

**Kirjeldav ja tõestav (kinnitav) andmeanalüüs**

**Kirjeldav analüüs** (*descriptive, exploratory data analysis*)

- Alustatakse andmetest (andmete kogumisest ilma kindla plaanita ja sidumata end kindla teooriaga).
- Järeldused on kirjeldavad ja ei anna nende kehtivuse tõenäosust.
- Tulemuste esitusviis on olulisem.
- Võib olla kiirem, odavam ja paindlikum.
- Tulemuslikkus sõltub:
  - üldistamisoskusest (eelnevatest teadmistest),
  - kasutatud andmete esinduslikkusest.

{Millest sõltub esinduslikkus?}

**Tõestav analüüs** (*inferential, deductive, confirmatory analysis*)

- Lähtutakse oletusest (hüpoteesist), mida üritatakse tõestada nullhüpoteesi ümberlukkamise abil.
- Tugineb rangetele teooriale.
- Tulemuseks on olulisustõenäosus, mis võimaldab hinnata, kui kindel võib järeldustes olla.
- Nõuab katseplaani järgimist.
- Tulemus kehtib vaid etteantud eelduste kehtimisel.
- Tulemuslikkus sõltub:
  - hüpoteeside püstitamise oskusest (eelnevatest teadmistest),
  - katse (vaatluse) planeerimise oskusest,
  - andmete esinduslikkusest.

Selles loengus käsitletakse kirjeldava analüüsi meetodeid.

**Kirjeldava uuringu põhilised viisid ja vahendid**

**Ühe andmekogumi kirjeldamine**  
keskmiste ja variatsiooninäitajate arvutamine ja esitamine, jaotuste iseloomustamine, kirjeldamine, vaatluste järjestamine mingi tunnuse järgi (ordinatsioon), tüüpiliste näidiste leidmine, vaatluste ja tunnuste klassifitseerimine.

**Andmekogumite võrdlemine**  
keskmiste ja varieeruvuse näitajate võrdlemine, jaotuste võrdlemine ja sarnasuste arvutamine.

**Seoste kirjeldamine** (tunnuste võrdlemine)  
**korrelatsioonikordajad**, seosed aja ja ruumiga (trend ja tsüklilisus), **vastavus** (*correspondence*) nominaalsete muutujate korral.

**Seose üldistamine ja formaliseerimine** (modelleerimine)  
regressioon jt statistilised **modelid**, **klassifitseerimine** (klassifikatsioonisüsteemi loomine), **näidiste** (tüüpide) komplekti (näidistebaasi) moodustamine.

**Tulemuste visualiseerimine ja visuaalne analüüs**  
graafikute ja kaartide võrdlemine.

Iga analüüsitüübi jaoks on omad vahendid

**Ühe andmekogumi kirjeldamine**

**Keskmsed**

- aritmeetiline keskmine.
- kaalutud keskmine (kui vaatlused ei ole võrdsed, esindavad nähtuse erinevat mahtu või on erineva usaldusväärsusega). Näiteks riigi metsasuse arvutamine maakondade metsasuste (osakaalu) andmetest.
- geomeetiline keskmine (kui olulised on kordsed erinevused).
- mood (asümmeetrilise jaotuse või nominaalse tunnuse korral).
- mediaan (tundmatu ebakorrapärase jaotustüübi korral).

**Varieeruvuse ehk väärtuste hajuvuse näitajad**

- dispersioon ehk keskmine ruuthälvetus
- standardhälve ehk ruutkeskmine
- varieeruvuse ulatus (teoreetiline ja praktiline)

**Sagedused**

- absoluutsagedused.
- suhteline sagedus ehk osakaal (mis on tervik, mis on ühik?).
- sagedusjaotus (mitme klassi puhul).

Väärtus muutuja arvteljel, millest ühele ja teisele pole jääb kindel osa vaatlustest.  
Mediaan on 50% kvantiil.  
Alumisest kvantiilist väiksemaid väärtusi on 25% ja ülemisest suuremaid väärtusi on 25%.  
Mis on kvantiili mõõtühik?

**Kvartiilid ehk veerandi kvantiilid**

**Arvuderea**  
1; 2; 3; 4; 5; 6; 7; 8; 9; 10; 11; 12  
alumine kvartiil on **3,5**; ülemine kvartiil **9,5**; keskmine on nende vahel = **6,5**

**Kas keskmine võib olla suurem kui ülemine kvartiil?  
Kas keskmine võib olla väiksem kui alumine kvartiil?**

Arvude 1; 2; 3; 4; 5; 6; 7; 10; 100; 200; 500; 10000  
alumine kvartiil on **3,5** ja ülemine kvartiil **150**  
keskmine = **903,1667** on kõrgem kui ülemine kvartiil

Arvude 1; 100; 101; 102; 103; 104; 105; 106; 107; 108; 109; 110  
alumine kvartiil on **101,5** ja ülemine kvartiil **107,5**  
keskmine = **96,3** on neist väiksem.

**Standardhälve ja standardviga**

**Üldkogumis**  
$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N}$$
  
Ühe vaatluse puhul dispersioon = 0

**Valimi järgi hinnates**  
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$
  
Ühe vaatluse järgi ei saa dispersiooni hinnata. Nulli jagamine nulliga

**Standardhälve** kirjeldab vaatluste varieeruvust tunnuse ühikutes.  
$$\sigma = \sqrt{\sigma^2}$$

**Standardviga** kirjeldab hinnangute oodatavat varieeruvust, s.o. hinnangu täpsust.  
Kõiksel analüüsil on kõik objektid mõõdetud ja viga ei ole.  
Tunnuse enda varieeruvus ja mõõtmisvead on erinevad asjad.  
**Keskmise standardviga valimi järgi hinnates**  
$$SE = \frac{s}{\sqrt{n}}$$

**Milleks ja kuidas klassifitseerida?**

**Milleks klassifitseerida?**

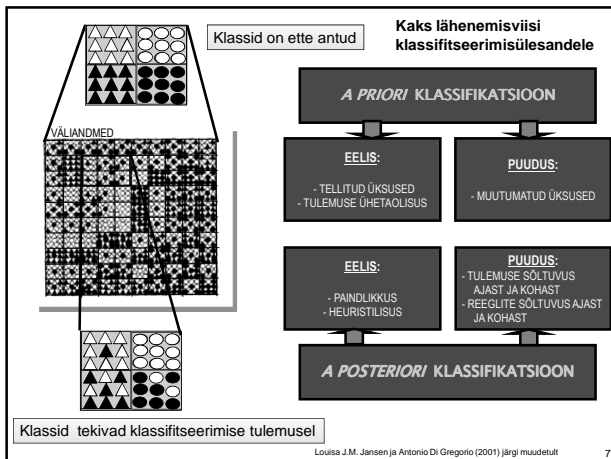
**Tunnustuslik vajadus.** Nähtuste lõputu mitmekesisus oleks midu haaramatu.

**Andmetöötuslik vajadus.** Vaatluseid on palju, kõiki ei jõua kõigega võrrelda, rühmitamise järel piisab rühmakuuluvuse määramisest ja rühmade võrdlemisest.

**Kuidas määrata tundmatu objekti klassikuuluvus?**  
Numbriliste väärtuste prognoosimisel samad võimalused

Klassid on ette antud

- Küsida eksperti (asjatundjalt või ekspertsüsteemilt)
- Võrrelda objekti ja klassikirjeldustega
- Kasutada otsustamise reegleid (määramistabelit)
- Võrrelda olemasolevate näidistega (muuseumisäilikutega, andmebaasikirjetega, juhtumitega)
- Arvutada mudelite või ekspertsüsteemide abil.



### Statistilise seose kirjeldamine

**Meetodid:**

- ❖ jaotuste võrdlemine (nominaalsed andmed),
- ❖ korrelatsioon ja astakkorrelatsioon (järjestatavad),
  - korrelatsioonimaatriks ja korrelogramm,
- ❖ regressioonimudelid, s.h. globaalselt ja lokaalselt,
  - trend (regressioon on aja- või ruumikoordinaatidega),
  - lokaalselt sobitatud regressioonid (ka libe keskmine),
- ❖ näidiste järgi järeldamine.

**Põhilised tulemuste esitusvormid:**

- ❖ tabel (näiteks korrelatsioonimaatriks),
- ❖ joonis (näiteks korrelogramm),
- ❖ valem (regressioonivõrrand),
- ❖ tekst,
- ❖ esitus multimeedia vahenditega.

### Kovariatsioon ja korrelatsioon

**Kovariatsioon** — hälvete korrutiste keskmine.

Kaks muutujat  $x$  ja  $y$

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

**Mille hälbed millest?** Kummagi muutuja väärtuste hälbed sama muutuja keskmisest.

**Mis on kovariatsiooni ühik?** Muutuja ühikute korrutis.

**Kas kovariatsiooni väärtus sõltub muutujate mõõtühikutest?** Jah, sõltub.

**Millal on kovariatsioon positiivne, millal negatiivne, millal null?** Null, kui kui muutujate vahel ei ole lineaarset seost; positiivne, kui muutujatel on kalduvus muutuda samas suunas; negatiivne, kui ühe muutuja suurenedes teine enamasti väheneb.

**Pearsoni lineaarne korrelatsioonikordaja** — standardhälvetega normeeritud kovariatsioon.

**Miks kovariatsiooni normeeritakse?** Võrreldavus.

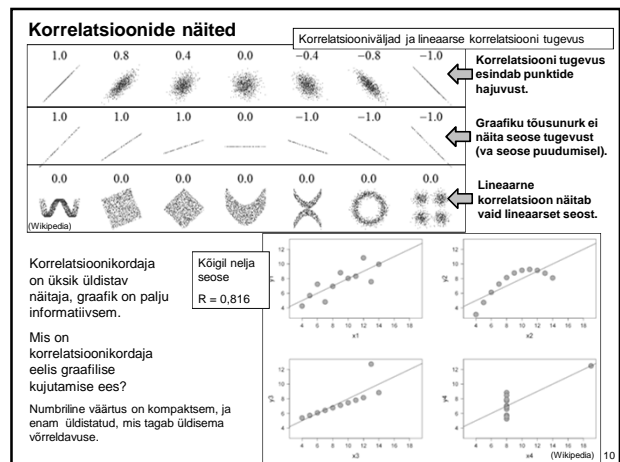
**Mida korrelatsioonikordaja näitab?** Vastastikust seost.

**Milliste tunnuste puhul saab korrelatsioonikordajat arvutada, milliste puhul on mõtet arvutada?** Arvutada saab numbritest, mõtet on arvutada, kui ruuthälvetel on mõtet.

**Kas korrelatsioonikordaja abil saab tõendada seose olemasolu (statistilist olulisust)?** Otse ei saa, vaatluste arvu on ka tarvis teada.

**Mis on korrelatsioonikordaja ühik?** Ühikuta.

**Determinatsioonikordaja** — korrelatsioonikordaja ruut  $R^2$ . Väljendab mudeliga (korrelatsioonikordaja puhul lineaarse mudeliga) seletatud dispersiooni suhet muutuja kogudispersiooni.



### Astakkorrelatsioonid

**Astak** — järjekorranumber.

Astakute kasutamisel ei ole muutujate jaotus oluline.

**Muutujad** peavad küll olema järjestatavad, aga ei pea olema pidevad ega normaaljaotusega. See on astakkorrelatsioonide eelis. Otsitav seos peab olema **monotoonne** (aina kasvav või kahanev).

Kui need tingimused ei ole täidetud, siis saab statistilisi seoseid uurida jaotuste võrdlemise testidega, näiteks  $\chi^2$  testiga.

**Spearmani astakkorrelatsioonikordaja  $\rho$  (rho)** — lineaarne korrelatsioonikordaja vaatluste järjekorranumbritest.

Omadused:  
 $-1 \leq \rho \leq 1$ ;  
 kui tunnuste vahel on kasvav seos, on  $\rho > 0$ ;  
 kui tunnuste vahel on kahanev seos, on  $\rho < 0$ ;  
 kui tunnuste vahel on funktsionaalne seos, siis  $|\rho| = 1$ ;  
 kui tunnused on sõltumatud, siis  $\rho = 0$ .

### Dice-Sørenseni sarnasuskordaja

Ühendi suhteline osa ehk Dice-Sørenseni sarnasuskordaja

$$QS = \frac{2C}{A+B}$$

Kokkulangev osa

Ühe vaatluse kogumaht ja teise vaatluse kogumaht

$$QS = 2 \cdot (0,2+0,07) / (0,3+0,2+0,12+0,08+0,5+0,07+0,2) = 0,37$$

Puude katvused ühes eraldises: Mä 30%, Ka 20%, Ku 12%, Le 8%

Katvused teises eraldises: Ka 50%, Ku 7%, muu 20%

Katvused kolmandas eraldises: Ka 40%, Ku 30% Le 10%

Sarnasuskordaja arvestab korraga mitut tunnust.

Kõigi vaatluspaaride puhul saab leida vaatlustevahelise kauguse. Saadakse

- 1) sarnasuste maatriks ja
- 2) kauguste maatriks.

Saab uurida, kas sarnasus sõltub kaugusest ja kui kaugele sarnasus ulatub.

### Sagedusjaotuste võrdlemine (2 x 2 sagedusjaotus)

#### Näidisülesanne

Vaadeldi 1090 rasvatihase reaktsiooni raudkulli topisele pesapaiga läheduses. 450st vaadeldud isalinnust reageeris agressiivselt 320, emaslinnudest näitas agressiivsust välja 420.

#### Ülesandes antud sagedused

Vaenulikkus	♀	♂	Kokku
+	420	320	?
-	?	?	?
<b>Kokku</b>	<b>?</b>	<b>450</b>	<b>1090</b>

#### Seose puudumise korral oodatavad sagedused

Vaenulikkus	♀	♂	Kokku
+	420	320	740
-	220	130	350
<b>Kokku</b>	<b>640</b>	<b>450</b>	<b>1090</b>

Vaenulikkuse osa  
 ♀ —  $420 / 640 = 65,625\%$   
 ♂ —  $320 / 450 = 71,11\%$   
 Paistab, et isased on vaenlase suhtes vaenulikumad.

Sagedustabel kirjeldab, aga ei tõesta seose olemasolu!  
 Ka juhulikke arve tabelisse seades saab mingid erinevused.

13

### Vigade maatriks

Kontrollandmed

	A	B	C	Kokku
A	120	23	12	155
B	11	505	23	539
C	35	45	400	480
Kokku	166	573	435	1174

Ridades ja veergudes on sama tunnus. Mõõtmisviis on erinev.  
 Diagonaalil on kokkulangevate vaatluste arv.

Vigade maatriksist saab arvutada klassifikatsioonide vastavuse suhtarve, näiteks:  
 Klassi A eristamise täpsus kasutaja jaoks (*user's accuracy*), **kontrollandmete suhtes** =  $120/166 \approx 72\%$   
 Klassi A eristamise täpsus klassifitseerija jaoks (*producer's accuracy*), **õpetusandmete suhtes** =  $120/155 \approx 77\%$

14

### Vastavuse indeks kapa (inglise k kappa ehk KHAT)

Vigade maatriks

	A	B	C	Σ
A	120	23	12	155
B	11	505	23	539
C	35	45	400	480
Σ	166	573	435	1174

#### Valemid

$$P_c = \frac{\sum_{i=1}^k n(i,i)}{n}$$

$$P_0 = \frac{\sum_{i=1}^k (n(+,i) \cdot n(i,+))}{n^2}$$

$$K = \frac{(P_c - P_0)}{(1 - P_0)}$$

Diagonaalil olevate pikslike oodatav osakaal **0 tähistab juhuslikkust** muutub vahemikus 0 ... 0.5;  
 $P_0 = (155 \cdot 166 + 573 \cdot 539 + 435 \cdot 480) / 1174^2 = 0.394$  **sõltub klasside arvust**, paljude ühtlase suurusega klasside puhul väiksem

Tegelikult diagonaalil olevate pikslike osakaal **C tähistab korrektsust** muutub vahemikus 0 ... 1  
 $P_c = (120 + 505 + 400) / 1174 = 0.873$

Kui klasse on palju, on diagonaali ruute suhteliselt vähem. Sellepärast õigete osakaal ( $P_0$ ) ise ei sobi klassifikatsioonide kooskõla hindamiseks.

Kapa koefitsient **K muutub vahemikus -1 ... 1**  
 $K = (0.873 - 0.394) / (1 - 0.394) = 0.79 = 79\%$   
**0 tähistab seose puudumist**  
**Millal on K väärtus -1?**

16

### Võiks meelde jääda

- Andmeanalüüs jaguneb kirjeldavaks ja tõestavaks.
- Kirjeldav analüüs
  - otsib tüüpilist (keskmist),
  - kirjeldab tunnuste varieeruvust,
  - iseloomustab seoseid tunnuste vahel ja sarnasust objektide vahel,
  - üldistab ja modelleerib.
- Kvantiil on koht muutuja arvteljel, millest ühel ja teisel pool on kindel osa vaatlustulemusi.
- Korrelatsioon on vastastikune seos tunnuste vahel.
- Korrelatsioonikordaja kirjeldab lineaarset seost ja kasutab ruuthälbeid.
- Mitteparameetrilised korrelatsioonikordajad ei sea tunnuste jaotustüübi eeldusi.
- Sarnasus on vastastikune seos objektide (vaatluste) vahel.
- Kapa kordaja mõõdab vastavust klassifitseerimistulemustes.
- Kapa muutumisvahemik on -1...+1. Null näitab juhuslikkuse korral oodatavat vastavust.