

Tässä näkee helposti puut metsältä mutta näetkö metsän puita?  
Siinä tieteenkin iso kysymys vai onko?



Tiede, Tilastot ja Media, Säätytalo 8.2.2016 Seppo Laaksonen

# Media ja tieteetkin vaikeuksissa tilastotieteen kanssa

Seppo Laaksonen

Helsingin yliopisto

**ISI:n** jäsen vuodesta 1995

Suomalaisen tiedeakatemian seminaari 8.2.2016

Säätytalo

**Statistical Science for a Better World**



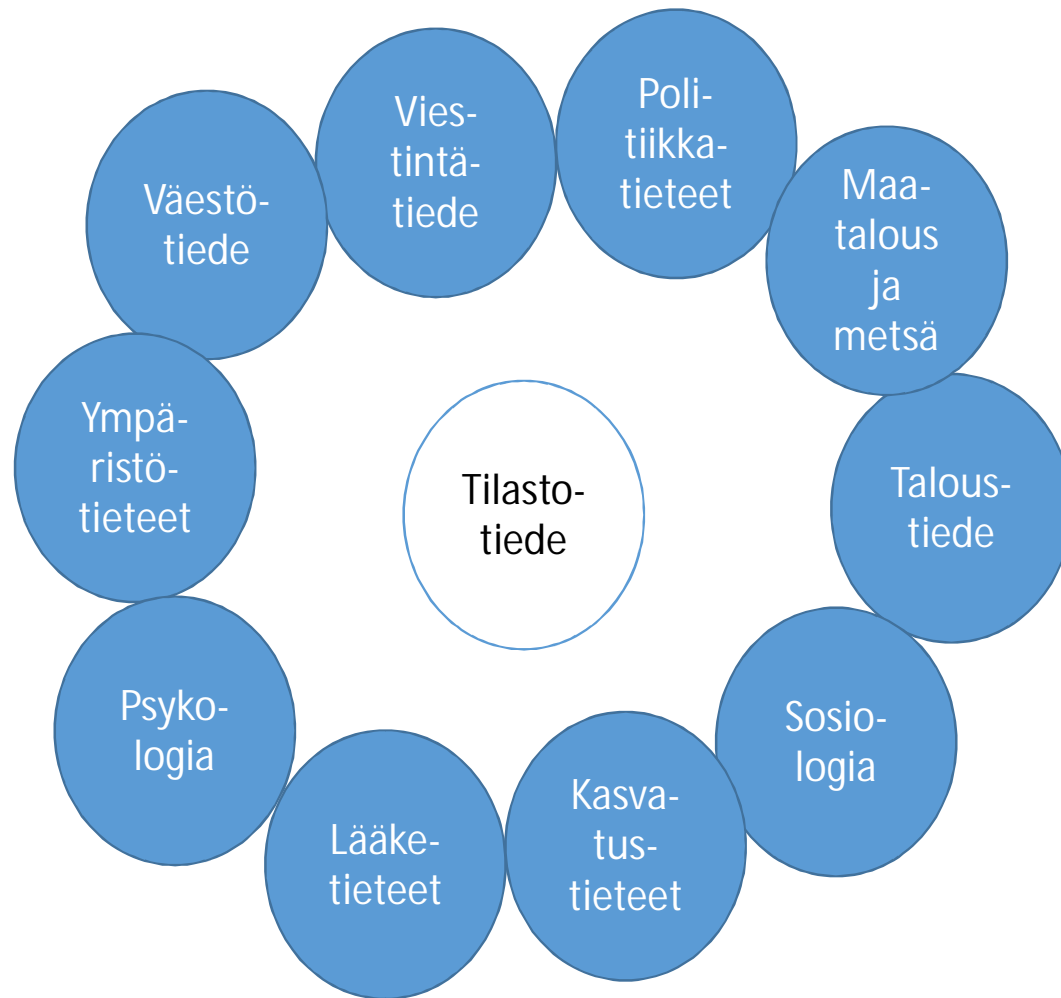
Tilastotiede siis pyrkii paremman maailman kehittämiseen kuten muutkin tieteet.

Media ei ainakaan tiedettä arvosta. Esimerkki:

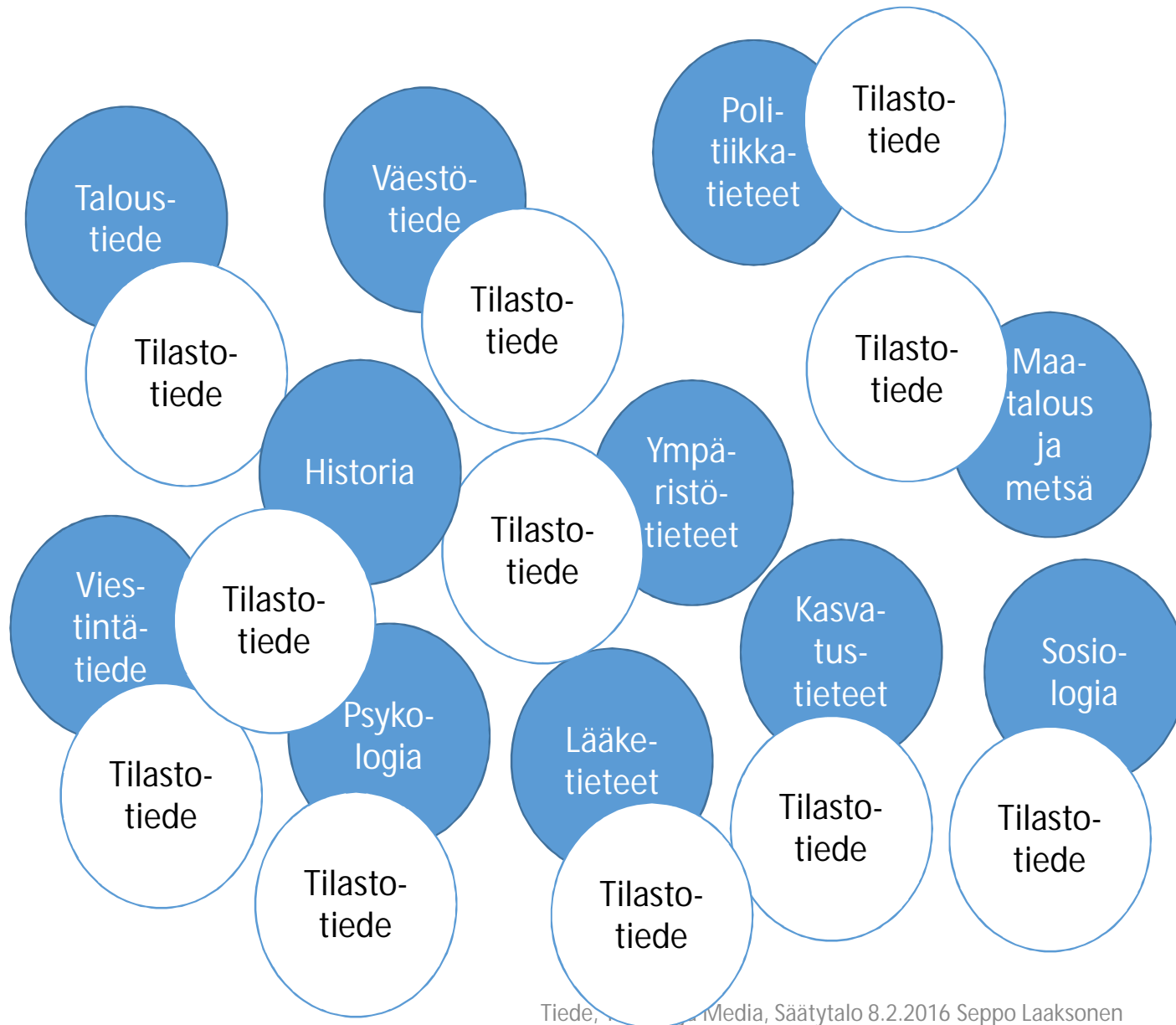
- Maailman Tilastotiedepäivä on nyt ollut 5 vuoden välein, viimeksi 20.10.2015.
- Yritin muutamaa viikkoa ennen houkutella valtamedioita sen huomioimiseen lähettämällä asiasta tietoa ja joitakin ajatuksia jutuiksi.
- Mitään ei tapahtunut paitsi yhdestä paikasta tuli selkeä hylkäys.

Jokin tässä tieteessä mättää. Toivottavasti ei liika kunnioitus ja siitä johtuva pelko.

Yritän tässä esityksessä häivyttää pelkoakin 'tilastopolitikoinnillani.'  
Tilastotiede on inhimillinen (humaani) tiede.



Tilastotieteen ensimmäinen professori Leo Törnqvist (Linus Thorvaldsin äidin isä) markkinoi tilastotiedettä oheisen kaltaisella **LETTUPANNULLA** jossa tilastotiede on keskiössä. Se on tietysti liioiteltu. Itse tekisin siitä seuraavan kaltaisen:



Tämän idea on, että yhdessä hyvä tulee, siinä kehittyvät molemmat tieteet ja useammatkin jos tiedon siirtoa tapahtuu. Siltä pohjalta tilastotieteen suuret saavutukset ovat syntyneet. Kun myöhemmin ryhdyttiin toimimaan liikaa omissa oloissa kehitys pysähtyi tai meni väärään suuntaan (monia viittauksia tähän löytyy).

Tilastotiede on kuitenkin pieni tiede eikä hyvin resurssoitu. Koulussakaan sitä ei juuri tunneta. Akatemian ja muiden käyttämissä luokituksissa ei sitä juuri löydy. Joudun valitsemaan erilaisissa lomakkeissa vähimmän huonon vaihtoehdon kun oikeata ei ole. Akatemian mainio Tieteen Tila –raportti jokin aika takaisin ei kerro mitään tilastotieteestä tai sitä on vaikea sieltä havaita. En usko että tämä päätieteeni olisi silti kovin korkealla arvoasteikoissa.

## Kaikki havaintotieteet tarvitsevat ja käyttävät tilastotiedettä

- Valitettavan usein keksien sen metodeja uudelleen vaikkapa kovan laskentakapasiteetin avulla kuten tietotekniikassa (big data), ja huonosti.
- Tutkijat taas pyytävät apua liian usein liian myöhään, jolloin ei kovin paljoa enää ole tehtävissä.
- Tutkimushankkeista päätettäessä epäilen arvioijien tuntevan tilastotieteen mahdollisuudet huonosti. Tulosten tultua käyttöön epäilen myös, onko varmasti käytetty oikeita menetelmiä ja oikein? Eli paraneeko maailma siis? Tämän päivän kuuma aihe Tiedevilppi on hyvä esimerkki.
- Eri tilastotieteen sovellusalueilla näyttää olevan eri termejä. Joskus on vaikea keksiä mitä mikäkin tarkoittaa vaikkapa terveystieteissä jos on enimmäkseen yhteiskuntatieteiden kanssa tekemisissä.

Vielä yksi yleisnäkökulma ennen kuin siirryn konkreettisiin esimerkkeihin:

- Nykyäänkin puhutaan keskittämisestä. Tilastotieteen ja muiden metoditieteiden (matematiikka ym) keskittämisellä jonnekin on vain huonoja seurauksia koska tietty annos tilastotiedettä tarvitaan melkein kaikilla aloilla.
- Ainakaan OPETUSTA ei voi keskittää muutamaa korkeakouluun. Koska opetus taas vaatii uutta tietoa tutkimuksesta, on myös tutkittava.
- Keskittäminen joihinkin tilastotieteen erityisalueisiin on hyväksyttävää. Tietysti minusta mieluiten Surveymetodiikkaan, sen laajassa merkityksessä, mutta kaikki tuskin ovat samaa mieltä.



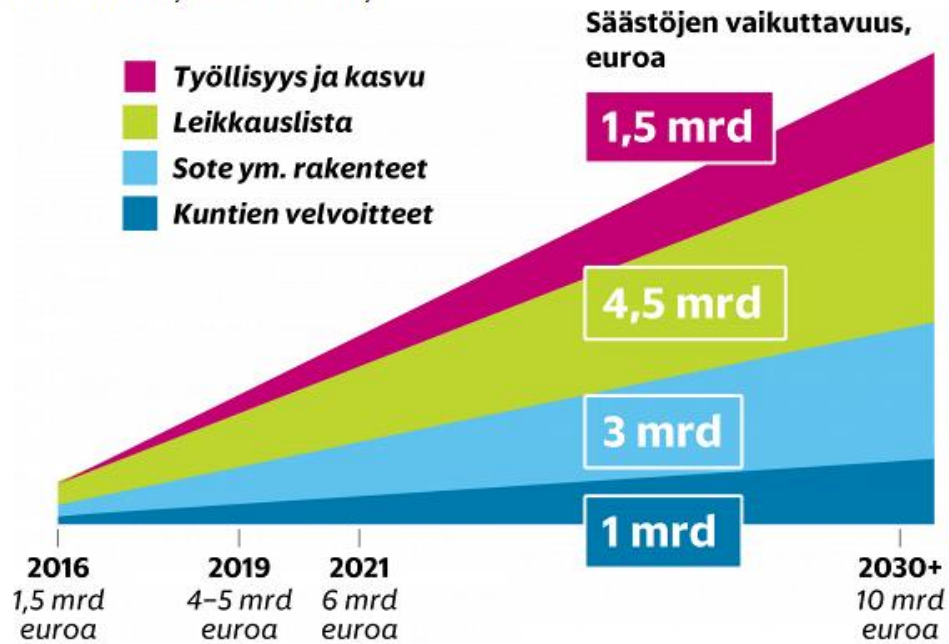
NYT SIIS ESIMERKKIEN SESSIO.

Mukana on paljon tilastografiikkaa mitä on hyvä käyttää mutta on vaara, että huonolla grafiikalla annetaan väärä kuva. Pelkät numerot taas voivat olla vaikeita ymmärtää.

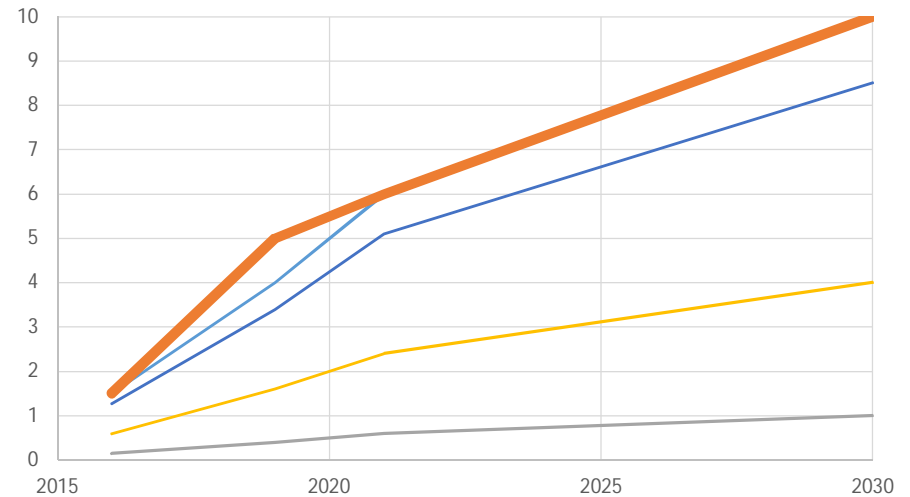
Helsingin Sanomien kuvio viime kesältä. En tiedä kenen tekemä? Hyvä esimerkki siitä, että on kaksi perusvirhettä. Tässä ensin se ettei tunnu olevan väliä miten x-akseli asetetaan eli akselin arvot ei vastaa vuosia.

### Sipilän viuhka esittelee hallituksen taloussuunnitelmaa

Hallitusneuvottelijat esittelivät keskiviikkona kaavion, joka kuvaa Suomen julkisen talouden korjaamista. Klikkaa kuvan lukuja nähdäksesi lisätietoja.



Sipilän viuhka HS:n tietojen ohjalta, lineaarinen

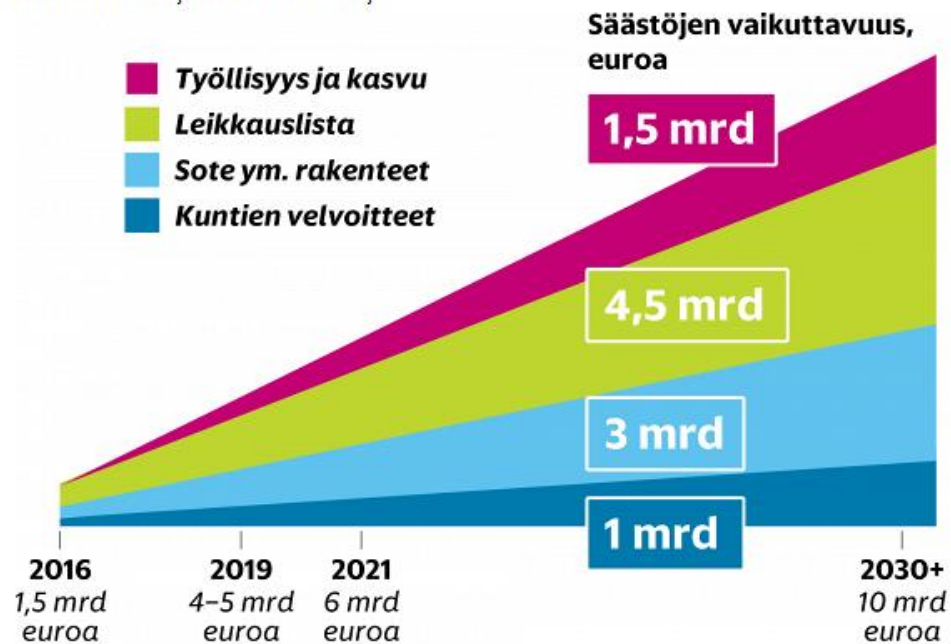


Tässä akseli loppuu 2030:een kun + on epäselvä

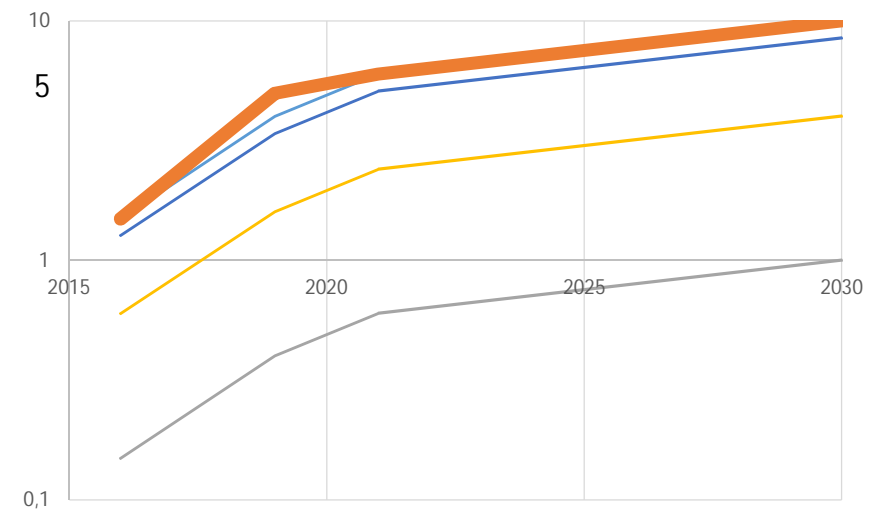
Helsingin Sanomien kuvio viime kesältä. Tässä toiseksi ettei välitetä suhteellisuudesta joka saataisiin helposti logaritmisella asteikolla y-akselille. Suhteellisuus vie tosin hienon viuhkan mutta kaikki rahat erityisesti ovat suhteellisia eli jatkuvat muuttujat ja siksi logaritmia pitäisi käyttää.

### Sipilän viuhka esittelee hallituksen taloussuunnitelmaa

Hallitusneuvottelijat esittelivät keskiviikkona kaavion, joka kuvaa Suomen julkisen talouden korjaamista. Klikkaa kuvan lukuja nähdäksesi lisätietoja.



Sipilän Viuhka HS:n tietojen pohjalta, log

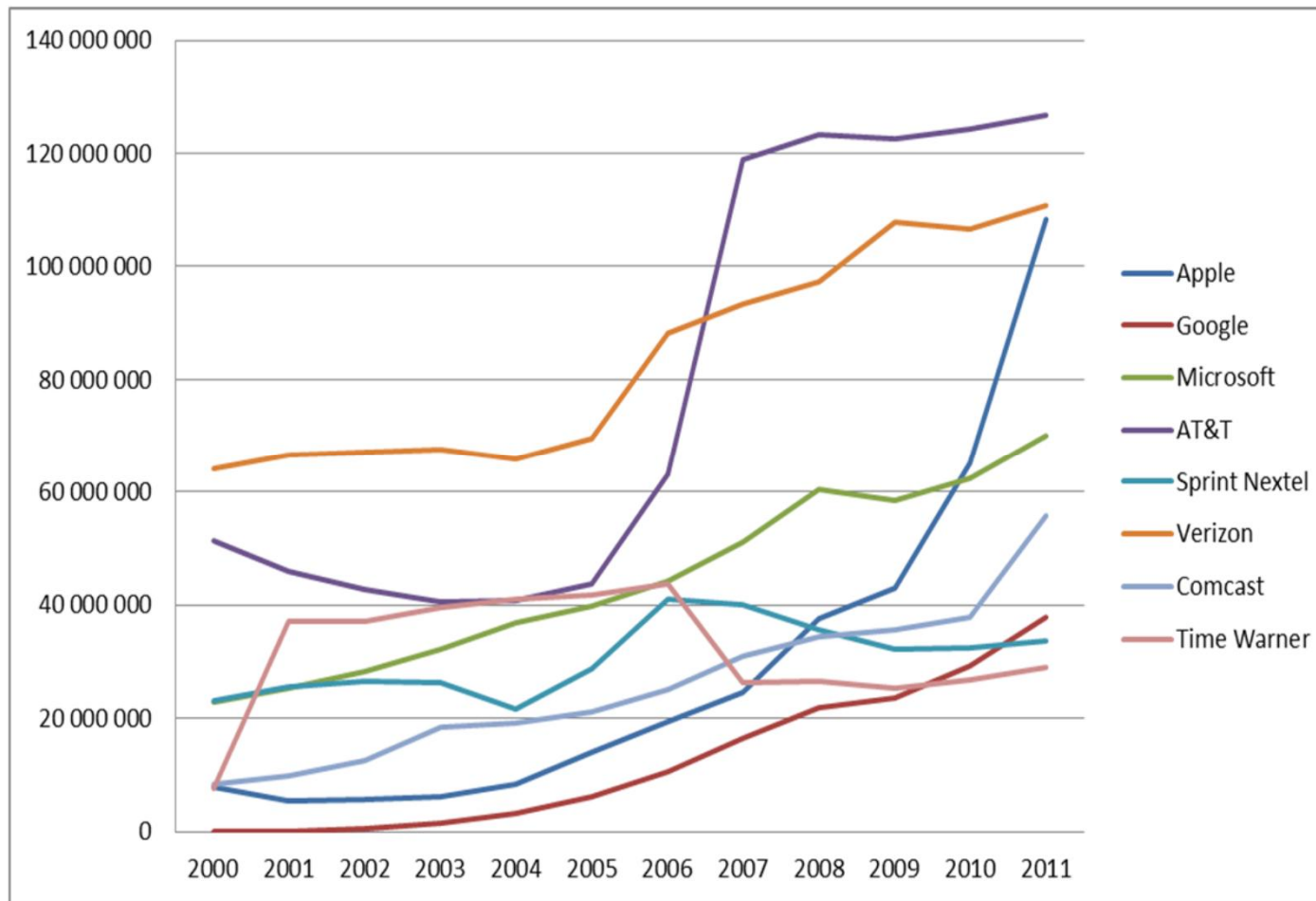


Jatkan vielä logaritmin kanssa koska se niin huonosti hallitaan sekä mediassa että useimmissa tieteissä vaikka nykyään on helpohkoa koska graafiset ohjelmistot ovat kehittyneet. Toista oli 10 vuottakin sitten.

Nytkin y-akselin logaritointi edistää asiaa.

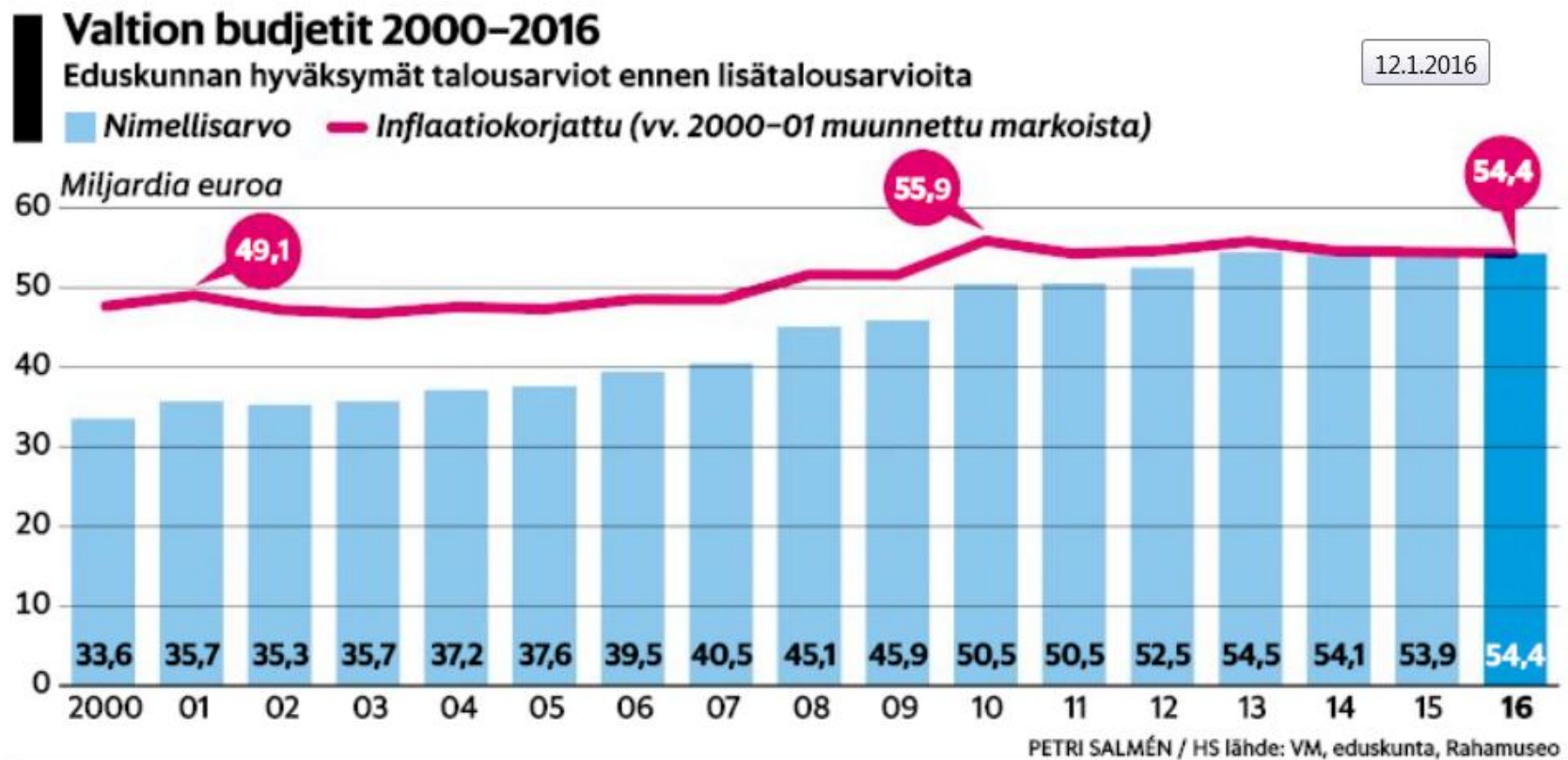
Seuraava esimerkki on vain puhdas esimerkki enkä osaa sanoa kuinka erilainen kuva tilanteesta tulisi jos pystyakseli olisi logaritminen. Tämä on tässä siis siksi että uskotte että logaritmista muunnosta ei ole kovin yleisesti tapana tehdä.

Tämä ei siis huono esimerkki. Olisi mukava nähdä y-akseli logaritmisena niin näkisi miten **muutos on suhteellisena**.



Etlan A sarja 48, 2014

Tämä HS:n kuvio on itse asiassa melko hyvä vaikkei logaritmia eli suhteellisuutta olekaan. Hyvä asia on Inflaatiokorjattu kuvio. Se ei ole logaritminen mutta auttaa muutoksen arvioimisessa paljon.



Jatkan vielä logaritmin kanssa ja usealla aikasarjalla.  
Nytkin y-akselin logaritointi edistää asiaa.

Tämä on tavanomainen huono kuvio lineaariselta pohjalta. Suurituloisin 1% on vieläpä poissa vaikka väitetään toisin.

HS 19.12.2014

# Tuloerot kasvoivat viime vuonna

**Suomalaisten mediaanitulo aleni 150 euroa**

STT-HS

**TULOEROT** kasvoivat viime vuonna toissa vuoteen verrattuna, koska suurituloisten tulot kasvoivat muita nopeammin, kertoo Tilastokeskus.

Tilastokeskuksen torstaina julkistaman tulojakotilaston mukaan väestön suurituloisimman kymmenesosan tulotaso kasvoi vuonna 2013 reaalisesti 2,6 prosenttia, kun pienituloisimman kymmenesosan tulotaso kasvoi 0,3 prosenttia.

Näiden välille jäävissä tulo-kymmenyksissä tulokehitys oli vielä heikompaa.

**SUOMALAISTEN** käytettävissä olevien tulojen mediaani vuonna 2013 oli 23 715 euroa yhden hengen asuntokunnalle laskettuna.

Kun tulonsaajat asetetaan tulojen mukaan suuruusjärjestykseen, mediaanituloon keskimäisen tulonsaajan tulo.

Mediaanitulo aleni edellisvuodesta 150 euroa.

Pienituloisin kymmenes suomalaisista ansaitsi nettotuloina 13 040 euroa vuodessa, eli 1 090 euroa kuussa.

**SUURITULOISIMMAN** kymmeneksen raja oli 40 955 euroa.

Suurituloisimman kymmeneksen tulo-osuus kasvoi 0,6 prosenttiyksikköä edellisvuodesta. Joukko sai tuloista 23 prosenttia eli lähes saman verran kuin pienituloisin neljäkymmentä prosenttia yhteensä.

Suurituloisimman prosentin raja oli 79 600 euroa. Se tarkoittaa 6 630 euron nettotuloja yhden hengen asuntokunnalle ja 10 000 euron nettotuloja kahden aikuisen asuntokunnalle.

**SUHTEELLISIA** tuloeroja kuvaava niin sanottu Gini-kerroin oli kuitenkin likimain samalla tasolla kuin kymmenen vuotta sitten.

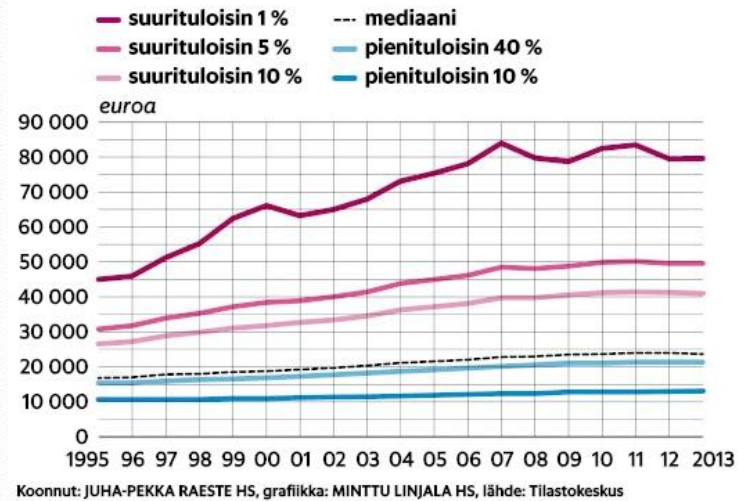
Sen lievä nousu vuonna 2013 johtui Tilastokeskuksen mukaan suurituloisimmille kohdentuvien myyntivoittojen kasvusta.

Vuoteen 1995 verrattuna Gini-kerroin on kasvanut Suomessa merkittävästi eli noin 5,4 prosenttiyksikköä.

Tuloerot ovat olleet Tilastokeskuksen mukaan suurimmillaan vuonna 2007.

## Näin tulotasot ovat eri ryhmissä kehittyneet

Tulorajat tulojakauman eri kohdissa, euroa kulutusyksikköä kohden vuoden 2013 rahassa



Koonnut: JUHA-PEKKA RAESTE HS, grafiikka: MINTTU LINJALA HS, lähde: Tilastokeskus



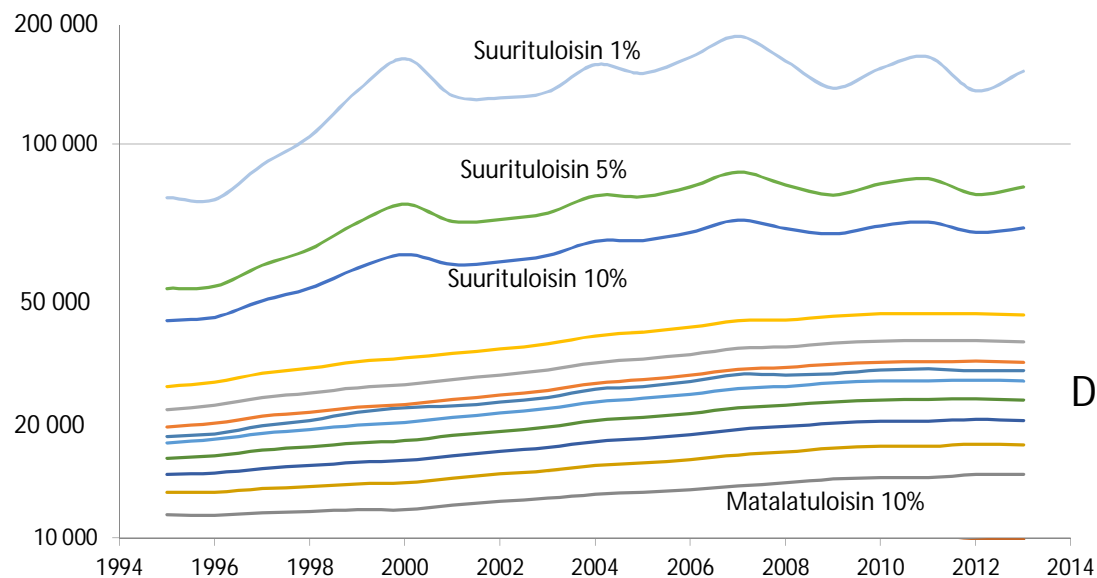
lähde, Ilmatorvi ja Media, Säätytalo 8.2.2016 Seppo Laaksonen



Keskiarvoja eri tuloryhmissä Tilastokeskuksen mukaan välillä 1995-2013

Tässä tekemässäni kuviossa kaikki tuloryhmät mahtuvat kuvioon ja ylinkin on mukana kuten ei ollut edellisessä. Nähdään että suurituloisissa on nousuja paremmissa taloussuhdanteissa. Matalampituloisten käyrät melko tasaisia. Inflaatiokorjauksenkin voisi tehdä.

Nyt pystyakseli on logaritminen.



Desiiliryhmät



Mikään tiede eikä tieteilijä ole erehtymätön. Koskee tietysti minua ja muitakin (tilasto)tieteilijöitä kuin Juha Alhoa  
HS 15.1.2016.

Kyse tässä on lasten tiedekysymyksestä jonka esitti Leena Koskinen, 4

## **Onko enemmän setiä vai enoja?**

[Tiede](#) 15.1.2016 2:00

Enoja oli ensin  
15.1. enemmän  
kuin setiä. Kolme  
päivää  
myöhemmin  
suunnilleen yhtä  
paljon.



Nasse-setä kuuluu vähemmistöön.

## Onko enemmän setiä vai enoja?

**Leena Koskinen, 4**

**ENOJA** on enemmän.

Lapsella on aina äiti ja isä. Oletetaan, että vanhemmat tulevat kaksilapsisista perheistä, joita voi lasten sukupuolen ja syntymäjärjestyksen mukaan lajiteltuna olla neljää tyyppiä: a) tyttö-tyttö, b) tyttö-poika, c) poika-tyttö tai d) poika-poika.

Jos vastasyntyneistä puolet olisi tyttöjä ja puolet poikia, niin kaikkia neljää perhetyyppiä olisi yhtä paljon.

Äiti voi olla peräisin perhees-

tä a, b tai c (d-perheessä ei ole tyttöä, josta voisi tulla äiti). Näistä perheistä kahdessa eli b:ssä ja c:ssä äidillä on veli. Näin ollen kahdella kolmesta äidistä olisi enoksi kelpaava veli.

Isät taas tulevat perheistä b, c, d. Näistä perheistä vain yhdessä isällä olisi veli. Niinpä vain joka kolmannella isällä olisi sedäksi kelpaava veli.

Näin laskien enoja olisi siis kaksi kertaa niin paljon kuin setiä.

Tosiasiasa todennäköisyys,

että vastasyntynyt on tyttö, on noin 0,49, ja poikien osuus on vastaavasti 0,51. Tällä ei kuitenkaan ole juuri vaikutusta. Näin tarkemmin laskien enoja olisi 1,97 kertaa niin paljon kuin setiä.

Lapsia voi myös olla perheessä enemmän kuin kaksi. Tämä kaventaa eroa, mutta melko vähän. Jos vanhemmat olisivat peräisin esimerkiksi kolmilapsisista perheistä, niin enoja olisi 1,78 kertaa niin paljon kuin setiä.

**Juha Alho**

tilastotieteen professori

Helsingin yliopisto

## Enoja ja setiä on suunnilleen yhtä paljon.

Lapsella on aina äiti ja isä. Oletetaan, että vanhemmat tulevat kaksilapsisista perheistä, joita voi lasten sukupuolen ja syntymäjärjestyksen mukaan lajiteltuna olla neljää tyyppiä: a) tyttö-tyttö, b) tyttö-poika, c) poika-tyttö tai d) poika-poika.

Jos vastasyntyneistä puolet olisi tyttöjä ja puolet poikia, niin kaikkia neljää perhetyyppiä olisi yhtä paljon.

Näistä perheistä kahdessa eli b:ssä ja c:ssä äidillä on veli. Näin ollen kahdella neljästä äidistä on enoksi kelpaava veli.

Isällä taas on veli vain perheessä d, mutta perheen molemmat pojat voivat tulla sediksi toistensa lapsille. Niinpä kahdella neljästä isästä on sedäksi kelpaava veli.

Näin laskien enoja ja setiä on yhtä paljon. Tosiasiassa todennäköisyys, että vastasyntynyt on tyttö, on noin 0,49. Poikien osuus on vastaavasti 0,51. Tällä ei kuitenkaan ole juuri vaikutusta.

Juha Alho

tilastotieteen professori

Helsingin yliopisto

**Oikaisu 18.1. klo 13.30: Toisin kuin vastauksessa alunperin sanottiin, enoja ei ole kahta kertaa enempää kuin setiä vaan molempia on suunnilleen yhtä paljon.**

Tiede, Tilastot ja Media, Säätytalo 8.2.2016 Seppo Laaksonen

Minulla on ollut kaksi enoa mutta viisi setää mutta tätä ei voi yleistää mihinkään. Sen kuitenkin uskallan yleistää että tätejä on enemmän kuin enoja ja setiä niin kauan kuin annetaan eri nimi Äidin Sisarelle ja Isän Sisarelle.

Medialla on taipumus tehdä yksilöistä haastatteluja. Tämä on ihan oikein koska ne kiinnostavat monia lukijoita. Ongelma syntyy jos annetaan jutun perusteella sellainen kuva että haastatellut ihmiset, yritykset tms edustavat jotakin määrättyä TAVOITEPERUSJOUKKOA johon tilastotieteilijä tulokset YLEISTÄÄ. Vaikkei siis tätä sanotakaan näin niin valistumaton lukija vetää siitä helposti sellaisen johtopäätöksen. [Moni samaistaa lööpin tavoiteperusjoukkoon, minä en.](#) Vielä pidemmälle saattavat mennä jotkin muut tahot kuten äskettäisessä TV:n terveyskeskustelussa jossa Vaihtoehtoterveysmenetelmien esittelijä sanoi suunnilleen näin:

**"Ihminen tulee paljon vakuuttuneemmaksi kun toinen ihminen kertoo kokemuksensa eikä tarvita tilastoa 10 000 ihmisestä." Jälkimmäisen datan pohjaltahan keskustelussa Tieteilijä perusteli uusia terveys suosituksia mitkä menevät varmaan täällä melko hyvin perille. Noudattaminen on eri asia.**

Esimerkki havainnollistaa yleistä ongelmaa.

Yleistäminen voidaan tehdä luotettavasti vain kun tiedetään mahdollisimman tarkasti mihin YLEISTETÄÄN. Tilastotieteessä se siis on tyypillisimmin TAVOITEPERUSJOUKKO, kuten

- (i) Gallupeissa tulevien vaalien äänestäjät (mitä ei tiedetä mutta yritetään tavoittaa ÄÄNESTYSKELPOISTEN kautta). Galluppien tulokset ovat kohtuullisia siksi, että ÄÄNESTÄMÄTTÖMISTÄ vain harvat osallistuvat galluppeihin.
- (ii) PISA-tutkimuksissa 15 -vuotiaat koulua käyvät.
- (iii) Muissa kyselyissä tietyn ikäiset maassa asuvat henkilöt. European Social Surveyssä 15 vuotta täyttäneet, useimmissa muissa on jokin yläikäraja kuten gallupeissa kai 79 vuotta.

Ymmärrän hyvin etteivät monia lukijoita kiinnosta kaikki luvut, eivät erityisestikään epävarmuusluvut kuten luottamusvälit, virhemarginaalit, keskivirheet tai p-arvot mitkä ovat melkein samoja käytännössä. Siksi toimittajan, tutkijan ja tulosten tulkitsijan tulisi osata kertoa tämä sopivalla tavalla eikä ainakaan siten että esittää tuloksia virheellisesti. Vaaditaan siis riittävä tilastotieteen osaaminen, ainakin niin paljon että osaisi kysyä asiantuntijalta.

Jos siis tavoiteperusjoukko on kirkas, tilanne on hyvä. Tavoiteperusjoukossa voi olla erilaisia osajoukkoja kuten, sukupuoli, ikä, tuloryhmä ja koulutus. Niihinkin voidaan tuloksia yleistää jos tietoa on riittävästi. Viime vuosikymmeninä on voimakkaasti yleistynyt ns. kvalitatiivinen tutkimus jossa esimerkiksi syvähaastatellaan jotakin ihmisryhmää. Tästäkin saatuja tuloksia voidaan YLEISTÄÄ jos tiedetään mitä ihmisryhmä edustaa koko tavoiteperusjoukossa. Tällöin siis joidenkin mielestä tylsien numeroiden taakse tulee lisää, kuin 'lihaa luun ympärille'. Tällainen kvalitatiivinen tutkimus on mitä kannatettavin. Median edustajat voisivat parhaimmillaan olla hyviä kvalitatiivisia tutkijoita mutta ilman tietoa kokonaisuudesta tilanne jää epämääräiseksi. Jotkut lukijat toki voidaan saada uskomaan mitä vain.



## Lisänäkökohta:

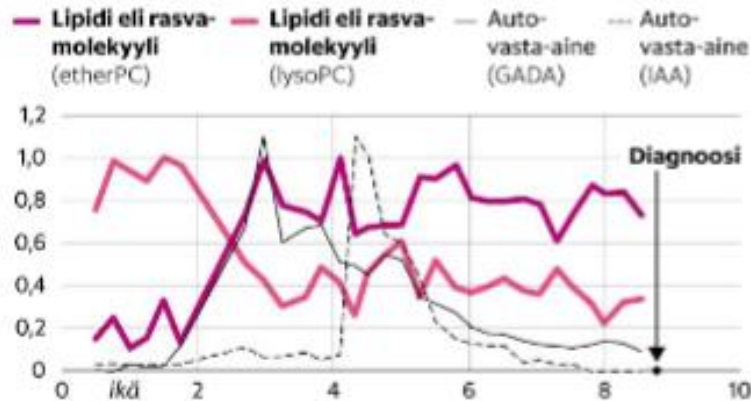
Jos media, kvalitatiivinen tutkija tai muu kertoo yksittäistä ihmistä, yritystä tai muuta toimijaa koskevan tarinan totuudenmukaisesti, se on aina arvokas. Toisaalta mikrotason tilastodatasta löytyvät myös yksilöt, sisältäen yleensä monia muuttujia niistä. Tietosuoja on molemmissa tilanteissa otettava vakavasti. Makrodatassa, jos 'makro' on iso, on tietosuoja helpompi hoitaa mutta ei tietenkään päästä kovin 'syväälle.'

Esimerkki Tiedevilpistä mistä HS teki laajan jutun 7.2.. Ansiokas juttu. Tänään 8.2. lehti jatkaa: "HS toi esiin sen, että terveenä pysyvän lapsen veritesteissä näkyy sama ilmiö kuin diabetekseen sairastuvan tytön tuloksissa. Tämän olennaisen eron havainnollistavaa kuvaparia ei näytetty alkuperäisessä tiedeartikkelissa. "

Maallikkona olen ollut monista kohdista ymmällään. Esimerkiksi jos yhden havainnon perusteella tehdään ratkaisevia johtopäätöksiä, se ei oikein mene minulle perille. Normaalisti pitää siis olla edustava aineisto ja riittävästi havaintoja mutta onko tämä alue jotenkin poikkeus. Sen voisi joku mulle selittää? Yhdestä havainnosta en yleistäisi vaikka jokainen havainto on arvokas. Tässä jutussa oli lisäksi kaksi kuviota joiden ymmärtäminen minulle oli melko mahdotonta. Entä teille? Kuvion pitäisi kertoa olennaiset asiat sellaisenaan ja edes akselit pitäisi kertoa. Vaaka-akseli on mutta pystyakseli pitää itse keksiä.

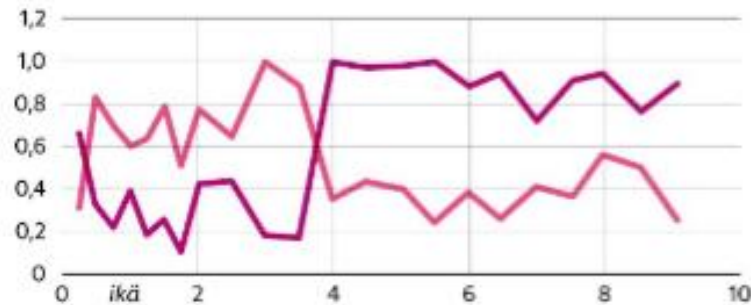
### Diabetesartikkelissa julkaistiin tämä kuva

Kuvalla haluttiin osoittaa, että lipidit eli rasvamolekyylit reagoivat voimakkaasti ennen ykköstyyppin diabeteksen puhkeamista. Tämän tytön diabetes puhkesi noin 9-vuotiaana.



### Terveenä pysyneellä verrokilla on samankaltainen ilmiö

Siksi yllä olevasta kuvasta ei voi vetää johtopäätöksiä. Tätä kuvaa ei julkistettu.



Koonnut: KATJA KUOKKANEN / HS, grafiikka: MINTTU LINJALA / HS

En pidä legendoista joiden perusteella kukin viiva on hidas löytää. Joissakin tieteissä legendat kelpaavat mutta mediassa niiden ei tulisi kelvata paitsi itsestään selvissä tapauksissa.

Tässä pitää vielä keksiä itse mitä viivat tarkoittavat enkä tiedä mitä tästä voisi päätellä. Käyrien välinen korrelaatio aika negatiivinen.

Valtakunnan päälehden tulisi olla esikuva. Jatkan samasta lehdestä mutta eri aiheella.

Tässä kolme lehtileikettä HS:n Vaaligalluppien julkistuksista. Pieniä eroja on.

HS 7.7.2012

**FAKTA**

### Näin tutkimus tehtiin

- TNS Gallup selvitti Helsingin Sanomien toimeksiannosta puolueiden kannatusta eduskuntavaaleissa.
- Haastateltavilta kysyttiin myös, kuinka varmoja he ovat puoluevalinnastaan.
- Haastatteluja tehtiin 2 487. Vastaajat edustavat Suomen äänestysikäistä väestöä Ahvenanmaata lukuun ottamatta.
- Tutkimus toteutettiin 4. kesäkuuta – 1. heinäkuuta 2012.
- Tutkimuksen virhemarginaali on suurimmillaan vajaat 2 prosenttiyksikköä suuntaansa.

**FAKTA** HS 7.4.2013

### Näin tutkittiin

- TNS Gallup toteutti HS:n pyynnöstä tutkimuksen kansalaisten suhtautumisesta vireillä oleviin kansalaisaloitteisiin.
- Tutkimus tehtiin puhelinhaastatteluna 19. 1.–27. 3. 2013. Haastattelujen määrä oli 1 004.
- Otos edustaa maan 15 vuotta täyttäneitä väestöä Ahvenanmaata lukuun ottamatta.
- Tutkimustulosten virhemarginaali on suurimmillaan kolme prosenttiyksikköä suuntaansa.

Heide, i

HS 22.1.2016

### Fakta

#### Näin tutkittiin

- TNS Gallup selvitti Helsingin Sanomille puolueiden kannatusta eduskuntavaaleissa.
- Tutkimusta varten haastateltiin 2 432 ihmistä 14.12.2015–14.1.2016.
- Otos edustaa manner-Suomen äänestysikäistä väestöä.
- Virhemarginaali on noin kaksi prosenttiyksikköä suuntaansa suurimmilla puolueilla.

sonen

HS on esittänyt kerran 19.3.2011 epävarmuuden laajemminkin . Syynä oli se, että opetin toimittajaa ja tein paremman graafisen esitystavan. Ks. lehdestä 4.2011 tai Kanava 8/2011. Myöhemmin lehti palasi entiselle tyylilleen, kun Raeste siirtyi toisiin hommiin lehden sisällä. Kanavaan laitoimme Ville Pernaan kanssa nimen **EPÄVARMUUSHAARUKKA** kun hänestä **EPÄVARMUUSMARGINAALI** ei ollut hyvä..

## ■ TAUSTA

### Virhemarginaali on satunnais- ja harhamarginaalin sekasikiö

**VIRHEMARGINAALI** on epäselvä käsite, joka syntyy satunnais- ja harhamarginaalin yhdistämisestä.

**SATUNNAISMARGINAALI** mittaa yhdessä tutkimustilanteessa tietyn otosmäärän sisällä olevaa todennäköisyyttä.

Jos esimerkiksi yhdessä kahden tuhannen hengen mielipidetiedustelussa 20 prosenttia vastaajista sanoo äänestävänsä kokoomusta, tutkija tietää taulukosta katsomalla, että kannatuksen satunnaismarginaali on 1,9

prosenttia suuntaansa. Toisin sanoen kokoomuksen kannatus vaihtelee 95 prosentin todennäköisyydellä 18,1–21,9 prosenttiyksikön välillä.

Silti keskimäinen ilmoitettu 20 prosentin luku on todennäköisin, ja noin kaksi kertaa todennäköisempi kuin luku 18,1 tai 21,9. Tässä tapauksessa todennäköisyys noudattaa pitkälti ns. normaalijakaumaa, jonka kansa tuntee Gaussin käyränä.

**HARHAMARGINAALI** selittää eroa, joka syntyy kun esimerkiksi

Taloustutkimus ja TNS Gallup tutkivat samalla viikolla puolueiden kannatusta, mutta tulokset ovat aivan erilaisia.

Harhamarginaali on systemaattinen, aineiston keruussa ja laskentametodissa oleva epävarmuus. Harhan mahdollisuus on erityisen suuri uusissa tilanteissa. Esimerkiksi nyt, kun perussuomalaisten gallup-kannatus on lyhyessä ajassa moninkertaistunut, tämä mahdollisuus kasvaa. Sitä ei kuitenkaan voida varmastikin arvioida.

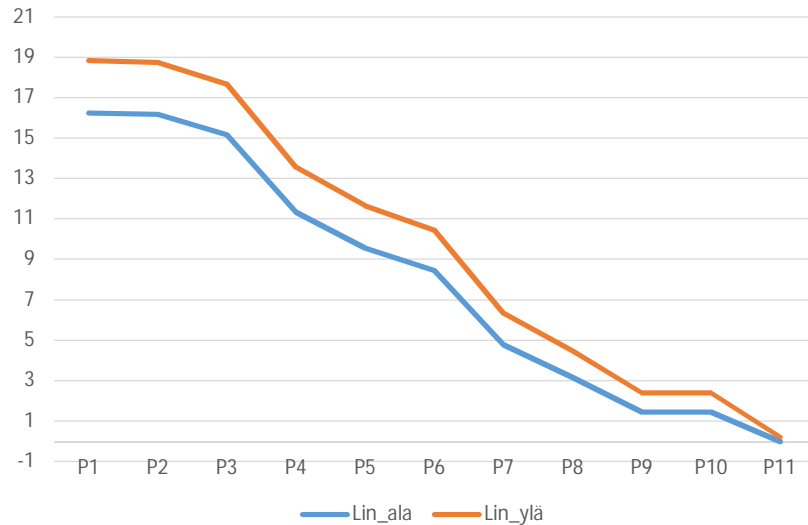
**Juha-Pekka Raeste**

Seuraavalla sivulla Virhemarginaali melko perinteisesti koska harhasta ei ollut tietoa. Tässä on siis oletus että vastaajat ovat satunnaisesti valikoituneet tiettyjen ositteiden sisällä. Sillä perusteella niitä voidaan laskea.

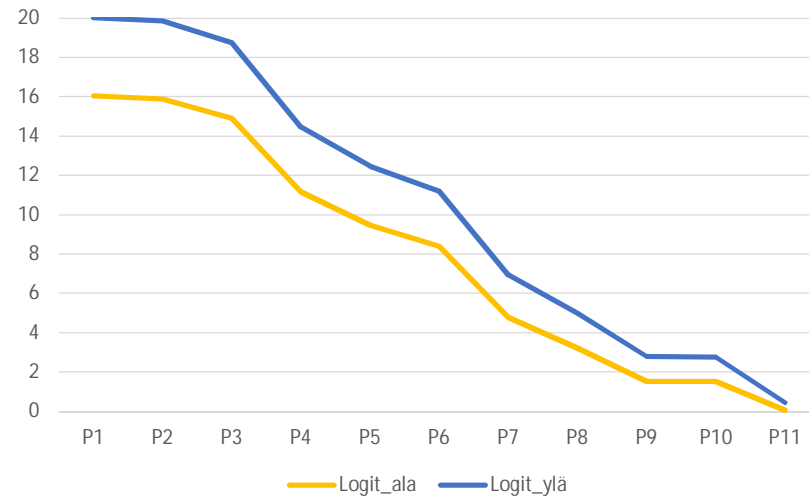
Vasemmassa kuviossa on perinteinen marginaali joka perustuu normaalijakaumaan ja toinen oikealla jossa on Logit-muunnos eli virhemarginaali ei ole symmetrinen kuten edellisessä. Pienissä arvoissa tämä voi viedä negatiiviseen kannatukseen.

Esimerkkien puolueet ovat vaalien 2019 puolueita.

Puolueiden kannatus ja lineaarinen virhemarginaali, P11:n alaraja negatiivinen



Puolueiden kannatus ja logit-virhemarginaali, ei negatiivisia



Toinen puoli virhemarginaalista on lineaarisessa on välillä 1,3 ja 0,2 %-yksikköä, logit-pohjaisessa siis yläpuolella näitä isompi, alapuolella melko sama. Molemmat ovat vähemmän kuin 2 %yksikköä suuntaansa suurimmilla puolueilla. Ohessa ei liene yhtä suuria puolueita. Jos kannatus on 20%, marginaali nousee hieman. Veikkaus: Gallup-firmoissa asetetaan hieman lisää marginaalia varmuuden vuoksi, koska on monia muita epävarmuuksia kuin vastaajien valikoituminen. Muistelen että joskus ennen virhemarginaalit ilmoitettiin pienemmiksi.



On mainio asia, että virhemarginaali on amerikkalaisen mallin mukaan otettu käyttöön vaikei hyvin. Olen ihmetellyt miksei sitä kaikissa muissakin otostulosteissa esitetä tai sitten siitä ei erityisemmin välitetä.

Katsoin äskettäin netistä yhtä selostetta:

" Tutkimusaineisto koottiin Gallup Kanavalla ja haastatteluja tehtiin yhteensä 1 149. Vastaajat edustavat Manner-Suomen 18–75 vuotta täyttäneitä väestöä. Tulosten virhemarginaali on suurimmillaan vajaan kolme prosenttiyksikköä suuntaansa."

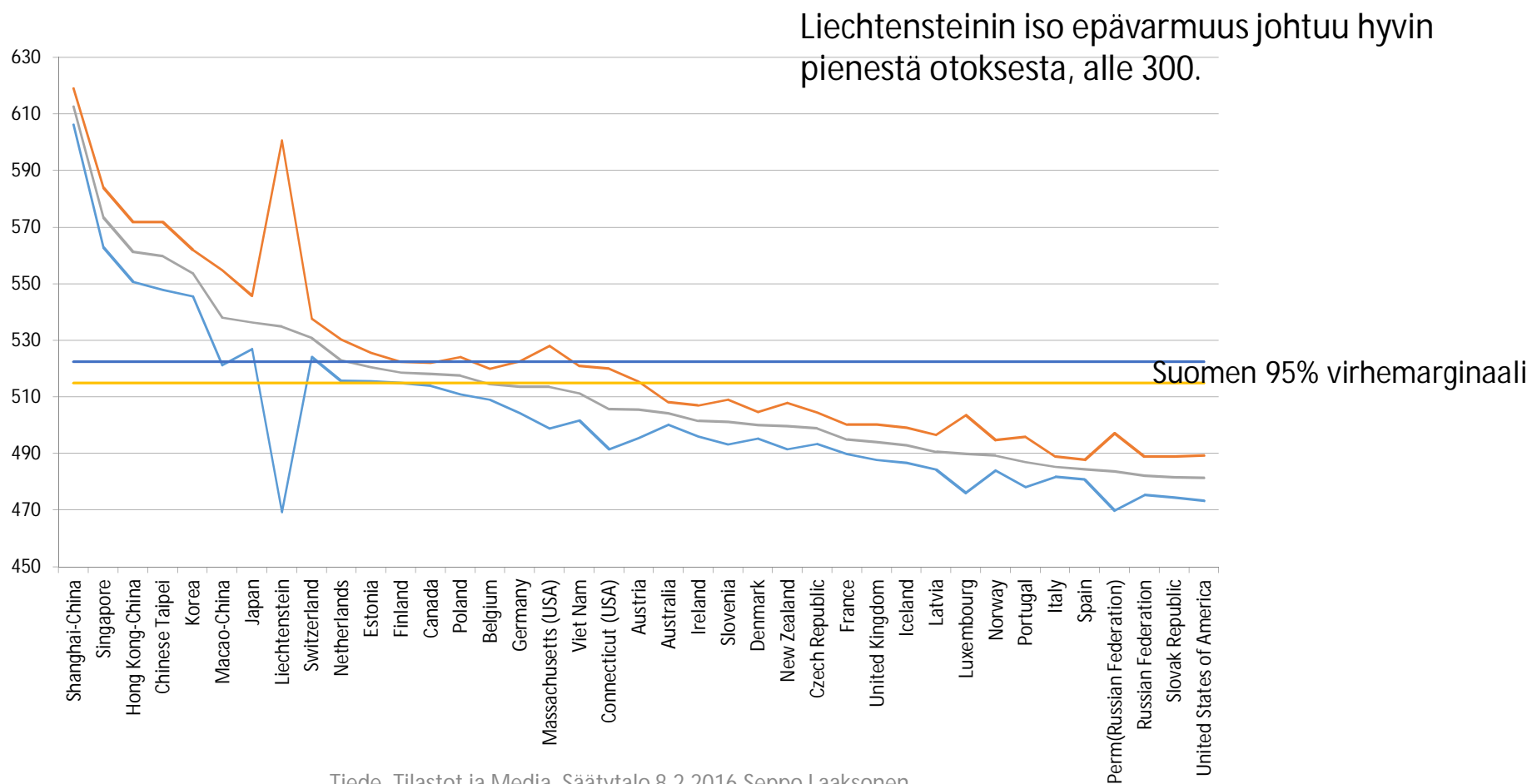
Tuloksia on oikeasti erilaisilla prosenttiosuuksilla eli virhemarginaali vaihtelee paljonkin. Otan nyt esimerkin jossa olivat kaikki virhemarginaalit (95% luottamusvälit) mediankin tiedossa, mutta käyttö oli suoraviivainen.

Pisa 2012: Matemaattis-tilastollisen lukutaidon kärki (maita ja talouksia on 65, tässä 39).

Ylempi viiva: 95% luottamusvälin yläraja; Alempi viiva: 95% luottamusvälin alaraja;

Keskellä keskiarvo jonka mukaan järjestys

Tehtävä: Kuinka MONES on Suomi? Yksinkertainen vastaus= 12:s mikä oli mediassa



Parempia vastauksia:

- Jos otetaan huomioon luottamusvälit kuten pitää, niin Suomea edellä on kahdeksan maata tai Taloudellista aluetta. Näistä neljä on maita eli Singapore, Korea, Japani ja Sveitsi.
- Suomen kanssa samalla tasolla on yhteensä 12 muuta maata tai Taloudellista aluetta. Näistä 9 on maita.
- Nyt voit laskea Suomen sijoituksen kahdella tavalla, joista yksi on
  - - Sijoitus on 9-22 (kaikki mukana)
  - Ja toinen (mukana vain maat)
  - - Sijoitus on 5-15.

Tämä on liian vaikeaa medialle ja kansalle, vai onko?

Hyvä sijoitus kaikissa tapauksissa mikä on ollut laskussa.

Seuraava mikrodatan esimerkki koskee yhä ajankohtaisempaa aihetta eli suhtautumista kolmeen ihmisryhmään joihin ei ole suhtauduttu varauksettoman myönteisesti Euroopassa, eli Juutalaisiin, Muslimeihin ja Romaneihin.

Tässä kysymys on esitetty 2014-2015 European Social Surveyssä, ensi kertaa. Vastaajien määrät ovat melko saman suuruisia kussakin maassa, yleensä yli 1500. Siksi virhemarginaalien erot ovat varsin pieniä ja on helpompi tulkita tuloksia.

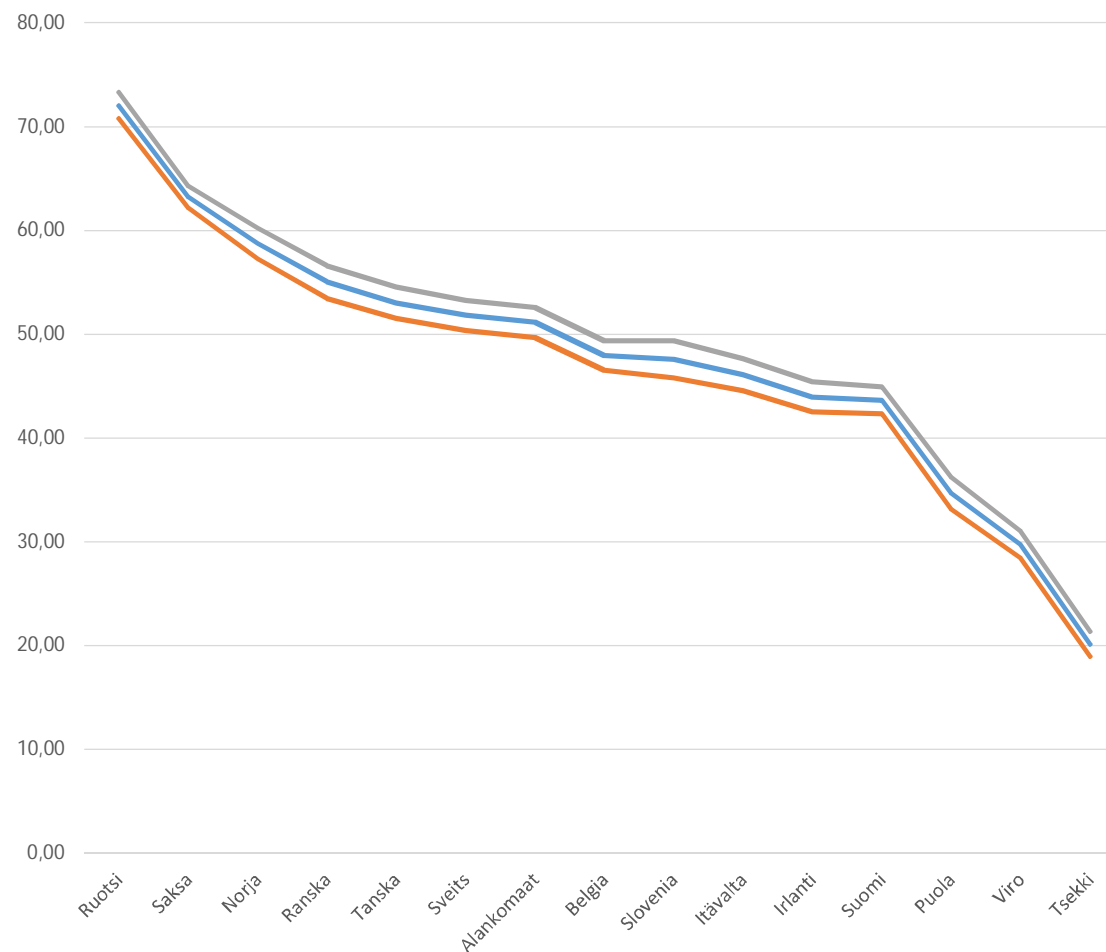
Ensimmäisessä kuviossa niiden ylä- ja alarajat ovat mukana, toisessa eivät koska kuvioista tulisi sotkuinen. Kun tiedät edellisestä luottamusvälit, niin kykenet melko hyvin tulkitsemaan erojen merkitsevyyden.

# Muslimimyönteisyys 15 Euroopan maassa 2014-2015

European Social Survey. Arvot: 0 = erittäin kielteinen, 100 = erittäin myönteinen)

Kiva puoli mediankin kannalta on, että erot ovat melkein aina merkitseviä lähimpien maiden välillä mutta Suomen ja Irlannin ero ei ole sitä.

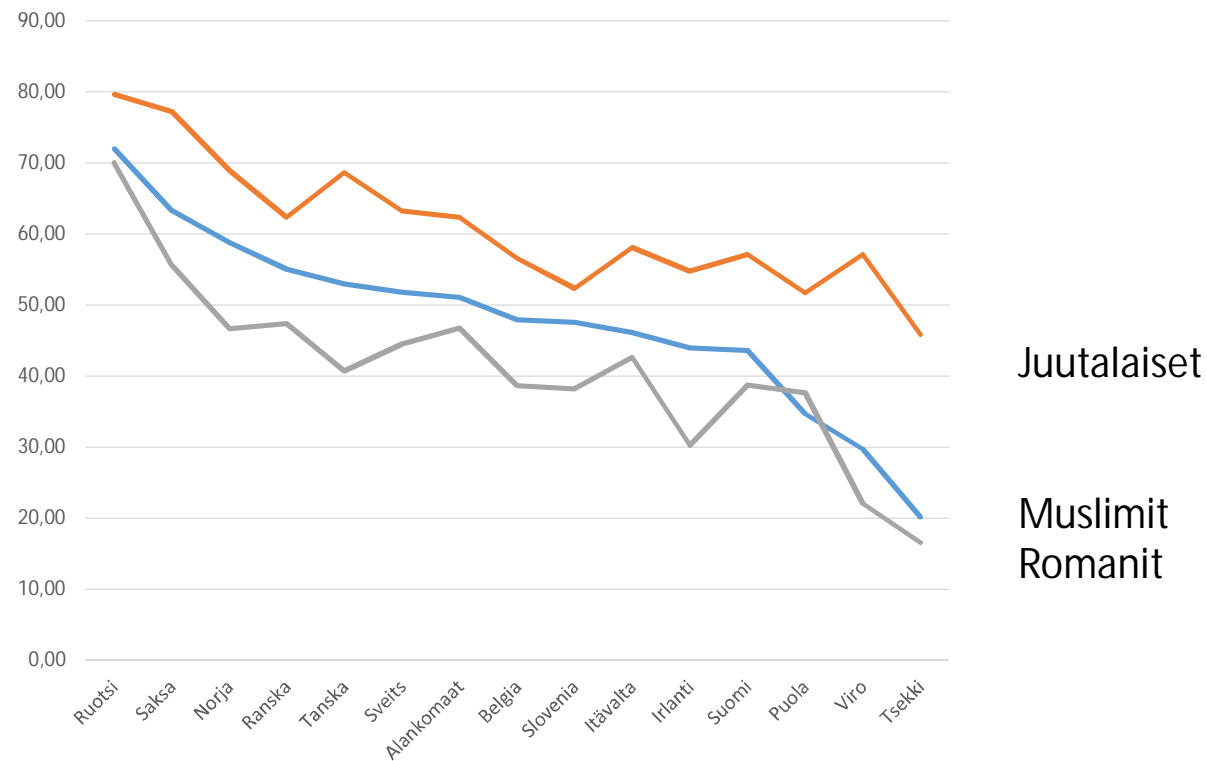
European Social Survey  
micro data



## Myönteisyys kolmea ihmisryhmää kohtaan 15 Euroopan maassa 2014-2015 European Social Survey

Juutalaisiin suhtaudutaan myönteisimmin kaikissa maissa, Romanit niukasti epäsuosituimpia mutta muutamissa maissa ei eroa.

Tiedot noin vuoden takaa. Onkohan myönteisyys laskenut?



# Keski-ikäiset valkoiset amerikkalaismiehet kuolevat kärsimykseen

Tiedemiehet havaitsivat yllättävän ilmiön. Keskiluokalla menee niin huonosti, että kuolleisuus kääntyi kasvuun. Tämä on ollut syksyn puheenaihe Yhdysvalloissa, kirjoittaa Saska Saarikoski Washingtonista.

SUNNUNTAI 22.11.2015 2:00 Päivitetty: 24.11.2015 14:52

Saska Saarikoski HELSINGIN SANOMAT

Tässä [sävähdyttävässä](#) esimerkissä kuolleisuus on laskettu yksinkertaisimmalla mahdollisella tavalla eli suoraan suhteuttamalla 100 000 asukkaaseen. Tämä pitäisi toimittajankin havaita ja ymmärtää. Seuraavilla sivuilla samoille valkoisille 'ei-hispaaneille' lasketut sarjat ikävakioituina USA:n alueilla. Tekijät mainittu. [En vastaa kummistakaan laskelmista.](#)

1. [Anne Case<sup>1</sup>](#) and
2. [Angus Deaton<sup>1</sup>](#)

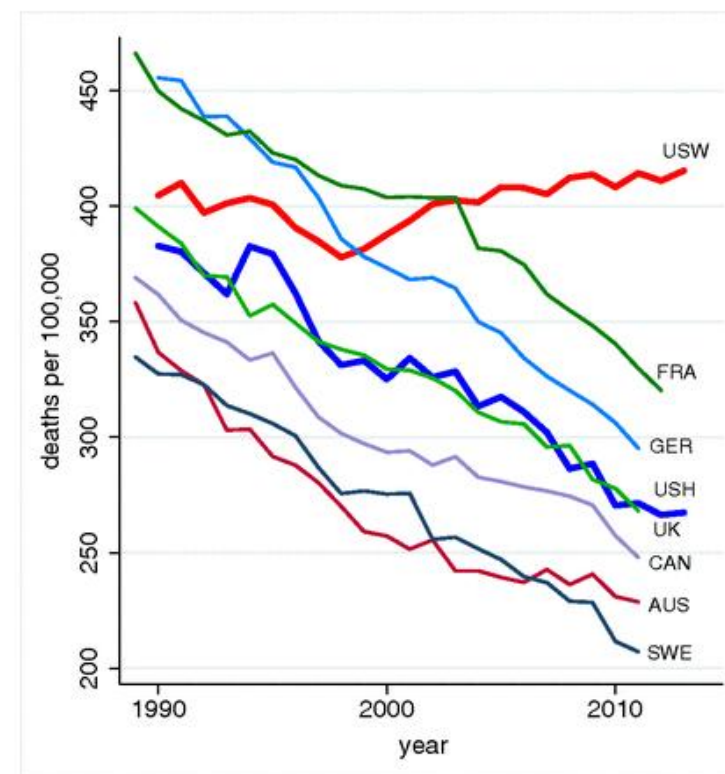


Fig. 1.

All-cause mortality, ages 45–54 for US White non-Hispanics (USW), US Hispanics (USH), and six comparison countries: France (FRA), Germany (GER), the United Kingdom (UK), Canada (CAN), Australia (AUS), and Sweden (SWE).

## Age-aggregation bias in mortality trends

Andrew Gelman ← † and Jonathan Auerbach

[Statistical Modeling, Causal Inference, and Social Science](#) 29.1.2016

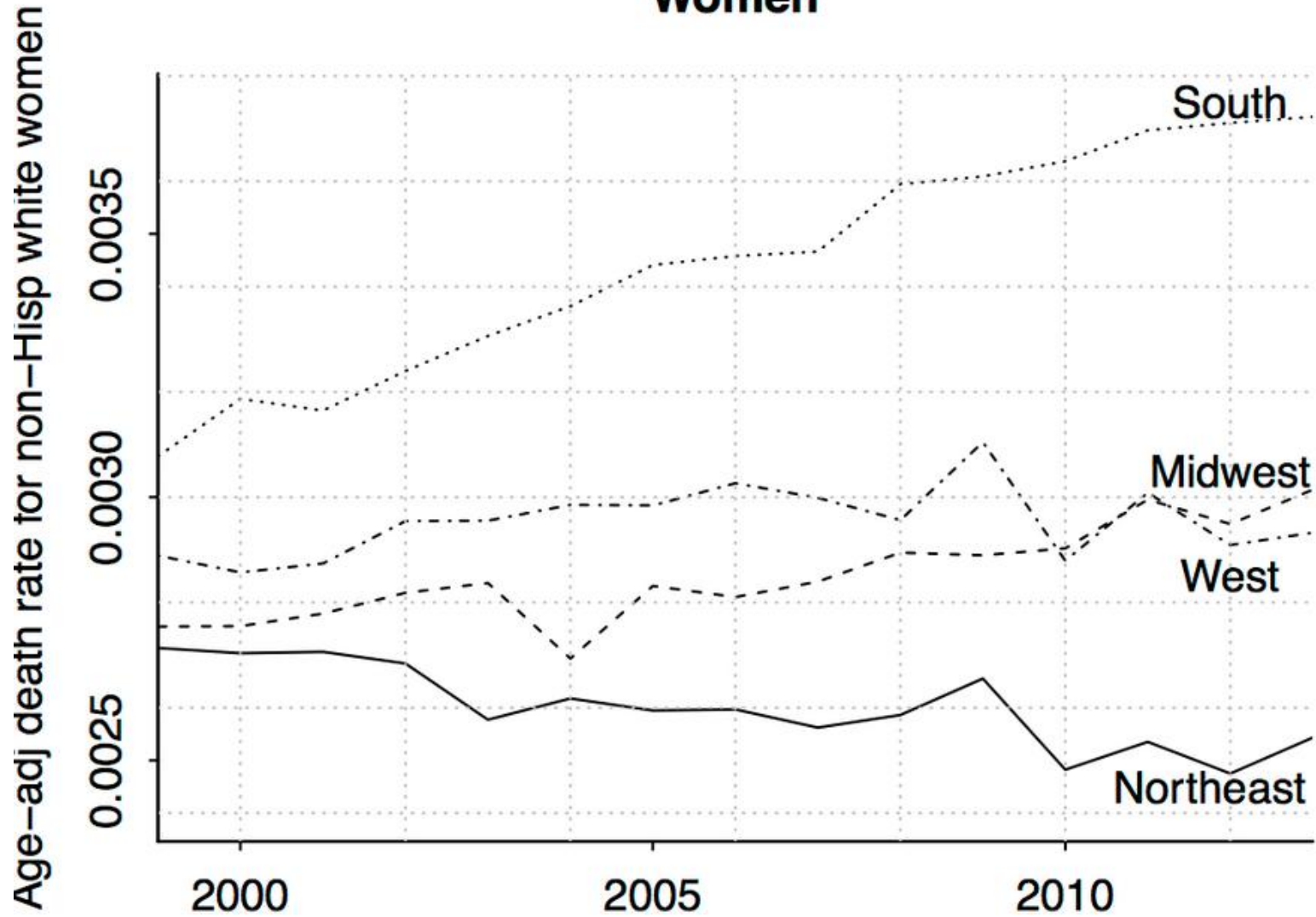
And here's the new analysis we did showing age-adjusted death rates for 45-54-year-old non-Hispanic white men and women.

Seuraavilla sivuilla kaksi kuviota tästä analyysistä.

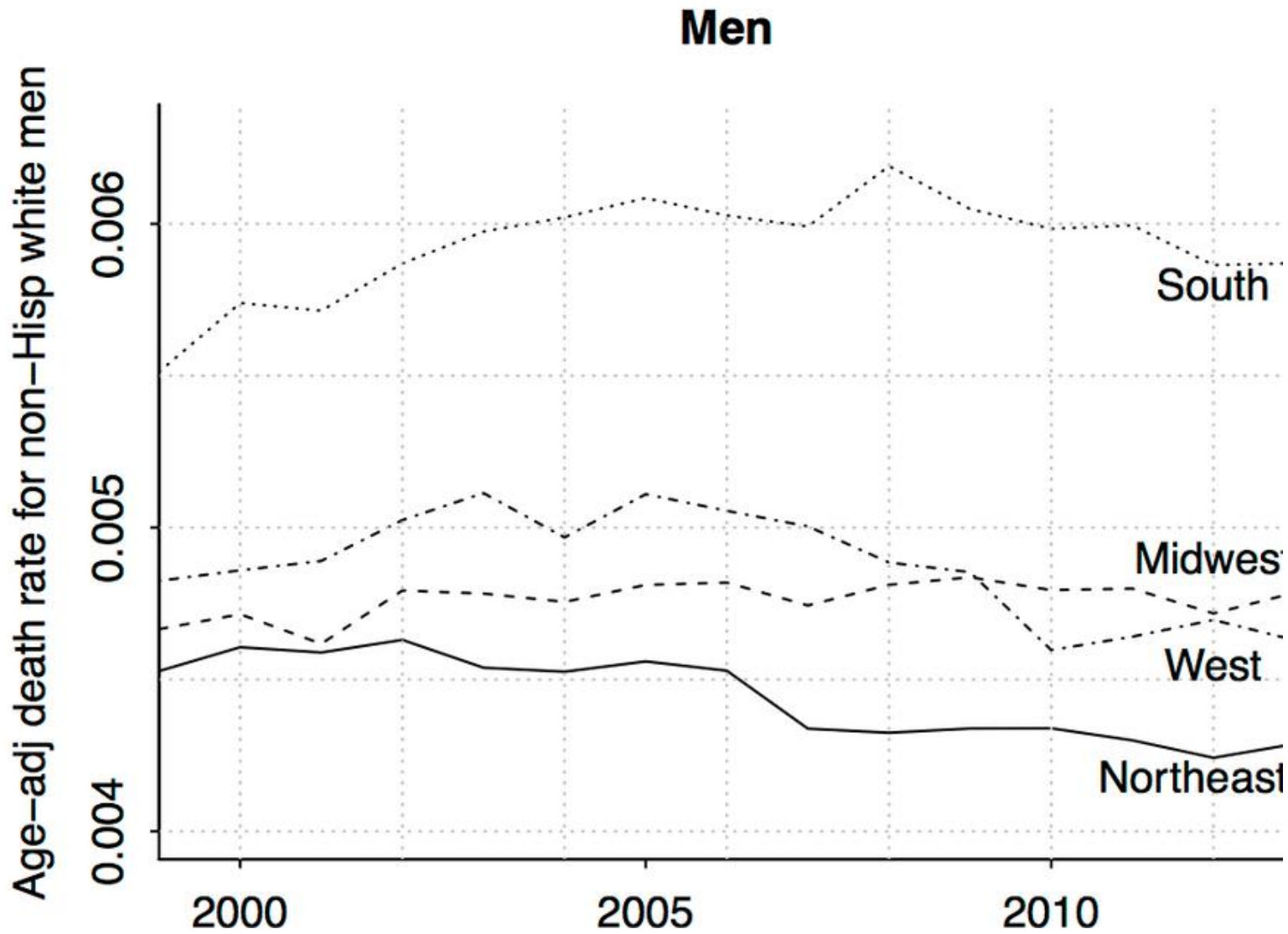
[http://andrewgelman.com/2016/01/19/death-trends-update-its-all-about-women-in-the-south/#.Vp6Jzo\\_\\_g8g.facebook](http://andrewgelman.com/2016/01/19/death-trends-update-its-all-about-women-in-the-south/#.Vp6Jzo__g8g.facebook)



# Women



Naisilla etelässä  
kuolleisuus  
lisääntynyt,  
muualla ei trendiä

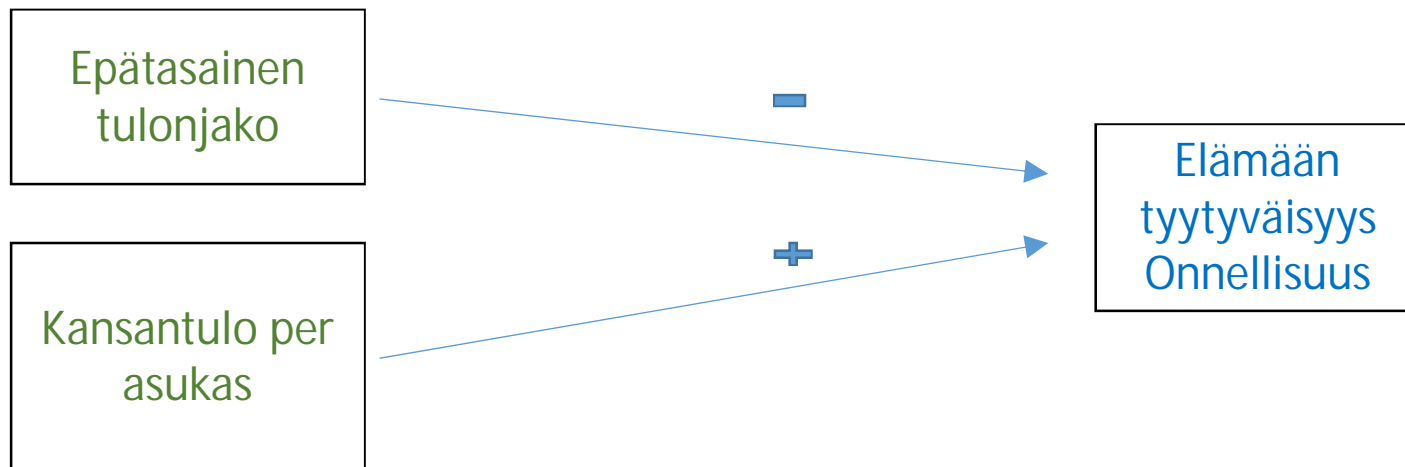


Miehillä etelässä kuolleisuus on ollut laskussa, muuallakin hieman.

Miesten luvut korkeampia kuten kaikkein kaikkialla

Lienette huomannut että en ole makrodatojen erityistutkija ja – suosija. Molempia tarvitaan, se on selvä. Median julkkistaloustieteilijät ovat enimmäkseen makroihmisiä jolloin voi samat yleisen tason 'totuudet' ilman konkreettisia lukuja, erityisesti mikrotason lukuja, esittää sujuvasti mihin tahansa tilanteeseen. Samalla tavalla esiintyvät politologit, usein Turusta. Juristit on kolmas median suosima ryhmä mutta heiltä ei lukufaktoja odotetakaan. Neljäs suosittu mediatiede on psykologia jonka edustajat taas upeasti yleistävät mikrotasolta asioita. Terveystieteilijöillä on ehkä samanlainen rooli. Monet tieteenalat eivät juuri mediassa esiinny. Lopuksi otan pienen esimerkin sekä makro- että mikrotasolta.

MAKRO: Aggregaatti eli laajalla maa-aineistolla on useissakin tutkimuksissa saatu kuvion kaksi tilastollisesti merkitsevää lineaarista yhteyttä. Pidän näitä karkeina tuloksina eli jos yhteyttä tutkittaisiin muitakin yhteyksiä ja kontrollimuuttujia käyttäen ja pienemmillä aggregaateilla (vaikkapa alueet, eivät maat), tulokset varmasti saisivat uutta syvyyttä. Esimerkiksi kansakunnan tulotason vaikutus onnellisuuteen ei ole ihan näin suoraviivainen. [Jätän teille tämän paremman tutkimisen.](#) Itse pysyn lähinnä mikrotasolla.



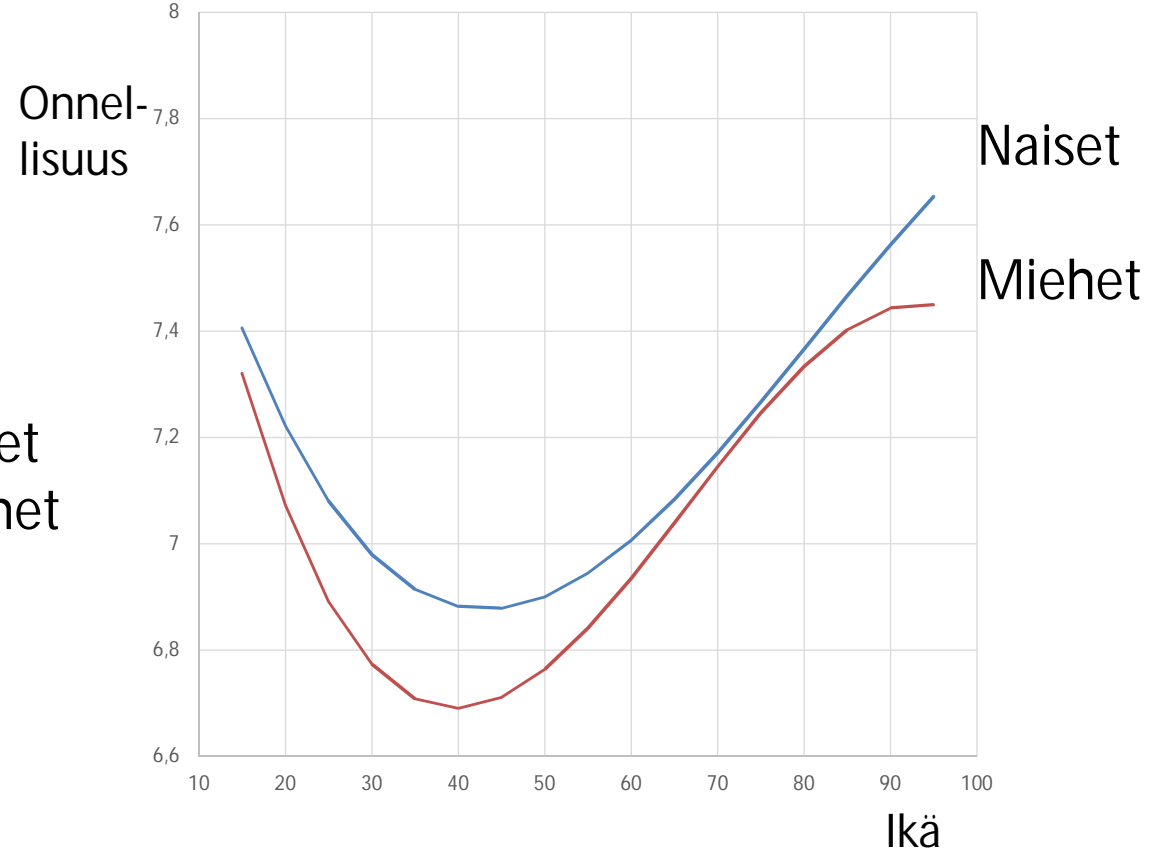
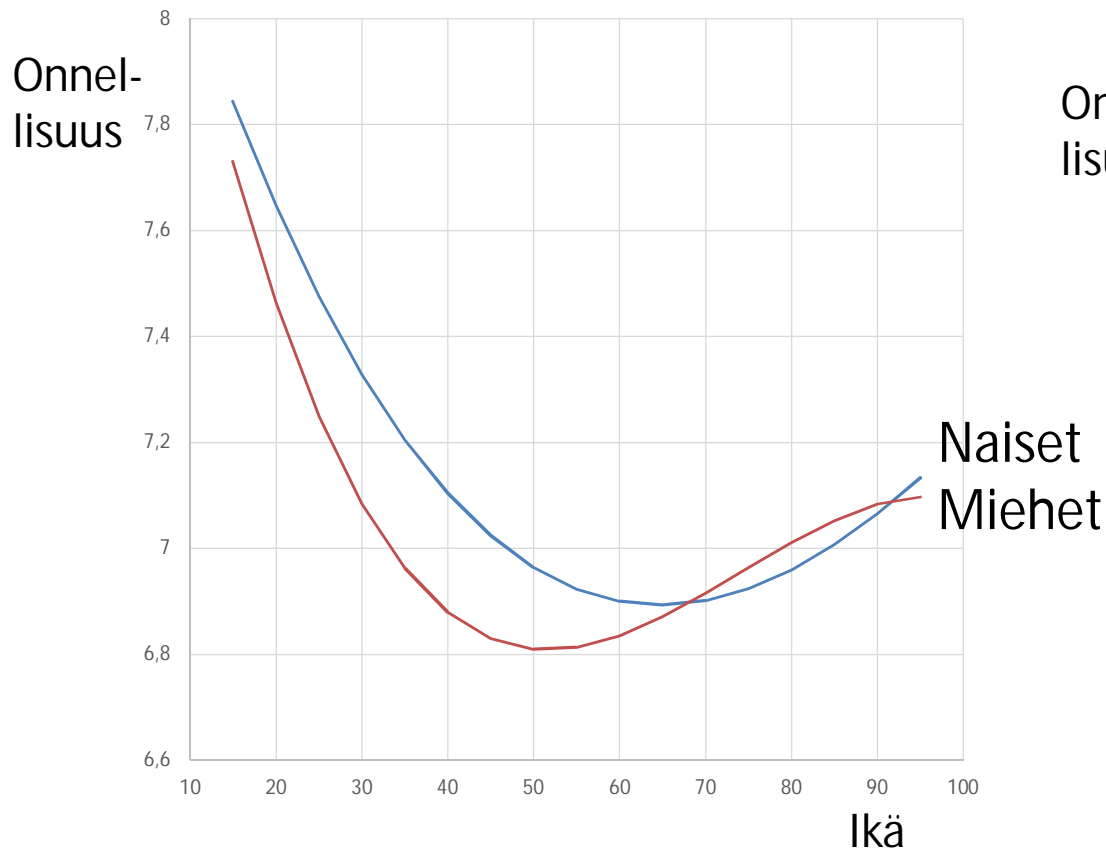
MIKRO: Lineaariset yhteydet aggregaatti- eli makroaineistolla ovat ehkä ainoa vaihtoehto soveltaa koska maa-aineistotkin ovat vain satojen suuruusluokkaa, mutta mikroaineistossa (ihmiset, yritykset) lineaarisuus huonommin toimii vaikka liikaa sitä käytetään. Eli pitää ensin tutkia minkälainen yhteys on. Otan loppuun esimerkin [Mikrotason Onnellisuuden mittaamisesta iän mukaan eli 15 vuotta täyttäneiden joukossa](#). Aineisto on noin 50 000 ihmistä 30 ESS-maasta (Venäjä, Israel ja Turkki ovat ESS:n Eurooppaa), vuosilta 2008-2012.

ESS= European Social Survey

Onnellisuus sukupuolittain iän mukaan 30 Euroopan maassa, kahdella epälineaarisella mallilla ( Asteikko [0, 10] ). Käyrät ovat ehdollisia kontroleihin

Kontrollit: Tulot (14 kategoriaa), Koulutus (6)

Täällä lisäksi Terveys (6)





Kukat ja Tiedekin tuovat onnellisuutta  
Onnellista oloa kaikille  
Kommentoikaa tarvittaessa sähköpostitse  
Kiitos

Tiede, Tilastot ja Media, Säätytalo 8.2.2016 Seppo Laaksonen