

TILM 3701 - Tilastotiede ja Data 2022

Koonneet Henri Nyberg¹ Roope Rihtamo²

2022-08-22

¹Turun Yliopisto, matematiikan ja tilastotieteen laitos, henri.nyberg@utu.fi

²Turun Yliopisto, matematiikan ja tilastotieteen laitos, roope.rihtamo@utu.fi

Contents

Kurssin rakenne	7
1 Johdantoa ja johdattelua tilastotieteeseen	9
1.1 Tilastotiede ja kurssin idea	9
1.2 Tilastotieteen asema tutkimusyhteisön ulkopuolella	11
1.3 Kurssin luonne tilastotieteen (ja datatieteen/data-analytiikan) opintojen esittelijänä	12
2 Tieteellinen tieto, tilastot ja arkitieto yhteiskunnassa	13
2.1 Mitä on tiede?	13
2.2 Tieteellinen menetelmä	13
2.3 Tilastojen yleisestä roolista yhteiskunnassa	13
2.4 Mitä on tutkimus?	13
2.5 Tieteellisen tutkimuksen vaiheet ja tulosten julkaiseminen	13
3 Tilastotiede tieteenalana	15
4 Sattuma ja satunnaisuus	17
4.1 Satunnaisilmiöt ja satunnaismuuttujat tilastotieteessä	17
4.2 Tilastotieteen suhde satunnaisuuteen ja todennäköisyyksiin . . .	17
4.3 Tilastolliset mallit, jakaumat ja parametrit	17
4.4 Odotusarvo ja varianssi	17
4.5 Joitain jakaumia	17
4.6 Sattuman rooli tieteenteossa: Vale-emävale-tilasto?	17

5	Tilastolliset aineistot, niiden kerääminen ja mittaaminen	19
5.1	Kertausta: Data eli aineisto	20
5.2	Otannan idea	23
5.3	Mittaaminen, mitta-asteikot ja tilastolliset muuttujat	23
5.4	Kontrolloidut kokeet ja suorat havainnot	23
5.5	Otantamenetelmät	23
5.6	Otantaesimerkkejä	23
5.7	Otannan haasteita vielä kootusti	23
6	Otokset ja otosjakaumat: tilastollisen päättelyn näkökulma	25
6.1	Satunnaisotos, yhteisjakauma ja tilastollinen malli	25
6.2	Otosjakauma: Estimaattori ja estimaatti	25
6.3	Otoskeskiarvo ja otosvarianssi (estimaattoreinta)	25
6.4	Suhteellisen frekvenssin otosjakauma	25
6.5	Muita tunnuslukuja	25
6.6	Luottamusvälit	25
6.7	Otoskoko	25
7	Tilastollinen riippuvuus ja korrelaatio	27
7.1	Muuttujien väliset riippuvuudet tilastollisen tutkimuksen kohteena	27
7.2	Kahden muuttujan havaintoaineiston kuvaaminen	27
7.3	Tunnusluvut	27
7.4	Satunnaismuuttujien kovarianssi ja korrelaatio	27
8	Regressioanalyysi	29
8.1	Johdatus regressioanalyysin ideaan	29
8.2	Yhden selittäjän lineaarinen regressiomalli	29
8.3	Muita regressiomalleja	29
9	Tilastotieteen rooli uuden tiedon tuottamisessa	31
9.1	Tilastollisen tutkimuksen yhteisiä elementtejä	31
9.2	Tutkimusprosessi	31

<i>CONTENTS</i>	5
10 Aineisto- ja tutkimustyytit ja koeasetelmat	33
10.1 Tutkimustyytit	33
10.2 Tutkimusstrategiat	33
10.3 Erilaisia aineistoja ja aineistolähteitä	33
11 Tilastollisesta ennustamisesta	35
11.1 Tilastollinen selittäminen vs. ennustaminen	35
11.2 Tilastolliseen ennustamiseen liittyviä huomioita	35
12 Tilastotieteen kehityksen nykytrendejä	37

Kurssin rakenne

- Tällä kurssilla tarkoituksena on melko yleisellä tasolla johdatella tilastotieteen ja aineistojen (datan) maailmaan pohtimalla myös näiden laajempia merkityksiä tieteellisen tutkimuksen hyvin keskeisinä osina.
- Kurssilla vältetään, mahdollisuuksien mukaan, kovin teknistä matemaattista esitystapaa, mutta tarvittavissa määrin tullaan myös käyttämään tilastotieteen perusopinnoissa tarvittavia matemaattisia merkintöjä ja määritelmiä. Esim. todennäköisyyslaskennan ja tilastollisen päättelyn perusteita ei käydä vielä riittävällä matemaattisella tarkkuudella lävitse, vaan nämä tarkastelut jäävät tätä kurssia seuraavien kurssien (TILM3553 Todennäköisyyslaskennan peruskurssi tai TILM3568 Todennäköisyyslaskenta sivuaineopiskelijoille sekä TILM3555 Tilastollisen päättelyn peruskurssi) asiaksi. Nämä kurssit, yhdessä alkuvaiheen pakollisten matematiikan kurssin lisäksi, muodostavat siis tämän kurssin johdannon kanssa lähtökohdan tilastotieteen opinnoille.
- Luennot eivät suoraan perustu yhteen kirjaan tai lähteeseen. Käytettyjä lähdemateriaaleja luetellaan alapuolella oheislukemiston myötä.
- Oheislukemistoa (sopivilta osin):
 - Mellin, I. (2004). Johdatus tilastotieteeseen: Tilastotieteen johdantokurssi (1.kirja). Yliopistopaino, Helsingin yliopisto.
 - Mellin, I. (2000). Johdatus tilastotieteeseen: Tilastotieteen jatkokurssi (2.kirja). Yliopistopaino, Helsingin yliopisto.
 - Mellin, I. (2006). Tilastolliset menetelmät. Luentomoniste, Aalto yliopisto (TKK).
 - Holopainen, M. ja P. Pulkkinen (2008). Tilastolliset menetelmät. Sanoma Pro Oy.
 - Pahkinen, E. ja R. Lehtonen (1989). Otanta-asetelmat ja tilastollinen analyysi. Gaudeamus, Helsinki.
 - Pahkinen, E. ja R. Lehtonen (2004). Practical Methods for Design and Analysis of Complex Surveys. 2. painos, Wiley.
 - Sund, R. (2003). Tilastotiede käytännön tutkimuksessa -kurssi. Helsingin yliopisto.

- Silver, N. (2014). Signaali ja kohina: Miksi monet ennusteet epäonnistuvat mutta jotkin eivät? Terra Cognita. (Suomentanut Kimmo Pietiläinen)
 - * Englanninkielinen teos: Silver, N. (2015). The Signal and the Noise: Why So Many Predictions Fail—but Some Don't. Penguin Books; Illustrated edition
 - Pesonen, M. (2017). Kurssimateriaali kurssille Aineistonhankinta ja tutkimusasetelmat, Turun yliopisto.
 - Vartia, Y. (1989). Tilastotieteen perusteet. Yliopistopaino, Helsinki. II painos.
- Muita taustamateriaaleja
 - Tilastokeskuksen tilastokoulu ([linkki](#))
 - Tilastotieteen sanasto suomi-englanti-suomi, ks. Juha Alho, Elja Arjas, Esa Läärä ja Pekka Pere (2021). Tilastotieteen sanasto. Suomen Tilastoseuran julkaisuja 8.

Suuret kiitokset Visa Kuntzelle ja Emil Lehdelle kommenteista ja avusta materiaalin työstämisessä. Kaikki jäljelle jääneet painovirheet ovat materiaalin kokoajien.

Chapter 1

Johdantoa ja johdattelua tilastotieteeseen

Ihmisellä on luontainen pyrkimys ymmärtää, mitä hänen ympärillään tapahtuu. Ymmärrys perustuu ihmisen tekemiin havaintoihin, joita luokittelemalla tai seuraamalla hän pyrkii löytämään säännönmukaisuuksia. Näiden säännönmukaisuuksien löytäminen vaatii loogisten johtopäätösten tekoa. Pelkän uteliaisuuden tyydyttämiseen ja älyllisen mielihyvän lisäksi ihminen pyrkii ennakoidaan tulevaa ja siten varautumaan tuleviin tapahtumiin... Edellä kuvattuja taitoja voi oppia.

Holopainen ja Pulkkinen, 2008

1.1 Tilastotiede ja kurssin idea

- Tämän tilastotieteen ensimmäisen kurssin ideana on (ainakin)
 - Esitellä ja johdatella **tilastolliseen ja tieteelliseen ajatteluun** ja sen hyödyntämiseen eri tyypisissä tutkimusongelmissa.
 - Esitellä tilastotieteen roolia **empiirisen tutkimusaineiston keräämisessä ja analyysissä** sekä tarkastella tieteentekemisen ja tilastotieteen suhdetta.
 - Pohtia **tilastotieteen olemusta tieteenalana** ja tarkastella tilastotieteen ja datatieteiden (data sciencen) samankaltaisuuksia ja eroja.
 - Pohtia **sattuman ja satunnaisuuden roolia** jokapäiväisessä elämässä ja erityisesti osana tieteellistä tutkimusprosessia.
 - Oppia tilastotieteen peruskäsitteitä ja (tilastollisen) tutkimuksenteon alkeita ja siihen liittyviä mahdollisia ongelmia esimerkiksi tilastollisten aineistojen keräämisessä.

- Oppia tilastollisten aineistojen **kuvaamisen ja käsittelyn** alkeita sekä tilasto(tieteellisen)llisen **mallintamisen ja koeasetelmien** peruskäsitteitä.
- Kurssilla käsitellään myös **tilastollisen päättelyn** peruskäsitteitä ja perusteita kuten
 - Mitä on **todennäköisyys** ja miten sen tulkitaan tilastotieteessä sekä laajemmin tieteessä. Erityisesti tilastotieteen osalta keskiössä on tämän kurssin osalta **satunnaismuuttujat** sekä niihin liitettävät käsitteet
 - * **Odotusarvo, varianssi** ja kahden (tai useamman) satunnaismuuttujan **korrelaatio**.
 - * Satunnaismuuttujien **todennäköisyysjakaumien** perusteita ja niiden yhteyksiä mm. normaalijakaumaan ja muutamiin muihin keskeisiin jakaumiin.
 - * Tilastollinen malli työkaluna satunnaismuuttujien formaalissa mallintamisessa ja päättelyssä. Tilastollisen malliin liittyy (usein) **parametreja** joihin tilastollinen päättely kohdistuu.
 - * Tilastollisten mallien **estimoinnin** perusidea, eli miten tilastollisen mallin parametreille muodostetaan arvot käytettävissä olevan aineiston pohjalta. Esimerkiksi: mitä tarkoittaa tilastollisen mallin parametrin **estimaattori** ja sen **harhattomuus**?
 - * Alustavia tarkasteluja tilastollisen mallin uskottavuuden käsitteelle ja **luottamusväleille** tilastollisen mallin estimoiduille parametreille.
- Toinen kurssin keskeisistä teemoista on tarkastella tieteellistä tutkimusprosessia teoriassa ja käytännössä. Tämä sisältää mm. seuraavia aiheita (joita siis käsitellään tällä kurssilla päällisin puolin ja varsin yleisestä näkökulmasta katsoen): tarkemmat yksityiskohdat jäävät tätä kurssia seuraavien tilastotieteen kurssien aihepiireiksi):
 - **Tutkimusongelman** asettaminen: mitä halutaan tutkia?
 - Tutkimusongelman täsmentäminen ja **tutkimusstrategian** laatiminen: millä keinoin asetettuun tutkimusongelmaan voidaan vastata?
 - **Tutkimusaineiston** (tai vain lyhyemmin **aineiston** eli **datan**) kerääminen
 - * **Aineiston ennakkoehdot**: mitkä ehdot tulee täyttyä, jotta asetettuun tutkimusongelmaan voidaan vastata?

- * **Otanta** (ja mittaaminen): miten tutkimusaineisto kerätään niin, että se täyttää aineiston ennakkoehdot? Erilaisissa tutkimuksissa käytetään erilaisia aineistoja kuten:
 - Survey- ja rekisteriaineistot
 - Havaintoarvojen välistä korrelaatiota esiintyy mm. aikasarja-aineistojen tai pitkittäisaineistojen tapauksessa
- **Aineiston kuvaaminen:** minkälaista aineistoa on kerätty ja vastaako se ennakkoehtoja?
- **Aineiston analyysin lähtökohtia**
 - Mitä tilastollista mallia/malleja käytetään?
 - Mitä tarkoitetaan mallien tuntemattomien parametrien arvojen estimoinnilla?
 - Tilastollinen päättely (estimointitulosten pohjalta)
- **Johtopäätelmien** tekeminen tilastollisen päättelyn pohjalta: saatiinko tutkimusongelmaan vastaus ja kuinka luotettava saatu vastaus on?

1.2 Tilastotieteen asema tutkimusyhteisön ulkopuolella

- Tilastotiede on oppiaineena usein varsin tuntematon toisen asteen opinnoista valmistuneelle, sillä sitä ei juurikaan opeteta lukioissa tai ammattikouluissa huolimatta sen keskeisestä ja kasvavasta roolista tiedemaailman kentillä.
- Tiedeyhteisön ulkopuolellakin **tilastotiedettä ja tilastotieteilijöitä arvostetaan laajalti.**
- **Tilastotiede onkin nostanut profiliaan viimeisten vuosikymmenien aikana** tietoteknisen kehityksen tuotua laajat tietoaineistot ja kehittyneet laskennalliset menetelmät lähes jokaisen kansalaisen saataville.
- Tämä “datavallankumous” näkyy tilastotieteilijöiden kysynnässä työmarkkinoilla: erilaisten aineistojen määrän lisääntyessä kasvaa myös kysyntä työntekijöistä, jotka osaavat ammatitaitoisesti käsitellä, tulkita ja mallintaa tilastollisia aineistoja.
- Ei siis liene ihmeäkään, että erilaisten “data”-alkuisten työpaikkojen, kuten **datatieteilijä** (eng. **data scientist**) tai **data-analyttikko** (**data-analyst**) määrä on kasvanut voimakkaasti jo pidempään. Kaikkia tieto- ja datainensiivisten ammattien tekijöitä yhdistää yksi tekijä: **heidän tulee hallita ja osata tilastotiedettä!** Karkeistettuna mitä paremmin ja enemmän (laajemmin), sen parempi palkka ja monipuolisemmat työtehtävät!

1.3 Kurssin luonne tilastotieteen (ja datatieteen/data-analytiikan) opintojen esittelijänä

Kurssin mittaan esitellään tilastotieteen perusteiden lisäksi **miten TY:ssa tilastotieteen opinnoissa syvennyttään** tällä kurssilla esiteltäviin menetelmiin, aineistotyyppeihin ja mallinnuskokonaisuuksiin.

Chapter 2

Tieteellinen tieto, tilastot ja arkitieto yhteiskunnassa

2.1 Mitä on tiede?

2.2 Tieteellinen menetelmä

2.3 Tilastojen yleisestä roolista yhteiskunnassa

2.4 Mitä on tutkimus?

2.5 Tieteellisen tutkimuksen vaiheet ja tulosten julkaiseminen

Chapter 3

Tilastotiede tieteenalana

Tässä luvussa hahmottelemme tilastotieteen piirteitä tieteenalana. Käymme läpi tilastotieteelle ominaisia piirteitä, jotka erottavat sen niin lähitieteistä, kuten matematiikasta ja tietojenkäsittelytieteestä, kuin myös sovellusaloista. Usein näkee tilastotieteen typistettävän vain työkaluksi eri sovellusalojen empiriseen tutkimukseen siitäkään huolimatta että tilastotieteellä on oma rikas teoriapohjansa sekä kiistaton asema omana tieteenalanaan. Tieteenalan määrittelemisen lyhyesti on aina hieman hankalaa. Tästä huolimatta seuraavassa yritämme osaltaan vastata seuraaviin kysymyksiin:

- Mitä tilastotiede on ja mitä se ei ole? Miksi tilastotiede ei ole vain sovellettua matematiikkaa tai matematiikalla höystettyä tietojenkäsittelyä?
- Mihin tilastotiedettä käytetään? Onko tilastotieteellä käyttöä ns. “akatemian” eli tutki- musyhteisön ulkopuolella?
- Tilastotieteelle tyypillistä kritiikkiä?

Chapter 4

Sattuma ja satunnaisuus

- 4.1 Satunnaisilmiöt ja satunnaismuuttujat tilastotieteessä
- 4.2 Tilastotieteen suhde satunnaisuuteen ja todennäköisyyksiin
- 4.3 Tilastolliset mallit, jakaumat ja parametrit
- 4.4 Odotusarvo ja varianssi
- 4.5 Joitain jakaumia
 - 4.5.1 Normaalijakauma
 - 4.5.2 Bernoulli-, binomi- ja Poisson-jakauma
- 4.6 Sattuman rooli tieteenteossa: Vale-emävaletilasto?

Chapter 5

Tilastolliset aineistot, niiden kerääminen ja mittaaminen

Edellisessä luvussa käsiteltiin tilastotieteen suhtautumista satunnaisilmiöihin. Tässä luvussa tarkastelemme lähemmin miten reaali maailman satunnaisilmiöistä kerätään tietoa ja miten niitä voidaan mitata. Tilastotieteen perusoppimäärä rakentuu ajatukselle ilmiöiden tutkimisesta rajallisen ja epävarman tiedon vallitessa. Käytännössä tämä tarkoittaa sitä, että tutkimuksen kohteena olevat rajalliset aineistot sisältävät niin systemaattista kuin satunnaisuudesta johtuvaa vaihtelua. Tilastollisten menetelmien avulla pyrimme erottamaan systemaattisen vaihtelun satunnaisesta sekä tekemään tilastollista päättelyä aineiston generoimasta mekanismista. Lyhyesti tämä tarkoittaa aineiston systemaattisen vaihtelun tilastollista mallintamista ja sen parametrien estimointia otoksesta, joka kattaa vain (pienen) osajoukon koko populaation (perusjoukon) tilastoyksiköistä.

Voidaksemme tehdä uskottavaa päättelyä “havainnoista parametreihin”, tulee otoksen olla riittävän **edustava**. Tämän luvun keskeisin oppi onkin, että miten otanta tulisi suorittaa, jotta havaintoaineisto olisi **edustava otos** populaatiosta, silloin kun aineisto kerätään otannalla. Vaikka aineiston hankinta vaatii yleensä runsaasti käytännön työtä, kannattaa se tehdä huolellisesti, sillä huonosti toteutetun otannan vuoksi tutkimusongelman kannalta keskeisiä johtopäätöksiä ei voida tehdä!

5.1 Kertausta: Data eli aineisto

- **Tilastollinen tutkimus** aloitetaan tutkimusaineiston keruun suunnittelulla.
- Kertauksen vuoksi: tilastollinen tutkimusaineisto (havaintoaineisto) koostuu tilastoyksiköiden populaatiosta havaituista tilastoyksiköiden muuttujien arvoista.
- Havaintoaineisto voidaan koota taulukoksi, johon listataan tilastoyksiköt riveille ja tilastomuuttujat sarakkeisiin. Jos havaintoaineisto koostuu n tilastoyksiköstä, joista jokaisesta on kerätty esim. m tilastomuuttujasta havainnot, niin havainnot voidaan kirjoittaa taulukon muotoon

Tässä siis rivillä i on i . **tilastoyksikön** havainto ja j sarakkeessa on j . tilastollisesta muuttujasta havaitut arvot $x_{i,j}$. Ts. yhdellä rivillä on yhden tilastoyksikön tiedot kaikista tilastomuuttujista ja yksi sarake on kaikkien tilastoyksiköiden tiedot yhdestä tilastomuuttujasta.

- Usein (varsinkin parhaillaan kiihtyvällä vauhdilla) kerättävät havaintoaineistot ovat niin suuria, ettei edellisenkaltaisesta havaintotaulukosta voida usein suoraan tarkastelemalla nähdä aineiston pääpiirteitä.
 - Tällöin on tarpeen luokitella aineistoa taulukon muodostamiseksi.
 - Luokittelussa on kysymys aineiston tiivistämisestä kohtuullisen kokoiseksi ja havainnollisempaan muotoon. Luokittelussa tilastomuuttujan arvot sijoitetaan eri luokkiin siten, että yhden tilastomuuttujan arvo voi kuulua vain yhteen luokkaan. Luokka ilmoitetaan yleensä luokkavälinä, kuten reaalityyppinä. Esimerkiksi henkilön ikä on tapana luokitella ikäjakauman kuvaamisessa 10-vuotislukuihin (15-24, 25-34, ...), vaikka periaatteessa ikä voitaisiin ilmoittaa minuutinkin tarkkuudella.
 - Luokkien lukumäärään vaikuttavat muun muassa tilastomuuttujan arvojen vaihteluväli ja havaintoaineiston laajuus. Luokittelussa pyritään siihen, että luokkien lukumäärä saadaan tarvittaessa luokkia yhdistämällä kohtuulliseksi ja että luokat valitaan tasavälisesti eli siten, että kahden peräkkäisen luokan alarajojen erotus on vakio. Kun aineistoa luokitellaan, aineiston luettavuus paranee mutta toisaalta osa tiedoista menetetään eivätkä yksittäiset havaintoarvot ole enää tiedossa.
 - Emme vielä tällä kurssilla etene tämän pidemmälle tilastografiikan esittämisessä ja siihen liittyvissä pohdinnoissa. Muun muassa tilastollisen päättelyn peruskurssi (TILM3555) vastaa näihin kysymyksiin tarkemmin. Graafiset menetelmät ovat joka tapauksessa erittäin

tärkeä osa aineiston havainnollistamista. Kuvat helpottavat aineiston tulkitsemista ja toimivat usein perusteltuna lähtökohtana monimutkaisempien tilastollisten mallien (ja algoritmien) sovittamiselle.

- Kvantitatiivisen tutkimuksen aineistoksi kelpaa periaatteessa kaikki havaintoihin perustuva informaatio, joka on **mittauksen** avulla muutettavissa numeeriseen muotoon.
 - Havaintoyksiköiden tilastollisten muuttujien numeerisia arvoja kutsutaan **havaintoarvoiksi** tai **havainnoiksi**.
 - Kaikki havaitut tilastolliset muuttujat eivät ole aina mielenkiintoisia. Tutkimuksen kannalta mielenkiintoisia muuttujia kutsutaan **tutkimusmuuttujiksi**, joiden lisäksi havaintoaineisto pitää mahdollisesti sisällään **taustamuuttujia**.
 - * Esimerkiksi, jos tutkimuksella halutaan tietoa suomalaisen aikuisväestön mielipiteistä, havaintoyksikköinä ovat aikuisväestöön kuuluvat henkilöt. Jos halutaan tietoa suomalaisista kunnista, havaintoyksikköinä ovat Suomen kunnat jne.
 - * Ensimmäisessä tapauksessa tilastollisina muuttujina on aikuisväestön mielipiteet, joita voidaan selvittää esimerkiksi kyselytutkimuksella. Toisaalta voidaan myös kerätä taustamuuttujiksi haastatelluista muita tietoja, kuten asuinpaikka, ikä ja ammatti.
 - Kaikkia mielenkiintoisia muuttujia ei kuitenkaan välttämättä voida havaita, eli niille ei voida määrittää numeerista arvoa.
 - Tällöin puhutaan nk. **latenteista muuttujista**, eli muuttujista joita ei suoraan havaita mutta joiden oletetaan vaikuttavan havaittavien muuttujien taustalla. Latenteja muuttujia voidaan rakentaa tilastollisten mallien avulla käyttäen hyödyksi niihin liittyviä havaittuja muuttujia.
 - Latenteja muuttujia ovat esimerkiksi elämänlaatu, onnellisuus, konservatiivisuus, yms.
- Tilastollinen tutkimus voi olla joko **kokonaistutkimus** tai **otanta-tutkimus**.
 - **Kokonaistutkimuksessa** tutkitaan kaikkia ajateltavissa olevia kohteita (kaikki perusjoukon alkiot tutkitaan).
 - * Esimerkiksi jos tutkitaan Suomen kuntia, niin kokonaistutkimuksessa tutkitaan kaikki kunnat.
 - * Tai jos tutkitaan jonkin lääkeaineen vaikutuksia ihmisiin, niin tutkitaan jokainen ihminen erikseen. Selvää on, että tällainen kokonaistutkimus olisi liian vaikeaa toteuttaa.
 - **Otantatutkimuksessa** tutkimus kohdistetaan johonkin (populaation/perusjoukon) osajoukkoon ja johtopäätelmiä populaatiosta/perusjoukosta tehdään otokseen perustuen.

- * Perusjoukosta otokseen poimittuja alkioita kutsutaan **otosyksiköiksi** ja niiden muodostama osajoukko, eli **otos**, on se osa perusjoukkoa, joka tutkitaan tutkimusaineiston keräämisen jälkeen.
 - * Lääketutkimusta tehdäänkin poikkeuksetta otantatutkimuksena (ja kontrolloituina kokeina, ks. alemmaa), jolloin lääkettä testataan vain osajoukolla koko ihmiskokoon otantaan ja tämän osajoukon alkioita ovat otosyksiköitä.
 - * Näin toimimalla, ja riittävän edustavalla otoksella, saadaan kuitenkin tarpeeksi tietoa lääkeaineen vaikutuksista ja tulokset voidaan yleistää populaatiotasolle ja lääke ottaa käyttöön.
 - * Otantatutkimus on halvempi kuin kokonaistutkimus ja tulokset saadaan nopeammin!
- Usein on kuitenkin niin, että koko populaation tutkiminen ei ole mahdollista tai kannattavaa. Tällöin tehtävä tutkimus on otanta-tutkimus ja tutkittavaksi valitaan perusjoukon osajoukko sopivaa **otantamenetelmää** (ks. alaluku 5.5) käyttäen.
 - Esimerkkinä aseiden patruunoita valmistava tehtailija, joka haluaisi tutkia toimivatko kaikki ammuksien tai kaikkien suomalaisten haastatteluun suomalaisten mielipiteitä kartoitettaessa. Myöskään valaisimien valmistaja tuskin tekee kokonaistutkimuksia valmistamiensa tuotteiden kestoajan selvittämiseksi.
 - Tämän vuoksi useimmiten keskitytään perusjoukkoa edustavan pienemmän, mieluiten satunnaisesti valitun osajoukon eli **otoksen** tutkimiseen.
 - Otantatutkimuksissa tiedot kerätään useimmiten haastattelulla, kirjallisella/sähköisellä kyselyllä tai suoraan tietorekistereistä. Tiedonkeruun toteuttaminen (eri sovelluksissa) määrää osaltaan käytettävän otantamenetelmän.
 - Teoriassa äärelliseen perusjoukkoon kohdistuvat kokonaistutkimukset voidaan aina tulkita otantatutkimuksiksi (perusjoukko tulkitaan otokseksi hypoteettisesta äärettömästä perusjoukosta)!
 - * Esimerkiksi Galilein tekemät painovoiman vaikutusta kappaleiden putoamisaikaan liittyneet mittaukset. Koetuloksia (mittauksia) voidaan pitää otoksena äärettömästä mahdollisten koe-tulosten joukosta. Tällöin ainoa mahdollisuus ilmiön tutkimiseen on käyttää otantaa.
 - Otantatutkimuksen tulokset voivat olla luotettavampia kuin kokonaistutkimuksen.
 - Otantatutkimuksessa voidaan panostaa enemmän huolelliseen ja tarkkaan mittaamiseen sekä valitun otoksen tavoittamiseen.

- Kokonaistutkimuksessa vastauskato ja tarkasteltavan populaation valintavirhe ovat mahdollisia siinä kuin otantatutkimuksessakin.
- Otantateoria on yksi tilastotieteen keskeisimpiä oppeja ja tarjoaa teoreettisen kehikon empiiristen tutkimusten tulosten yleistämiseen. Tarkastellaan siis tarkemmin otannan ideaa ja toteuttamista seuraavassa alaluvussa.

5.2 Otannan idea

5.3 Mittaaminen, mitta-asteikot ja tilastolliset muuttujat

5.4 Kontrolloidut kokeet ja suorat havainnot

5.5 Otantamenetelmät

5.5.1 Yksinkertainen satunnaisotanta

5.5.2 Systemaattinen otanta

5.5.3 Ositettu otanta

5.5.4 Ryväotanta

5.6 Otantaesimerkkejä

5.7 Otannan haasteita vielä kootusti

Chapter 6

Otokset ja otosjakaumat: tilastollisen päättelyn näkökulma

- 6.1 Satunnaisotos, yhteisjakauma ja tilastollinen malli
- 6.2 Otosjakauma: Estimaattori ja estimaatti
- 6.3 Otokeskiarvo ja otosvarianssi (estimaattoreinta)
- 6.4 Suhteellisen frekvenssin otosjakauma
- 6.5 Muita tunnuslukuja
- 6.6 Luottamusvälit
- 6.7 Otokoko

Chapter 7

Tilastollinen riippuvuus ja korrelaatio

- 7.1 Muuttujien väliset riippuvuudet tilastollisen tutkimuksen kohteena
- 7.2 Kahden muuttujan havaintoaineiston kuvaaminen
- 7.3 Tunnusluvut
- 7.4 Satunnaismuuttujien kovarianssi ja korrelaatio

Chapter 8

Regressioanalyysi

8.1 Johdatus regressioanalyysin ideaan

8.2 Yhden selittäjän lineaarinen regressiomalli

8.3 Muita regressiomalleja

Chapter 9

Tilastotieteen rooli uuden tiedon tuottamisessa

9.1 Tilastollisen tutkimuksen yhteisiä elementtejä

9.2 Tutkimusprosessi

Chapter 10

Aineisto- ja tutkimustyyppit ja koeasetelmat

10.1 Tutkimustyyppit

10.2 Tutkimusstrategiat

10.3 Erilaisia aineistoja ja aineistolähteitä

10.3.1 Rekisteriaineistot

10.3.2 Aikasarjat ja paneeliaineistot

10.3.3 Survey eli haastattelu- tai kyselytutkimus

Chapter 11

Tilastollisesta ennustamisesta

11.1 Tilastollinen selittäminen vs. ennustaminen

11.2 Tilastolliseen ennustamiseen liittyviä huomioita

Chapter 12

Tilastotieteen kehityksen nykytrendejä