



**Mitmemõõtmelise statistika
koolitus**

**Korrespondents-, klaster- ja
diskriminantanalüüs ning
mitmemõõtmeline skaleerimine**

Eesti Maaülikool
20.-24. jaanuar 2014
Tanel Kaart



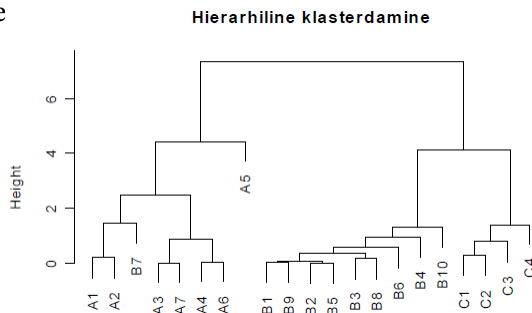
**Klasteranalüüs
[cluster analysis]**

Klasteranalüüsi eesmärgiks on kas tunnuste või uuritavate objektide rühmitamine. Kaks peamist algoritmi:

- hierarhiline klasterdamine;
- k-keskmiste klasterdamine.

Hierarhiline klasterdamine

- Loodusuurija sõitis Asustamata Saarele ja avastas seal suure hulga seni teadusele tundmatuid putukaid.
- Ta mõõtis oma lühikese saarelvibimise jooksul tervel hunnikul seninägematutel putukatel igasuguseid näitajaid (igatsorti pikkuseid ja mustrielementide arvu ja paljut muudki).
- Järgmiseks soovis loodusuurija määratleda, mitmesse alamliiki leitud putukad võiksid kuuluda.
- Et saada esimest ligikaudset lähendit, kust oma uurimistööga pihta hakata, soovis ta leida sarnaste putukate rühmad – kes oleksid siis alamliikide kandidaatideks.
- Selleks sõitis ta oma andmed klasteranalüüsi teostavasse programmi, mis joonistas järgmise pildi:



Märt Mölsi loengukonseptist

Hierarhiline klasterdamine

... on hästi kasutatav siis, kui meil on suhteliselt vähe objekte või kui on oodata, et klastrid suhteliselt selgelt üksteisest eristuvad.

Hierarhiline klasteranalüüs põhineb väga lihtsal algoritmil:

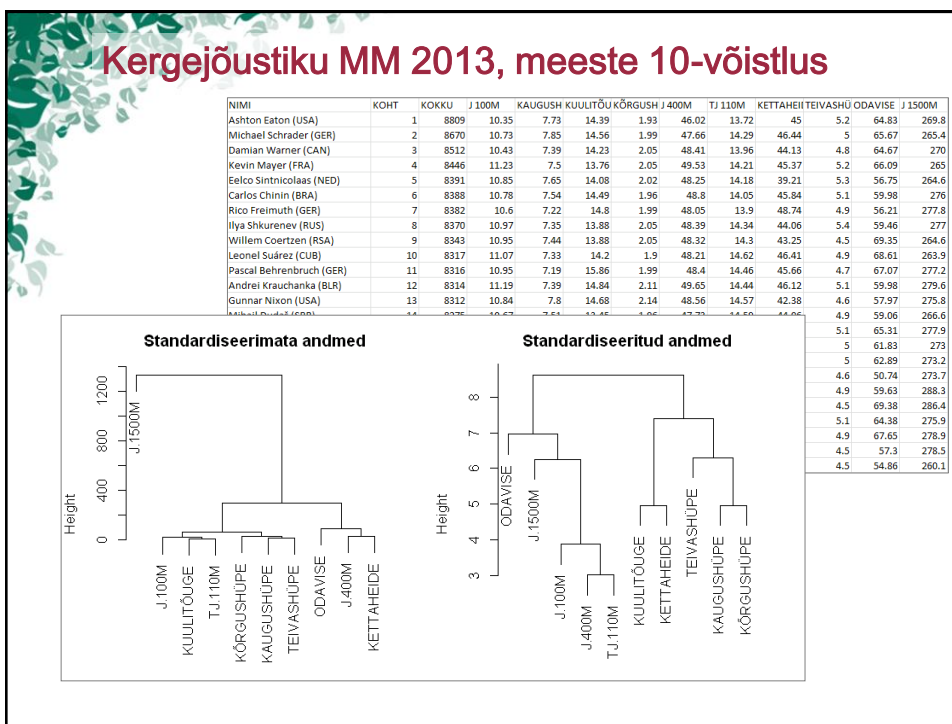
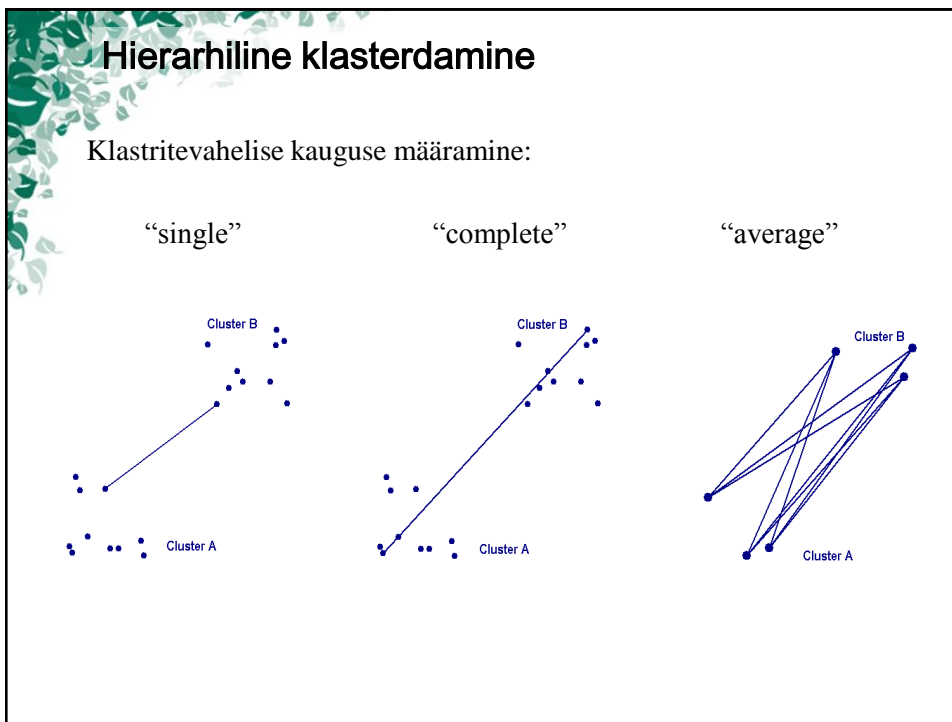
samm-sammult hakatakse omavahel kokku panema kõige sarnasemaid objekte. Näiteks, kui leidub kaks täpselt ühesuguste tulemustega objekti, siis liidetakse nad esimesel sammul üheks klatriks, peale seda võrreldakse kõiki üksikobjekte ja juba tekkinud klastreid ja liidetakse jälle kõige sarnasemad omavahel jne.

Vaatluste omavahelise kauguse määramine:

Eukleidese kaugus:

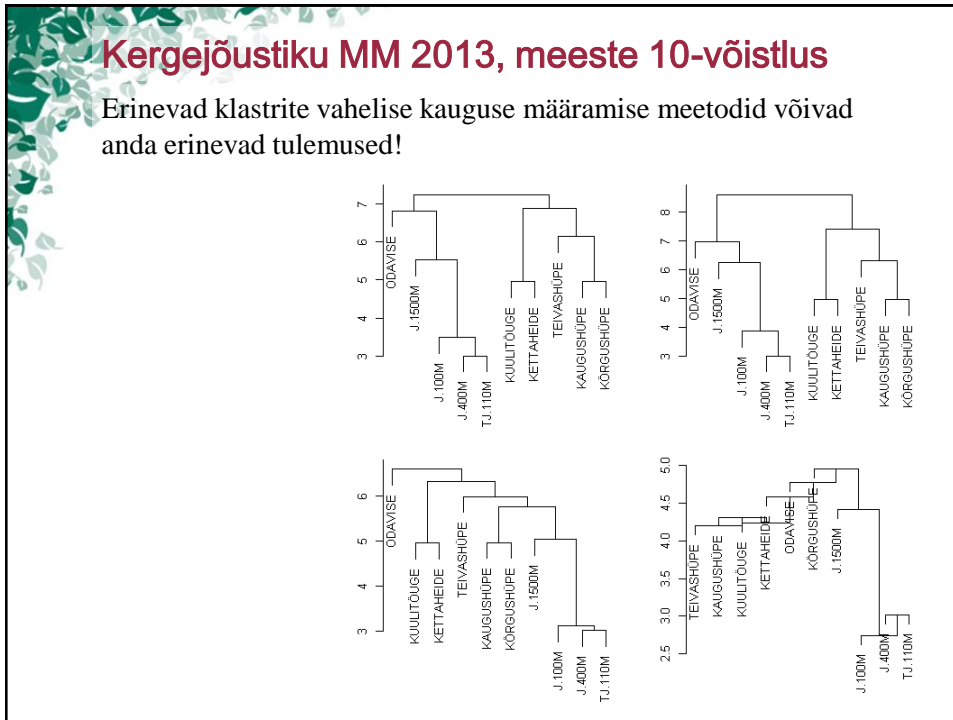
$$\text{Manhattani kaugus: } d(x_1, y_1), (x_2, y_2) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

$$\dots \quad d(x_1, y_1), (x_2, y_2) = |x_1 - x_2| + |y_1 - y_2|$$



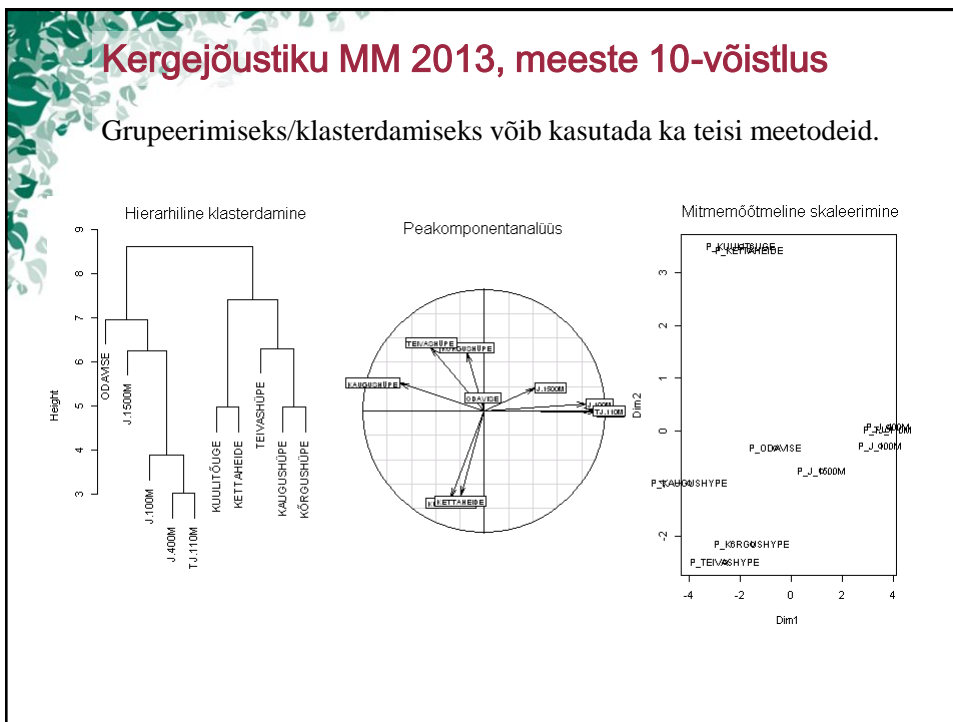
Kergejõustiku MM 2013, meeste 10-võistlus

Erinevad klasterite vahelise kauguse määramise meetodid võivad anda erinevad tulemused!



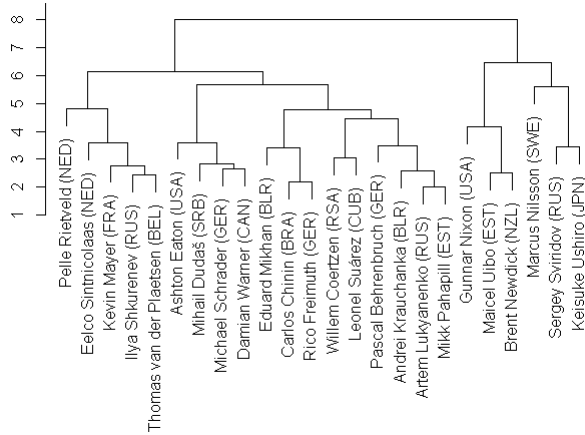
Kergejõustiku MM 2013, meeste 10-võistlus

Grupeerimiseks/klasterdamiseks võib kasutada ka teisi meetodeid.



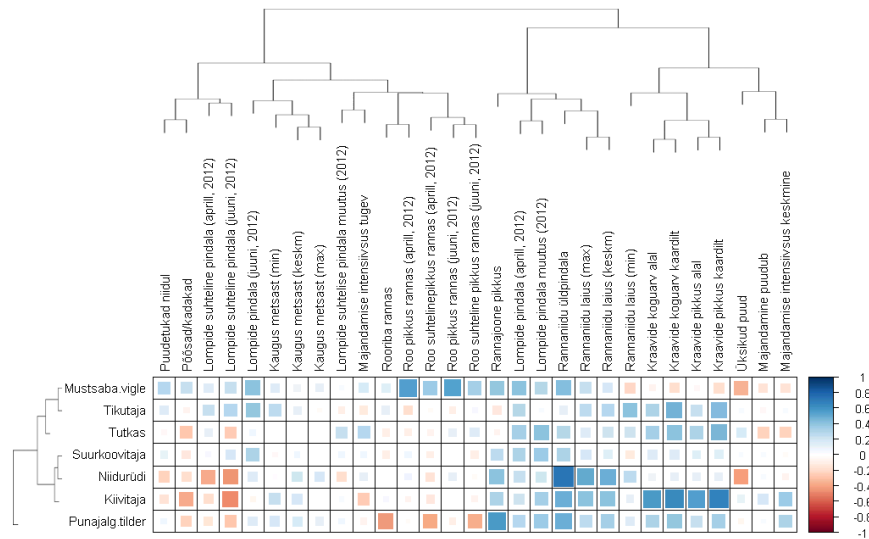
Kergejõustiku MM 2013, meeste 10-võistlus

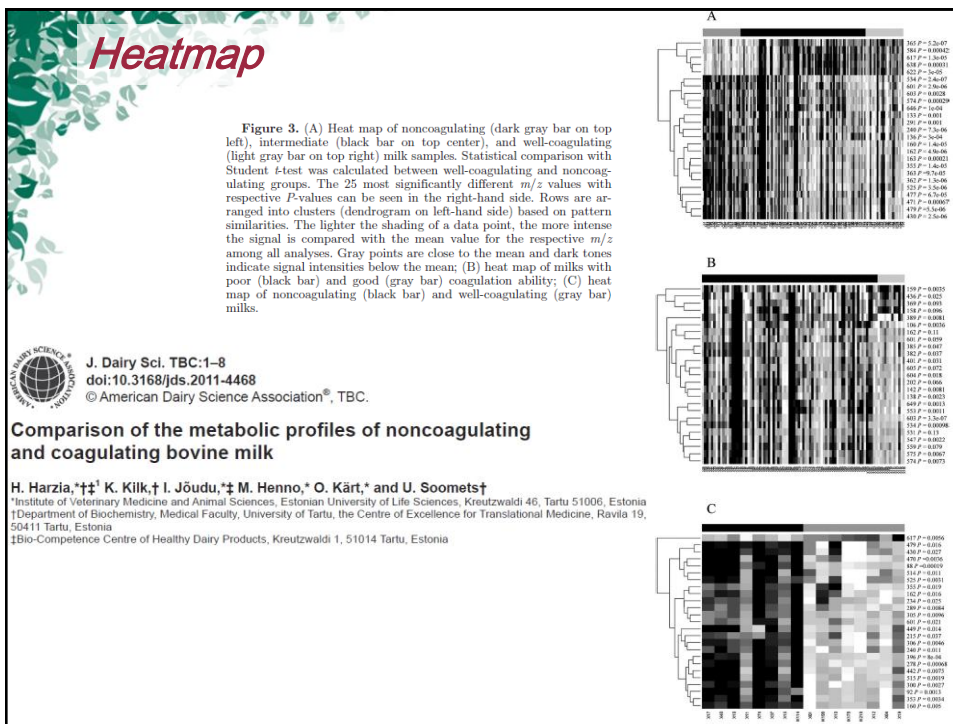
Klasterdada võib ka andmebaasi ridu.



Klasteranalüüs tunnuste järjestamise meetodina

Eesmärgiks tuua visuaalselt selgemalt välja andmetes peituvad sarnasusmustrid.





k-keskmiste klasterdamine

... sobib grupeerimise meetodiks siis, kui objekte on nii palju, et hierarhilise klasteranalüüsi tulemus muutub ebaülevaatlikuks, kuid ka siis kui me oskame meile sobivat klastrite arvu ligilähedaselt ennustada ning ühtlasi soovime saada ka tekkivate klastrite kirjelduse nende tunnuste osas, mis on grupeerimise aluseks.

Algoritm:

- kõigepealt tuleb määrata klastrite arv, siis
- jagada objektid esialgsetesse klastritesse,
- arvutada välja klastrite keskpunktid ning
- hakata võrdlema igat objekti kõigi klastrite keskpunktidega; kui osutub, et objekti kaugus mõne muu klatri keskpunktist on väiksem kui selle klatri keskpunktist, milles ta parasjagu asub, siis tuleb objekt teise klattrisse ümber tõsta;
- peale objekti ümber tõstmist tuleb pöörduda uuesti sammu 3 juurde ja jätkata protsessi niikaua kui kõik objektid on klattris, mille keskpunktile nad kõige lähemal asuvad.

Näide: puude klasterdamine

ID	KUIV	PE	H100	KKT	ARENGUK D	H	A	
410057	0	MA	15.9	MS	K	14	12	55
410082	0	KU	23.6	JM	V	28	23	90
410111	0	KS	16.3	TR	K	14	13	60
410151	0	MA	14.8	MS	L	12	11	55
410173	0	MA	17.5	KN	L	6	4	29
410179	0	MA	22	MS	K	26	22	100
410202	0	MA	17.2	MS	K	16	15	70
410230	0	MA	24	JM	V	26	23	85
410248	1	MA	19.5	KM	K	24	20	110
410251	1	MA	17.5	SN	N	4	5	26

- 29539 puud,
- eesmärk jagada puud klastritesse diameetri (D) ja vanuse (A) alusel.

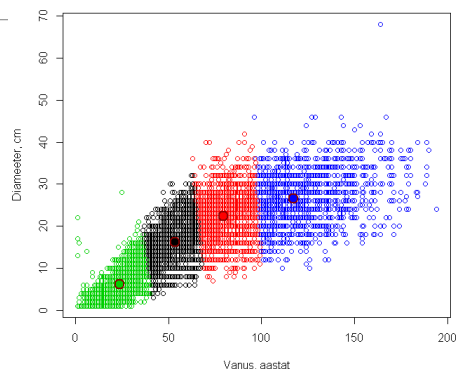
Näide: puude klasterdamine

```
> table(puu_kmean$cluster) # klastrite suurused
 1     2     3     4
8964 7493 5941 3026
> # klasterdamise aluseks olnud tunnuste keskmised klastrite kaupa
> puu_kmean$centers
      [,1]      [,2]
1  53.52443 16.21586
2  79.33298 22.34472
3  23.54385  6.14543
4 117.20787 26.56345
```

Miks tundub, et klastrid on moodustatud eelkõige puu vanuse järgi?

Vastus – puude vanust märkivad arvud (eelkõige nende varieeruvus) on suuremad, kui diameetrit märkivad arvud.

Soovides, et klasterdamisalgorithm peaks puu diameetrit sama oluliseks kui vanust, tuleb klasterdamiseks kasutatavad tunnused eelnevalt standardiseerida.

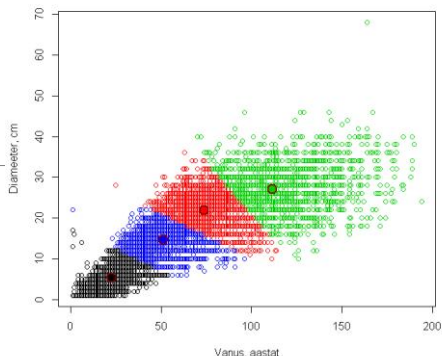


Näide: puude klasterdamine

Klasterdamise tulemused peale vanuse ja diameetri standardiseerimist

```
> table(puu_kmean2$cluster) # klastrite suurused
 1     2     3     4
5346 8295 3836 7947
>
> # klasterdamise aluseks olnud tunnuste keskmised klastrite kaupa
> puu_kmean2$centers
      [,1]      [,2]
1 -1.2838316 -1.4539809
2  0.3980414  0.6267682
3  1.6430013  1.2736169
4 -0.3449028 -0.2908839
>
> # Mis tüüpi puud mingitesse klastritesse kuuluvad?
> table(puu_kmean2$cluster, puudi$ARENGUKL)
      A     K     L     N     S     V     Y
1     1  148 2480 2566     3  110   38
2     0 5003     1     0     0 1592 1699
3     0  895     0     1     3  697 2240
4     0 5825 1529     0     1  358  234
```

Kuidas on puud klasterdunud sõltuvalt arenguklassist (A – lage, N – noorendikud, L – latimets, K – keskealised, V – valmiv, Y – küps, S – selguseta)?



Näide: Sirje Värvi veiste klasterdamine geneetiliste markerite alusel



MARKER-BASED GENETIC CHARACTERIZATION OF THE ESTONIAN DAIRY CATTLE BREEDS

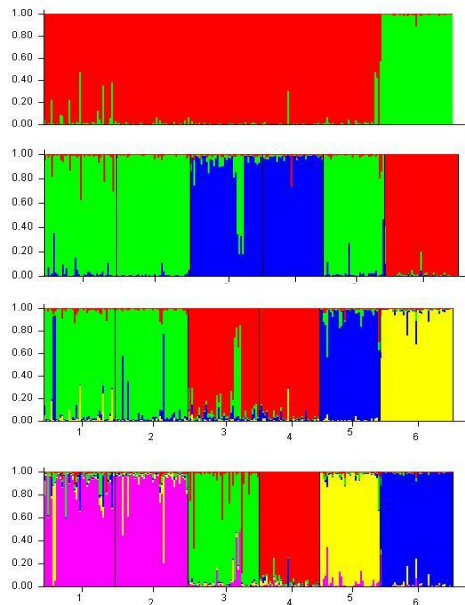
EESTI PIIMAVEISTEÕUGUDE ISELOOMUSTAMINE GENEETILISTE MARKERITE ALUSEL

SIRJE VÄRV

A Thesis for applying for the degree of Doctor of Philosophy in Animal Sciences

Väljõud filosoofidoktori kraadi taotlemiseks loomakasvatuse erialal

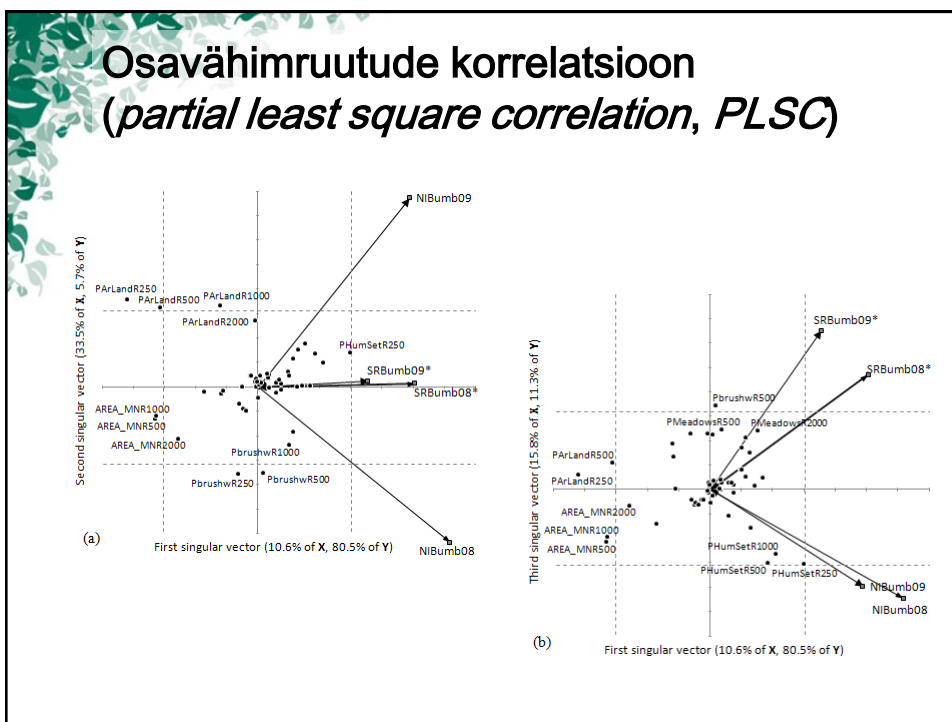
Tartu 2012

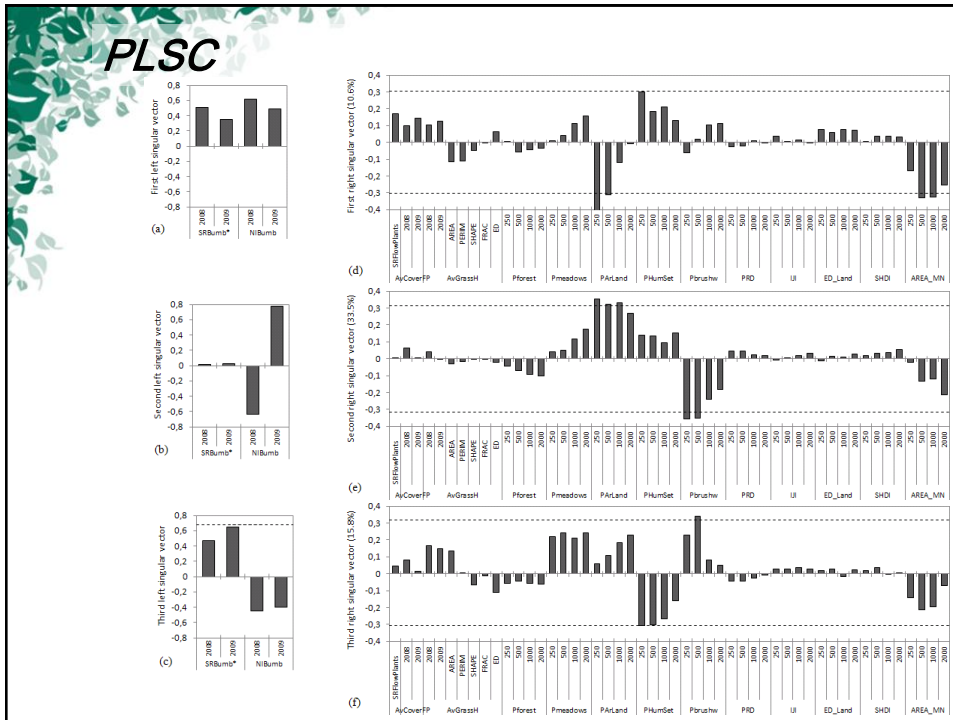


1 – EN,
2 – WFC,
3 – ER,
4 – EHF,
5 – SwRP,
6 – DkJer

Osavähimruutude regressioon ja korrelatsioon

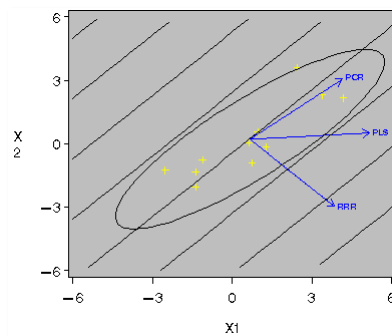
[Partial Least Square Regression and Corelation, Reduced Rank Regression, Principal Component Regression, ...]





Mustrid (*patterns*)?

- Peakomponentide regressioon (*principal components regression, PCR*)
- Osavähimruutude korrelatsioon ja regressioon (*partial least squares correlation, PLSC*), kanooniline korrelatsioon (*canonical correlation*)
- *Reduced rank regression (RRR)*



SAS Online Doc

Klasteranalüüs

Klasteranalüüsi eesmärgiks on kas tunnuste või uuritavate objektide rühmitamine. Kaks peamist algoritmi:

- hierarhiline klasterdamine;
- k-keskmiste klasterdamine.

Kasutusala: kus iganes (geneetika, ökoloogia, sotsioloogia, meditsiin, majandus, ...)

Klasteranalüüs

- Loodusuurija sõitis Asustamata Saarele ja avastas seal suure hulga seni teadusele tundmatuid putukaid.
- Ta mõõtis oma lühikese saarelvibimise jooksul tervel hunnikul seninägematutel putukatel igasuguseid näitajaid (igatsorti pikkuseid ja mustrielementide arvu ja paljut muudki).
- Järgmiseks soovis loodusuurija määratlada, mitmesse alamliiki leitud putukad võiksid kuuluda.
- Et saada esimest ligikaudset lähendit, kust oma uurimistööga pihta hakata, soovis ta leida sarnaste putukate rühmad – kes oleksid siis alamliikide kandidaatideks.
- Selleks sõitis ta oma andmed klasteranalüüsi teostavasse programmi, mis joonistas järgmise pildi:

Märt Mölsi loengukonspektist

Hierarhiline klasterdamine

Height

A1 A2 A3 A7 A4 A6 B1 B9 B2 B5 B3 B8 B6 B4 B10 C1 C2 C3 C4

Klasteranalüüs

Hierarhiline klasterdamine

... on hästi kasutatav siis, kui meil on suhteliselt vähe objekte või kui on oodata, et klastrid suhteliselt selgelt üksteisest eristuvad.

Hierarhiline klasteranalüüs põhineb väga lihtsal algoritmil: samm-sammult hakatakse omavahel kokku panema kõige sarnasemaid objekte. Näiteks, kui leidub kaks täpselt ühesuguste tulemustega objekti, siis liidetakse nad esimesel sammul üheks klastriks, peale seda võrreldakse kõiki üksikobjekte ja juba tekkinud klastreid ja liidetakse jälle kõige sarnasemad omavahel jne.

Vaatluste omavahelise kauguse määramine:

$$\text{Eukleidese kaugus: } d(x_1, y_1), (x_2, y_2) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

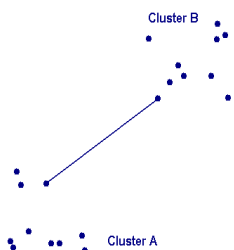
$$\text{Manhattani kaugus: } d(x_1, y_1), (x_2, y_2) = |x_1 - x_2| + |y_1 - y_2|$$

...

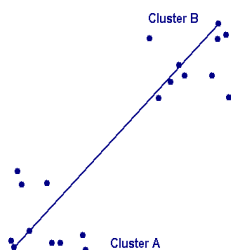
Klasteranalüüs

Klastritevahelise kauguse määramine:

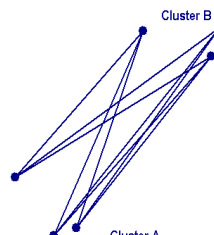
“single”



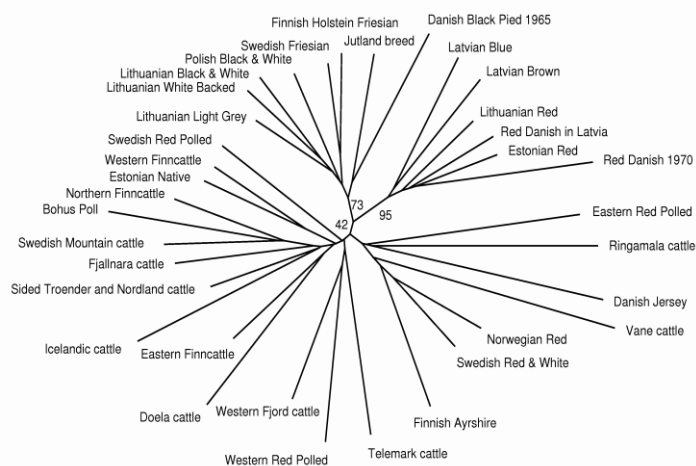
“complete”



“average”



Hierarhiline klasterdamine – näide



Tapio et al 2006. Prioritization for Conservation of Northern European Cattle Breeds Based on Analysis of Microsatellite Data. *Conservation Biology* 20, 1768-1779.

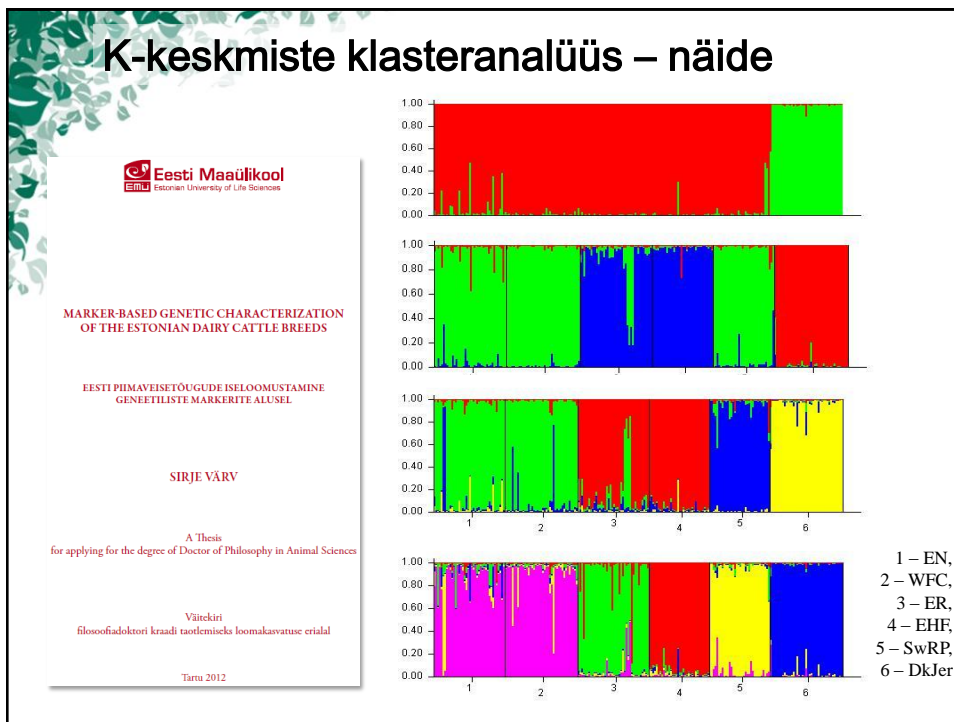
Klasteranalüüs

k-keskmiste klasterdamine

... sobib grupeerimise meetodiks siis, kui objekte on nii palju, et hierarhiline klasteranalüüsi tulemus muutub ebaülevaatlikuks, kuid ka siis kui me oskame meile sobivat klastrite arvu ligilähedaselt ennustada ning ühtlasi soovime saada ka tekkivate klastrite kirjelduse nende tunnuste osas, mis on grupeerimise aluseks.

Algoritm:

- kõigepealt tuleb määrata klastrite arv, siis
- jagada objektid esialgsesse klastritesse,
- arvutada välja klastrite keskpunktid ning
- hakata võrdlema igat objekti kõigi klastrite keskpunktidega; kui osutub, et objekti kaugus mõne muu klatri keskpunktist on väiksem kui selle klatri keskpunktist, milles ta parasjagu asub, siis tuleb objekt teise klastrisse ümber tõsta;
- peale objekti ümbertõstmist tuleb pöörduda uuesti sammu 3 juurde ja jätkata protsessi niikaua kui kõik objektid on klastris, mille keskpunktile nad kõige lähemal asuvad.



Diskriminantanalüüs

Diskriminantanalüüsi eesmärgiks on objektide rühmitamine nendel mõõdetud tunnuste alusel.

Seejuures on objektide klassidesse kuulumine enne analüüsi teada (erinevalt peakomponent või klasteranalüüsist).

Diskriminantanalüüs

Objektidel mõõdetud tunnuste alusel koostueeritakse nn diskreemineeriv funktsioon, mis eristaks grupe võimalikult selgelt:

$$d = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

Näiteks

2,3×“jalgade pikkus“ + 6,7×“saba pikkus“ – 2,8×“kere pikkus“ + 5,1×“noka pikkus“

Kui saadud väärtus <10,2, siis on ilmselt tegu isase isendiga.

Täpsemalt öeldes hinnatakse tunnuste vektori \mathbf{x} tihedusfunktsioon $f_t(\mathbf{x})$ igas grupis t ning arvutatakse iga objekti mingisse gruppi kuulumise tõenäosus Bayesi valemist kujul

$$P(t | \mathbf{x}) = \frac{q_t f_t(\mathbf{x})}{\sum_u q_u f_u(\mathbf{x})},$$

misjärel määratakse iga objekt tõenäoliseimasse gruppi

(suurus q_t eelnevas valemis on objekti gruppi t kuulumise alg tõenäosus).

Diskriminantanalüüs – näide

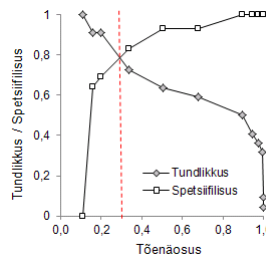
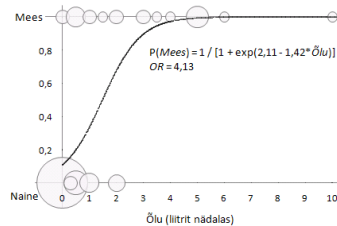
Jaanika Hõimra

Morfomeetriliste tunnuste varieeruvuse sõltuvus keskkonnatingimustes kortslehe (*Alchemilla L.*) viiel mikroliigil eksperimendi tingimustes

Tabel 8. Klassifitseeriv diskriminantanalüüs 18 tunnuse alusel (2001). Ridades on antud empiirilisel määratud liigid ja tulpades prognoos.

	%	KAREDA- KARVANE	VÄIKE	TERAVA- HÖLMINE	KÜÜT
KAREDAKARVANE	93,2	262	19	0	0
VÄIKE	99,7	1	286	0	0
TERAVAHÖLMINE	91,5	0	0	259	24
KÜÜT	87,4	0	0	36	250
KOKKU	92,9	263	305	295	274

Diskriminantanalüüs vs logistiline regressioon



Logistilise regressioonanalüüsi tulemus:

nii mehed kui ka naised identifitseeritakse õigesti 80%-lise tõenäosusega.

Et diskriminantanalüüs loeb objekti kuuluvaks suurima tõenäosusega gruppi, ei pruugi saadav klassifitseerimiseeskiri olla optimaalseim.

Diskriminantanalüüsi tulemus:

meestest identifitseeritakse õigesti 59,1%, naistest 92,9%.

Number of Observations and Percent Classified into SUGU			
From SUGU	0	1	Total
0	39	3	42
	92.86	7.14	100.00
1	9	13	22
	40.91	59.09	100.00
Total	48	16	64
	75.00	25.00	100.00

Korrespondentsanalüüs

(*correspondence analysis*)

Korrespondentsanalüüs võimaldab graafiliselt kirjeldada sagedustabelite kujul esitatud tunnuste vahelisi seoseid.

Korrespondentsanalüüs

Andmestik (juuste ja silmade värv Caith'is, Šotimaal):

Silmavärv/juuksevärv	Blond	Punapea	Šataän	Brünett	Süsimust
Sinine	326	38	241	110	3
Hele	688	116	584	188	4
Keskmine	343	84	909	412	26
Tumedad	98	48	403	681	85

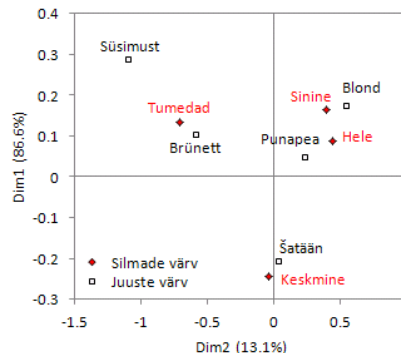
Korrespondentsanalüüsi esitus:

```

The CORRESP Procedure
Inertia and Chi-Square Decomposition
Singular Value Principal Inertia Chi-Square Percent Cumulative Percent
0.44637 0.13924 1073.33 86.56 86.56 *****
0.17346 0.02009 162.08 13.07 99.63 *****
0.02932 0.00086 4.63 0.37 100.00 *****
Total 0.23019 1240.04 100.00
Degrees of Freedom = 12

Row Coordinates
Dim1 Dim2
Sinine 0.4003 0.1854
Hele 0.4467 0.2085
Keskmine -0.0336 -0.2450
Tumedad -0.7027 0.1339

Column Coordinates
Dim1 Dim2
Blond 0.5440 0.1738
Punapea 0.2333 0.0483
Šataän 0.0420 -0.2083
Brünett -0.5887 0.1940
Süsimust -1.0944 0.2864
    
```



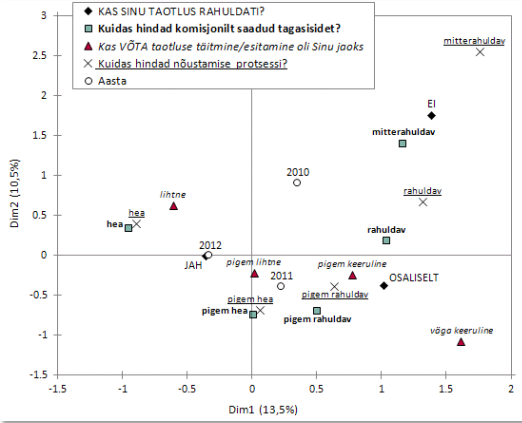
Mitmene korrespondentsanalüüs

- Mitmemõõtmeline korrespondentsanalüüs (*multiple correspondence analysis, MCA*) on nn andmete tihendamise meetod, mis võimaldab hulga mitteamvuliste tunnuste (näiteks küsitluse tulemuste) mitmemõõtmeliste sagedustabelite kujul esitatavad seosed projitseerida madalamamõõtmelisse (enamasti kahemõõtmelisse) ruumi.
- Analüüsi tulemused esitatakse enamasti joonise kujul, kus teatud väärtuste (küsimuste vastuste) lähedikkude paiknemine viitab nende sagedasemale koosinemisele.
- Seeläbi annab mitmemõõtmeline korrespondentsanalüüs võimaluse uurida korraka mitmete küsimuste ja vastuste omavahelisi seoseid ning tuvastada võimalikke mustreid vastanute seas.
- Projitseerimise tarvis arvutatavad nn dimensioonid on edasi kasutatavad juba arvuks tunnuseid eeldavates analüüsides.

Mitmene korrespondentsanalüüs

VÕTA taotluse esitanud tudengite ankeedivastuste mitme-mõõtmelise korrespondentsanalüüsi tulemused.

- Horisontaalsihis eristuvad tudengid, kelle VÕTA taotlus rahuldati ning kes hindasid nii taotluse esitamise protsessi lihtsaks kui ka nõustamise protsessi komisjonilt saadud tagasisidet heaks, ülejäänutest.
- Kõik VÕTA protsessiga seonduvad positiivsed vastused paiknevad joonise vasakus pooles ja negatiivsed vastused paremas pooles, neutraalsed vastused paiknevad vertikaaltelje lähedal.
- Vertikaalsihis eristuvad kindlad vastused (äärmused) kahtlevatest (neutraalsetest) vastustest.
- Seega kipuvad hinnangud taotluse esitamise ja nõustamise protsessile ning komisjonilt saadud tagasisidele olema sarnased, peegeldudes suurel määral ka esitatud taotluse tulemuslikkuses.
- Aastate lõikes on 2010. aastal olnud enam VÕTA protsessiga mitterahulolevaid tudengeid, 2011. aastal VÕTA protsessi rahuldavalt suhtuvaid tudengeid ja 2012. aastal VÕTA protsessi positiivselt suhtuvaid tudengeid (aastaarvud paiknevad joonise vastavais piirkondades).



Mitmene korrespondentsanalüüs

Raaperi, K., Bougeard, S., Aleksejev, A., Orro, T., Viltrop, A., 2012. ASSOCIATION OF HERD BRV AND BHV-1 SEROPREVALENCE WITH RESPIRATORY DISEASE AND REPRODUCTIVE PERFORMANCE IN ADULT DAIRY CATTLE.

Table 2. Descriptive characteristics of the variables included in the models (100 herds in Model 1 and 77 herds in Model 2)

Variable	Definition of the categories of the variable	Number of herds	
		Model 1	Model 2
Nasal discharge (red nose?) in cows and/or pregnant heifers	0 - not present at all or was shown only as single cases at some point during the last two years	82	
NASCOW	1 - present in more than just single cases at some point during the last two years	18	
	0 - not present at all or was shown only as single cases at some point during the last two years	80	
BHV/heif=0	1 - present in more than just single cases at some point during the last two years	20	
	0 - not present at all or was shown only as single cases at some point during the last two years	88	
BHV/heif=1	1 - present in more than just single cases at some point during the last two years	12	
	0 - less than two respiratory disease symptoms were present in more than a single case at some time during the last two years	81	
BRSV=0	1 - at least two respiratory disease symptoms were present in more than a single case at some time during the last two years	19	
	0 - < 1.3% in a herd (median for cut-off value)	37	
BRSV=1	1 - > 1.3% in a herd	40	
	0 - < 1.9 in a herd (median for cut-off value)	39	
BRSV=2	1 - 20-99 cows	40	25
	2 - 100-199 cows	19	17
BRSV=3	3 - 200-399 cows	23	24
	4 - > 400 cows	18	14
BRSV=4	0 - no	77	
	1 - yes	23	
BRSV=5	0 - no	38	
	1 - yes	72	
BRSV=6	0 - no	56	
	1 - yes	44	
BRSV=7	0 - no	61	
	1 - yes	59	
BRSV=8	0 - together	51	
	1 - in separate building from 6 months until pregnancy	49	
BRSV=9	1 - herd	23	
	2 - loose	30	
BRSV=10	3 - some period of life lived, some period loose	47	
	1 - herd	71	
BRSV=11	2 - loose	29	
	0 - no	47	
BRSV=12	1 - yes	53	
	0 - no	77	57
BRSV=13	1 - yes (at least one animal tested positive)	23	30
	0 - negative	46	34
BRSV=14	1 - 1 - 50%	40	31
	2 - >50%	14	10
BRSV=15	0 - negative	37	25
	1 - 1 - 50%	28	24
BRSV=16	2 - >50%	35	38
	0 - negative	56	40
BRSV=17	1 - positive	41	37

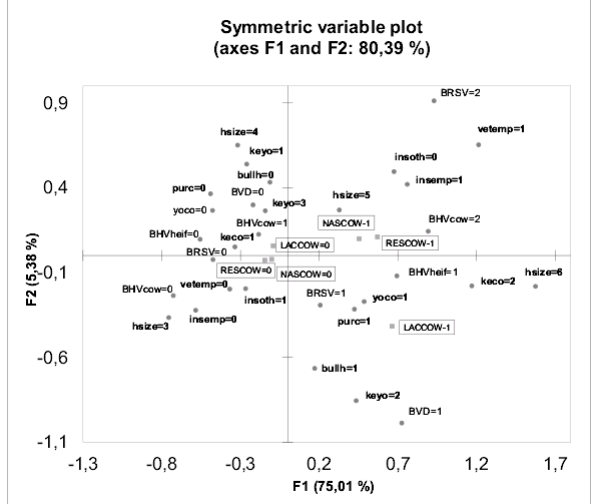


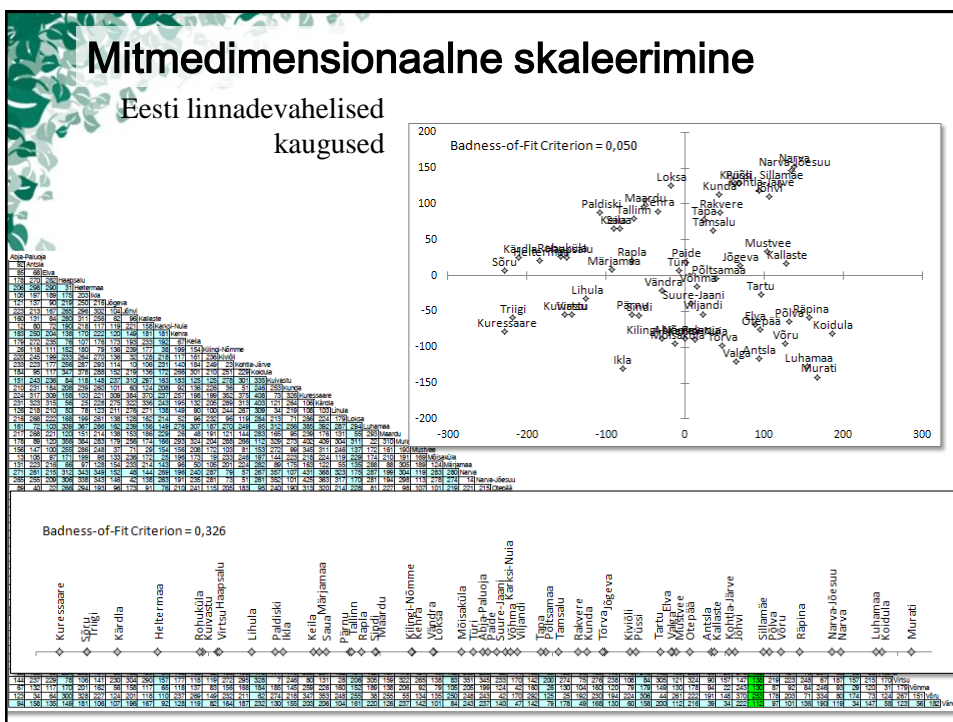
Figure 1 Graphical display of MCA for high incidence of respiratory disease symptoms in cows and pregnant heifers (100 herds)

Mitmemõõtmeline skaleerimine

(multidimensional scaling)

Mitmedimensionaalse skaleerimise eesmärgiks on objektide erinevuste või sarnasuste teisendamine distantsiks mitmemõõtmelises ruumis ja saadu esitamine kaardina objektide omavahelisest paiknemisest.

Kasutusala: kus iganes (geneetikas, ökoloogias, majandusteaduses, sotsioloogias, ...)



Mitmedimensionaalne skaleerimine

Egiptuse vaaraode 18. dünastia sugupuu algus

