

Tilastotieteen johdantokurssi -luentokalvot

Yliopistonlehtori Juho Kopra

Itä-Suomen yliopisto, Tietojenkäsittelytieteen laitos

24.10.2023

- 1 Tilastotiede
- 2 Tilastotiedettä soveltavat tieteet
- 3 Miksi tilastotiedettä kannattaa opiskella?
- 4 Sattuma ja satunnaisuus
- 5 Tilastollisen päättelyn perusidea
- 6 R-kieli tilastotieteen työkaluna
- 7 Otanta
- 8 Mittaaminen
- 9 Muuttujat
- 10 Aineisto ja havaintomatriisi
- 15 Summamerkintä ja keskiarvon laskeminen
- 16 Hajontaluvut
- 17 Varianssin ja keskihajonnan laskeminen
- 18 Kvantiilit
- 19 Lineaarimuunnos
- 20 Riippuvuuden tarkastelu
- 21 Pearsonin tulomomentti-korrelaatiokerroin
- 22 Spearmanin järjestyskorrelaatiokerroin
- 23 Ristiintaulukko eli kontingenssitaulu
- 24 Tilastollinen testaaminen
- 25 Yhden otoksen Z-testi

Tilastotiede

“Statistics is the grammar of science.” - Karl Pearson

Mitä tilastotiede on?

“Statistics is a branch of applied mathematics that involves the collection, description, analysis, and inference of conclusions from quantitative data.” (Investopedia)

- Tilastotiede on tieteenala, joka käsittelee numeerisen aineiston hankintaa, kuvailua, analysointia, tulkintaa ja esittämistä.
- Käydään seuraavaksi läpi tilastotiede -sanon ja statistics -sanon alkuperää
- Tarkastellaan tutkimusprosessia ja tilastotieteen merkitystä suhteessa tutkimukseen.
- Tarkastellaan yleisimpiä tutkimustyypppejä

- Kantasana tila, eli maatila
- Mainittu ensimmäisen kerran teoksessa Suomen Suuriruhtinaskunnan nykyinen tilasto (1848)
- Kyse on ollut maatilojen ja alueiden luetteloinnista ja niitä koskevien tietojen kuvailusta sanallisesti ja keskiarvon avulla.
- Vedetään yhteen tietoja, saadaan aikaan tilasto (esim. montako asukasta maailmassa on peninkulmaa kohden)

- Tilastotiede on englanniksi statistics, joka tulee sanasta state (valtio) vuodelta 1791. Statistics on tarkoittanut valtion hallinnolliseen käyttöön kerättyä yleensä numeromuotoista aineistoa, jota voidaan esittää myös kuvien avulla.
- Suomeksi statistics kääntyy myös muotoon tunnusluvut. Tunnusluvuilla on paljon käyttöä tilastotieteessä.
- Suomessa viralliset tilastot mm. valtion käyttöön tuottaa Tilastokeskus (www.stat.fi). Ks. myös www.findikaattori.fi.

Tiede on (wikipedia.org)

“todellisuuden ilmiöiden ja niiden välisten suhteiden järjestelmällistä ja arvostelevaa tutkimista”

“sekä sen avulla saatua tietojen jäsentynyttä kokonaisuutta”

Eli kun opiskelemme tilastotiedettä, niin opiskelemme tiedettä. Olemme oppimassa tutkimisen taitoja ja samalla tietojen kokonaisuutta. Tilastotiede on luonteeltaan menetelmätiede, eli oppi menetelmistä, joilla voidaan tehdä tutkimusta koskien yleensä numeerisia aineistoja.

- 1 Ongelman asettaminen
 - 2 Ongelman täsmentäminen ja tutkimusstrategian laatiminen
 - 3 Aineiston kerääminen
 - 4 Aineiston kuvaaminen
 - 5 Aineiston analyysi
 - 6 Johtopäätösten teko
 - 7 Tutkielman tai raportin laatiminen
 - 8 Tutkimustulosten julkaiseminen
- Kohdat 1 ja 2: sovellusalan osaaminen ja tutkimuskirjallisuuden tuntemus
 - kohdat 3-6: Puhuttaessa määrällisestä tutkimuksesta, vaiheet 3-6 ovat oleellisesti TILASTOTIEDETTÄ riippumatta siitä, mistä sovellusalasta tutkimus on peräisin
 - Kohdat 7 ja 8: kirjoittamisen taito

- ③ Aineiston keräämisen tekniikat: otantateoria, koesuunnittelun teoria, mittaaminen
- ④ Aineiston kuvaamisen tekniikat: laadullisen aineiston luokittelu, numeeristen aineistojen muodostaminen, näiden tilastollinen kuvailu jakaumien ja tunnuslukujen avulla, graafinen esitys
- ⑤ Johtopäätösten teon tekniikat: tilastollisten testien teoria, tilastollinen päättely, tulosten tulkinta

Tilastotieteen menetelmät ovat hyvin tärkeitä soveltavien tieteiden aloilla. Ilman tilastomenetelmien tuntemista ei usein voi toteuttaa tutkimusta soveltavilla aloilla, kuten lääketiede, terveystieteet, hoitotiede, farmasia, biolääketiede, metsätiede, ympäristötiede, kauppatiede tai tietojenkäsittelytiede.

- Tilastotieteen tutkimus voidaan jakaa teoreettiseen ja soveltavaan
 - **teoreettinen tilastotiede** kehittää tilastomenetelmiä matematiikan avulla
 - **soveltava tilastotiede** käyttää tilastomenetelmiä jonkin toisen tieteenalan tutkimuksessa aineiston analyysissä
- Soveltavat tutkimukset voidaan jaotella sen mukaan, miten aineisto kerätään
 - **kokeellinen tutkimus**: tutkija kontrolloi ainakin joitakin muuttujia
 - **havainnoiva tutkimus**: tutkija seuraa ilmiötä tehden mittauksia, kuitenkin määräämättä muuttujien arvoja
 - **poikkileikkaustutkimus**: mitataan tietyssä ajankohtana tutkittavan ilmiön tilaa (satunnaisotanta)
 - **pitkittäistutkimus**: samoille yksilöille toistetaan mittauksia eri ajanhetkinä (esim. lasten pituuden seuranta)
 - Tutkimustyyppejä on paljon muitakin, erityisesti terveystieteissä, joita ei tässä käsitellä.

- **kokeellinen tutkimus:** tutkija kontrolloi ainakin joitakin muuttujia
 - esim. halutaan tutkia, miten eri tekijät vaikuttavat mansikan taimien kasvuun. Tutkija kontrolloi mullan laatua, lannoitusta, valon määrää ja kastelua.
- **havainnoiva tutkimus:** tutkija seuraa ilmiötä tehden mittauksia, kuitenkin määräämättä muuttujien arvoja
 - **poikkileikkaustutkimus:** mitataan tiettyä ajankohtana tutkittavan ilmiön tilaa (satunnaisotanta)
 - esim. FININTERVEYS selvittää suomalaisen aikuisväestön terveydentilaa. Mitataan pituus, paino, BMI, verenpaine, kysytään ravitsemusta, tupakointia, alkoholinkäyttöä, liikuntaa jne.
 - haastattelututkimus eli gallup: esim. vaaligallup
 - **pitkittäistutkimus:** samoille yksilöille toistetaan mittauksia eri ajanhetkinä (esim. lasten pituuden seuranta)

Tilastotiedettä soveltavat tieteet

- Seuraavaa kuutta kalvoa (tämä ml.) ei käsitellä luennolla, mutta silmäile ne läpi ja lue oman alasi osalta tarkemmin.
- Tilastotiede on tieteenala, joka kehittää menetelmiä, joiden avulla tutkitaan ilmiöitä ja niiden säännönmukaisuuksia numeerisen aineiston perusteella.
- Tilastotiede auttaa tekemään johtopäätöksiä ja ennusteita ilmiöistä, jotka ovat epävarmoja tai satunnaisia. Esimerkiksi tilastotiedettä voidaan käyttää arvioimaan Suomen väestön ominaisuuksia, kuten työttömien lukumääriä maakunnittain.

- Luonnontieteissä tilastotieteen menetelmiä käytetään mm. seuraaviin tarkoituksiin:
 - Hypoteesien testaaminen kokeellisilla tai havainnollisilla aineistoilla
 - Mallintaminen ja simulointi fysikaalisista, kemiallisista tai biologisista ilmiöistä
 - Mittausten ja kokeiden suunnittelu ja optimointi
 - Virheiden ja epävarmuuksien arviointi ja hallinta
 - Monimuuttuja-analyysi ja koneoppiminen suurten ja monimutkaisten aineistojen käsittelyssä

- Yhteiskuntatieteissä tilastotieteen menetelmiä käytetään mm. seuraaviin tarkoituksiin:
 - Kyselytutkimusten suunnittelu, toteutus ja analysointi
 - Tilastollinen päättely väestöllisistä tai yhteiskunnallisista ilmiöistä
 - Regressioanalyysi ja muut riippuvuussuhteiden tutkimisen menetelmät
 - Faktori- ja klusterianalyysi ja muut latenttien rakenteiden paljastamisen menetelmät
 - Aikasarja-analyysi ja ennustaminen taloudellisista tai sosiaalisista muuttujista

- Lääketieteessä tilastotieteen menetelmiä käytetään mm. seuraaviin tarkoituksiin:
 - Kliinisten tutkimusten suunnittelu, toteutus ja analysointi
 - Lääkekehitys ja lääkevaikutusten arviointi
 - Epidemiologia ja tartuntatautien leviämisen mallintaminen
 - Biostatistiikka ja geneettisten tai molekyylitason aineistojen analysointi
 - Etiikka ja tilastollinen merkitsevyys lääketieteellisessä päätöksenteossa

- Taloustieteessä tilastotieteen menetelmiä käytetään mm. seuraaviin tarkoituksiin:
 - Makro- ja mikrotaloudellisten ilmiöiden kuvaaminen, selittäminen ja ennustaminen
 - Taloudellisten teorioiden testaaminen empiirisillä aineistoilla
 - Taloudellisten mallien estimointi ja validointi
 - Ekonometria ja taloudellisten riippuvuussuhteiden tutkiminen
 - Tilinpito ja rahoitusmarkkinoiden analysointi

- Tekniikassa tilastotieteen menetelmiä käytetään mm. seuraaviin tarkoituksiin:
 - Laadunvalvonta ja prosessien parantaminen
 - Suunnittelukokeet ja tuotekehitys
 - Luotettavuus- ja riskianalyysi
 - Signaalinkäsittely ja kuvantaminen
 - Teollisuusmatematiikka ja optimointi

Miksi tilastotiedettä kannattaa opiskella?

Miksi tilastotiedettä kannattaa opiskella?

- Tilastotiedettä kannattaa opiskella yliopistossa sivuaineena, koska se:
 - Antaa sinulle vahvan metodologisen osaamisen ja kriittisen ajattelun taidon
 - Parantaa sinun mahdollisuuksiasi työllistyä ja edetä urallasi
 - Laajentaa sinun näkökulmaasi ja ymmärrystäsi eri alojen ilmiöistä ja ongelmista
 - Mahdollistaa sinulle monitieteisen yhteistyön ja verkostoitumisen

- Tilastotiede sopii hyvin sivuaineeksi minkä alan tutkintoon tahansa
- Voit valita tilastotieteen perus- tai aineopinnot tai edetä pidemmälle tilastollisen koneoppimisen suuntaan
- Voit opiskella tilastotiedettä eri ohjelmistoilla, kuten R tai SPSS
- Voit saada tilastollista neuvontaa opintojesi aikana

Sattuma ja satunnaisuus

- Arkipuheessa, kun jotain tapahtuu sattumalta, on se jotain mitä ei voinut arvata ennalta.
- **Satunnaisuus** on keskeinen termi tilastotieteessä ja todennäköisyyslaskennassa.
- Kun tieteessä jotain tapahtumaa pidetään täysin satunnaisena, se tarkoittaa, että kyseistä tapahtumaa ei voida mitenkään ennustaa.
 - Esim. nopanheiton silmäluku on satunnainen.
- Satunnaisuus ei tarkoita, että kaikki mahdolliset arvot olisivat yhtä todennäköisiä.
- Satunnaisen tapahtuman eri arvojen yleisyyttä voidaan jäsentää **todennäköisyyden** avulla.

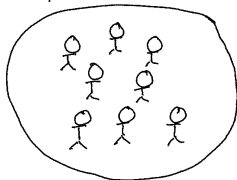
Tilastotieteessä satunnaisuutta käytetään tilanteesta riippuen hyödyksi tai sitä pyritään kontrolloimaan.

- tilastollinen mallintaminen pyrkii löytämään aineistosta ei-satunnaisen (ns. signaalin) ja erottamaan tästä satunnaisvaihtelun
- tilastollinen hypoteesintestaus selvittää, johtuuko aineistossa havaittu vaihtelu sattumasta, vai löytyykö näyttöä tutkittavan hypoteesin puolesta tai vastaan
- satunnaisotanta mahdollistaa tulosten yleistämisen perusjoukkoon
- koesuunnittelussa satunnaistamisen avulla voidaan erottaa lääkkeen todellinen vaikutus lumevaikutuksesta eli placebosta

Tilastollisen päättelyn perusidea

Tilastollisen päättelyn perusidea

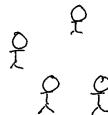
Pemysjoukko eli
populaatio



päätelmät
yleistetään
koskemaan
koko pemsjoukkoa



Satunnaisotos



ID	PITUUS	PAINO
1	158,1	58,4
2	172,3	65,6
3	189,5	79,0
4	166,4	75,3

R-kieli tilastotieteen työkaluna

- Vaikka tilastotiede on pitkälti matematiikkaa ja sen hyödyntämistä, käytännön työ on järkevintä tehdä tietokoneella. Kaikki kiinnostavat aineistot ovat sen verran suuria, että minkä tahansa analyysin tekeminen käsin on turhan aikaavieppää.
- Tällä kurssilla käytämme R-kieltä ja RStudiota työkaluina. Emme syvenny kieleen kovin syvällisesti, vaan tutustumme siihen pikkuhiljaa. Jo tässä vaiheessa on kuitenkin hyvä puhua hieman ohjelmoinnista, jotta voimme totutella antamaan tietokoneelle komentoja koodin muodossa.

Ohjelmoinnin perusperiaatteet

- Ohjelmointikielissä, kuten R-kieli, on tietyt perussanat, jotka on tiedettävä, jotta kieltä pystyy käyttämään. Samoin kielessä on oma ns. kielioppinsa, jota kutsutaan syntaksiksi. Sekä käytettävien sanojen että syntaksin on oltava täysin oikein, jotta tietokone suostuu tekemään mitään. Pienetkin piste- tai pilkkuvirheet tai väärä tai puuttuva kirjain aiheuttaa yleensä virheilmoituksen.
- Ohjelmointikielissä tietokoneelle annetaan ohjeita (siksi kai sanakin ohjelmointi?), joita tietokone noudattaa annetussa järjestyksessä. Siksi on väliä, että monivaiheiset komennot annetaan oikeassa järjestyksessä!
- Tässä vaiheessa ei ole tarpeen opetella kaikkea R:ään liittyvää ulkoa, vaan voit katsoa materiaalista oikeat komennot. Mikäli haluat oppia R:ää syvemmin, ei kuitenkaan ole haitaksi jos joitain keskeisimpiä työkaluja jää jo muistiin.

```
x <- 1.5 # sijoittaa luvun 1.5 objektiin x
x # tulostaa x:n tiedot konsoliin
y <- 1:5 # sijoittaa luvut 1-5 objektiin y
y2 <- y*2 # y2 on kaksi kertaa y:n arvot
?sample # avaa sample-funktion ohjesivun
sample(y,2) # poimii 2 havainnon satunnaisotoksen y-objektista
```

- Peruslaskutoimitukset: $+$ $-$ $*$ $/$ $^$
- Vertailuoperaattorit: $==$ $!=$ $<=$ $>=$

- R ja RStudio ovat eri ohjelmia. RStudio on graafinen ympäristö R-kielen käyttöön.
- R:ää voi käyttää jopa komentoriviltä tai jonkin muun graafisen ympäristön kautta. Windows-tietokoneilla R:n mukana asentuu R Gui, jota ei suositella käytettävän.

1. Editori

```

1 print("Hello, world!")
2
3
4 80 * (1 - 0.35)
5
6 x <- 3
7 y <- 5
8 z <- x + y
9
10 z <- "x + y"
11
12 z
13
14 # Sum of x and y
15 z <- x + y
16 z
17
18 # VECTORS
19
20 x <- c(1, 2, 7.4, 15, 0.2)
21
22 x <- seq(1, 10)
23 x
24
25 seq(from = 0, to = 1, by = 0.2)
26
27 3:9
28
29 Trep
30 rep(1, times = 5)
31 rep(c(2, 3), times = 5)
32 rep(c(2, 3), each = 5)
33
34 # VECTOR ARITHMETIC
35
36 x <- c(1, 2, 3, 6, 10)
37
38 [Rep Level]

```

2. Konsoli

```

R version 3.5.3 (2019-03-11) -- "Great Truth"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

>

```

3. Työtila

Environment is empty

4. Kuvaajat Paketit Manuaali

Name	Description	Version
lme4	Smoothing Metabolomics Analyses	0.2.1
askpass	Safe Password Entry for R, Git, and SSH	0.2.1
assertthat	Easy Pre and Post Assertions	0.2.1
backports	Reimplementations of Functions Introduced Since R-3.0	0.1.2
base64enc	Tools for base64 encoding	0.1-3
BH	Boost C++ Header Files	1.72.0-3
Biobase	Biobase: Base Functions for Bioconductor	2.46.0
BioGenomics	54 generic functions for Bioconductor	1.30.10
BiocManager	Access the Bioconductor Project Package Repository	1.30.10
BiocVersion	Set the appropriate version of Bioconductor packages	3.8.0
bitops	Bitwise Operations	1.0-6
brew	Templating Framework for Report Generation	1.0-6
callr	Call R from R	3.4.2
caTools	Tools: moving window statistics, GFI, Base64, ROC AUC, etc.	1.14
cli	Helpers for Developing Command Line Interfaces	2.0.1
clipr	Read and Write from the System Clipboard	0.7.0
clisymbols	Unescape Symbols in the R Prompt	1.2.0
colorspace	A toolbox for Manipulating and Assessing Colors and Palettes	1.4-1
commonmark	High-Performance CommonMark and Github Markdown Rendering in R	1.7
covr	Test Coverage for Packages	3.4.0
cowplot	Streamlined Plot Theme and Plot Annotations for 'ggplot2'	1.0.0
crayon	Colored Terminal Output	1.3.4
crossstalk	Inter-Widget Interactivity for HTML Widgets	1.0.0
curl	A Modern and Flexible Web Client for R	4.3
desc	Manipulate DESCRIPTION Files	1.2.0
devtools	Tools to Make Developing R Packages Easier	2.2.2

Pakettien asentaminen ja lataaminen

- R koostuu perusosasta (base) sekä lisäosista eli paketeista, joilla R:n toiminnallisuutta voidaan laajentaa. Asennetaan nyt remotes -niminen paketti, jota tarvitaan toisen paketin asentamista varten.

```
install.packages("remotes")
```

- Tyypillisesti saat asennettua tarvitsemasi paketit samalla tavalla kuin yllä asennettiin remotes.
- Jos haluamaasi pakettia ei ole lähetetty R:n yleiseen pakettiarkistoon (Comprehensive R Archive Network, CRAN) vaan se on githubissa, niin sen voi asentaa seuraavasti. Asennetaan siis kurssilla tarvittava paketti datas4uef.

```
remotes::install_github("jukop/datas4uef")
```

- Kun mitä tahansa tietoja halutaan työstää R:ssä, on se ladattava muistiin jollekin nimelle, jonka käyttäjä voi itse valita. Nimen ei tule kuitenkaan olla mikään varattu sana R-kielessä, joten jos tiedät jo joitakin varattuja sanoja, niin vältä niiden käyttöä objektien niminä.
- Esimerkiksi voin luoda objektin nimeltä `x`, jossa on luvut 1, 5 ja 6:
 - `x <- c(1, 5, 6)`
- Edellä käytimme sijoitusoperaattoria `<-` sekä yhdistimme luvut funktiolla `c` (`c` kuten `combine`)
 - Sijoitusoperaattorin merkkien välissä ei saa olla välilyöntiä!
 - Funktiokutsun tunnistaa siitä, että sen jälkeen on sulut ja sulkujen väliin tulee funktioiden argumentit eroteltuna pilkuilla. Tässä tapauksessa argumentit olivat arvot 1, 5 ja 6.

- Käsittelemme kulloisenkin aiheen yhteydessä miten asioita tehdään R:llä. Tällä tavoin R tulee ainakin vähän tutuksi, vaikkei sitä tarvitsekaan tällä kurssilla osata tentissä itse käyttää (mutta sitä saa käyttää). Osassa harjoitustehtäviä kuitenkin tarvitaan R:n käyttöä.
- Jos haluat paneutua heti enemmän R:n käyttöön, voit hyödyntää R-kieli -kurssin materiaalia:
https://vm3751.kaj.pouta.csc.fi/shiny/_book/r-kurssi.html

- Kurssin alkuvaiheessa esiintyy paljon tilastotieteen termejä, jotka on syytä opetella tuntemaan.
- Mikäli luet englanninkielistä tutkimus- tai oppikirjallisuutta, niin suomalainen tilastotieteen sanasto kannattaa ottaa avuksi: <https://sanasto.tilastoseura.fi/>
 - Luentokalvoilla en esittele englanninkielistä termistöä, mutta [verkkoluentomonisteessa](#) on usein mainittu myös englanninkielinen termi.
- Uusimmille menetelmille ei välttämättä ole vakiintunutta suomenkielistä nimeä. Siksi englanninkielisiä termejä käytetään edelleen paljon.

Otanta

- Otanta on tilastotieteen osa-alue, joka tarkastelee kuinka havaintoyksiköt kannattaa poimia perusjoukosta
- Otanta on tärkeää, koska sen avulla voidaan tehdä tilastollista päättelyä koko perusjoukosta otoksen perusteella
- Otannassa pyritään poimimaan edustava otos perusjoukosta.
- Yleisimpiä otantamenetelmiä ovat yksinkertainen satunnaisotanta palauttaen ja palauttamatta
 - Riippuen tutkimuksen tavoitteista voidaan käyttää myös monimutkaisempia otantamenetelmiä, jotka voivat olla tehokkaampia tiettyyn kysymykseen vastattaessa.

- Perusjoukko on tilastollisen tutkimuksen kohdejoukko eli se joukko yksiköitä, joista halutaan saada tietoa. Esimerkiksi jos halutaan tutkia suomalaisten mielipiteitä EU:sta, niin perusjoukkona ovat kaikki Suomen kansalaiset.
- Otos on perusjoukon osajoukko eli se joukko yksiköitä, joilta kerätään aineisto tilastollista analyysia varten. Esimerkiksi jos halutaan tutkia suomalaisten mielipiteitä EU:sta, niin otoksena voi olla 1000 satunnaisesti valittua Suomen kansalaista.

- Tyypillinen vaatimus hyvälle otokselle on että se on poimittu nk. satunnaisotantaa käyttäen
 - Ts. tutkija tai mikään taustatekijä ei määrää sitä, kuka tulee poimituksi otokseen. Jos näin olisi, niin on vaikeaa tehdä päätelmiä perusjoukkoa koskien.
 - Jos otos ei ole satunnaisotos, siitä käytetään nimeä näyte.
- Otoksen tulee olla edustava eli sen tulee heijastaa perusjoukon ominaisuuksia mahdollisimman hyvin.
- Suurempi otos mahdollistaa tarkemman käsityksen saamisen perusjoukosta, mutta yleensä aineiston kerääminen on kallista. Siksi otoksen koko on painottelua tulosten tarkkuuden ja tutkimuksen kustannusten kanssa.

- Yksinkertainen satunnaisotanta on yleisin ja yksinkertaisin otantamenetelmä
- Siinä poiminta tehdään suoraan arpomalla alkiot perusjoukosta
- Jokaisella alkiolla on sama todennäköisyys tulla valituksi otokseen

Palauttaen vai palauttamatta?

- Yksinkertainen satunnaisotanta voidaan tehdä joko palauttaen tai palauttamatta valittu alkio takaisin perusjoukkoon
- Tehtäessä otanta palauttaen sama alkio voi tulla otokseen monta kertaa
- Tehtäessä otanta palauttamatta jokainen alkio voidaan valita vain kerran

Yksinkertainen satunnaisotanta palauttamatta

```
# luodaan vektori porukka, josta poimitaan otos  
porukka <- c("Erkki", "Jaska", "Niina")  
# porukka -vektorissa on kolme alkiota, joten  
# voidaan poimia kokoa 1, 2 tai 3 oleva otos  
sample(porukka, size=2)
```

```
## [1] "Jaska" "Niina"
```

```
# uusi funktiokutsu poimii uuden otoksen  
sample(porukka, size=2)
```

```
## [1] "Erkki" "Jaska"
```

Otanta R-kielellä jatk.

Yksinkertainen satunnaisotanta palauttaen

```
sample(porukka, size=4, replace=TRUE)
```

```
## [1] "Erkki" "Erkki" "Niina" "Jaska"
```

- Asettamalla nk. alkuluku (seed) ennen koodin suorittamista arpoo aina saman otoksen
 - Vertaa kaverin kanssa tai suorita useita kertoja peräkkäin

```
set.seed(42)  
sample(porukka, size=2)
```

```
## [1] "Erkki" "Niina"
```

Mittaaminen

- Jotta tutkittavasta ilmiöstä saadaan muodostettua numeerinen/määrällinen aineisto, on tehtävä mittauksia.
- Mittaamisen kohteena ovat tilastoyksiköt eli ne yksiköt, joista halutaan saada tietoa. Tilastoyksikkö voi olla esimerkiksi henkilö, organisaatio, tapahtuma tai asiakirja.
- Mittaamisen avulla voidaan määrittää tutkittavien tilastoyksiköiden ominaisuuksia numeerisesti.
- Mittaamisen tuloksena syntyy muuttujia, jotka kuvaavat tilastoyksiköiden ominaisuuksia. Muuttujan arvo kullekin tilastoyksikölle kertoo tilastoyksikölle tehdyn mittaustuloksen.

Esimerkkejä mittaamisesta

- pituus mitattuna mittanauhalla, punnittu paino
- henkilön oma käsitys painostaan (eri kuin edellä)
- ajan mittaaminen sekuntikellolla
- liikenneonnettomuuksien määrä Valtatie 9:llä vuoden aikana
- formulakuljettajan sijoitus kilpailussa
- lämpötila $^{\circ}C$ (vrt. $^{\circ}F$)
- silmien väri (laatuero)
- auton merkki, lintulaji
- älykkyyssosamäärä

Mittaustapoja on ainakin

- mittalaitteella mittaaminen (esim. vaaka)
- henkilöltä kysyminen (kasvotusten, puhelimitse, netissä tai kirjeellä)
- tietokannasta virallisen tilaston tms. mukaan, rekisteriaineistot
- havainnoimalla ja kirjaamalla havainnot ylös
- kyselylomakkeella

Muuttujien tyypit

- Muuttujat voidaan luokitella eri tavoin niiden ominaisuuksien perusteella. Yksi tapa on luokitella muuttujat niiden arvojen tyyppin mukaan numeerisiin ja luokkamuuttujiin.
- **Numeeriset muuttujat** ovat sellaisia, joiden arvot ovat lukuja. Numeeriset muuttujat voidaan jakaa edelleen jatkuviin ja diskreetteihin muuttujiin. Jatkuvilla muuttujilla on teoriassa äärettömän monta mahdollista arvoa tietyllä välillä, esimerkiksi pituus tai paino. Diskreeteillä muuttujilla on vain rajallinen määrä mahdollisia arvoja, esimerkiksi lasten lukumäärä tai silmien väri.
- **Luokkamuuttujat** ovat sellaisia, joiden arvot ovat sanoja tai merkkejä. Luokkamuuttujat voidaan jakaa edelleen dikotomisiin ja moniarvoisiin muuttujiin. Dikotomisilla muuttujilla on vain kaksi mahdollista arvoa, esimerkiksi sukupuoli tai kyllä/ei-vastaus. Moniarvoisilla muuttujilla on useampia mahdollisia arvoja, esimerkiksi poliittinen kanta tai ammatti.

Muuttuja kuuluu johonkin seuraavista mitta-asteikoista:

- 1 Laatueroasteikko
- 2 Järjestysasteikko
- 3 Välimatka-asteikko
- 4 Suhdeasteikko

① Laatueroasteikko

- toisensa poissulkevat luokat
- esim. mies/nainen, Volvo/Audi/BMW, talitiainen/sinitäinen/...
- voidaan sanoa, mihin luokkaan havainto kuuluu, mutta suuruusjärjestys ei ole mielekäs

② Järjestysasteikko

- myös toisensa poisulkevat luokat
- eri arvojen suuruusjärjestystä voidaan vertailla
- esim. Likert-asteikko täysin eri mieltä/jokseenkin eri mieltä/ei samaa eikä eri mieltä/jokseenkin samaa mieltä/täysin samaa mieltä
- eri arvojen etäisyyttä ei voida ilmaista numeerisesti

3 Välimatka-asteikko

- yleensä jatkuva, mutta mittaustarkkuus voi tehdä myös diskreetin
- kahden lukuarvon erotus määrittää etäisyyden
- käytössä on mittayksikkö, esim. $^{\circ}C$
 - $25^{\circ}C - 20^{\circ}C = 5^{\circ}C$
- jakolaskua ei pidetä mielekkäänä

4 Suhdeasteikko

- välimatka-asteikon ominaisuudet ja lisäksi:
- muuttujan arvon saadessa arvon 0, mitattava ominaisuus häviää
- esim. lämpötila kelvinasteina ($^{\circ}K$) mittaa lämpöliikkeen määrää ($0^{\circ}K =$ lämpöliike lakkaa), etäisyys sentteinä (cm), paino (kg)
- jakolaskulla on mielekäs tulkinta, esim. 75 cm pitkä lapsi on 1.5-kertaa pidempi kuin 50 cm pitkä lapsi $\frac{75\text{ cm}}{50\text{ cm}} = 1.5$

- R-kieli olettaa, että kaikki laskutoimitukset ovat sallittuja
 - Ts. jos muuttuja on numeerinen, niin R olettaa suhdeasteikon
- Toisaalta käytännön kannalta ei yleensä ole väliä onko muuttuja suhdeasteikollinen vai välimatka-asteikollinen (jakolasku on harvinainen vaatimus perusmenetelmissä)
 - Siispä keskitytään siihen onko muuttuja numeerinen (vähintään välimatka-asteikollinen) vai järjestysasteikollinen tai laatueroasteikollinen
- Esimerkiksi olkoon muuttuja, joka on koodattu lukuarvoin 1, 2 ja 3 siten että 1 = "Eri mieltä", 2="Samapa tuo" ja 3="Samaa mieltä".

```
x <- c(1,2,2,3)
```

- laatueroasteikolliseksi

```
x_nom <- factor(x, levels=c(1,2,3),  
               labels=c("Eri mieltä", "Samapa tuo", "Samaa mieltä"))  
x_nom
```

```
## [1] Eri mieltä   Samapa tuo   Samapa tuo   Samaa mieltä  
## Levels: Eri mieltä Samapa tuo Samaa mieltä
```

- järjestysasteikolliseksi

```
x_ord <- factor(x, levels=c(1,2,3),  
               labels=c("Eri mieltä", "Samapa tuo", "Samaa mieltä"),  
               ordered=TRUE)  
x_ord
```

```
## [1] Eri mieltä   Samapa tuo   Samapa tuo   Samaa mieltä  
## Levels: Eri mieltä < Samapa tuo < Samaa mieltä
```

Muuttujat

- Muuttujalla viitataan mitattavan kohteen ominaisuuteen, joka vaihtelee yksiköstä tai mittauksesta toiseen. Esimerkiksi henkilön pituus, paino ja poliittinen kanta ovat muuttujia.
- Muuttujan arvo on muuttujan mittaustulos tietyssä yksikössä. Esimerkiksi henkilön pituuden arvo voi olla 170 cm ja poliittisen kannan arvo voi olla vasemmisto.
- Muuttujien arvot syntyvät mittaamisen tuloksena. Arvojen vaihtelusta syntyy muuttujan jakauma, jota voidaan kuvata tilastollisilla tunnusluvuilla ja tilastollisen grafiikan eli kuvaajien avulla.

Aineisto ja havaintomatriisi

- Tarkastellaan seuraavaksi aineiston ja havaintomatriisin käsitteitä.
- Seuraaviin perussääntöihin on olemassa poikkeuksia, sillä kaikkia tietoja ei voi esittää sujuvasti taulukkomuodossa.
- Taulukkomuotoinen esitystapa kuitenkin toimii tutkimustarkoituksiin erittäin usein.

Aineisto ja havaintomatriisi

- Kun samoista tilastoyksiköistä koostetaan useita muuttujia yhteen saadaan aineisto.
- Yleensä aineisto esitetään havaintomatriisin muodossa.
 - Havaintomatriisin kukin sarake sisältää yhden muuttujan arvot.
 - Havaintomatriisin kukin rivi sisältää mittaustiedot yhdeltä tilastoyksiköltä.
 - Taulukon solu sisältää muuttujan arvon.

sukupuoli	ikä	pituus	paino	pääaine
mies	27	194	80	TTRA2
mies	26	170	67	TTRA2
nainen	23	165	47	TILTK
mies	17	170	61	TK1K
nainen	25	168	50	TILTK

Aineistosta tarkemmin

- Aineistossa on hyvä olla yksi sarake, joka yksilöi tilastoyksiköt.
 - Edelliseltä kalvolta tämä puuttui.
- Jos aineistossa on useita mittauksia samalta tilastoyksiköltä, niin silloin aineistossa voi olla useita rivejä.
 - R-kieltä käytettäessä tämä on suositeltavampi tapa muodostaa aineisto vs. että aineistossa olisi useita sarakkeita eri mittauskertoja varten.

- Leikitään ajatuksella, että aineistossa on kaksi mittausta jokaiselta henkilöltä vuoden välein

Toistomittausaineisto:

id	sukupuoli	ikä	paino
1	mies	27	80
1	mies	28	82
2	mies	26	67
2	mies	27	68
3	nainen	23	47
3	nainen	24	47
4	mies	17	61
4	mies	18	59
5	nainen	25	50
5	nainen	26	51

- Tarkemmin ja yleisemmin aineiston käyttöönottoa opetellaan kurssilla R-kieli, 2op.
- Tällä kurssilla tarvittavat aineistot saa käyttöön seuraavasti:
 - Lataa aiemmin asentamasi paketti `datas4uef` komennolla `library(datas4uef)`
 - Sitten komentoa `data` käyttäen ja tietämällä aineiston nimen saat aineiston käyttöön. Esim. `data(hlotsim_dat)`.

```
library(datas4uef)
data(hlotsim_dat)
```

Objektien tyypit

- Objekteja on eri tyyppisiä, joista tärkeimmät ovat `numeric`, `character` ja `factor`
 - `numeric` ilmaisee välimatka- ja suhdeasteikollisia muuttujia.
 - `factor` ilmaisee luokitteluasteikollisia ja järjestysasteikollisia muuttujia
 - `character` on tekstimuotoinen tieto, joka voidaan muuttaa `numeric` tai `factor`-tyyppiseksi.
- Aineiston objektityyppi on `data.frame` tai `tibble` (uudempi)
- Objektin tyyppin voi tarkastaa funktiolla `typeof`, esim. `typeof(hlotsim_dat)`
 - Jos objektin tyyppi on `data.frame` eli aineisto niin yhden muuttujan poiminta onnistuu `$`-merkin avulla (`aineiston_nimi$muuttuja`):

```
hlotsim_dat$ikä
```

```
## [1] 27 26 23 17 25 28 20 21 33 32 25 25 21 15 19 25 27 29 20 19 20  
## [26] 24 16 26 20 26 23 24 20 31 27 24 25 26 24 28
```

- Aineistot ovat siis `data.frame`-tyyppisiä tai vastaavia.
- Aineiston sarakkeissa olevat muuttujat ovat R:ssä vektoreita.
 - Vektoreita voi luoda itse `c()` -komennolla, esim. `c(1,5,6)` luo vektorin, jossa on luvut 1, 5, ja 6.
 - Aineiston sarakkeiden ei tarvitse olla keskenään samaa tyyppiä.
- Uuden sarakkeen lisääminen aineistolle `dat` onnistuu
`dat$uusi_sarake <- c(1,2,3)`
 - Huomaa, että sijoitusoperaattorin oikealla puolella on oltava sama määrä lukuja kuin `data.frame`:ssa on rivejä. Rivimäärän saat komennolla `nrow(dat)`
- Saraketta voi muokata ylikirjoittamalla vanhan sarakkeen tiedot. Esim. kerrotaan sarakkeen arvot kahdella
`dat$vanha_sarake <- 2*dat$vanha_sarake`
 - Tässä on oltava huolellinen, sillä muutoksia ei voi peruuttaa.

Jakaumat

- Jakaumia on kahdenlaisia:
 - Empiirinen jakauma eli aineistosta laskettu jakauma
 - Teoreettinen jakauma eli todennäköisyysjakauma

Empiirinen jakauma kuvaa, mitä arvoja muuttuja saa aineistossa ja miten yleisiä mitkäkin arvot ovat. Se voidaan esittää taulukkona tai kuvaajan avulla. Sen tietoa voidaan myös kuvailla tunnusluvulla.

Todennäköisyysjakauma ilmaisee satunnaismuuttujan arvojen todennäköisyydet. Esim. normaalijakauma (tuttu lukiosta) on todennäköisyysjakauma. Todennäköisyysjakaumiin palataan myöhemmin.

- Tosimaailmaa koskevassa tutkimuksessa ei tyypillisesti tunneta todellista todennäköisyysjakamaa, vaan sitä pyritään ns. estimoimaan (eli tekemään päätelmiä siitä) aineiston avulla.

Diskreetille aineistolle (laatuero- ja järjestysasteikko) voidaan laskea havaintojen lukumäärät kutakin mahdollista arvoa kohti.

Aineisto:

id	sukupuoli
1	mies
2	nainen
3	nainen
4	nainen
5	mies
6	nainen
7	mies

Frekvenssijakauma:

sukupuoli	lkm
mies	3
nainen	4

Prosenttijakauma:

sukupuoli	%
mies	42.9
nainen	57.1

Frekvenssi- ja prosenttijakauma R-kielellä

```
# frekvenssijakauma  
tab <- table(dat$sukupuoli)
```

```
tab
```

```
##  
##   mies nainen  
##      3      4
```

```
# prosenttijakauma  
pros_jakauma <- round(prop.table(tab)*100,1)
```

```
pros_jakauma
```

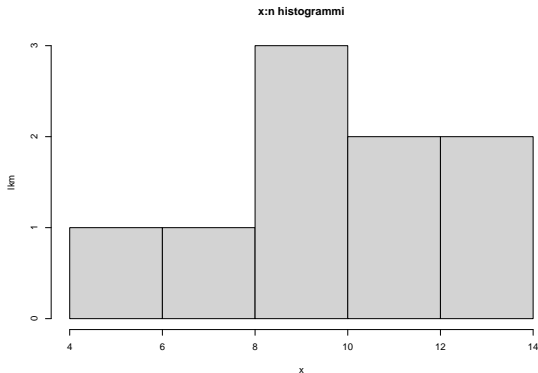
```
##  
##   mies nainen  
##  42.9   57.1
```

- Jatkuvalle muuttujalle arvoja ei kannata taulukoida, sillä taulukosta saattaa tulla liian pitkä luettavaksi.
- Yksi vaihtoehto on luokitella jatkuvan muuttujan aineisto ja esittää se taulukon tai kuvan avulla.
 - Tässä saatetaan kuitenkin menettää arvokasta informaatiota.
- Taulukon sijaan voidaan myös käyttää tunnuslukuja (engl. statistics) kyseisen jakauman kuvailemiseksi.

Aineisto:

id	x
1	5.12
2	6.47
3	8.38
4	8.72
5	9.39
6	10.41
7	10.92
8	12.07
9	13.79

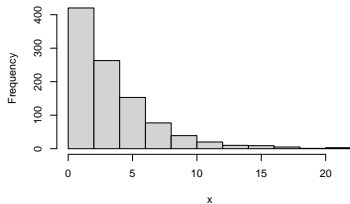
- Luokitellaan väleille 4-6, 6-8, 8-10, 10-12, 12-14.
 - Luokan alaraja ei kuulu luokkaan, yläraja kuuluu (puoliavoin väli).



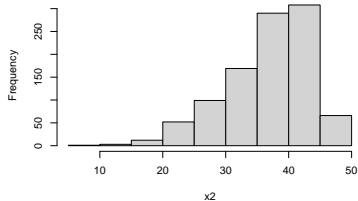
Histogrammista voidaan tulkita ainakin seuraavat seikat (ks. seuraavia kalvoja):

- vinous oikealle tai vasemmalle
- symmetrisyys (vastakkainen vinoudelle)
- monihuippuisuus
- valitun pylväiden määrän vaikutus kuvaajaan
- oudokki eli poikkeava havainto

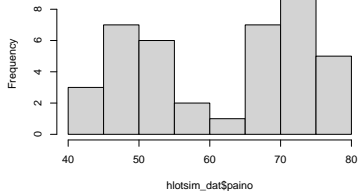
Histogram of x



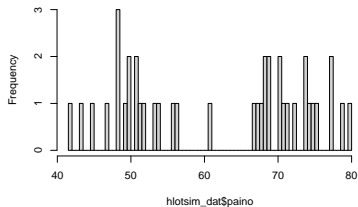
Histogram of x2



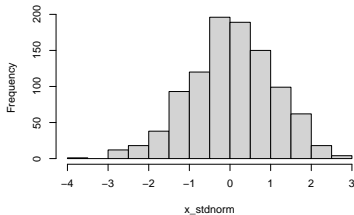
Histogram of hlotsim_dat\$paino



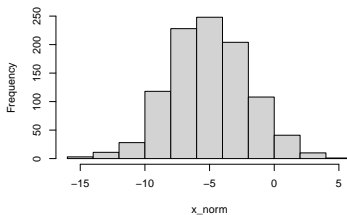
Histogram of hlotsim_dat\$paino



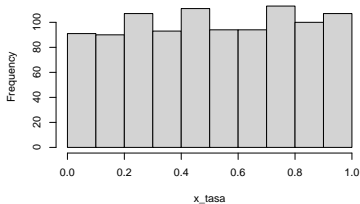
Histogram of x_stdnorm



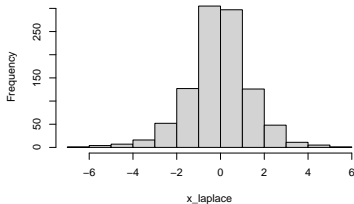
Histogram of x_norm



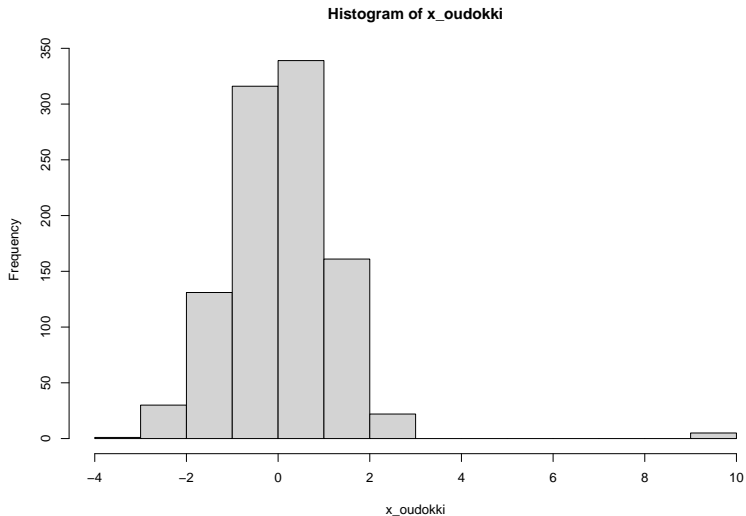
Histogram of x_tasa



Histogram of x_laplace



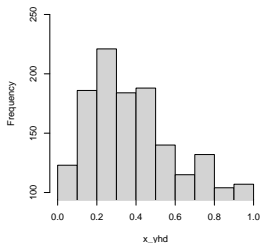
Poikkeava havainto eli oudokki



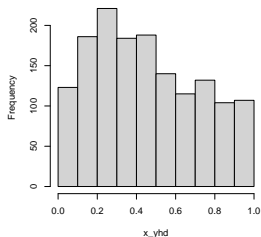
Ei kannata

- valita y-akselia siten että nolla ei ole x-akselin risteymässä (osa histogrammia leikkautuu pois)
- valita vaikeasti ymmärrettävää pylväsjakoa
- käyttää epätasavälistä pylväsjakoa

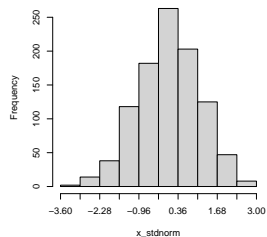
Histogram of x_yhd



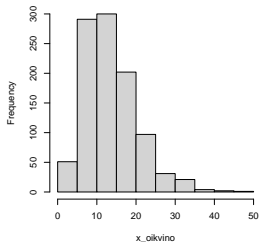
Histogram of x_yhd



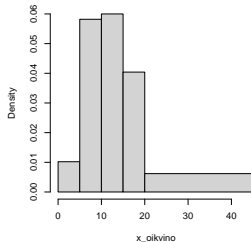
Histogram of x_stdnorm



Histogram of x_oikvino



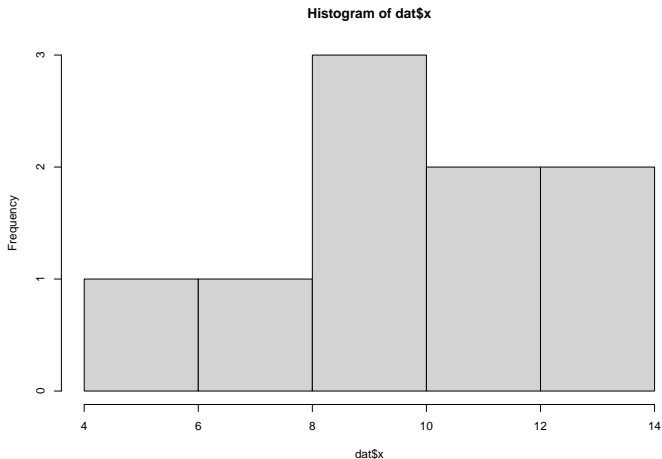
Histogram of x_oikvino



Histogrammi R-kielillä

Histogrammi aineistosta `dat` muuttujalle `x` oletusasetuksin

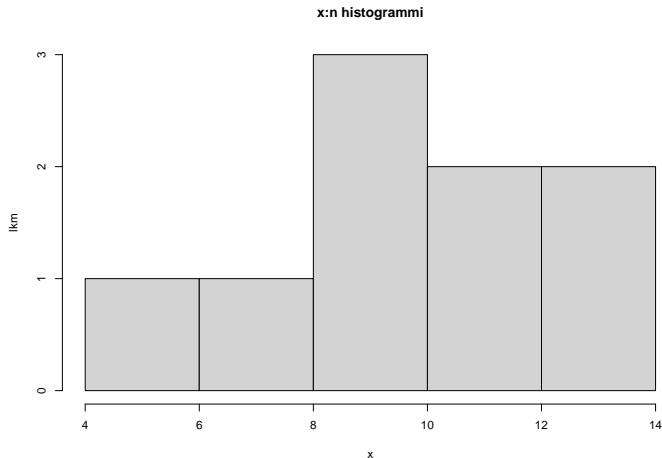
```
hist(dat$x)
```



Muutetaan x- ja y-akseleiden selitteet sekä otsikko:

```
# histogrammi
```

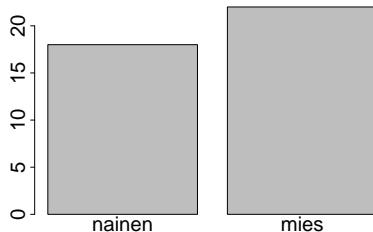
```
hist(dat$x,xlab="x",ylab="lkm",main="x:n histogrammi")
```



Kuvaajat

- Tutustutaan seuraaviin kuvaajiin tai kuvioihin:
 - pylväskuvaaja eli palkkikuvio
 - piirakkakuvio eli sektorikuvio
 - laatikkokuvio
 - histogrammi (käsiteltiin aiemmin)
 - sirontakuvio eli hajontakuvio
- Opetellaan tulkitsemaan ja piirtämään kyseiset kuviot
- Muita kuvioita voi opetella halutessaan osion viimeisen kalvon linkkien avulla

- luokitteluasteikollisen muuttujan frekvenssijakauman tai prosenttijakauman esittämiseen



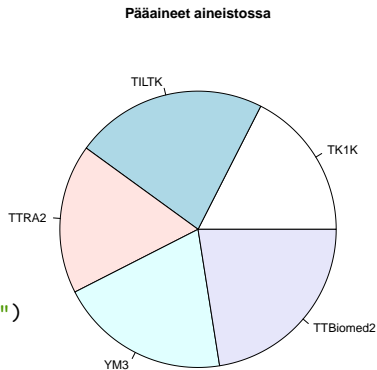
```
barplot(table(hlotsim_dat$sukupuoli))
```

```
# suhteelliset osuudet
```

```
barplot(prop.table(table(hlotsim_dat$sukupuoli)))
```

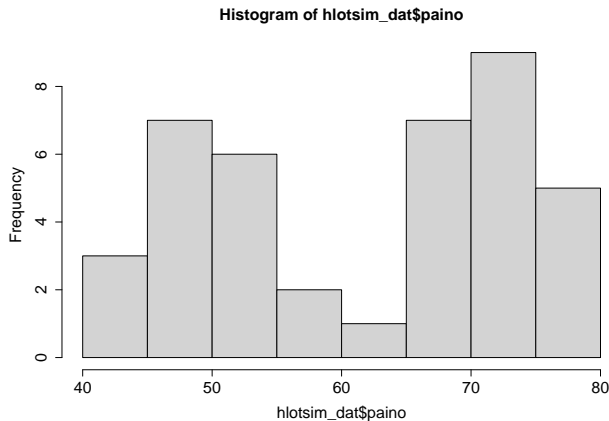
- luokitteluasteikollisen muuttujan lukumäärien tai osuuksien/prosenttien kuvaamiseen
- ei sovellu tieteelliseen esittämiseen
 - käyttöohje: Älä käytä!

```
lkm <- table(hlotsim_dat$pääaine)  
pie(lkm,main="Pääaineet aineistossa")
```

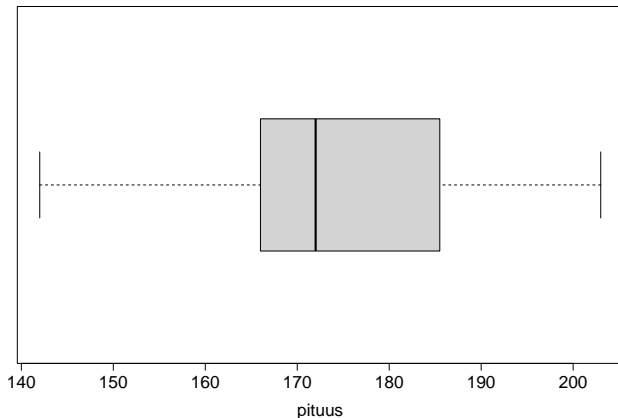


Histogrammi

- vähintään välimatka-asteikollinen muuttuja (diskreetti tai jatkuva)

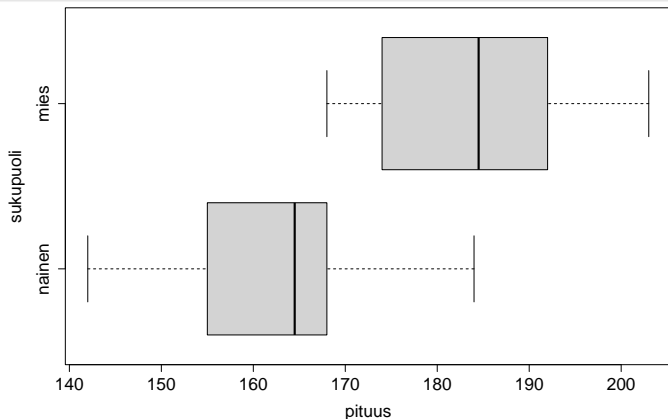


- vähintään välimatka-asteikollinen muuttuja (diskreetti tai jatkuva)
- voidaan piirtää pysty- tai vaakasuuntaisena



```
boxplot(hlotsim_dat$pituus,xlab="pituus",horizontal=TRUE)
```

Laatikkokuvio ryhmittäin

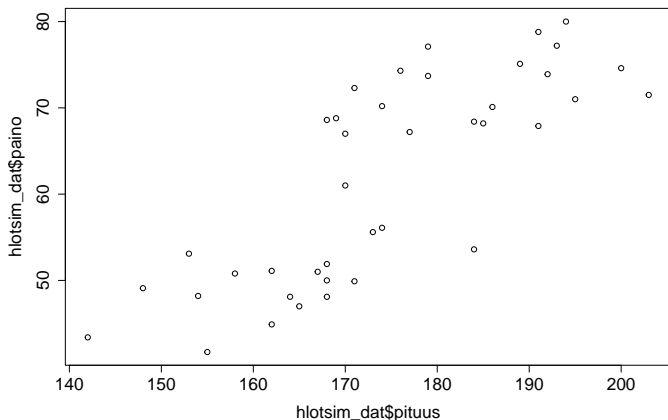


- Ao. koodissa `pituus ~ sukupuoli` on nk. formula, joka ilmaisee että tehdään haluttu asia (boxplot) pituudelle muuttujan sukupuoli suhteen
 - vaatii että aineisto on pitkässä muodossa, eli pituus oltava yhdessä sarakkeessa (ei eri sarakkeita eri sukupuolille)

```
boxplot(pituus ~ sukupuoli, data=hlotsim_dat, horizontal=TRUE)
```

Sirontakuvio eli hajontakuvio

- samoilta tilastoyksiköiltä kaksi vähintään välimatka-asteikollista muuttujaa (mittausparit)



```
plot(hlotsim_dat$pituus, hlotsim_dat$paino)
```

- Kuvaajia on lukuisia muitakin kuin edellä on esitetty. Ks. <https://r-graph-gallery.com/>
- Nykyaikainen tapa tehdä kuvaajia R-kielellä (ei käytössä tällä kurssilla): <https://ggplot2-book.org/>

Tilastolliset tunnusluvut

- Edellä esitettyjen kuvaajien avulla saa nopeasti käsityksen aineiston luonteesta ja jakaumista.
- Tarkempaa analyysiä varten on järkevää laskea tilastollisia tunnuslukuja. Tunnuslukuja tarvitaan myös, jotta osataan piirtää kuvaajia.
- Tunnusluvut ovat keskeisiä myös myöhemmin, kun siirytään tilastollisen päättelyn pariin.

- Tilastolliset tunnusluvut ovat lukuja, jotka kuvaavat tilastollista aineistoa tiivistetysti ja havainnollisesti.
 - Tunnusluvut siis kuvaavat jakauman ominaisuuksia.
- Tilastollisia tunnuslukuja ovat esimerkiksi keskiarvo, mediaani, moodi, varianssi, keskihajonta ja korrelaatio.
 - Näistä keskilukuja ovat keskiarvo, mediaani, moodi.
 - Hajontalukuja ovat varianssi ja keskihajonta.
- Muuttujan mitta-asteikosta riippuu, mitä tunnuslukuja muuttujalle voidaan laskea.
- Tilastollisia tunnuslukuja voidaan käyttää aineiston analysointiin, vertailuun ja esittämiseen.

Keskiarvo eli otoskeskiarvo

- Keskiarvo (aritmeettinen keskiarvo) kuvaa aineiston keskimääräistä arvoa.
- Vaatii vähintään välimatka-asteikollisen muuttujan, koska perustuu etäisyyksien laskemiseen (yhteen- ja vähennyslaskuun).
- Keskiarvo lasketaan jakamalla aineiston arvojen summa aineiston lukumäärällä.
- Keskiarvon merkintä on \bar{x} ja laskukaava

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

- Termiä otoskeskiarvo käytetään tähdentämään että tarkoitetaan aineistosta eli otoksesta laskettua keskiarvoa.
- Keskiarvon vastine populaatiossa on odotusarvo μ .

- Esimerkki: Jos luokan oppilaiden pituudet ovat 160, 165, 170, 175 ja 180 cm, niin niiden keskiarvo on $\frac{160+165+170+175+180}{5} = 170$ cm.
- Keskiarvo on herkkä poikkeaville arvoille (oudokeille).
 - Esim. lisätään viimeksimainittuun aineistoon hlö, jonka pituus on 250 cm. Nyt keskiarvo on 183.3 cm.

- Mediaani on tilastollinen tunnusluku, joka kuvaa aineiston keskimmäistä arvoa. Se jakaa aineiston kahteen yhtä suureen osaan: puolet arvoista on mediaania pienempiä ja puolet suurempia.
- Voidaan laskea vähintään järjestysasteikolliselle muuttujalle.
 - Käytetään usein välimatka-asteikolliselle muuttujalle sen robustisuuden vuoksi.
- Mediaani lasketaan järjestämällä aineiston arvot suuruusjärjestykseen ja valitsemalla keskimmainen arvo. Jos arvoja on parillinen määrä, mediaani on kahden keskimmäisen arvon keskiarvo.
- Mediaanin merkintä on M tai $Q2$.

- Esimerkki: Jos luokan oppilaiden pituudet ovat 160, 165, 170, 175 ja 180 cm, niin mediaani on 170 cm. Jos oppilaiden pituudet ovat 160, 165, 170, 175, 180 ja 185 cm, niin mediaani on $\frac{170+175}{2}=172.5$ cm.
- Poikkeavat havainnot eivät vaikuta herkästi mediaanin arvoon ts. mediaani on robusti tunnusluku.
 - Esim. lisätään viimeksimainittuun aineistoon hlö, jonka pituus on 250 cm. Nyt mediaani on 175 cm (aiemmin 172.5).

- Moodi kertoo aineiston arvon, joka on yleisin.
- Moodi soveltuu laatueroasteikollisten muuttujien kuvailuun.
 - Esim. katsastusasemalla vieraili autoja päivässä seuraavasti.
Moodiksi saadaan “Volkswagen”, koska sen frekvenssi on suurin.

autot

##	Audi	Tesla	Volkswagen
##	8	5	13

- Numeerisille muuttujille voidaan laskea moodiluokka esim. histogrammista. Moodiluokka on silloin se väli, jonka pylväs on korkein.
- Teoreettisille jakaumille on helpompi määrittää moodi jakauman kuvaajan korkeimman kohdan perusteella.

Keskiarvo, mediaani ja moodi R:ssä

Keskiarvo, mediaani ja moodi R:ssä

- Lasketaan aineiston hlotsim_dat muuttujan pituus keskiarvo ja mediaani käyttäen funktioita mean(x) ja median(x)

```
mean(hlotsim_dat$pituus)
```

```
## [1] 174.3
```

```
median(hlotsim_dat$pituus)
```

```
## [1] 172
```

- Moodin laskemiseksi ei ole valmista R-funktiota, vaan se tulee päätellä frekvenssitaulukon perusteella. Tehdään frekvenssitaulukko muuttujalle pääaine
 - Kaksi yhtä suurta suurinta frekvenssiä, eli moodeja ovat TILTK ja TTBiomed2

```
table(hlotsim_dat$pääaine)
```

```
##
```

```
##      TK1K      TILTK      TTRA2      YM3 TTBiomed2
```

```
##          7          9          7          8          9
```

Summamerkintä ja keskiarvon laskeminen

- Summamerkintä on matemaattinen symboli, jolla ilmaistaan lukujen summaa lyhyesti ja kätevästi.
- Summamerkinnässä käytetään kreikkalaista kirjainta sigma (\sum), jonka alapuolella on summauksen alkuarvo ja yläpuolella loppuarvo. Summattavat luvut ilmaistaan yleensä indeksoidulla muuttujalla, jonka indeksi kasvaa yhdellä joka askeleella.
- Esimerkiksi summamerkitä

$$\sum_{i=1}^5 i$$

tarkoittaa lukujen 1, 2, 3, 4 ja 5 summaa eli 15.

- Sama summamerkitä voi saada myös muodon $\sum_{i=1}^5 i$

- Keskiarvon laskemiseen voidaan käyttää summamerkintää apuna. Jos lukujoukon luvut ovat x_1, x_2, \dots, x_n , niin niiden keskiarvo on $\frac{1}{n} \sum_{i=1}^n x_i$.
- Esimerkiksi lukujen 2, 4, 6 ja 8 keskiarvo on

$$\frac{1}{4} \sum_{i=1}^4 x_i = \frac{1}{4}(2 + 4 + 6 + 8) = \frac{20}{4} = 5$$

Hajontaluvut

- Hajontaluvut ovat tilastollisia tunnuslukuja, jotka kuvaavat aineiston hajontaa eli vaihtelua.
- Hajontalukuja ovat esimerkiksi varianssi, keskihajonta, kvartiilit, kvartiiliväli ja variaatiokerroin.
- Hajontalukuja voidaan laskea numeeriselle aineistolle
 - Kvartiilit voidaan laskea myös järjestysasteikolliselle (harvoin käytetty)
- Hajontalukuja voidaan käyttää aineiston vaihtelun arvioimiseen ja vertailuun.

- Varianssi kuvaa aineiston arvojen tyypillistä neliöpoikkeamaa keskiarvosta.
- Varianssi lasketaan neliöimällä aineiston arvojen poikkeamat keskiarvosta ja jakamalla havaintojen määrä vähennettynä yhdellä.
- Otosvarianssin merkintä on s^2 ja sitä vastaava populaation varianssi on σ^2 .
- Otosvarianssin laskukaava on

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2,$$

missä n on aineiston havaintojen lkm, x_i käy läpi aineiston havainnot ja \bar{x} on keskiarvo.

- Esimerkki: Jos luokan oppilaiden pituudet ovat 160, 165, 170, 175 ja 180 cm, niin varianssi on

$$\frac{(160-170)^2 + (165-170)^2 + (170-170)^2 + (175-170)^2 + (180-170)^2}{5-1} = 62.5$$

cm².

- Huomaa mittayksikön muutos! cm \rightarrow cm² Siksi varianssin lukuarvoa on vaikea tulkita.

- Varianssin tulkinnan vaikeudesta päästään eroon ottamalla varianssista neliöjuuri: $\sqrt{x \text{ cm}^2} = \sqrt{x} \text{ cm}$. Näin saadaan keskihajonta.
- Keskihajonta on kuvaa aineiston arvojen keskimääräistä poikkeamaa keskiarvosta.
 - Mittayksikkö on sama kuin alkuperäisillä havainnoilla.
- Keskihajonnan merkintä on s ja populaation keskihajonta on σ .
- Esimerkki: Jos luokan oppilaiden pituudet ovat 160, 165, 170, 175 ja 180 cm, niin keskihajonta on $\sqrt{62.5} \approx 7.91 \text{ cm}$.

- Kvartiilit ovat lukuja, jotka jakavat aineiston neljään yhtä suureen osaan. Ne ovat mediaanin kaltaisia tunnuslukuja.
- Perusidea kvartiilien laskemisessa on: järjestetään aineiston arvot suuruusjärjestykseen ja valitaan neljännesvälien kohdalta arvot. Kvartiilien tarkkaa laskemista varten on useita tapoja (R:ssä on 9 eri tapaa) ja eri ohjelmistot käyttävät eri tapoja. Käsitellään laskemista myöhemmin.
- Kvartiileja merkitään Q_1 , Q_2 ja Q_3 .
- Esimerkki: Jos luokan oppilaiden pituudet ovat 160, 165, 170, 175 ja 180 cm, niin ensimmäinen kvartiili on laskutavasta riippuen 162.5 cm tai 165 cm, toinen kvartiili eli mediaani on 170 cm ja kolmas kvartiili on 177.5 cm tai 175 cm.

- Kvartiiliväli kertoo millä välillä aineiston keskimäinen 50 prosenttia havainnoista sijaitsee.
 - Kvartiilivälistä voidaan laskea kvartiilivälin pituus, joka on kolmannen ja ensimmäisen kvartiilin erotus.
- Kvartiilivälin merkintä on $(Q1, Q3)$ ja kvartiilivälin pituus on IQR (interquartile range).
- Esimerkki: Jos luokan oppilaiden pituudet ovat 160, 165, 170, 175 ja 180 cm, niin ensimmäinen kvartiili on 162.5 cm ja kolmas kvartiili on 177.5 cm. Kvartiiliväli on $(162.5, 177.5)$ ja $IQR = 177.5 - 162.5 = 15$ cm.

Varianssin ja keskihajonnan laskeminen

- Varianssi on aina ei-negatiivinen eli $s^2 \geq 0$. Mitä suurempi varianssi on, sitä enemmän havainnoissa on hajontaa eli vaihtelua. Jos varianssi on nolla, kaikki luvut ovat yhtä suuria kuin keskiarvo.
- Keskihajonta kuvaa siis havaintojen vaihtelua keskiarvon ympärillä. Mitä pienempi keskihajonta on, sitä tiiviimmin luvut ovat ryhmittynyt keskiarvon ympärille. Mitä suurempi keskihajonta on, sitä laajemmin luvut ovat levittäytyneet keskiarvosta poispäin.

- Esimerkiksi lukujen 2, 4, 6 ja 8 varianssin ja keskihajonnan laskemiseksi aloitetaan laskemalla niiden keskiarvo:

$$\bar{x} = \frac{1}{4} \sum_{i=1}^4 x_i = \frac{1}{4}(2 + 4 + 6 + 8) = \frac{20}{4} = 5.$$

- Sitten lasketaan jokaisen luvun ja keskiarvon erotuksen neliö:

$$(2 - 5)^2 = (-3)^2 = 9, (4 - 5)^2 = (-1)^2 = 1,$$

$$(6 - 5)^2 = 1^2 = 1, (8 - 5)^2 = 3^2 = 9$$

- Summataan neliöt yhteen ja jaetaan lukujen määrä vähennettynä yhdellä $n - 1$:

$$\frac{9 + 1 + 1 + 9}{4 - 1} = \frac{20}{3} \approx 6.67$$

- Keskihajonta on neliöjuuri (R:ssä `sqrt(x)`) varianssista:

$$\sqrt{\frac{20}{3}} = \frac{2\sqrt{5}}{\sqrt{3}} \approx 2.58$$

- Lasketaan aineiston `hlotsim_dat` muuttujan pituus varianssi ja keskihajonta

```
var(hlotsim_dat$pituus)
```

```
## [1] 211.959
```

```
sd(hlotsim_dat$pituus)
```

```
## [1] 14.55881
```


Kvantiilit

- Kvantiilit (eli fraktiilit) ovat hajontalukuja, jotka jakavat aineiston tiettyyn määrään yhtä suuria osia. Ne ovat yleistys mediaanille ja kvartiileille, jotka ovat itsessään myös kvantiileja.
- Kvantiilit lasketaan järjestämällä aineiston arvot suuruusjärjestykseen ja valitsemalla tiettyjen osuuksien kohdalta arvot. Jos osuus ei osu tarkalleen yhteen arvoon, voidaan käyttää interpolaatiota eli laskennallista väliarvoa.
 - Kvantiilien määrittämiseksi on useita eri menetelmiä (samat kuin kvartiileille). Esimerkiksi R-kielessä kvantiilit voidaan laskea yhdeksällä eri tavalla.
- Kvantiilien merkintöjä ovat Q_k , missä k on osuus desimaalimuodossa.
- Esimerkki: Jos luokan oppilaiden pituudet ovat 160 cm, 165 cm, 170 cm, 175 cm ja 180 cm, niin 0.25-kvantiili eli ensimmäinen kvartiili on 162.5 cm ja 0.75-kvantiili eli kolmas kvartiili on 177.5 cm.

Olkoon meillä suuruusjärjestykseen asetettu aineisto ($n=20$), jossa pienin arvo 1 ja suurin arvo 19.

##	1.	2.	3.	4.	5.	6.	7.	8.	9.	10.	11.	12.	13.	14.	15.	16.	17.	18.	19.	20.
##	1	4	4	7	8	9	10	10	11	11	13	14	14	14	14	15	15	18	18	19

Käytetään seuraavilla kalvoille merkintöjä

- $x_{(j)}$ = “j:nneksi pienin arvo”
- $\lfloor m \rfloor$: luku m pyöristettynä alaspäin lähimpään kokonaislukuun

p -kvantiili voidaan laskea funktiolla

$$Q_i(p) = (1 - \gamma)x_{(j)} + \gamma x_{(j+1)}$$

missä

- γ määrittää miten paljon painotetaan lukuja $x_{(j)}$ ja $x_{(j+1)}$
- i vastaa R-funktion `quantile(x, probs, type)` `type`-argumenttia. Oletuksena R:ssä `type=7`.

p -kvantiili on siis painotettu summa kahdesta luvusta, jotka määräytyvät luvun $j = \lfloor np + m \rfloor$ perusteella. Luvun j laskemiseksi tarvittava luku m riippuu valitusta arvosta i (`type`). Oletusarvolla `type=7` jolloin luku m saadaan kaavalla $m = 1 - p$.

Kvantiilien laskeminen käsin (type=7)

- Halutaan laskea p -kvantiili. Esim. jos lasketaan alakvartiili, niin $p = 0.25$. Olkoon n aineiston havaintojen lukumäärä, tässä $n = 20$.
- Lasketaan apuparametri γ . Jos halutaan laskea R:n käyttämää oletusparametrisointia käyttäen (type=7), niin $m = 1 - p = 1 - 0.25 = 0.75$.
- Koska $n = 20$ ja $p = 0.25$ niin nyt $j = \lfloor 20 \cdot 0.25 + 0.75 \rfloor = \lfloor 5.75 \rfloor = 5$.
- Nyt γ saadaan siten että lasketaan ensin $g = np + m - j$, mistä saadaan $g = 20 \cdot 0.25 + 0.75 - 5 = 0.75$. Jos type = 4,5,6,7,8 tai 9, niin $\gamma = g$.

Eli tässä tapauksessa

$$Q_7(0.25) = (1 - 0.75)x_{(5)} + 0.75x_{(6)} = 0.25 \cdot 8 + 0.75 \cdot 9 = 8.75$$

Kvantiilien laskeminen R-kielellä

```
# halutut luvut vektorissa (ei tarvitse olla järjestetty)
luvut <- c(1,4,4,7,8,9,10,10,11,11,13,14,14,14,14,15,15,18,18,19)
# 0.25-kvantiili eli alakvartiili
quantile(luvut,probs=0.25) #oletus type=7
```

```
## 25%
## 8.75
```

```
quantile(luvut,probs=0.25,type=2)
```

```
## 25%
## 8.5
```

```
quantile(luvut,probs=0.25,type=6) # SPSS
```

```
## 25%
## 8.25
```

```
# useita kvantiileja samassa komennossa
quantile(luvut,probs=c(0,0.25,0.50,0.75,1))
```

```
##      0%      25%      50%      75%     100%
##  1.00   8.75  12.00  14.25  19.00
```

Viiden numeron yhteenveto

- Viiden numeron yhteenveto: minimi, alakvartiili, mediaani, yläkvartiili ja maksimi
 - Tarvitaan laatikkokuvion piirtämiseen (tähän palataan myöhemmin)
- Voidaan laskea nopeasti R-funktiolla `summary` (tulostaa myös keskiarvon) tai `fivenum`
 - Lasketaan muuttujalle pituus

```
summary(hlotsim_dat$pituus)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    142.0   166.5   172.0   174.3   185.2   203.0
```

```
fivenum(hlotsim_dat$pituus)
```

```
## [1] 142.0 166.0 172.0 185.5 203.0
```

Lineaarimuunnos

Olkoon a ja b reaalityyji ja X on muuttuja.

Lineaarimuunnos havainnoille x_1, x_2, \dots, x_n on

$$y_i = a + bx_i$$

jossa i käy läpi arvot $1, \dots, n$.

- Esim. jos muuttuja X saa lukuarvot 160, 165, 170, 175, 180 ja on valittu $a = -70$ ja $b = 0.50$ niin Y -muuttujan arvot ovat: 10, 12.5, 15, 17.5, 20
- a -parametri siirtää havaintojen sijaintia: negatiivinen siirtää vasemmalle, positiivinen oikealle.
- b -parametri skaalaa havaintojen etäisyyttä origosta: $b < 1$ siirtää lähemmäs nollaa, $b > 1$ kauemmas

Lineaarimuunnoksen vaikutus keskiarvoon

Lineaarimuunnos vaikuttaa keskiarvoon samalla tapaa kuin yksittäisiin arvoihin:

$$\bar{y} = a + b\bar{x}$$

Ts. jos tiedämme että X-muuttujan keskiarvo on 170 niin voimme laskea että Y-muuttujan keskiarvo on $-70 + 0.50 * 170 = 15$

- Kokeillaan R:n avulla

```
x <- c(160, 165, 170, 175, 180)
mean(x)
```

```
## [1] 170
```

```
y <- -70+0.50*x
mean(y)
```

```
## [1] 15
```

Lineaarimuunnoksen vaikutus varianssiin

Olkoon s_x^2 X-muuttujan varianssi ja Y-muuttuja on saatu lineaarimuutoksella kuten edellisellä kalvolla. Nyt

$$s_y^2 = b^2 s_x^2$$

Ts. b -kerroin vaikuttaa varianssiin siten että alkuperäisen muuttujan varianssi kerrotaan b^2 :lla, jolloin saadaan uuden muuttujan varianssi. Kokeillaan R:n avulla

```
x <- c(160, 165, 170, 175, 180)
var(x)
```

```
## [1] 62.5
```

```
y <- -70+0.50*x
var(y)
```

```
## [1] 15.625
```

```
0.50^2 * var(x)
```

```
## [1] 15.625
```

Lineaarimuunnoksen vaikutus keskihajontaan

Olkoon s_x X-muuttujan keskihajonta ja Y-muuttuja on saatu lineaarimuutoksella kuten edellisellä kalvolla. Nyt

$$s_y = |b|s_x$$

Ts. uusi keskihajonta saadaan kertomalla alkuperäisen muuttujan keskihajontaa b :n itseisarvolla. Kokeillaan R:n avulla

```
x <- c(160, 165, 170, 175, 180)
sd(x)
```

```
## [1] 7.905694
```

```
y <- -70+0.50*x
sd(y)
```

```
## [1] 3.952847
```

```
0.50 * sd(x)
```

```
## [1] 3.952847
```

Riippuvuuden tarkastelu

Kahden muuttujan välistä riippuvutta voidaan tarkastella riippuen muuttujien luonteesta:

- Määrälliset muuttujat
 - Hajontakuvion avulla
 - Pearsonin tai Spearmanin korrelaatiokertoimen avulla
- Järjestysasteikolliset muuttujat
 - Spearmanin korrelaatiokertoimen avulla
- Luokitteluasteikolliset muuttujat
 - Ristiintaulukoinnin avulla (testaaminen jää peruskurssille)

Pearsonin tulomomenttikorrelaatiokerroin

- Soveltuu jatkuville ja väh. välimatka-asteikollisille muuttujille
- Mittaa lineaarista riippuvuutta eli miten vahvasti hajontakuviossa pisteparvi on keskittynyt suoralle
- Arvot välillä $-1 \leq r_{xy} \leq 1$
- Määritellään tätä varten otoskovarianssi, joka on otosvarianssin läheinen sukulainen

Olkoon muuttujista X ja Y havaintoja x_1, x_2, \dots, x_n ja y_1, y_2, \dots, y_n s.e. (x_i, y_i) on mitattu samalta tilastoyksiköltä. Otoskovarianssi s_{xy} on

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

Jos lasketaan muuttujan X kovarianssi itsensä kanssa, niin saadaan X :n varianssi s_x^2 .

Pearsonin (tulomomentti)korrelaatiokerroin saadaan jakamalla otoskovarianssi X :n ja Y :n otoskeskihajonnoilla:

$$r_{xy} = \frac{s_{xy}}{s_x \cdot s_y} \quad (1)$$

$$= \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x \cdot s_y} \quad (2)$$

$$= \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})}{s_x} \frac{(y_i - \bar{y})}{s_y} \quad (3)$$

- $r_{xy} = 0$: ei lainkaan lineaarista riippuvuutta
- $r_{xy} = 1$: täydellinen positiivinen lineaarinen riippuvuus (arvot nousevalla suoralla)
- $r_{xy} = -1$: täydellinen negatiivinen lineaarinen riippuvuus (arvot laskevalla suoralla)

Muuta huomioitavaa

- korrelaatiolla ei ole mittayksikköä (ei ole esim. kg)
- on herkkä oudokeille kuten keskiarvo ja keskihajonta

Pearsonin korrelaatio käsin laskien

- Lasketaan ensin keskiarvot ja keskihajonnat
 - $\bar{x} = 8.175$ ja $s_x = 1.5107945$,
 - $\bar{y} = 4.25$ ja $s_y = 4.5$

x	y	\bar{x}	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})/s_x$	$(y - \bar{y})/s_y$	$\frac{(x - \bar{x})}{s_x} \frac{(y - \bar{y})}{s_y}$
6.3	-2						
8.2	4						
8.2	8						
10.0	7						

$$r_{xy} = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})}{s_x} \frac{(y_i - \bar{y})}{s_y} =$$

Olkoon, lineaarimuunnokset havainnoille $(x_1, y_1), \dots, (x_n, y_n)$

$$x'_i = a + bx_i$$

ja

$$y'_i = c + dy_i$$

jossa i käy läpi arvot $1, \dots, n$ ja a, b, c, d ovat reaalityyppisiä lukuja. Nyt

$$r_{x'y'} = r_{xy}$$

kun b ja d ovat etumerkeiltään samat, ja

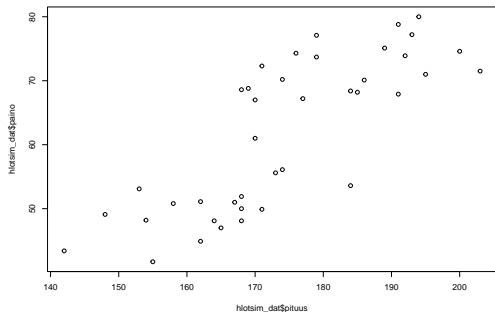
$$r_{x'y'} = -r_{xy}$$

kun b ja d ovat etumerkeiltään erit.

Pearsonin korrelaatio R-kielellä

- Pearsonin korrelaatio R:ssä lasketaan funktiolla `cor(x,y)` vektoreille `x` ja `y`

```
library(datas4uef)
data("hlotsim_dat")
plot(hlotsim_dat$pituus, hlotsim_dat$paino)
```



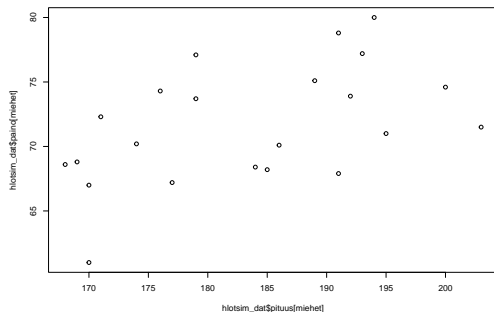
```
cor(hlotsim_dat$pituus, hlotsim_dat$paino)
```

```
## [1] 0.7930337
```

```
# rajoitetaan miehiin
```

```
miehet <- hlotsim_dat$sukupuoli=="mies"
```

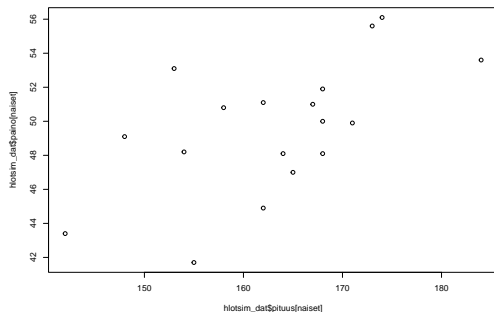
```
plot(hlotsim_dat$pituus[miehet],hlotsim_dat$paino[miehet])
```



```
cor(hlotsim_dat$pituus[miehet],hlotsim_dat$paino[miehet])
```

```
## [1] 0.4974266
```

```
# rajoitetaan naisiin  
naiset <- hlotsim_dat$sukupuoli=="nainen"  
plot(hlotsim_dat$pituus[naiset],hlotsim_dat$paino[naiset])
```



```
cor(hlotsim_dat$pituus[naiset],hlotsim_dat$paino[naiset])
```

```
## [1] 0.5958413
```


Spearmanin järjestyskorrelaatiokerroin

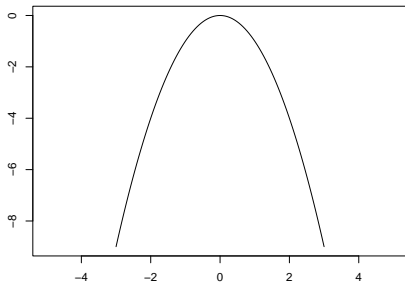
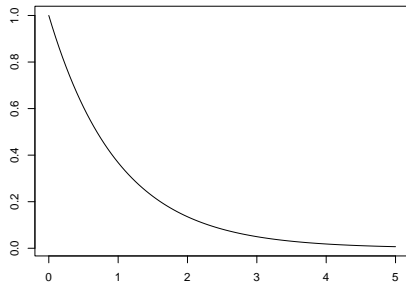
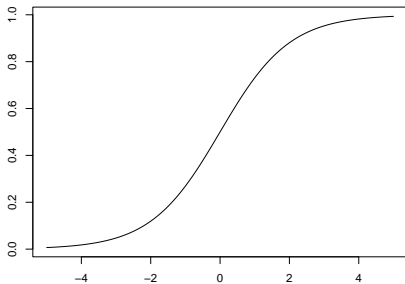
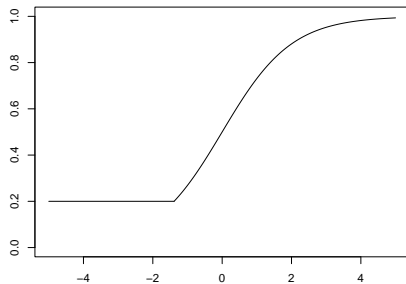
Spearmanin järjestyskorrelaatiokerroin

- Soveltuu jatkuville ja väh. järjestysasteikollisille muuttujille
- Mittaa monotonista riippuvuutta, eli miten yhteensopivia kahden muuttujan havaintoarvojen suuruusjärjestykset ovat.
 - Tyypillisin käyttötilanne on silti välimatka-asteikollisille muuttujille, kun lineaarinen yhteys ei ole uskottava, mutta monotoninen on.
- Arvot välillä $-1 \leq r_{xy} \leq 1$
- Määritellään tätä varten monotonisesti kasvava (ja vähenevä) funktio

Olkoon f funktio s.e. $y = f(x)$.

- Funktio f on monotonisesti kasvava, jos kaikille x -arvoille pätee:
 - kun x :n arvoa kasvatetaan, niin y :n arvo kasvaa tai vähintään pysyy samana
 - eli matemaattisemmin ilmaistuna jos $x_1 < x_2$ niin $y_1 = f(x_1) \leq f(x_2) = y_2$
- Vastaavasti funktio f on monotonisesti vähenevä, jos kaikille x -arvoille pätee:
 - kun x :n arvoa kasvatetaan, niin y :n arvo vähenee tai korkeintaan pysyy samana
 - eli matemaattisemmin ilmaistuna jos $x_1 < x_2$ niin $y_1 = f(x_1) \geq f(x_2) = y_2$

Mitkä seuraavista ovat monotonisia funktioita?



- Oletetaan jälleen, että muuttujista X ja Y on havaintoja x_1, x_2, \dots, x_n ja y_1, y_2, \dots, y_n s.e. (x_i, y_i) on mitattu samalta tilastoyksiköltä, ja että X ja Y ovat vähintään järjestysasteikollisia.
- Muodostetaan X :n ja Y :n havaintoja vastaavat järjestysluvut kummallekin erikseen s.e. pienin arvo saa järjestysluvuksi 1, seuraavaksi pienin 2 jne. Mikäli jotain arvoa on enemmän kuin yksi, eli tulee "tasapeli", niin tällaisille annetaan kyseisten järjestyslukujen keskiarvo. Esim. jos havainnot ovat 6.3, 8.2, 8.2 ja 10.0, niin järjestyslukuiksi tulee 1, 2.5, 2.5, 4, missä $2.5 = \frac{2+3}{2}$.
 - R-kielessä $\text{rank}(x)$ määrittää x -vektorin järjestysluvut.
- Jos muuttujat ovat valmiiksi järjestysasteikollisia, niin järjestysluvut ovat valmiina, eikä niitä tarvitse erikseen määrittää.

Yleisesti pätee seuraava:

$$r_s = \frac{s_{R(X), R(Y)}}{s_{R(X)} \cdot s_{R(Y)}},$$

missä $R(X)$ ja $R(Y)$ ovat X - ja Y -havaintojen järjestysluvut.

Ts. kun X :n ja Y :n havainnot korvataan järjestysluvuillaan, niin niiden Spearmanin korrelaatio voidaan laskea järjestysluvuista samalla kaavalla kuin Pearsonin korrelaatio.

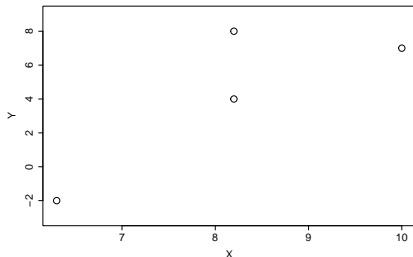
Mikäli järjestysluvut ovat kokonaislukuja (ei tasapelejä) niin voidaan käyttää seuraavaa kaavaa:

$$r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)},$$

missä

$d_i = R(x_i) - R(y_i)$ on järjestyslukujen erotus havaintoparille (x_i, y_i) .
 n = havaintoparien (rivien) määrä.

X	Y	R(X)	R(Y)
6.3	-2		
8.2	4		
8.2	8		
10.0	7		



R-kielillä:

```
x <- c(6.3, 8.2, 8.2, 10.0)
y <- c(-2, 4, 8, 7)
cor(x,y,method="spearman")

## [1] 0.6324555
```


Ristiintaulukko eli kontingenssitaulu

- Ristiintaulukoinnin ja khi toiseen -riippumattomuustestin (Tilastotieteen peruskurssilla) avulla voidaan selvittää riippuvatko kaksi luokitteluasteikollista muuttujaa toisistaan tilastollisesti.
 - Perusidea on, että jos muuttujat eivät riipu toisistaan, niin toisen arvon tietäminen ei anna mitään tietoa (todennäköisyyksimielessä) toisen arvosta.
 - Jos riippuvuutta on, niin ristiintaulukolla voidaan tarkastella millaista riippuvuus on: mitkä havaintoarvoparit ovat yleisempiä kuin, jos muuttujat olisivat riippumattomia.
- Oppikirjoissa näkee usein esimerkkejä, joissa muuttujat on johdettu välimatka-asteikollisista muuttujista luokittelemalla ne. Tämä ei ole välttämättä hyvä idea tutkimuskäyttöön, sillä ko. tilanteeseen on parempia menetelmiä.

Kuopion sosiaali- ja terveystointien yhdistämisen vaikutuksista tehdyssä tutkimuksessa kyselylomakkeessa on kysytty henkilön sukupuoli ja takaavatko vastaanottotilat yksityisyyden.

Taulukko 1: Vastaajien kokemus vastaanottotilojen yksityisyyden suojasta Kuopiossa.

Mielipide	Mies	Nainen	Yht.
Huonosti	7	144	151
Hyvin	40	292	332
Yht.	47	436	483

Taulukko 2: Vastaajien kokemus vastaanottotilojen yksityisyyden suojasta Kuopiossa (ehdolliset prosenttijakaumat).

Mielipide	Mies	Nainen	Yht.
Huonosti	14.9 %	33.0 %	31.3 %
Hyvin	85.1 %	67.0 %	68.7 %
Yht.	100.0 %	100.0 %	100.0 %

Ristiintaulukon muodostaminen aineistosta

Sukupuoli	Tiedekunta
M	FILO
N	LUMET
N	LUMET
M	TT
N	TT

##

##

FILO TT LUMET

##

M 1 1 0

##

N 0 1 2

Taulukon indeksointi ja merkinnät

- Ristiintaulukko (ja kaikki muut taulukot matematiikassa) indeksoidaan s.e. indeksi $i = 1, 2, \dots, r$ määrittää millä rivillä taulukossa ollaan. Indeksillä $j = 1, 2, \dots, c$ määrittää mikä sarake on kyseessä.
- Koska ristiintaulukko sisältää lukumääriä eli frekvenssejä, niin merkitään niitä f_{ij} rivillä i sarakkeessa j .

f_{11}	f_{12}	f_{13}
f_{21}	f_{22}	f_{23}

- Oletetaan jatkossa, että muuttuja X :n arvot ovat riveillä ja muuttujan Y arvot sarakkeissa.

- Usein on tarpeen laskea frekvenssien summia riveittäin tai sarakeittain
 - Rivisumma riville i on

$$f_{i.} = \sum_{j=1}^c f_{ij},$$

missä i on valitun rivin indeksi (luku).

- Sarakesumma sarakkeelle j on

$$f_{.j} = \sum_{i=1}^c f_{ij},$$

missä j on valitun sarakkeen indeksi (luku).

Marginaalijakauma eli reunajakauma

- X :n marginaalijakauma saadaan laskemalla kaikki rivisummat $f_{i\cdot}$ ja raportoimalla ne yhdessä niihin liittyvien X :n arvojen kanssa.
- Vastaavasti Y :n marginaalijakauma saadaan laskemalla kaikki sarakesummat $f_{\cdot j}$ ja raportoimalla ne yhdessä niihin liittyvien Y :n arvojen kanssa.

Ehdollinen (empiirinen) jakauma

- Kiinnittämällä muuttujan X arvo, ja tarkastelemalla muuttujan Y jakaumaa saadaan Y :n ehdollinen jakauma ehdolla X .
 - Vastaavasti vaihtamalla X :n ja Y :n roolit keskenään saadaan X :n ehdollinen jakauma Y :n suhteen (eli ehdolla Y).
- Taulukon rivit ovat sarakemuuttujan Y ehdolliset jakaumat. Samoin sarakkeet ovat rivimuuttujan X ehdolliset jakaumat.
- Ehdollisista (frekvenssi)jakaumista voidaan edelleen määrittää ehdolliset prosenttijakaumat.
 - Riippuvuuden tarkastelu empiirisesti tapahtuu ehdollisia prosenttijakaumia vertailemalla. Suuret erot ehdollisissa prosenttijakaumissa antavat viitteitä riippuvuudesta - varsinkin, jos aineisto ei ole kovin pieni!

- Ristiintaulukon muodostaminen
- Marginaalijakaumien laskeminen vaiheuttain + funktiolla `add.margins`
- Ehdollisten jakaumien laskeminen
- Ehdollisten prosenttijakaumien laskeminen

- Havaintomatriisista ristiintaulukko saadaan tehtyä table-funktiolla

```
data("hlotsim_dat")  
rt <- table(hlotsim_dat$sukupuoli, hlotsim_dat$pääaine)  
rt
```

```
##  
##           TK1K  TILTK  TTRA2  YM3  TTBiomed2  
##   nainen      2      7      4      4          1  
##   mies       5      2      3      4          8
```

- Jos aineistoa ei ole, mutta lukumäärät ovat tiedossa, niin voidaan luoda taulukko matrix-komennolla (aineisto, joka oli aiemmin)

```
rt2 <- matrix(c(7, 144,  
               40,292),nrow=2,ncol=2,byrow=T)  
dimnames(rt2) <- list("Mielipide"=c("Huonosti","Hyvin"),  
                      "Sukupuoli"=c("Mies","Nainen"))  
rt2
```

```
##           Sukupuoli  
## Mielipide  Mies Nainen  
## Huonosti    7   144  
## Hyvin      40   292
```

Marginaalijakaumien laskeminen vaiheittain

- Marginaalijakaumat eli reunajakaumat

```
rowSums(rt2)
```

```
## Huonosti    Hyvin  
##      151      332
```

```
colSums(rt2)
```

```
##   Mies Nainen  
##    47    436
```

```
sum(rt2)
```

```
## [1] 483
```

- marginaaliprosenttijakaumat

```
prop.table(rowSums(rt2))
```

```
## Huonosti    Hyvin  
## 0.3126294 0.6873706
```

```
prop.table(colSums(rt2))
```

Marginaalijakaumien laskeminen funktiolla `addmargins`

- Tai `addmargins` -funktiolla
 - **Huom!** `addmargins`-funktio antaa taulukon, josta ei tule enää laskea rivi-/sarakesummaa eikä osuuksia R:ssä! Sum-sarake ja rivi eivät ole erityisasemassa.

```
addmargins(rt2)
```

```
##              Sukupuoli
## Mielipide   Mies Nainen Sum
## Huonosti    7    144 151
## Hyvin       40    292 332
## Sum         47    436 483
```

Ehdollisten jakaumien laskeminen

```
prop.table(rt2,margin=1) #riveittäin, eli rivien osuuksien summa on 1
```

```
##           Sukupuoli
## Mieli-pide      Mies      Nainen
## Huonosti 0.04635762 0.9536424
## Hyvin      0.12048193 0.8795181
```

```
prop.table(rt2,margin=2) #sarakeittain
```

```
##           Sukupuoli
## Mieli-pide      Mies      Nainen
## Huonosti 0.1489362 0.3302752
## Hyvin      0.8510638 0.6697248
```

Ehdollisten prosenttijakaumien laskeminen

```
100*prop.table(rt2,margin=1) #riveittäin
```

```
##           Sukupuoli
## Mieliipide      Mies   Nainen
## Huonosti    4.635762 95.36424
## Hyvin       12.048193 87.95181
```

```
100*prop.table(rt2,margin=2) #sarakeittain
```

```
##           Sukupuoli
## Mieliipide      Mies   Nainen
## Huonosti    14.89362 33.02752
## Hyvin       85.10638 66.97248
```


Tilastollinen testaaminen

- Aiemmin laskimme aineistosta empiirisiä tunnuslukuja, kuten keskiarvo tai varianssi.
- Aineiston keskiarvo antaa viitteitä siitä, minkä suuruinen sitä vastaava populaation parametri, siis odotusarvo μ on. Tähän liittyy kuitenkin epävarmuutta.
 - Pienellä aineistolla epävarmuus odotusarvosta on suurempaa kuin jos aineistoa olisi enemmän.
- Tilastollisella testillä voidaan selvittää esimerkiksi se, onko populaation odotusarvo tietyn arvon suuruinen.
 - Tilastollisia testejä on muidenkin suureiden tarkasteluun, esim. voidaan tutkia ovatko varianssit samansuuruiset kahdessa ryhmässä tai voidaan tutkia onko pisteparveen sovitettun suoran kulmakerroin nolasta poikkeava.

- Testaamista varten on asetettava hypoteeseja, joiden voimassaoloa testataan.
- Tutkija asettaa nk. nollahypoteesin, jota merkitään H_0 :lla, esim. $H_0 : \mu = \mu_0$.
- Usein nollahypoteesi on muotoa, että jokin parametri on yhtä suurta kuin nolla.
- Nollahypoteesi voidaan muotoilla tieteenalan nykytietämyksen perusteella.
- Vastahypoteesi H_1 valitaan siten, että se on vastakkainen nollahypoteesille, ts. jos nollahypoteesin ehto ei ole voimassa, niin vastahypoteesin ehto astuu voimaan.
 - Esim. $H_0 : \mu = 0$ ja $H_1 : \mu \neq 0$.

- Jos hypoteesipari on muotoa $H_0 : \mu = \mu_0$ ja $H_1 : \mu \leq \mu_0$ niin vastahypoteesi on yksisuuntainen.
- Kaksisuuntainen vastahypoteesi esiintyi edellisellä kalvolla eli $H_0 : \mu = 0$ ja $H_1 : \mu \neq 0$
 - Kaksisuuntaisuus korostuu erityisesti, kun muotoillaan H_1 uudelleen: $H_1 : \mu < 0$ tai $\mu > 0$. Tässä on siis kaksi suuntaa!
- Lähtökohtaisesti käytetään kaksisuuntaista hypoteesiparia.
- Yksisuuntaisen vastahypoteesin käyttö vaatii erityisen vahvan perustelun: Onko esim. mahdollista, että testatava suure ei voi muuttua ylös tai alas päin? Esim. voisi ajatella, että ihmisen pituus ei voi vähentyä nuorilla ihmisillä.
- Yksi- ja kaksisuuntaisten vastahypoteesien tapauksessa p -arvon laskeminen tapahtuu hieman eri tavalla.

Oikeanpuoleinen ja vasemmanpuoleinen testaus

- Jos vastahypoteesi on yksisuuntainen, niin tutkija voi asettaa hypoteesiparin joko oikean- tai vasemmanpuoleiseksi.
 - Oikeanpuoleinen testaus: $H_0 : \mu = \mu_0$ ja $H_1 : \mu > \mu_0$
 - Vasemmanpuoleinen testaus: $H_0 : \mu = \mu_0$ ja $H_1 : \mu < \mu_0$
- Oikeanpuoleisessa μ -parametri ei voi alittaa arvoa μ_0 .
- Vasemmanpuoleisessa μ -parametri ei voi ylittää arvoa μ_0 .

- 1 Hypoteesien asettaminen, H_0 ja H_1
- 2 Sopivan testin valinta
- 3 Testisuureen havaitun arvon laskeminen
- 4 p -arvon laskeminen
- 5 Päätelmien tekeminen
- 6 Jatkotarkastelut (tarvittaessa)

Hylkäämisvirhe ja hyväksymisvirhe

- Hylkäämisvirhe (type I error) tapahtuu, kun H_0 hylätään, vaikka H_0 on totta.
- Hyväksymisvirhe (type II error) tapahtuu, kun H_0 hyväksytään, vaikka H_0 on epätotta.

	H_0 hyväksytään	H_0 hylätään
H_0 voimassa	oikea johtopäätös	hylkäämisvirhe eli lajin 1 virhe
H_0 ei voimassa	hyväksymisvirhe eli lajin 2 virhe	oikea johtopäätös

Huom! Aineistosta ja/tai testin tuloksesta ei voida 100 %:n varmuudella tietää, onko H_0 totta vai epätotta. Aina, kun H_0 hyväksytään tai hylätään, otetaan riski siitä, että ollaan tehty väärä päätelmä. Mitä suurempi aineisto on, sitä pienempiä ovat sekä hylkäämisvirheen että hyväksymisvirheen riskit.

Todennäköisyys tehdä hylkäämisvirhe, ehdolla että H_0 on oikeasti voimassa on: $\alpha = P(H_0 \text{ hylätään} | H_0 \text{ on totta})$.

Todennäköisyys tehdä hyväksymisvirhe, ehdolla että H_1 on oikeasti voimassa on: $\beta = P(H_0 \text{ hyväksytään} | H_0 \text{ on epätotta})$.

Merkitsevyystaso α on tutkijan ennen tutkimuksen aloittamista valitsema raja, jota pienemmät p -arvot ovat tilastollisesti merkitseviä (tästä lisää myöhemmin).

Parametrissa β saadaan nk. testin voimakkuus eli $1 - \beta$. Testin voimakkuutta tarvitaan silloin, kun halutaan selvittää miten suuri otoskoko tarvitaan, kun hyväksymisvirheen halutaan olla alle tietyn rajan (tätä ei käsitellä tällä kurssilla tarkemmin).

Esittelemme seuraavat testit:

- Yhden otoksen Z-testi (normaalijakauma)
- Yhden otoksen t-testi (t-jakauma)
- Kahden otoksen t-testi (t-jakauma)
- Levenen testi (ei syvennyttä)
- Parittainen t-testi (yksi otos, parittaiset mittaukset, t-jakauma)

Yhden otoksen Z-testi

- Mikäli populaation varianssi σ^2 on tunnettu, niin voidaan käyttää Z-testiä. Käytännön tilanteissa on erittäin harvinaista, että σ^2 olisi tunnettu, mutta tällä testillä on opettavaista teoriaa.

Valitaan merkitsevyystaso, yleensä $\alpha = 0.05$.

Asetetaan hypoteesit (voidaan asettaa myös yksisuuntaisena)

H_0 : Odotusarvo μ on yhtä suuri kuin μ_0 , ts. $\mu = \mu_0$

H_1 : Odotusarvo μ on eri suuri kuin μ_0 , ts. $\mu \neq \mu_0$

Testisuureksi Z saadaan (vain 5 op laskevat näitä käsin)

$$Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

missä $\sigma = \sqrt{\sigma^2}$ on populaation keskihajonta.

Luvuista laskettua havaittua testisuuretta merkitään: z_{obs} .

Koska Z saadaan otoskeskiarvosta ja tunnetuista suureista (ei-satunnaismuuttujia) σ ja n standardoimalla, niin Z noudattaa standardinormaalijakaumaa

$$Z \sim N(0, 1)$$

Näin ollen p -arvo voidaan laskea normaalijakauman avulla.

Esimerkki 8.1: Phyla-asteroidin etäisyys maasta

Tähtitieteilijä käyttää etäisyyksien mittaukseen laitetta, jonka mittaustulosten odotusarvo μ on sama kuin varsinainen oikea etäisyys mitattavaan taivaankappaleeseen. Mittalaitteen mittausten keskihajonta $\sigma = 0.5$ valovuotta. Nykyteorian mukaan etäisyys maasta asteroidiin *Phyla* on 14.4 valovuotta. Tähtitieteilijä valitsee merkitsevyystasoksi $\alpha = 0.05$. Tähtitieteilijä tekee kuusi toisistaan riippumatonta mittausta tästä etäisyydestä ja saa tulokset

```
x <- c(15.1, 14.8, 14.0, 15.2, 14.7, 14.5)
mean(x)
```

```
## [1] 14.71667
```

Näiden keskiarvo on 15 valovuotta.

Väittämä: *Maan ja Phyla:n välinen etäisyys on 14.4 valovuotta.*

Hypoteesit muotoillaan:

$H_0 : \mu = 14.4$ eli todellinen etäisyys on 14.4 valovuotta.

$H_1 : \mu \neq 14.4$ eli todellinen etäisyys ei ole 14.4 valovuotta.

Esimerkki 8.1 (jatkuu)

- Z-testin testisuure ja p-arvo saadaan laskettua kirjaston BSDA funktiolla `z.test`
 - On tärkeää raportoida p-arvon lisäksi myös ne suureet, joista p-arvo voidaan tarkastaa (hyvä tieteellinen käytäntö). Eli tässä tapauksessa testisuure $z_{obs} = 1.55$

```
library(BSDA)
z.test(x,mu=14.4,sigma.x=0.5)
```

```
##
## One-sample z-Test
##
## data:  x
## z = 1.5513, p-value = 0.1208
## alternative hypothesis: true mean is not equal to 14.4
## 95 percent confidence interval:
##  14.31659 15.11674
## sample estimates:
## mean of x
##  14.71667
```


Aiemmin mainittiin että

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}.$$

Tässä esimerkissä $\bar{x} = 1.5513$, $\mu_0 = 14.4$, $\sigma = 0.5$, ja $n = 6$.

Saadaan

$$z_{obs} = \frac{14.71\bar{6} - 14.4}{0.5/\sqrt{6}} = 1.5513$$

Esimerkki 8.1 (jatkuu)

- p-arvo on todennäköisyys, että nollahypoteesin ollessa voimassa saadaan vähintään yhtä poikkeuksellinen testisuureen arvo.
- p-arvo lasketaan testisuureen jakaumasta (eli tässä normaalijakaumasta).
- Koska vastahypoteesi on kaksisuuntainen, niin p-arvo on todennäköisyys:

$$\begin{aligned} p &= P(Z \leq -|z_{obs}| \text{ tai } |z_{obs}| \leq Z) \\ &= 2(1 - P(|z_{obs}| \leq Z)) = 2(1 - \Phi(z_{obs})) \end{aligned}$$

- R:n avulla voidaan laskea

```
2*(1-pnorm(1.5513))
```

```
## [1] 0.1208298
```

- Koska $p \geq 0.05$ niin nollahypoteesi jää voimaan.
- Tulkitaan että todellinen etäisyys on 14.4 valovuotta.
 - Täsmällisemmin sanottuna, meillä ei ole riittävästi näyttöä sitä vastaan, että todellinen etäisyys on 14.4 valovuotta.
- Jos otos olisi ollut suurempi olisi voitu saada pienempi p -arvo

Jos olisi ollut yksisuuntainen hypoteesi

- Leikitään, että olisi mahdotonta, että etäisyys olisi alle 14.4 valovuotta.
- Hypoteesit olisivat

$H_0 : \mu = 14.4$ eli todellinen etäisyys on 14.4 valovuotta.

$H_1 : \mu > 14.4$ eli todellinen etäisyys on yli 14.4 valovuotta.

Nyt p -arvo laskettaisiin testisuureen oikeasta hännästä

$$\begin{aligned} p &= P(|z_{obs}| \leq Z) \\ &= 1 - P(|z_{obs}| \leq Z) = 1 - \Phi(z_{obs}) \end{aligned}$$

- R:n avulla voidaan laskea

```
1-pnorm(1.5513)
```

```
## [1] 0.0604149
```

Testisuureksi saadaan

$$z_{obs} = \frac{14.71\bar{6} - 14.4}{0.5/\sqrt{20}} = 2.8324$$

ja p -arvoksi

```
2*(1-pnorm(2.8324))
```

```
## [1] 0.004620001
```

- Hahmotellaan toinen tapa perustuen testisuureen **kriittiseen rajaan**.
- Kriittinen raja testisuurelle Z löytyy, kun etsitään z_{obs} s.e.

$$p = P(Z \leq -|z_{obs}| \text{ tai } |z_{obs}| \leq Z) < \alpha$$

eli, jos $\alpha = 0.05$ ja vastahypoteesi on kaksisuuntainen, niin

$$p = P(Z \leq -|z_{obs}| \text{ tai } |z_{obs}| \leq Z) = 2(1 - \Phi(z_{obs})) < 0.05$$

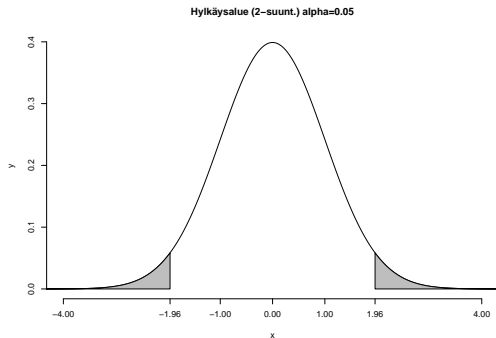
Ratkaistaan $\Phi(z_{obs}) \geq 1 - 0.05/2 = 0.975$. Selvitetään millä arvoilla epäyhtälö toteutuu. Tähän voidaan käyttää joko normaalijakauman taulukkoa "takaperin" tai laskea normaalijakauman 97.5 %:n kvantiiliarvo $z_{0.975}$, joka saadaan R:ssä

```
qnorm(0.975)
```

```
## [1] 1.959964
```

Ts. jos $z_{obs} > 1.96$ tai $z_{obs} < -1.96$ niin silloin $p < 0.05$ ja nollahypoteesi hylätään.

Hylkäysalueen kuvaaja



- Jos havaittu testisuure on hylkäysalueella, niin p -arvo on alle merkitsevyysrajan.
 - Eli jos $z_{\text{obs}} < -1.96$ tai $z_{\text{obs}} > 1.96$ niin $p < \alpha$.

Näin ollen p -arvo saadaan laskemalla tapahtuman

$$p = P(Z < z_{\text{obs}} \cup Z > z_{\text{obs}})$$

todennäköisyys.

- On hyvä idea piirtää vastaava kuvaaja aina, kun testin testisuureen ja p-arvon laskee käsin.

Yhden otoksen t-testi

- Tutkitaan, onko populaation odotusarvo sama, kuin tutkijan asettama arvo μ_0
- Populaation varianssi σ^2 ei ole tunnettu, vaan sen sijaan käytetään aineistosta estimoitua otosvarianssia

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}.$$

- Yhden otoksen t-testiä voidaan käyttää, kun
 - muuttuja on vähintään välimatka-asteikollinen,
 - havainnot ovat riippumattomia toisistaan,
 - muuttuja noudattaa populaatiossa likimain normaalijakaumaa (oletus muodosta).

- Hypoteesit asetetaan samoin kuin Z-testissä eli ovat muotoa

H_0 : Odotusarvo μ on yhtä suuri kuin μ_0 , ts. $\mu = \mu_0$

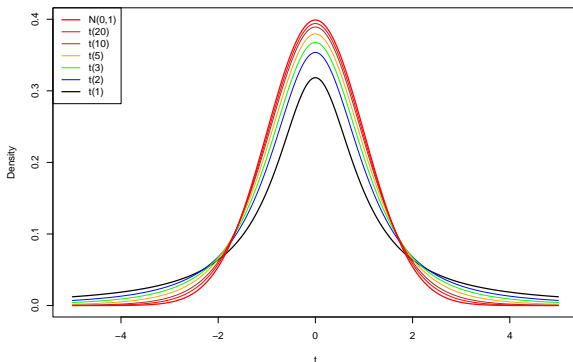
H_1 : Odotusarvo μ on eri suuri kuin μ_0 , ts. $\mu \neq \mu_0$

- Esim. Phyla-asteroidin tapauksessa

H_0 : $\mu = 14.4$ eli todellinen etäisyys on 14.4 valovuotta.

H_1 : $\mu \neq 14.4$ eli todellinen etäisyys ei ole 14.4 valovuotta.

- Jos otosvarianssi ei ole tunnettu, vaan se täytyy laskea aineistosta, niin testisuure ei noudata normaalijakaumaa, vaan sen vaihtelu on suurempaa.
 - Tämä johtuu siitä, että otosvarianssiin liittyy epävarmuutta, joka lisää testisuureen vaihtelua.
- Tällöin testisuure noudattaa t -jakaumaa vapausastein $n - 1$, missä n on havaintojen määrä.



- t-jakauman tiheysfunktio ja kertymäfunktio ovat monimutkaisia, eikä niitä kannata käyttää käsin laskemiseen.
- Sen sijaan laskemiseen käytetään joko R:ssä funktiota $pt(q, df)$ tai t-jakauman kertymäfunktion taulukkoa.
 - q on kvantiili eli testisuureen arvo ja df on vapausasteluku.

```
pnorm(-1.96) # normaalijakauma
```

```
## [1] 0.0249979
```

```
pt(-1.96, df=30)
```

```
## [1] 0.02967116
```

```
pt(-1.96, df=10)
```

```
## [1] 0.03921812
```

```
pt(-1.96, df=1)
```

```
## [1] 0.1501714
```

Riippumattomien normaalijakautuneiden satunnaismuuttujien otosvarianssin otantajakauma on

$$(n-1)\frac{s^2}{\sigma^2} \sim \chi^2(n-1)$$

eli se noudattaa khi toiseen -jakaumaa vapausastein $n-1$.

- Emme perehdy χ^2 -jakaumaan tällä kurssilla.

Miksi päädytään t-jakaumaan? 2/2

Oletetaan että $Z \sim N(0, 1)$ ja $V \sim \chi^2(\nu)$ ja Z ja V ovat riippumattomia.

Erään matemaattisen teorian mukaan standardinormaalijakautuneen sm:jan ja χ^2 -jakautuneen sm:jan neliöjuuren jakolasku noudattaa t-jakaumaa

$$\frac{Z}{\sqrt{V/\nu}} \sim t(\nu)$$

vapausastein ν .

- Tätä t-jakautuneisuuden perustelua eitarvitse opiskella tarkemmin. Nyt olette nähneet idean, mihin lasku perustuu ja siihen voi tarvittaessa perehtyä itsenäisesti tarkemmin.

Testisuure on

$$T = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

ja se noudattaa t -jakaumaa vapausastein $n - 1$ eli

$$T \sim t(n - 1)$$

p -arvo saadaan t -jakauman kertymäfunktion avulla, eli käytännössä joko kertymäfunktion taulukon avulla tai helpommin R:ssä funktiolla `pt(q, df)`.

Esimerkki 8.2

Esimerkki 8.1 olettaen, että emme tunne varianssia σ^2 .

Otoskeskihajonnaksi saadaan $s = 0.436$

```
x <- c(15.1, 14.8, 14.0, 15.2, 14.7, 14.5)
sd(x)
```

```
## [1] 0.4355074
```

ja testisuureen havaituksi arvoksi

$$t_{obs} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{14.71\bar{6} - 14.4}{0.4355\dots/\sqrt{6}} \approx 1.781$$

p -arvoksi saadaan käyttäen t -jakaumaa

$$P\left(\frac{\bar{x} - \mu_0}{s/\sqrt{n}} \leq -1.78 \cup 1.78 \leq \frac{\bar{x} - \mu_0}{s/\sqrt{n}}\right) = 0.135.$$

Esimerkki 8.2 (jatkuu)

Samat laskut R-koodissa

```
Tobs <- (mean(x)-14.4)/(sd(x)/sqrt(6))  
Tobs
```

```
## [1] 1.781076
```

```
2*(1-pt(Tobs,df=6-1))
```

```
## [1] 0.1350077
```

Esimerkki 8.2 (jatkuu)

Tämä voidaan tehdä myös käyttäen R:n valmista funktiota `t.test`

```
x <- c(15.1, 14.8, 14.0, 15.2, 14.7, 14.5)
t.test(x,mu = 14.4)
```

```
##
##  One Sample t-test
##
## data:  x
## t = 1.7811, df = 5, p-value = 0.135
## alternative hypothesis: true mean is not equal to 14.4
## 95 percent confidence interval:
##  14.25963 15.17370
## sample estimates:
## mean of x
##  14.71667
```

Esimerkki 8.2 (jatkuu)

- Koska $p > 0.05 = \alpha$ niin todetaan että nollahypoteesi jää voimaan.
- Jos olisimme havainneet $p < 0.05$, niin täytyisi vielä tehdä jatkotarkastelut.

- Jatkotarkasteluilla tarkoitetaan tarkempaa päätelmää esim. eron suunnasta silloin, kun tilastollisessa testissä on löytynyt eroa, eli nollahypoteesi on hylätty.
- Edellisen esimerkin mukaisessa tilanteessa, jos nollahypoteesi hylätään, niin tiedetään että Phyla-asteroidi ei ole 14.4 valovuoden päässä.
 - Jatkotarkastelu tälle selvittää, onko aineiston perusteella Phyla-asteroidi alle 14.4 valovuoden päässä vai yli 14.4. valovuoden päässä. Verrataan otoskeskiarvoa lukuun 14.4 valovuotta.

```
mean(x)
```

```
## [1] 14.71667
```

- Nyt $\bar{x} = 14.72 > 14.4$, joten Phyla-asteroidin täytyy olla yli 14.4 valovuoden päässä.

Riippumattomien otosten (kahden otoksen) t-testi

Riippumattomien otosten (kahden otoksen) t-testi

- Tutkitaan ovatko kahden populaation odotusarvot samat
 - Onko $\mu_1 = \mu_2$
- Populaatiot ovat erillisiä, eli eivät riipu toisistaan
- Jos varianssit yhtä suuret: voidaan laskea käsin (5op)
- Jos varianssit erisuuret: tarvitaan Welchin korjaus
 - Welchin korjausta emme laske käsin tällä kurssilla, vaan käytämme siihen R:ää

- Riippumattomien otosten t-testiä voidaan käyttää, kun
 - kiinnostuksen kohteena oleva (vaste)muuttuja on vähintään välimatka-asteikollinen,
 - vastemuuttujasta on mittauksia kahdessa erillisessä ryhmässä,
 - aineisto on hankittu satunnaisotannalla populaatioistaan (\Rightarrow riippumattomuus),
 - muuttujan arvojen vaihtelu on yhtä suurta molemmissa populaatioissa (ryhmittäiset varianssit yhtä suuria, $\sigma_1^2 = \sigma_2^2$), ja
 - ryhmäkeskiarvojen otantajakauma noudattaa kummassakin populaatiossa jotakuinkin normaalijakaumaa,
 - pienillä otoksilla, havaintojen on oltava normaalijakautuneita ryhmittäin
 - suurilla otoksilla keskeisen raja-arvolauseen mukaan riippumattomuuden takia otantajakauma on likimain normaalijakautunut

$H_0: \mu_1 = \mu_2$, odotusarvot ovat yhtäsuuret

$H_1: \mu_1 \neq \mu_2$, odotusarvot ovat erisuuret

• Alustavat tarkastelut

```
head(pcb.data)
```

```
##   aika  pcb  
## 1 2018 11.2  
## 2 2018 10.4  
## 3 2018 10.8  
## 4 2018 11.6  
## 5 2018 12.5  
## 6 2018 10.1
```

```
nrow(pcb.data)
```

```
## [1] 20
```

```
# vuosittaiset keskiarvot
```

```
aggregate(pcb.data$pcb, list(pcb.data$aika), FUN=mean)
```

```
##   Group.1      x  
## 1    2018 11.18  
## 2    2020 11.72
```

```
# vuosittaiset keskihajonnat
```

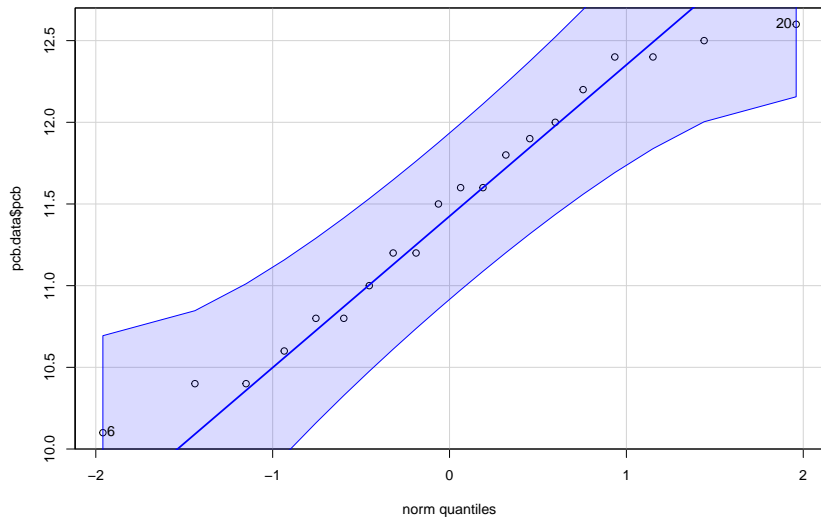
```
aggregate(pcb.data$pcb, list(pcb.data$aika), FUN=sd)
```

```
##   Group.1      x  
## 1    2018 0.7955431  
## 2    2020 0.6860515
```

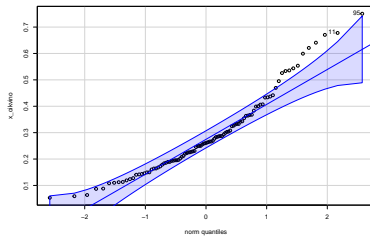
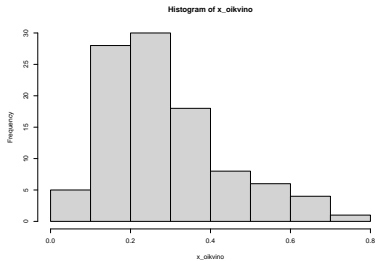
- Vastemuuttujan jakaumaa voidaan tarkatella nk. kvantiili-kvantiili -kuvion avulla (engl. quantile-quantile plot, eli qqplot)
- Kvantiili-kvantiili -kuvaajan idea on, että jos havainnot ovat normaalijakautuneita, niin havaintojen kvantiilien on oltava lähellä normaalijakauman kvantiileja
- Perus R:ssä on myös mahdollista tehdä qqplot komennolla `qqplot(x)`, mutta sen tulkinta on vaikeaa, koska on hankala sanoa, milloin pisteet poikkeavat tarpeeksi suoralta
- `car` -paketin `qqPlot`-funktio antaa myös ns. luottamusvälit, eli rajat joiden sisällä 95 %:a havaintojen kvantiileista tulisi olla

Normaalisuustarkastelut: qqPlot

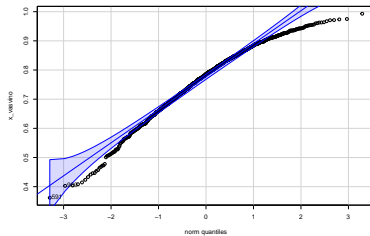
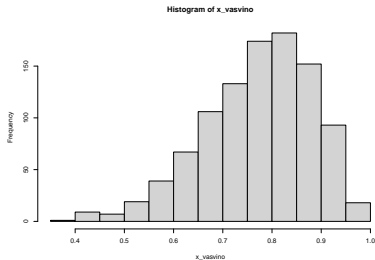
```
library(car)  
qqPlot(pcb.data$pcb)
```



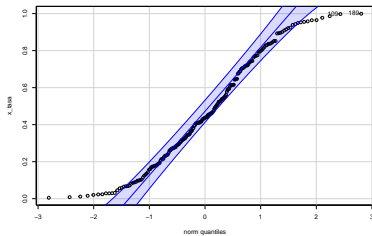
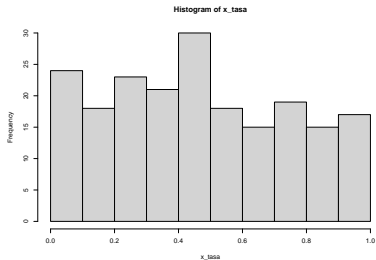
Poikkeama normaalisuudesta (oikealle vino)



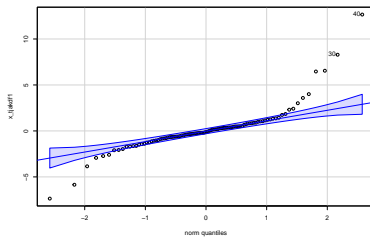
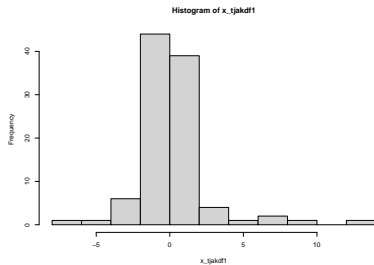
Poikkeama normaalisuudesta (vasemmalle vino)



Poikkeama normaalisuudesta (tasajakauma)



Poikkeama normaalisuudesta: $t(df=2)$



Riippumattomien otosten t-testi R:llä

- jos varianssit voidaan olettaa yhtä suuriksi, niin `var.equal=TRUE`
- tämä testi olettaa että populaatiot ovat erilliset, eli ei eivät muodosta pareja, siksi `paired=FALSE`
- ao. koodissa mittaukset on muuttujassa (sarakkeessa) `pcb` ja aika määrittelee populaatiot (ryhmän)

```
##?t.test  
t.test(pcb ~ aika, var.equal=TRUE, paired=FALSE, data=pcb.data)
```

```
##  
## Two Sample t-test  
##  
## data:  pcb by aika  
## t = -1.6255, df = 18, p-value = 0.1214  
## alternative hypothesis: true difference in means between group 2018 and group 2020 is not equal to 0  
## 95 percent confidence interval:  
## -1.2379222  0.1579222  
## sample estimates:  
## mean in group 2018 mean in group 2020  
##          11.18          11.72
```

Riippumattomien otosten t-testi R:llä

- edellä yksi muuttuja sisälsi tiedot mittauksista ja toinen muuttuja tiedot ryhmästä (kummastako populaatiosta mittaus oli)
- mikäli kahden ryhmän havainnot on eri sarakkeissa tai eri vektoreissa, niin t.test-funktioon annetaan ko. sarakkeet näin

```
pcb2018 <- pcb.data$pcb[pcb.data$aika==2018]
pcb2020 <- pcb.data$pcb[pcb.data$aika==2020]

t.test(pcb2018,pcb2020,var.equal=TRUE, paired=FALSE, data=pcb.data)
```

```
##
## Two Sample t-test
##
## data:  pcb2018 and pcb2020
## t = -1.6255, df = 18, p-value = 0.1214
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -1.2379222  0.1579222
## sample estimates:
## mean of x mean of y
##      11.18      11.72
```

Riippumattomien otosten t-testi R:llä

- Mikäli variansseja ei voi olettaa samoiksi, niin matematiikan puolella tapahtuu ns. Welchin korjaus (R tulostaa Welch Two Sample t-test)
 - R:ssä riittää muuttaa `var.equal=FALSE` vrt. aiempaan
- p-arvoksi tulee eri luku, josta voi tulla eri päätelmät nollahypoteesin hylkäämiseen liittyen

```
t.test(pcb ~ aika, var.equal=FALSE, paired=FALSE, data=pcb.data)
```

```
##  
## Welch Two Sample t-test  
##  
## data:  pcb by aika  
## t = -1.6255, df = 17.619, p-value = 0.1218  
## alternative hypothesis: true difference in means between group 2018  
## 95 percent confidence interval:  
## -1.2390048  0.1590048  
## sample estimates:  
## mean in group 2018 mean in group 2020  
##           11.18           11.72
```

Testisuureen ja p-arvon laskeminen käsin (5op)

- Oletetaan että varianssit on samat, eli ei käytetä Welchin korjausta

$$t_{obs} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

jossa

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

Jos H_0 on totta, niin tällöin

$$t_{obs} \sim t(n_1 + n_2 - 2)$$

eli p-arvo voidaan laskea t-jakauman avulla, kun testisuure t_{obs} on saatu laskettua ja tiedetään otosten koot n_1 ja n_2 .

Esimerkki 8.3 käsin (5op)

- Keskiarvot ovat $\bar{x}_1 = 11.18$ ja $\bar{x}_2 = 11.72$
- Otosvarianssit ovat $s_1^2 = 0.6328889$ ja $s_2^2 = 0.4706667$
- Yhdistetty eli poolattu varianssi (huomaa vapausasteiden käyttö painoina)

$$s^2 = \frac{(10 - 1)0.6328889 + (10 - 1)0.4706667}{10 + 10 - 2} = 0.5517778$$

- testisuure saadaan kaavalla

$$t_{obs} = \frac{11.18 - 11.72}{\sqrt{0.5517778 \left(\frac{1}{10} + \frac{1}{10} \right)}} = -1.6255363$$

ja vapausasteet $df = n_1 + n_2 - 2 = 10 + 10 - 2 = 18$.

Nyt $t_{obs} \sim t(18)$ eli p-arvo saadaan funktion `pt(1.6255363,df=18)` avulla (muista itseisarvo, R:ssä `abs(x)` tai pudota miinus pois)

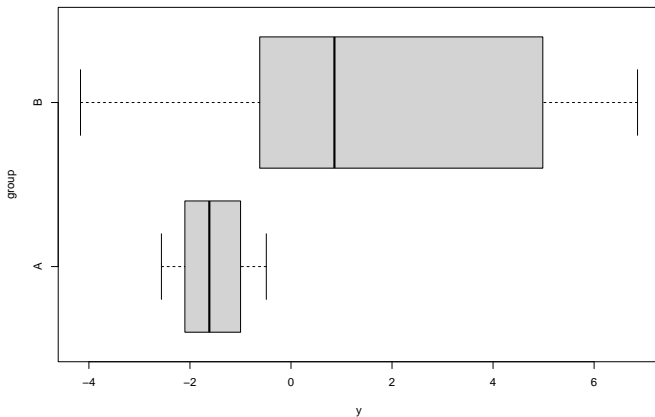
```
2*(1-pt(abs(-1.6255363),df=18))
```

```
## [1] 0.1214283
```

Levenen testi

- Selvittää, ovatko varianssit samat kahdessa tai useammassa eri otoksessa
- Käytetään yleisesti tarkasteltaessa t-testin ja varianssianalyysin oletuksia.
- Perustuu varianssianalyysiin, joka käsitellään Tilastotieteen peruskurssilla.
 - Siksi emme syvenny teoriaan yksityiskohtaisesti
- Samavarianssisuus eli varianssin homoskedastisuus
 - Jos varianssit ovat erisuuret niin aineisto on heteroskedastinen

```
boxplot(y~group,data=dat,range=3,horizontal = T)
```



- Voidaan käyttää kahden tai useamman ryhmän samavarianssisuuden testaamiseen
- Ajatus: keskihajonnoista on hankala sanoa, onko näyttöä tarpeeksi siitä että varianssit eroavat.

- 1 Tarkasteltavat otokset eivät riipu toisistaan.
- 2 Populaatio noudattaa likimain normaalijakaumaa.

Tällä kurssilla

$$H_0 : \sigma_1^2 = \sigma_2^2$$

H_1 : Ryhmien varianssit ovat erisuuria.

Yleisesti, kun ryhmiä 3 tai enemmän

$$H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_t^2$$

H_1 : Kaikki varianssit eivät ole yhtäsuuria.

- 1 Lasketaan havaintojen absoluuttinen etäisyys (itseisarvo) otoskeskiarvoon ryhmittäin.
- 2 Toteutetaan varianssianalyysi näin saadulle aineistolle.

- car-paketin `leveneTest`-funktioilla
- Tarvittaessa asenna car-paketti komennolla
`install.packages("car")`
- ongelmatilanteissa Rcourse-paketin `leveneTest`-funktioilla
- Ota käyttöön R-kieli -kurssin ohjeiden mukaan

```
library(car)
leveneTest(y~group,data=dat)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##           Df F value  Pr(>F)
## group    1  11.081 0.00497 **
##           14
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
library(Rcourse)
leveneTest(y~group,data=dat)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##           Df F value  Pr(>F)
## group    1  11.081 0.00497 **
##           14
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- jos p-arvo < 0.05 niin H_0 hylätään (ja silloin varianssit ovat eri suuret)
- jos p-arvo on ≥ 0.05 niin H_0 saa tukea (ja silloin varianssit voidaan olettaa samoiksi)

Parittainen t-testi

- Käytössä on yksi otos, josta samoista tilastoyksiköistä on kaksi mittausta
 - Soveltuu myös tilanteeseen, jossa esim. näytteet on jaettu kahteen osaan, joille on tehty eri käsittely
 - Oleellista on että otokset riippuvat toisistaan tilastoyksiköiden tasolla, ts. aineiston mittauksista on luontevaa muodostaa pareja
- Tutkitaan onko odotusarvo eri mittausten välillä eri suuruinen
 - Esim. onko samojen kalojen PCB-pitoisuus yhtä suurta nyt kuin kaksi vuotta sitten

- mittaukset muodostavat pareja, ts. samoilta havaintoyksiköiltä tehdyt mittaukset riippuvat toisistaan
- mitatut muuttujat ovat vähintään välimatka-asteikollisia
- erotukset noudattavat jotakuinkin normaalijakaumaa

H_0 : Mittausten erotusten odotusarvo ei poikkea nolasta. ($\mu_D = 0$)

H_1 : Mittausten erotusten odotusarvo poikkeaa nolasta. ($\mu_D \neq 0$)

- Jos kahden mittauskerran välinen ero johtuu sattumasta, niin silloin $\mu_D = 0$.
 - Tämä johtuu siitä että, jos odotusarvo on molemmilla mittauskerroilla sama niin
$$E(X_1 - X_2) = E(X_1) - E(X_2) = \mu - \mu = 0$$

Olkoon meillä aineisto siten, että mittauksista tehdyt havainnot muodostavat pareja $(x_{11}, x_{12}), (x_{21}, x_{22}), \dots, (x_{n1}, x_{n2})$.

- Mittauksen x_{ij} ensimmäinen indeksi i kertoo monesko mittauspari on kyseessä ja toinen indeksi j kertoo onko kyseessä mittauskerta 1 vai 2.

Lasketaan ensin havainnoista parittaiset erotukset $x_{Di} = x_{i1} - x_{i2}$.

Erotukset voidaan laskea myös $x'_{Di} = x_{i2} - x_{i1}$, kunhan tulkintoja tehtäessä muistetaan kumminko päin erotukset laskettiin.

Kotilo	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Ennen	5	3	2	3	2	1	4	3	5	3	2	3	2	1	4	3
Jälkeen	2	0	1	2	5	1	1	5	1	1	0	1	2	3	0	1
Erotus	3	3	1	1	-3	0	3	-2	4	2	2	2	0	-2	4	2

- Tästä saadaan erotusten keskiarvo

$$\bar{x}_D = \frac{1}{n} \sum_{i=1}^n x_{i1} - x_{i2}$$

- Sekä myös erotusten otosvarianssi

$$s_D^2 = \frac{1}{n-1} \sum_{i=1}^n (x_{Di} - \bar{x}_D)^2$$

- Keskiarvo ja varianssi lasketaan siis tavalliseen tapaan, mutta erotuksia käyttäen.

Testisuureen otantajakauma on

$$\frac{\bar{x}_D - \mu_D}{\sqrt{\frac{s_D^2}{n}}} \sim t(n-1)$$

missä n on erotusten lukumäärä, \bar{x}_D on erotusten keskiarvo ja s_D^2 on erotusten varianssi.

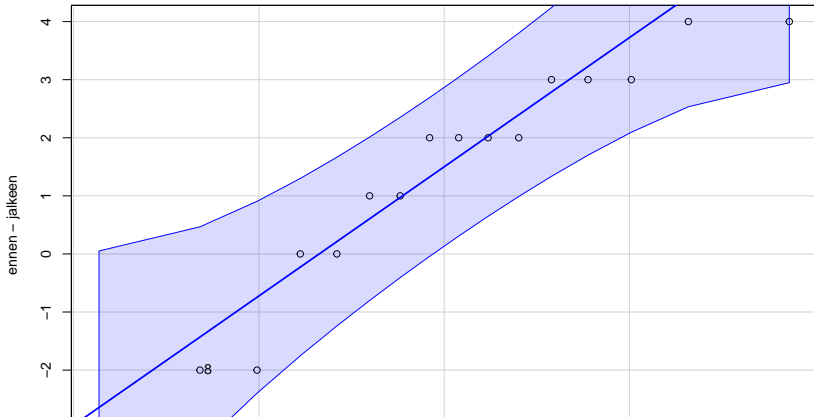
- Huomaa, että havaintoja on n kpl (kaavassa n), eli yhtä monta kuin erotuksia. **Ei siis yhtä monta kuin aineistossa lukuja!**

Parittainen t-testi R:llä (tapa 1)

Jos mittaukset ennen ja jälkeen ovat eri sarakkeissa

```
ennen <- c(5, 3, 2, 3, 2, 1, 4, 3, 5, 3, 2, 3, 2, 1, 4, 3)
jalkeen <- c(2, 0, 1, 2, 5, 1, 1, 5, 1, 1, 0, 1, 2, 3, 0, 1)
```

```
library(car)
qqPlot(ennen-jalkeen)
```



Parittainen t-testi R:llä (tapa 1 jatkuu)

```
t.test(ennen, jalkeen, paired=TRUE)

##
## Paired t-test
##
## data:  ennen and jalkeen
## t = 2.3313, df = 15, p-value = 0.0341
## alternative hypothesis: true mean difference is not equal to 0
## 95 percent confidence interval:
##  0.1071375 2.3928625
## sample estimates:
## mean difference
##           1.25
# Erotusten keskiarvot
mean(ennen-jalkeen)

## [1] 1.25
```

Parittainen t-testi R:llä (tapa 2)

Jos mittaukset ennen ja jälkeen ovat samassa sarakkeessa ja toinen muuttuja kertoo mittauksen ajankohdan.

- Tässä oletetaan, että 1. ennen-mittaus vastaa 1. jälkeen-mittausta! Jos järjestys sotketaan, muuttuvat myös tulokset.
- Aineisto on tässä ns. pitkässä muodossa.

```
t.test(loiset ~ mittauskerta, data=dat, paired=TRUE)
```

```
##  
## Paired t-test  
##  
## data: loiset by mittauskerta  
## t = 2.3313, df = 15, p-value = 0.0341  
## alternative hypothesis: true mean difference is not equal to 0  
## 95 percent confidence interval:  
##  0.1071375 2.3928625  
## sample estimates:  
## mean difference  
##          1.25
```

```
# Ryhmäkeskiarvot kertovat muutoksen suunnan  
aggregate(dat$loiset, list(dat$mittauskerta), FUN=mean)
```

```
## Group.1      x  
## 1   ennen 2.875  
## 2  jälkeen 1.625
```

- $p\text{-arvo} = 0.0341 < 0.05$, joten nollahypoteesi hylätään
- Siispä mittausten erotusten odotusarvo poikkeaa nollasta.
($\mu_D \neq 0$)

Edellä aineisto saatiin pitkään muotoon ao. koodilla:

```
mit <- rep(c("ennen", "jälkeen"),  
          c(length(ennen), length(jalkeen)))  
dat <- data.frame(loiset = c(ennen, jalkeen),  
                  mittauskerta = as.factor(mit))
```

- Koska testisuureen havaittu arvo $2.3313 > 0$, niin erotusten keskiarvolle on voimassa $1.25 > 0$. Tällöin $\mu_D > 0$.
- Erotuksen laskettiin ennen - jälkeen (tämän voi tarkastaa aggregate-tulosteesta), joten on oltava $ennen - jälkeen > 0 \Rightarrow ennen > jälkeen$.
 - Siispä matolääkekylvytys on vähentänyt loisten määrää kotiloilla.

Testisuureen laskeminen käsin (5op)

- erotusten keskiarvo $\bar{x}_D = 1.25$
- varianssi $s_D^2 = 4.6$
- erotusten lkm $n = 16$.

Testisuureen havaituksi arvoksi saadaan ($\mu_D = 0$, kun H_0 on totta),

$$\frac{1.25}{\sqrt{\frac{4.6}{16}}} = 2.331262 \dots \approx 2.33$$

Koska testisuure noudattaa jakaumaa $t(n - 1)$, niin p-arvo saadaan kaksisuuntaisen hypoteesin tapauksessa t-jakauman kertymäfunktion (tai taulukon) avulla. Esim.

```
2*(1-pt(2.331262,df=16-1))
```

```
## [1] 0.03409656
```

- p-arvo = $0.034 < 0.05$, joten H_0 hylätään.
- Kotiloiden lukumäärän odotusarvo on muuttunut.
 - Tutkitaan jatkotarkasteluilla muutoksen suunta. (tämä tehtiin jo pari kalvoa aiemmin)

- Harjoitustehtävissä on sekä tehtäviä, missä tarvitaan parittaista t-testiä että kahden otoksen testiä.
 - Onko kaksi otosta, vai muodostaako mittaukset pareja samoilta tilastoyksiköiltä?
- Miten aineisto on kerätty, onko yksi otos vai kaksi otosta?
 - Pohdi näiden perusteella, tuleeko kyseiseen tilanteeseen käyttää kahden otoksen t-testiä vai parittaista t-testiä.
- Parittaisen t-testin tilanteessa ei käytetä Levenen testiä, eikä Welchin korjausta
- Kahden otoksen t-testissä Levenen testillä tutkitaan, ovatko varianssit yhtäsuuret vai erisuuret.
 - Erisuuruuden tapauksessa käytetään Welchin korjausta.

- Parittainen t-testi on erikoistapaus yhden otoksen t-testistä, jossa havaintoina ovat parin havaintojen erotukset ja testataan onko $\mu_D = 0$ mikä on matemaattisesti yhtäpitävää hypoteesin $\mu = 0$ kanssa.

Alla demonstroitu R:llä miten saadaan sama testisuure ja p-arvo (vrt. seuraavan kalvon tulosteeseen).

```
# parittainen t-testi
t.test(ennen, jalkeen, paired=TRUE)

##
## Paired t-test
##
## data:  ennen and jalkeen
## t = 2.3313, df = 15, p-value = 0.0341
## alternative hypothesis: true mean difference is not equal to 0
## 95 percent confidence interval:
##  0.1071375 2.3928625
## sample estimates:
## mean difference
##           1.25
```

Lisähuomautus (jatkuu)

```
# yhden otoksen t-testin parittaisille erotuksille  
y <- ennen-jälkeen  
t.test(y, mu=0)
```

```
##  
## One Sample t-test  
##  
## data: y  
## t = 2.3313, df = 15, p-value = 0.0341  
## alternative hypothesis: true mean is not equal to 0  
## 95 percent confidence interval:  
## 0.1071375 2.3928625  
## sample estimates:  
## mean of x  
## 1.25
```

- On myös huomauttamisen arvoista, että koska parittaisessa testissä on vain yksi otos, niin Levenen testiä ei tule käyttää tällaisessa tilanteessa.
 - Ei nimittäin ole kahta otosta (kahta ryhmää), joiden välillä varianssia voitaisiin testata. On vain yksi otos, joka on mitattu kaksi kertaa.