf{y} \\ \mathbf{Z}^T\mathbf{y} \end{pmatrix}.$$

Kuna praktikas sageli dispersioonimaatrikseid võrrandis (20) või dispersioonikomponente võrrandis (21) teada pole, tuleb need enne efektide prognoosimist andmetest hinnata, misjärel tehakse segamudeli võrrandis asendused $\mathbf{G} = \hat{\mathbf{G}}$ ja $\mathbf{R} = \hat{\mathbf{R}}$ (ehk $\sigma_u^2 = \hat{\sigma}_u^2$ ja $\sigma_e^2 = \hat{\sigma}_e^2$). On tõestatud, et juhul, kui uuritav tunnus on sümmeetrilise jaotusega (näiteks normaaljaotusega), on selliselt saadud fikseeritud efektide hinnangud ning juhuslike efektide prognoosid nihketa hinnangud BLUE($\boldsymbol{\beta}$)-le ja BLUP(\mathbf{u})-le.

Segamudeli võrrandis on juhuslikele efektidele vastavad read/veerud lineaarselt sõltumatud ega vaja hinnangu ühesuseks erinevalt fikseeritud efektidest lisakitsenduste rakendamist (lisakitsendused võivad olla endiselt vajalikud fikseeritud efektide hindamisel).

Andmaks pisut ettekujutust, mis vahe on faktori mõjude hinnanguil sõltuvalt sellest, kas käsitleda faktorit fikseerituna või juhuslikuna, vaatleme lihtsaimat dispersioonanalüüsi mudelit, mis sisaldab vaid vabaliiget μ ja üht diskreetset faktorit u :

$$y_{ij} = \mu + u_i + e_{ij}. \quad (22)$$

Käsitledes faktorit u fikseerituna, võime parima lineaarse nihketa hinnangu [valem (13)] faktori u tasemele u_i kirjutada kujul

$$\hat{u}_i = \bar{y}_i - \mu, \quad (23)$$

kus $\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$ ja n_i tähistab vaatluste arvu faktori i . tasemel – seega on vaid üht faktorit sisaldava dispersioonanalüüsi mudeli korral iga faktori mõju hinnanguks temale vastavate uuritava tunnuse väärtuste aritmeetilise keskmise erinevus üldkeskmisest μ .

Lugedes aga u juhuslikuks, avaldub tema i . taseme mõju u_i parim lineaarne nihketa prognoos (19) seosena

$$\tilde{u}_i = \frac{n_i \sigma_u^2}{\sigma_e^2 + n_i \sigma_u^2} (\bar{y}_i - \mu). \quad (24)$$

Valemite (23) ja (24) võrdlemisel hakkab silma, et juhuslike efektide hinnangud on alati väiksemad, võrreldes fikseeritud efektide hinnangutega (sest $n_i \sigma_u^2 / (\sigma_e^2 + n_i \sigma_u^2) \leq 1$). Seejuures kehtib intuiitiivselt mõistetav loogika – mida tugevam on faktori mõju, seda sarnasemad on fikseeritud ja juhusliku mudeli eeldusel leitud efektide hinnangud (seda suurem on dispersioonikomponendi σ_u^2 väärtus võrreldes jääkvarieeruvusega σ_e^2). Samuti nähtub valemist (23) ja (24), et faktori i . tasemel sooritatud mõõtmiste arvu n_i suurenedes väheneb erinevus selle taseme mõju hinnangute \hat{u}_i ja \tilde{u}_i (leituna vastavalt kas fikseeritud või juhuslikust mudelist) vahel.

Näide 12. Käsitleme nüüd mudeliga (1) kirjeldatud tallede võõrutusmassi näiteanalüüsis jäära mõju juhuslikuna (andmestikus oleva kahe jäära ja nende kokku üheksa järglase alusel hinnatavad isa mõjud esindavad vaid juhuslikku väljavõtet kõigist isalt järglasele päranduda võivatest geenikomplektidest). Mudeli fikseeritud ja juhuslikele efektidele vastavad plaanimaatriksid on

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 1 & 91 \\ 1 & 1 & 0 & 1 & 0 & 0 & 90 \\ 1 & 1 & 0 & 0 & 1 & 0 & 86 \\ 1 & 0 & 1 & 0 & 1 & 0 & 101 \\ 1 & 0 & 1 & 0 & 1 & 0 & 86 \\ 1 & 1 & 0 & 0 & 1 & 0 & 117 \\ 1 & 1 & 0 & 0 & 1 & 0 & 152 \\ 1 & 0 & 1 & 1 & 0 & 0 & 110 \\ 1 & 0 & 1 & 0 & 1 & 0 & 107 \end{pmatrix} \quad \text{ja} \quad \mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

ning hinnatavate parameetrite vektorid esituvad kujul

$$\boldsymbol{\beta} = \begin{pmatrix} \mu \\ S_1 \\ S_2 \\ P_1 \\ P_2 \\ P_3 \\ b_1 \end{pmatrix} \quad \text{ja} \quad \mathbf{u} = \begin{pmatrix} I_{1025} \\ I_{1027} \end{pmatrix}.$$

Üheste lahendite saamise huvides võrdsustame nulliga fikseeritud faktorite viimaste tasemete mõjud ($S_2 = 0$, $P_3 = 0$) ning eemaldame neile vastavad veerud plaanimaatriksist \mathbf{X} , saades maatriksthete läbiviimisel kasutatava täisastakuga plaanimaatriksi

$$\mathbf{X}^* = \begin{pmatrix} 1 & 1 & 0 & 0 & 91 \\ 1 & 1 & 1 & 0 & 90 \\ 1 & 1 & 0 & 1 & 86 \\ 1 & 0 & 0 & 1 & 101 \\ 1 & 0 & 0 & 1 & 86 \\ 1 & 1 & 0 & 1 & 117 \\ 1 & 1 & 0 & 1 & 152 \\ 1 & 0 & 1 & 0 & 110 \\ 1 & 0 & 0 & 1 & 107 \end{pmatrix}.$$

Et juhuslike efektide \mathbf{u} ja \mathbf{e} dispersioonimaatriksid on diagonaalsel kujul (11), on mudeli parameetrid hinnatavad segamudeli võrrandist (21). Viimase rakendamiseks vajaliku dispersioonikomponentide suhte σ_e^2 / σ_u^2 saame avaldada tallede võõrutusmassi päritavusest h^2 , lugedes viimase teadaolevaks:

$$\frac{\sigma_e^2}{\sigma_u^2} = \frac{4 - h^2}{h^2}.$$

Võttes $h^2 = 0,2$, saame $\sigma_e^2 / \sigma_u^2 = 19$.

Seega avaldub segamudeli võrrandi (21) kordajate maatriksi parem alumine nurk $\mathbf{Z}^T \mathbf{Z} + \sigma_e^2 / \sigma_u^2 \mathbf{I}$ kujul

$$\begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} + 19 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 5+19 & 0 \\ 0 & 4+19 \end{pmatrix} = \begin{pmatrix} 24 & 0 \\ 0 & 23 \end{pmatrix}$$

ning kogu segamudeli võrrand on peale plaanimaatriksite korrutamistehete sooritamist järgmine:

$$\begin{pmatrix} 9 & 5 & 2 & 6 & 940 & 5 & 4 \\ 5 & 5 & 1 & 3 & 536 & 3 & 2 \\ 2 & 1 & 2 & 0 & 200 & 1 & 1 \\ 6 & 3 & 0 & 6 & 649 & 3 & 3 \\ 940 & 536 & 200 & 649 & 101716 & 454 & 486 \\ \hline 5 & 3 & 1 & 3 & 454 & 24 & 0 \\ 4 & 2 & 1 & 3 & 486 & 0 & 23 \end{pmatrix} \begin{pmatrix} \mu \\ S_1 \\ P_1 \\ P_2 \\ b_1 \\ \hline I_{1025} \\ I_{1027} \end{pmatrix} = \begin{pmatrix} 254 \\ 146 \\ 70 \\ 162 \\ 26746 \\ \hline 136 \\ 118 \end{pmatrix}.$$

Selle võrrandi lahendina saame parameetrite hinnangute vektori

$$\begin{pmatrix} \hat{\mu} \\ \hat{S}_1 \\ \hat{P}_1 \\ \hat{P}_2 \\ \hat{b}_1 \\ \hline \hat{I}_{1025} \\ \hat{I}_{1027} \end{pmatrix} = \begin{pmatrix} 14,12 \\ 3,51 \\ 14,37 \\ 5,98 \\ 0,05 \\ \hline -0,04 \\ 0,04 \end{pmatrix}.$$

Võrreldes fikseeritud mudelist saadud hinnangutega (18) on ootuspäraselt vähenenud jääradevaheline erinevus – on ju jäära mõjuga seletatav võõrutusmasside erinevus 19 korda väiksem juhusliku vea arvele jäänud varieeruvusest ($\sigma_e^2/\sigma_u^2 = 19$), mis, nagu nähtus valemeist (23) ja (24), toob kaasa väiksemad juhuslike isamõjude hinnangud. Loomulikult muutusid tänu mudeli muutmisele mingil määral ka fikseeritud efektide hinnangud, aga faktorite tasemete järjestus jäi endiseks – jäärtallede võõrutusmass on suurem; suurima võõrutusmassiga on üksiktalled, seejärel kaksik- ja kolmiktalled; võõrutusvanuse kasvades suureneb ka tallede võõrutusmass.

Eelnevas näites rakendatud mudelit, kus uuritava tunnuse geneetilist determineeritust väljendava juhusliku faktorina käsitletakse isa ning ülejäänud faktorid on fikseeritud ja kirjeldavad mitte aditiiv-geneetilisi mõjusid, nimetatakse põllumajandusloomade aretuses **isa mudeliks** (inglise keeles *sire model*). Lähemalt isa mudelist ja ka teistest segamudeli rakendustest geneetiliste parameetrite hindamisel järgmises peatükis.

8 DISPERSIOONIKOMPONENTIDE HINDAMINE *

Nii fikseeritud efektide hindamine üldistatud vähimruutude meetodil (17), juhuslike efektide parim lineaarne nihketa prognoos (19) kui ka Hendersoni segamudeli võrrandi (20) või (21) lahendamine eeldab dispersioonikomponentide kaudu defineeritud dispersioonimaatriksite teadmist. Nagu nähtus eelneva peatüki näitest, piisab sageli ka dispersioonikomponentide suhte või mõne populatsiooni geneetilise parameetri väärtuse teadmisest. Samas aga ei ole meil kuidagi võimalik teada dispersioonikomponentide suhet kõigi ette tulla võivate tunnuste või mudelite korral ning ka meie poolt analüüsitava populatsiooni geneetiline struktuur ei pruugi vastata teadaolevate geneetiliste parameetritega kirjeldatud populatsioonile. Valede dispersiooniparameetrite kasutamisega kaasnevad aga ka valed mõjuhindangud, aretusväärtused ja seeläbi ka valed otsused tuleviku strateegiate osas.

Dispersioonikomponentide hindamiseks on välja töötatud hulk erinevaid meetodeid. Nende mitmesus on tingitud sellest, et vähegi keerulisema mudeli puhul ei ole võimalik välja kirjutada dispersioonikomponentide parimate omadustega analüütilisi lahendeid, kasutada tuleb kas teatud lihtsustavaid eeldusi või siis ligilähedasi matemaatilisi arvutusmeetodeid. Meetodite matemaatiline keerukus on ka põhjuseks, miks üldiste lineaarsete segamudelite rakenduste tudeerimisel dispersioonikomponentide hindamine kõrvale jäetakse. Samas baseerub peatükis 3 tutvustatud peamiste populatsioonigeneetiliste parameetrite hindamine enamasti just dispersioonikomponentide hindamisel.

Andmete tasakaalulisuse (vt pt 6.2.4) ja vaatluste sõltumatuse eeldusel on dispersioonikomponentide hinnangud leitavad tavalise dispersioonanalüüsi abil. Mittetasakaalulise andmestiku ja keerulisema kovariatsioonistruktuuri korral on kõige universaalsemaks osutunud REML-meetod.

8.1 Dispersioonanalüüsi meetod (ANOVA-meetod)

Võtame vaatluse alla ühefaktorilise dispersioonanalüüsi mudeli (22). Eeldame et faktor u on juhuslik ning et efektid u_i ja e_{ij} on sõltumatud ja normaaljaotusega. Samuti eeldame, et faktoril u on a taset, igal tasemel on sooritatud täpselt n mõõtmist ning kokku on mõõtmisi $N = an$. Dispersioonanalüüsi arvutused koondatakse nn dispersioonanalüüsi tabelisse, mis vaatlusaluse mudeli tarvis on toodud tabelis 1.

Tabel 1. Dispersioonanalüüs tabel mudeli (22) korral

Varieeruvuse allikas	Ruutude summa *	Vabadusastmete arv	Keskruut	Keskruudu ooteväärtus
Faktor u	$SS(u) = \sum_{i=1}^a \frac{y_i^2}{n} - \frac{y_{..}^2}{N}$	$a - 1$	$MS(u) = SS(u)/a - 1$	$n\sigma_u^2 + \sigma_e^2$
Jääk	$SS(e) = \sum_{i=1}^a \sum_{j=1}^n y_{ij}^2 - \sum_{i=1}^a \frac{y_i^2}{n}$	$N - a$	$MS(e) = SS(e)/N - a$	σ_e^2

* siin $y_i = \sum_{j=1}^n y_{ij}$ ja $y_{..} = \sum_{i=1}^a \sum_{j=1}^n y_{ij}$

ANOVA-hinnangud dispersioonikomponentidele σ_u^2 ja σ_e^2 saadakse, võrdsustades keskruudud nende ooteväärtustega:

$$\hat{\sigma}_u^2 = \frac{1}{n} [MS(u) - MS(e)]$$

ja

$$\hat{\sigma}_e^2 = MS(e).$$

Mittetasakaaluliste andmete ja ühefaktorilise dispersioonanalüüsi mudeli (22) korral on juhusliku faktori poolt määratud osa uuritava tunnuse varieeruvusest, σ_u^2 , ANOVA-meetodil hinnatav valemist

$$\hat{\sigma}_u^2 = \frac{1}{d} [MS(u) - MS(e)],$$

kus $d = \frac{1}{a-1} \left(N - \frac{1}{N} \sum_{i=1}^a n_i^2 \right)$ ja n_i tähistab mõõtmiste arvu faktori u tasemel i .

8.2 REML-meetod

REML-meetod sobib dispersioonikomponentide hindamiseks ükskõik millise dispersioonistruktuuriga mudelite korral, kusjuures tasakaalulise traditsioonilise segamudeli korral on dispersioonikomponentide REML-hinnangud võrdsed ANOVA-hinnangutega. Lühend REML tähendab **kitsendatud e jääkide e vähendatud suurima tõepära meetodit** (*REstricted e REsidual e REDuced Maximum Likelihood method*) ja baseerub nagu suurima tõepära meetodid ikka mingi teoreetilise jaotuse tõepärafunktsioonil (vt ka pt 2.4.1). Üldise lineaarse mudeli (4) korral on selleks teoreetiliseks jaotuseks normaaljaotus $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$. Kitsendatud suurima tõepära meetodi idee seisneb selles, et kuigi üldine lineaarne mudel sisaldab kaht tüüpi hinnatavaid parameetreid – fikseeritud faktorite efekte ja juhuslike faktorite dispersiooniparameetreid – püütakse viimaste hindamiseks teisendada esialgset mudelit kujule, kus juhuslike faktorite mõjudest tingitud varieeruvus oleks sama, aga fikseeritud faktorite mõju võrduks nulliga. Sellisel juhul on kogu arvutusprotsess suunatud dispersioonikomponentide hindamisele, mis peaks garanteerima täpsemad hinnangud võrreldes tavalise suurima tõepära meetodiga (ML, *Maximum Likelihood method*), mis hindab samaaegselt nii fikseeritud faktorite mõjusid kui ka dispersioonikomponente. Mudeli (4) teisendust võib esitada korrutisena

$$\mathbf{K}^T \mathbf{y} = \mathbf{K}^T \mathbf{X} \boldsymbol{\beta} + \mathbf{K}^T \mathbf{Z} \mathbf{u}$$

kus maatriks \mathbf{K} on defineeritud selliselt, et $\mathbf{K}^T \mathbf{X} = \mathbf{0}$, mistõttu

$$\mathbf{K}^T \mathbf{y} \sim N(\mathbf{0}, \mathbf{K}^T \mathbf{V} \mathbf{K}).$$

Teisendatud mudeli vasakut poolt $\mathbf{K}^T \mathbf{y}$ võib mõista ka kui fikseeritud efektide poolt kirjeldamata jäänud osa uuritavast tunnusest \mathbf{y} e erinevust (jääki) fikseeritud efektide poolt kirjeldatu ja uuritava tunnuse \mathbf{y} tegelike väärtuste vahel (siit ka nimetus *REsidual ML method*).

REML-meetodi korral maksimeeritakse teisendatud uuritava tunnuse logaritmiline tõepärafunktsioon

$$l(\mathbf{V}; \mathbf{y}) = -\frac{1}{2} \ln |2\pi \mathbf{K}^T \mathbf{V} \mathbf{K}| - \frac{1}{2} \mathbf{y}^T \mathbf{K} (\mathbf{K}^T \mathbf{V} \mathbf{K})^{-1} \mathbf{K}^T \mathbf{y} \quad (25)$$

dispersiooniparameetrite suhtes, võttes selleks funktsioonist (25) esimest järku osatuletise hinnatava parameetri suhtes, võrdsustades tulemuse nulliga ning lahendades – tulemuseks ongi dispersiooni-parameetri REML-hinnang. Üldjuhul ei ole dispersioonikomponentide hinnangute analüütiline leidmine võimalik, logaritmilise tõepärafunktsiooni maksimeerimise tulemusena saadakse keeruline maatriksvõrdus

$$\text{tr} \left[(\mathbf{K}^T \mathbf{V} \mathbf{K})^{-1} \mathbf{K}^T \mathbf{Z}_i \mathbf{Z}_i^T \mathbf{K} \right] = \mathbf{y}^T \mathbf{K} (\mathbf{K}^T \mathbf{V} \mathbf{K})^{-1} \mathbf{K}^T \mathbf{Z}_i \mathbf{Z}_i^T \mathbf{K} (\mathbf{K}^T \mathbf{V} \mathbf{K})^{-1} \mathbf{K}^T \mathbf{y},$$

kus mudeli dispersiooniparameetrid sõltuvad üksteisest ja ka juhuslike faktorite realiseerunud väärtustest (maatriks \mathbf{Z}_i koosneb i -le juhuslikule faktorile vastavatest plaanimaatriksi \mathbf{Z} veergudest). Üksikute parameetrite hinnangute leidmiseks kasutatakse ligikaudseid järkjärgulise lähendamise e iteratsioonimeetodeid: algselt (nn nullsammul) antakse mudeli parameetritele mingid algväärtused, neist lähtuvalt leitakse 1. sammul kõigi parameetrite uued hinnangud, seejärel lähtutakse 2. sammul parameetrite hinnangute arvutamisel juba 1. sammul leitudest jne – igal järgneval sammul kasutatakse eelmisel sammul leitud hinnanguid; protsess koondub parameetrite lõplikeks hinnanguteks, kui mingil sammul leitud hinnangute erinevus eelneval sammul leitudest on piisavalt väike. Algoritme, mida sellisel ligikaudsel arvutamisel kasutatakse, on jälle rohkem kui üks – mõni neist koondub väiksema, mõni suurema arvu sammude järel; teisalt ei pruugi enamasti kiiremini koonduvad algoritmid mõnel juhul üleüldse koonduda – kõik see teeb dispersioonikomponentide hindamise veel keerulisemaks.

9 ÜLESANDED

Järgnev tabel sisaldab andmeid seitsme lehma piimatoodangu, isa, vanuse ja karja kohta.

Loom	Isa	Vanus	Kari	Toodang
2	1	7	1	7000
4	3	7	1	6600
6	5	6	1	6800
7	5	4	1	7200
8	3	3	2	8000
9	3	2	2	8300
10	5	2	2	8900

- Lugedes kõik faktorid fikseerituteks, pange kirja üldine lineaarne mudel hindamaks isa, vanuse ja karja mõju piimatoodangule, seejuures käsitlege isa ja karja diskreetsete faktoritena ning vanust pidava faktorina
 - objektiivisi;
 - maatrikskujul, pannes kirja ka plaanimaatriksi \mathbf{X} .
- Pange kirja mudeli dispersioonimaatriksid \mathbf{G} , \mathbf{R} ja \mathbf{V} , lugedes isa mõju juhuslikuks.
- Lugedes kõik faktorid fikseerituteks, hinnake nende mõjud vähimruutude meetodil (13), kasutades SAS-i reparametriseerimistingimusi.
- Leidke isadele ja karjadele vastavad vähimruutkeskmised.
- Lugedes isa mõjud juhuslikuks ja võttes piimatoodangu päritavuseks $h^2 = 0,25$, leidke Hendersoni segamudeli võrrandi (21) abil fikseeritud efektide hinnangud ja juhusliku faktori 'Isa' realiseerunud väärtuste prognoosid.