

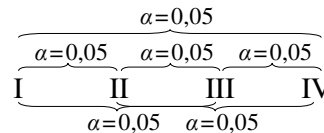
Biomeetria

Enam kui kahe populatsiooni keskväertuste võrdlemine – dispersioonanalüüs

Mitmene võrdlus



Võrdleme näiteks 4 gruppi, lubades iga üksikvõrdluse puhul eksimist 5% tõenäosusega.



Tõenäosus, et üksikvõrdlusel viga ei tehta, on $1-\alpha=0,95$.

Tõenäosus, et kuuel üksikvõrdlusel kokku ei eksita, on $(1-\alpha)^6=0,95^6\approx 0,735$.

Mistõttu tõenäosus teha üks (või mitu) vale otsus(t) 4 grupi paarikaupa võrdlemisel on $1-0,735=0,265$ (eksimise tõenäosus on üle 25%!).

Bonferroni meetod: piiramaks k üksikvõrdluse puhul ühe või enama vea tegemise tõenäosust olulisuse nivoo α , tuleb kõigil üksikvõrdlustel võtta olulisuse nivoo α/k .

Näiteks 4 grupi võrdlemisel, garanteerimaks kuue võrdluse peale kokku eksimist mitte üle 5%-lise tõenäosusega, tuleb üksikvõrdlustel võtta olulisuse nivoo $\alpha^* = \alpha/k = 0,05/6 \approx 0,0083$.

Enam kui kahe grupi keskmiste võrdlus



$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_1: \text{leiduvad sellised grupid } i, j, \text{ et } \mu_i \neq \mu_j$$

Eeldustel, et

- ▣ uuritav (sõltuv) tunnus on normaaljaotusega ja
- ▣ uuritava tunnuse varieeruvus võrreldavais gruppides on ühesugune, on rakendatavaks analüüsimeetodiks dispersioonanalüüs.

Dispersioonanalüüsil jagatakse tunnused vastavalt nende rollile kaheks: tunnus, mille keskmisi võrrelda soovitakse, on uuritav tunnus e funktsioon-tunnus (lehma piimatoodang, forelli kasvukiirus, talle mass, sea pekipaksus, jne); (diskreetne või mittearvuline) tunnus, mille väärtuste alusel võrreldavad grupid moodustatakse, on faktortunnus (tõug, lüpsiseade, laudatüüp jne).

Dispersioonanalüüsi tulemuste tõlgendamisel räägitaksegi enamasti faktortunnuse mõjust uuritavale tunnusele.

Näiteks, tõu või lüpsiseadme või laudatüübi vm mõju piimatoodangule, kasvanduse mõju forellide kasvukiirusele, omaniku mõju talle massile, genotüübi (teatud geenikombinatsioonide) mõju sigade pekipaksusele jne.

Dispersioonanalüüsi mudel

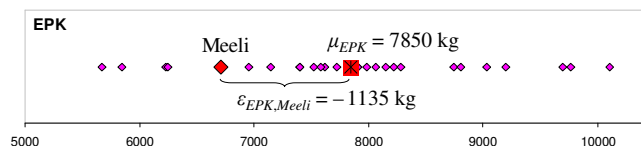


✓ Iga rühma i (kus $i=1, \dots, k$) iseloomustab keskmine uuritava tunnuse väärtus μ_i , mistõttu mõõtmistulemused saab esitada mudeliga


$$y_{ij} = \mu_i + \varepsilon_{ij},$$

kus y_{ij} on uuritava tunnuse väärtus i . rühma kuuluval j . objektil ja ε_{ij} on juhuslik mõju (objekti omapära).

Näiteks EPK-tõugu lehm Meeli 1. laktatsiooni piimatoodang 6715 kg on väljendatav kui uuritud EPK-tõugu lehmade 1. laktatsiooni keskmise toodangu $\mu_{EPK} = 7850$ kg ja Meeli tõusisese erinevuse $\varepsilon_{EPK, Meeli} = -1135$ kg summa.



Dispersioonanalüüsi mudel



✓ Faktortunnuse mõju uurimiseks esitatakse mudel kujul

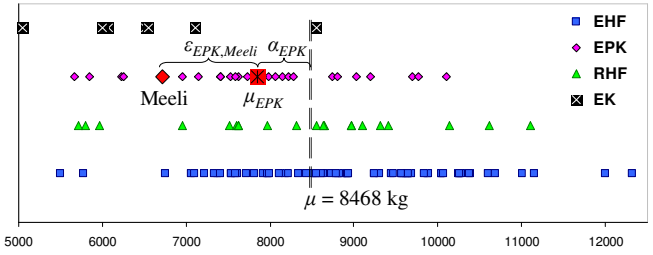
$$y_{ij} = \mu + \alpha_i + \varepsilon_{ij},$$

kus μ tähistab üldkeskmist ja α_i on faktori i . taseme poolt põhjustatud kõrvalekalle üldkeskmisest (i . taseme mõju), $\mu_i = \mu + \alpha_i$.


$H_0: \mu_1 = \mu_2 = \dots = \mu_k$ $H_0: \alpha_1 = \alpha_2 = \dots = \alpha_k = 0$
 $H_1: \text{leiduvad grupid } i, j, \text{ et } \mu_i \neq \mu_j$ $H_1: \text{leidub grupp } i, \text{ et } \alpha_i \neq 0$

Näiteks EPK-tõugu lehm Meeli 1. laktatsiooni piimatoodang 6715 kg on väljendatav kui kõigi uuritud lehmade keskmise 1. laktatsiooni piimatoodangu $\mu = 8468$ kg, EPK-tõu mõju

(EPK-tõugu lehmade 1. laktatsiooni keskmise tootangu erinevus üldkeskmisest) $\alpha_{EPK} = -618$ kg ja Meeli tõusisese erinevuse $\varepsilon_{EPK, Meeli} = -1135$ kg summa.



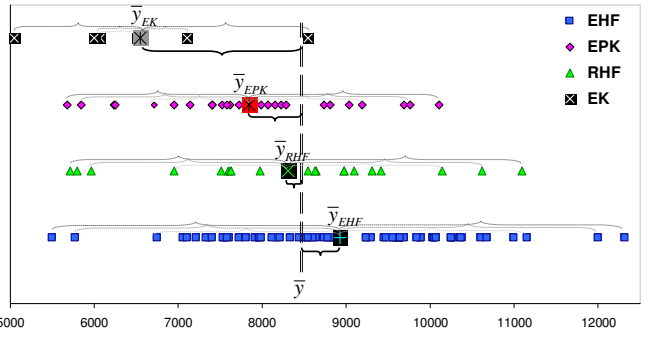
Dispersioonanalüüsi tööpõhimõte



Dispersioonanalüüsi tööpõhimõte seisneb uuritava tunnuse rühmadesisese (nn juhusliku) varieeruvuse $SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$ ja rühmadevahelise (faktori mõjust tingitud) varieeruvuse $SSA = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2$ võrdlemises – kui rühmadevaheline erinevus on suurem kui rühmadesisene erinevus, on tegu ilmse tõendiga faktortunnuse mõju olemasolu kohta.

Siit ka analüüsi nimetus – dispersioonanalüüs [analysis of variance, ANOVA].

Näide.
 $i = EK, EPK, RHF, EHF$



Dispersioonanalüüsi tabel



Dispersioonanalüüsiga seotud arvutused koondatakse tavaliselt alljärgnevasse nn. dispersioonanalüüsi tabelisse.

Varieeruvuse allikas	Hälvete ruutude summa	Vabadusastmeid	Keskruut	F -suhe	Olulisustõenäosus
Faktor	$SSA = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2$	$k - 1$	$MSA = \frac{SSA}{k - 1}$	$F = \frac{MSA}{MSE}$	p
Viga	$SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$	$n - k$	$MSE = \frac{SSE}{n - k}$		
Kokku	$SS = SSA + SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2$	$n - 1$			

Juhul, kui faktortunnuse mõjule vastav keskmine gruppide vaheline varieeruvus MSA on suurem, kui uuritavate objektide omapäralt vastav keskmine gruppide sisene varieeruvus MSE , on F -statistiku väärtus ühest suurem.

Piisavalt suure F -suhte väärtuse korral võib lugeda tõestatuks sisuka hüpoteesi – leiduvad vähemalt 2 teineteisest selgelt eristuvat gruppi.

Dispersioonanalüüs



Näide. Uuritakse ühes katsefarmis peetava 121 lehma 1. laktatsiooni piimatoodangu sõltuvust tõust (EHF, RHF, EPK, EK). Kontrollitav hüpoteeside paar on kujul:

$$\begin{aligned} H_0 : \mu_{EHF} = \mu_{RHF} = \mu_{EPK} = \mu_{EK} & \Leftrightarrow H_0 : \text{tõul ei ole mõju} \\ H_1 : \text{leiduvad tõugrupid } i, j, \text{ et } \mu_i \neq \mu_j & H_1 : \text{tõul on mõju} \end{aligned}$$

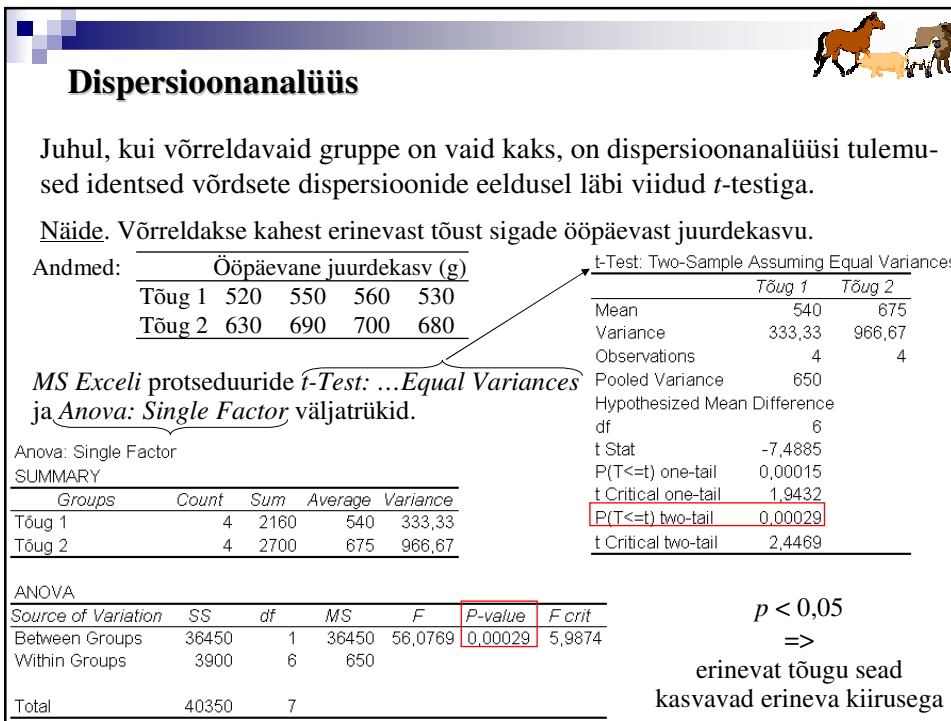
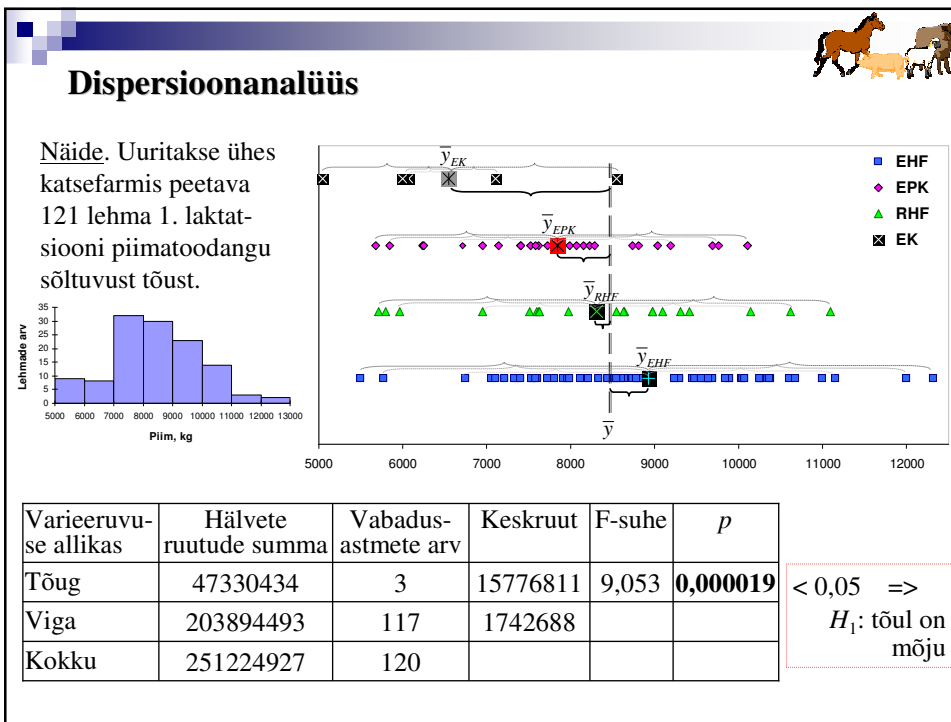
Dispersioonanalüüsi mudel on kujul $y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$, kus μ on kõigi farmi lehmade keskmine 1. laktatsiooni piimatoodang, α_i on i . tõu keskmine erinevus sellest (i . tõu mõju, $i = EHF, RHF, EPK, EK$) ning y_{ij} ja ε_{ij} on vastavalt i . tõugu j . lehma mõõdetud piimatoodang ja selle erinevus tõu keskmisest (lehma “omapära”, $j = 1, \dots, n_i$, n_i on lehmade arv i . tõus).

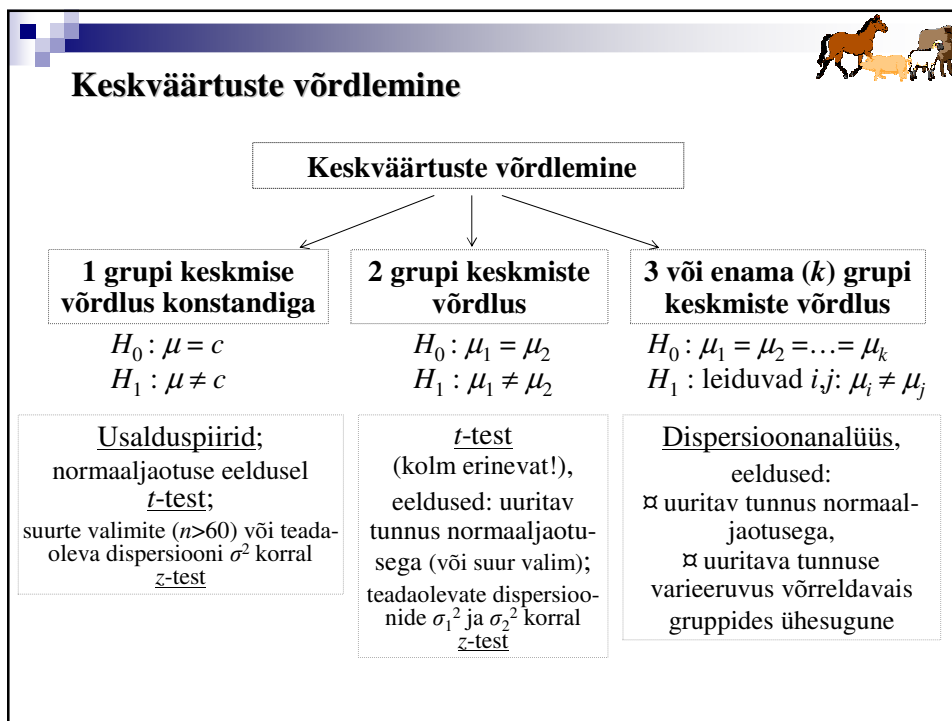
Tõug	n_i	$\bar{y}_i = \hat{\mu}_i$	$s_i^2 = \hat{\sigma}_i^2$
EHF	68	8929,9	1776652
RHF	20	8311,5	2265768
EPK	27	7850,1	1335053
EK	6	6549,3	1419571
Kokku	121	8468,1	2093541

Tõugude mõjud on kõrvaloleva keskmiste toodangute tabeli alusel leitavad kujul


$$\begin{aligned} \alpha_{EHF} &= 460,8 \text{ kg}; \alpha_{RHF} = -156,6 \text{ kg}; \\ \alpha_{EPK} &= -618,0 \text{ kg ja } \alpha_{EK} = -1918,8 \text{ kg}. \end{aligned}$$

Nende mõjude erinevuse kontrollimiseks tuleb läbi viia dispersioonanalüüs [viimase eeldused dispersioonide võrdsuse ja normaaljaotuse (vt ka järgmine lk) osas on enamvähem täidetud].





Mitmefaktoriline dispersioonanalüüs




Kui vaatlusobjekte saab rühmitada mitme tunnuse (faktortunnuse) järgi, võib osutada mõttekaks analüüsida korraga mitme faktortunnuse mõju (näiteks igal lehmil võib olla fikseeritud tema tõug ja farm, igal kalal tema sugu ja püügikoht).

Dispersioonanalüüsi mudel, mis hõlmab kahe faktortunnuse mõjusid, on kujul:

$$y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk},$$

kus μ tähistab üldkeskmist,
 α_i -d ja β_j -d märgivad uuritava tunnuse keskmise muutust vastavalt esimese ja teise faktori väärtuste muutumisele (α_i on esimese faktori i . taseme mõju ja β_j on teise faktori j . taseme mõju),
 y_{ijk} ning ε_{ijk} on vastavalt esimese faktori i . tasemel ja teise faktori j . tasemel sooritatud k . mõõtmise väärtus ning selle erinevus sama väärtuste kombinatsiooni keskmisest (vaatluse omapära, mudeli viga).

Mitmefaktoriline dispersioonanalüüs

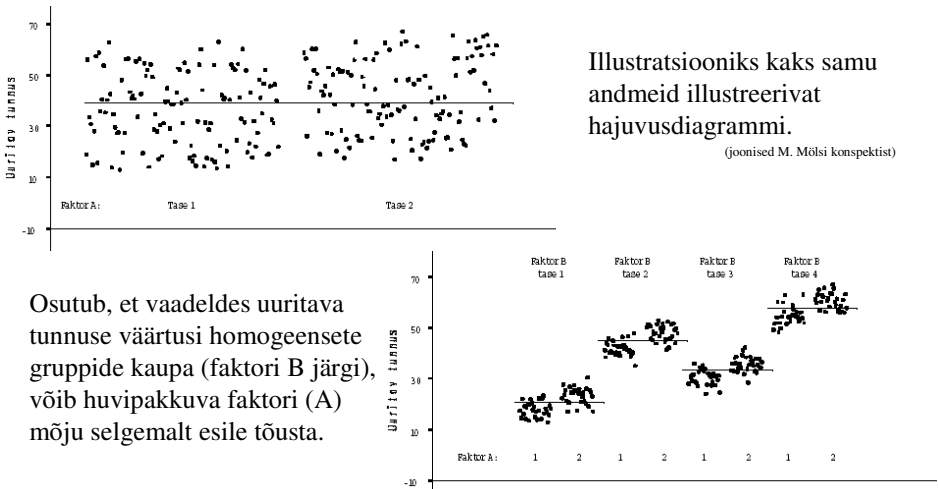


Miks seda vaja on?

1) Hinnangute ja otsustuste täpsus võib paraneda.


Illustratsiooniks kaks samu andmeid illustreerivat hajuvusdiagrammi.

(joonised M. Mölsi konseptist)



Osutub, et vaadeldes uuritava tunnuse väärtusi homogeensete gruppide kaupa (faktori B järgi), võib huvipakkuva faktori (A) mõju selgemalt esile tõusta.

Mitmefaktoriline dispersioonanalüüs



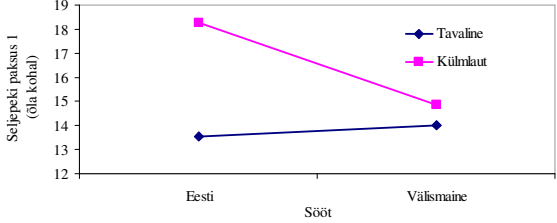
Miks seda vaja on?

2) Võimalik selgemalt väljendada uuritava tunnuse ja faktorite vahelisi seoseid.

3) Interaktsioonid e koosmõjud – uuritava tunnuse väärtused muutuvad ühe faktori tasemete vahel erinevalt, sõltuvalt teise faktori väärtustest.

4) Ilma ei pruugi mudel olla korrektne (jääkliige ei pruugi olla normaaljaotusega).

Sigade seljapeki paksus 1 sõltuvalt söödast ja lauda tüübist



Kahefaktoriline faktoritevahelist interaktsiooni arvestav mudel esitatakse kujul

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk},$$

kus γ_{ij} märgib esimese faktori i . taseme ja teise faktori j . taseme koosmõju.

