

$P(X) >$
 $f(x) = x \sigma$
 $Z = \frac{x - \mu}{\sigma}$
 π

Kansi: Hanna Svennevig

TILASTOLLISEN PÄÄTELYN KÄYTÄNTÖ *Tilastotiedettä soveltajille*

Pentti Manninen & Matti Ylén

TILASTOLLISEN PÄÄTELYN KÄYTÄNTÖ

Tilastotiedettä soveltajille

Pentti Manninen
Matti Ylén

$P(X) >$
 $Z = \frac{x - \mu}{\sigma}$
 π

Sisältö

1	JOHDANTO	1
1.1	Tilastomenetelmät luotettavan tutkimuksen perustana	1
1.1.1	Otos vs. näyte	1
1.1.2	Tilastollinen päättely ja tieteellisyden kriteerit	2
2	TEOREETTINEN JAKAUMA	3
2.1	Satunnaismuuttuja	3
2.1.1	Merkintätavat	3
2.2	Tilastollinen päättely ja teoreettinen jakauma	3
2.3	Empiirinen jakauma vs. teoreettinen jakauma	4
2.4	Teoreettisen jakauman peruskäsitteet	6
2.4.1	Todennäköisyys (p)	6
2.4.1.1	Satunnaismalli vs. deterministinen malli	6
2.4.1.2	Klassinen todennäköisyys	7
2.4.2	Tiheysfunktio $f(X)$	7
2.4.2.1	Diskreetti vs. jatkuva jakauma	8
2.4.2.2	Tasajakauma	8
2.4.2.3	Jatkuvan tasajakauman tiheysfunktio	9
2.4.2.4	Todennäköisyyksien määrääminen tiheysfunktion avulla	10
2.4.2.5	Tasajakauman odotusarvo ja varianssi	11
2.4.3	Kertymäfunktio $F(x)$	11
2.4.3.1	Tasajakauman kertymäfunktio	13
2.4.3.2	Kertymäfunktion ominaisuuksia	13
2.4.4	Diskreetti tasajakauma	15
2.5	Teoreettiset tunnusluvut	16
2.5.1	Odotusarvo ja varianssi	16
2.5.2	Muita tunnuslukuja	17
2.6	Normaalijakauma	20
2.6.1	Keskeinen raja-arvolause (central limit theorem)	20
2.6.2	Muuttujatyypit jotka noudattavat likimäärin normaalijakaumaa	20
2.6.2.1	Virhemuuttujat	20
2.6.2.2	Ominaisuudet	21
2.6.2.3	Otoskeskiarvon otantajakauma	21
2.6.3	Normaalijakauman tiheysfunktio	21
2.6.4	Normaalijakauman ominaisuuksista	23

2.6.5	Parametrien μ ja σ^2 vaikutus tiheysfunktioon	24
2.6.6	Taulukon korvaaminen Excelin funktioilla	24
2.6.7	Todennäköisyyksien määrittäminen (Z-jakaumassa)	25
3	LINEAARINEN MUUNNOS	27
3.1	Yleistä	27
3.1.1	Lineaarisen muunnoksen vaikutus muuttujan keskiarvoon ja varianssiin	27
3.2	Muuttujan standardoiminen	28
3.3	Lineaarisen muunnoksen vaikutus muuttujan fraktiilin määrittämiseen	30
3.3.1	Lineaarisen muunnoksen vaikutus muuttujan todennäköisyyksien määrittämiseen	30
4	OTANTAJAKAUMA	32
4.1	Yleistä	32
4.2	Otoskeskiarvon otantajakauma	32
4.3	Prosenttiluvun otantajakauma	34
4.4	Binomijakauma	36
4.4.1	Binomijakauma ja otosprosenttiluvun otantajakauma	37
4.4.2	Binomijakauman pistetodennäköisyysfunktio	37
4.4.3	Binomijakauman muodon riippuvuus parametreista n ja p	39
4.4.4	Binomijakauman kertymäfunktio	39
4.4.5	Binomijakauman odotusarvo ja varianssi	40
4.4.6	Binomi- ja normaalijakauman välinen yhteys	41
4.4.7	Excelin funktiot, jotka liittyvät binomijakumaan	41
4.5	t-jakauma	42
4.5.1	t-jakauman yhteydet muihin jakaumiin	42
4.5.2	t-jakauman fraktiilit	43
4.5.3	t-jakauman tiheysfunktio	44
4.5.4	Excelin t-jakaumaan liittyvät funktiot	44
4.6	χ^2-jakauma — eli otosvarianssin s^2 otantajakauma	45
4.6.1	χ^2 -jakauman tiheysfunktio	45
4.6.2	χ^2 -jakauman odotusarvo ja varianssi	46
4.6.3	χ^2 -jakauman kriittiset arvot	47
4.6.4	Otosvarianssin s^2 ja χ^2 -jakauman välinen yhteys	48
4.6.5	Excelin funktiot, jotka liittyvät χ^2 -jakaumaan	49

4.7	F-jakauma — eli varianssitestien otantajakauma	50
4.7.1	F-jakauman kriittiset arvot	51
4.7.2	Excelin funktiot, jotka liittyvät F-jakaumaan	53
5	ESTIMOINTI	54
5.1	Johdanto	54
5.2	Piste-estimointi	54
5.2.1	Estimaattorin ”hyvyys”	54
5.3	Väliestimointi	56
6	LUOTTAMUSVÄLIEN LASKEMINEN	57
6.1	Luottamusvälin määrittäminen populaation keskiarvolle μ, kun σ^2 tunnetaan	57
6.2	Populaation keskiarvon μ luottamusvälin määrittäminen, kun σ^2 on tuntematon	59
6.3	Luottamusvälin määrittäminen populaation prosenttiluvulle θ	61
6.4	Luottamusvälin määrittäminen populaation varianssille σ^2	62
6.5	Luottamusvälin määrittäminen populaation korrelaatiokertoimelle ρ	63
7	TILASTOLLINEN HYPOTEESIN TESTAUS	66
7.1	Tilastollisen testauksen lähtökohta	66
7.2	Hypoteesit	66
7.2.1	Esimerkkejä tilastollisista hypoteeseista	67
7.3	Tilastollisen testin mahdolliset virheet	68
7.4	Testisuure	69
7.5	Kriittinen alue	70
7.6	Esimerkki tilastollisesta testauksesta	71
7.7	Tilastollinen testaus käytännössä	72
7.8	Testin oletukset eli vaatimukset	73

8	TILASTOLLISET TESTIT	73
8.1	Testien ryhmittely	73
8.2	Tixelin testipohja	76
8.3	Prosenttilukutestit	77
8.3.1	Yhden otoksen prosenttilukutesti	77
8.3.1.1	Pieni otoskoko	79
8.3.2	Kahden riippumattoman otoksen prosenttilukutesti	80
8.3.3	Kahden riippuvan otoksen prosenttilukutesti	81
8.4	Keskiarvotestit	84
8.4.1	Yhden otoksen keskiarvotesti	84
8.4.2	Kahden riippumattoman otoksen keskiarvotesti	86
8.4.3	Perusjoukon varianssit tunnetaan	87
8.5	Kahden riippuvan otoksen keskiarvotesti	88
8.5.1	Erotusmuuttujan keskihajonta tunnuslukujen avulla	89
8.5.2	Excelin keskiarvotestit	90
8.5.3	Yksisuuntainen varianssianalyysi	91
8.6	Riippuvuuslukutestit	91
8.6.1	Korrelaatiokerroin	91
8.6.2	Korrelaatiokertoimen kriittiset arvot	93
8.6.3	Korrelaatiokerroin nollihypoteesin mukaan eri suuri kuin nolla	94
8.6.4	Kahden riippumattoman korrelaatiokertoimen yhtäsuuruuden testaus	95
8.6.5	χ^2 -riippumattomuustesti	96
8.6.5.1	Ristiintaulukon riippuvuusluvut	99
8.7	χ^2 -yhteensopivuustesti	100
8.8	Kahden varianssin yhtäsuuruuden testi	101
9	TESTIN VOIMAKKUUSFUNKTIO (POWER)	105
9.1	Hyväksymisvirheen todennäköisyys β	105
9.2	Testin voimakkuuteen vaikuttavia tekijöitä	108
9.2.1	α :n vaikutus testin voimakkuuteen	108
9.2.2	μ :n arvo	109
9.2.3	Otoskoko (n)	110
9.2.4	σ^2 :n arvo	111
9.2.5	Käytetty testityyppi	112
9.3	Testin voimakkuusfunktion arvon laskeminen tietyille vaihtoehdoiselle hypoteesille	112

10	OTOSKOON MÄÄRÄÄMINEN	114
10.1	Otoskoon määrääminen testin voimakkuusfunktioita ($1-\beta$) käyttäen	114
10.1.1	Keskiarvo	115
10.1.2	Prosenttiluku	117
10.2	Otoskoon määrääminen luottamusvälin avulla	118
10.2.1	Keskiarvo	118
10.2.2	Prosenttiluku	119
10.2.3	Korrelaatiokerroin	121
11	OTANTAMENETELMÄT	123
11.1	Otantamenetelmien käytön perusta	123
11.2	Otantamenetelmät	124
11.3	Yksinkertainen satunnaisotanta	124
11.4	Systemaattinen otanta	126
11.5	Estimointi yksinkertaisessa satunnaisotannassa	127
11.5.1	Luottamusväli	128
11.6	Ositettu otanta	128
11.6.1	Estimointi ositetussa otannassa	129
11.6.2	Otoksen kiintiöinti	130
11.7	Ryväotanta	132
11.8	Otoskoon määrääminen äärellisen perusjoukon tapauksessa	132
12	LIITTEET	134

LIITE 1: Standardoidun normaalijakauman kertymäfunktion arvoja

LIITE 2: t-jakauman kriittisiä arvoja

LIITE 3: Binomijakauman kertymäfunktion arvoja

LIITE 4: Binomijakauman kriittisiä arvoja

LIITE 5: χ^2 -jakauman kriittisiä arvoja

LIITE 6: F-jakauman kriittisiä arvoja

LIITE 7: Pearsonin korrelaatiokertoimen kriittisten arvojen itseisarvoja

LIITE 8: Asteikkovälin lineaarinen muunnos

1 Johdanto

1.1 Tilastomenetelmät luotettavan tutkimuksen perustana

Tilastotieteelliset menetelmät koostuvat erilaisista tilastoaineiston *keräämis-*, *käsittely-*, *analysointi* ja *päätelymenetelmistä*. Näillä kaikilla on merkityksensä luotettavan tilastollisen tutkimuksen aikaansaamiseksi. Tässä kirjassa käsitellään erityisesti tilastollista testausta ja päätöksentekoa. Tällöin olemme tekemisissä etenkin päätelymenetelmien ja niihin liittyvän todennäköisyysmatematiikan kanssa. Ts. tarkastelemme niitä tilastotieteen asettamia periaatteita ja kriteereitä, joiden avulla voimme päättää, millaisia johtopäätöksiä olemme oikeutettuja tekemään tilastoaineiston pohjalta. Tilastollinen päätöksenteko voidaan jakaa kahteen osaan, *estimointiin* ja *hypoteesien testaukseen*. Estimoinnissa pyritään muodostamaan *estimaatteja* eli arvioita perusjoukon parametrien (tunnuslukujen) todellisille arvoille. Hypoteesien testauksessa on tavoitteena tehdä oikeita johtopäätöksiä ennakkoon asetettujen tilastollisten hypoteesien paikkansapitävyydestä. Estimointia ja tilastollisten hypoteesien testausta käsitellään myöhemmin erillisissä luvuissaan.

Toisaalta on todettava, että muutkin menetelmät vaikuttavat siihen millaisiin johtopäätöksiin tilastoaineisto meidät oikeuttaa. Etenkin aineiston keräämismenetelmän vaikutus on tässä suhteessa huomattava. Keräämismenetelmät muodostuvat etupäässä erilaisista otantamenetelmistä, joita esitellään tarkemmin omassa luvussaan. On kuitenkin syytä jo alustavasti tässä yhteydessä painottaa joitakin tähän liittyviä seikkoja. Tilastomenetelmiä hyödyntävä tutkimus perustuu useimmiten erilaisista perusjoukoista poimittaviin otoksiin, jolloin otokset eivät edusta perusjoukkojaan täydellisesti. Tästä nimenomaisesta syystä tilastolliseen tutkimukseen (eli perusjoukon tunnuslukujen todellisten arvojen mittaamiseen) muodostuu virhettä, jonka suuruuteen ja hallittavuuteen vaikuttaa osaltaan tilastoaineiston keräämismenettely, ts. se miten otos poimitaan. Näin ollen tilastollisen tutkimuksen luotettavuuden eräs perustava lähtökohta on käytetyn otantamenetelmän pätevyys¹. Käytännössä tämä useimmiten tarkoittaa (tutkittavasta perusjoukosta riippuen) tietyllä tavalla toteutettua satunnaisotantaa.

1.1.1 Otos vs. näyte

Mikäli asianmukaisesti muodostettu satunnaistaminen ei otoksen poiminnassa toteudu, ei voida puhua otoksesta, vaan *näytteestä*. Näyte ei edusta perusjoukkoa yhtä luotettavasti kuin otos. Näytteeseen sisältyvät tilastoyksiköt ovat usein jollakin kontrolloimattomalla tavalla valikoituneita. Näin ollen tilastollisen päätelyn kannalta on syytä pitää mielessä, että näytteen pohjalta muodostetut ti-

¹ Muita keskeisiä tutkimuksen luotettavuuden edellytyksiä ovat käytetyn mittaustmenetelmän sopivuus ja tarkkuus, suotuisat mittaolosuhteet, mittajan pätevyys jne.

lastolliset johtopäätökset ovat usein otokseen verrattuna huomattavasti epäluotettavampia. Näytteestä on esimerkiksi kysymys silloin, kun graduaan tekevä opiskelija jakaa tutkimuslomakkeitaan jollain massaluennolla ja pyytää lomakkeen ottaneita palauttamaan sen täytettynä tiettyyn paikkaan. Pätevästi tehdyn otospohjaisen tutkimuksen yhteydessä voidaan tulosten tarkkuus arvioida suhteellisen luotettavasti. Käytännössä tämä tarkoittaa usein sitä, että kun voidaan olettaa otoksen edustavan hyvin perusjoukkoa, on mahdollista määrätä havaittujen otostunnuslukujen (esim. keskiarvojen) keskivirheet. Keskivirheen avulla saadaan arvioita otantavirheen suuruudesta eli siitä kuinka paljon otostunnusluvun arvo mahdollisesti poikkeaa tutkittavasta populaation tunnusluvusta.

Vaikkakin käytännön tutkimuksen yhteydessä otantamenetelmien tinkimätön käyttö onkin joskus problemaattista, on kuitenkin varsin oleellista tuoda esiin otantamenetelmien merkitys luotettavan tilastollisen tutkimuksen aikaansaamiseksi. Tutkimuksen näytepohjaisuus on mm. jo sinällään yksi merkittävä syy siihen, miksi samaa aihetta koskevat eri tutkimuksen päätyvät usein huomattavastikin erilaisiin tai keskenään ristiriitaisiin tutkimustuloksiin.

1.1.2 Tilastollinen päättely ja tieteellisyyden kriteerit

Kuten yllä on pyritty esittämään, oleellinen kysymys tilastollisessa päättelyssä on tietenkin tutkittavaa asiaa koskevan arvion tai päätöksen luotettavuus / varmuus. Niin tieteellisten, mielipide- kuin taloustutkimustenkin keskeinen lähtökohta on se, että kyettäisiin saamaan mahdollisimman luotettavia arvioita asioiden vallitsevasta tilasta, asioiden vaikuttavuudesta johonkin, asioiden muutoksista sekä arvioita asioiden keskinäisistä suhteista ja eroista. Jotta tutkimus olisi luotettavaa, tarvitaan tieteellisiä kriteereitä, joiden mukaan tutkimus olisi suoritettava. Mikäli tutkimus on luonteeltaan tilastollinen, muodostaa tilastollinen päätöksenteko oleellisen osan näistä luotettavuuden kriteereistä.

Yhteenvedona voidaan esittää, että tämän kirjan keskeisellä aiheella, *tilastollisella päättelyllä*, tarkoitetaan tilastotieteellisten menetelmien avulla toteutettua käytäntöä, jossa tarkasteltavasta tilastoaineistosta edetään ko. perusjoukkoa koskeviin (mahdollisimman luotettaviin) päätelmiin ja induktiivisiin yleistyksiin.

Tämän kirjan päätehtävänä on näin ollen johdattaa yhtäältä ymmärtämään sitä, mihin tilastollinen päätöksenteko perustuu ja toisaalta hallitsemaan sitä, mitä tilastollinen päätöksenteko käytännössä tarkoittaa.

2 Teoreettinen jakauma

2.1 Satunnaismuuttuja

Erilaisissa satunnaiskokeissa, kuten nopanheitossa, saatavia arvoja voidaan kutsua *satunnaismuuttujan* arvoiksi. Satunnaisilmiön tulokseen on siis liitetty jokin reaaliluku — nopanheiton tapauksessa kokonaisluku. Tällöin satunnaismuuttuja saa arvoja väliltä 1–6. Myös kaikki erityistieteellisessä tutkimuksessa käytetyt muuttujat, kuten esim. kyselylomakkeiden muuttujat, ovat luonteeltaan satunnaismuuttujia. Tällöin satunnaismuuttujan mitta-asteikko kuvaa tietyn empiirisen satunnaisilmiön alkeistapauksiin liitettyä arvojoukkoa. Oman ryhmänsä muodostavat otostunnusluvut, kuten esim. otoskeskiarvo ja –prosenttiluku, joiden arvot vaihtelevat otoksesta toiseen. Tällöin satunnaiskoe on otoksen poiminta. Satunnaismuuttujiin liittyy periaatteessa aina jokin teoreettinen jakauma, ts. todennäköisyysjakauma, joka kuvaa niiden arvojen esiintymisen todennäköisyyttä. Usein satunnaismuuttuja kirjoitetaan ilman ’satunnais’ –osaa tai sitä voidaan kutsua, kuten tässä kirjassa, empiiriseksi muuttujaksi.

2.1.1 Merkintätavat

Tässä kirjassa satunnaismuuttujaa merkitään kursivoituilla isoilla kirjaimilla X , Y , Z jne. ja niitä vastaavia satunnaismuuttujan arvoja pienillä kirjaimilla x , y , z jne. Matemaattisia muuttujia merkitään erotukseksi puolestaan isoilla kursivoimattomilla kirjaimilla X , Y , Z jne. Matemaattiset muuttujat eivät liity mihinkään satunnaisilmiöön ja niitä käytetään sen tähden matemaattisesti tarkoissa lausekkeissa, kuten esim. tiheysfunktioiden lausekkeissa, joihin ei voi sisältyä satunnaismuuttujia.

Otostunnuslukuja merkitään kursivoimattomilla kirjaimilla asiayhteydestä selviävällä tavalla. Teoreettisia tunnuslukuja sekä eräitä muita matemaattisia käsitteitä merkitään kursivoimattomilla kreikkalaisilla aakkosilla niin ikään asiayhteydestä tarkemmin selviävällä tavalla.

2.2 Tilastollinen päättely ja teoreettinen jakauma

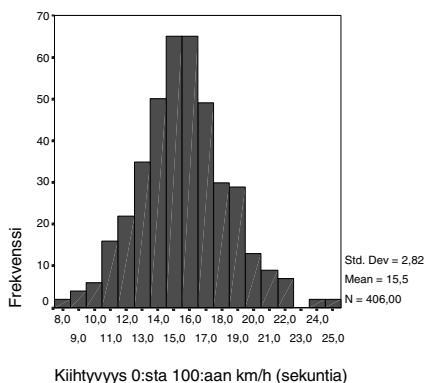
Tässä kirjassa käsitellään tilastollisen päättelyn teoriaa ja käytäntöä lähinnä parametristen testien kannalta. Toisin kuin non-parametriset testit, parametriset testit perustuvat tiettyihin teoreettisiin jakaumiin sekä perusjoukon parametrien todellisia arvoja estimoiviin otostunnuslukuihin eli estimaattoreihin. Myös näiden estimaattoreiden satunnaiskäyttäytyminen hallitaan teoreettisten jakaumien avulla. Tästä syystä on välttämätöntä käsitellä joitakin keskeisimpiä teoreettisia jakaumia, joita käytetään empiiristen muuttujien tai johdettujen tilastollisten testisuureiden matemaattisina malleina.

Empiiriset muuttujat eivät kuitenkaan noudata tarkasti vastaavia teoreettisia jakaumia, mutta niitä voidaan käyttää tästä likimääräisyydestä huolimatta. Likimääräisyyden ajatellaan johtuvan satunnaisvaihtelusta, jonka lähteinä nähdään erilaisia toisistaan riippumattomia tekijöitä, kuten otanta- ja mittausvirhe. Parametristen testien yhteydessä satunnaismuuttujien oletetaan perusjoukon tasolla noudattavat varsin tarkasti niitä mallintavia teoreettisia jakaumia. Helposti voisi ajatella, että koska empiiriset ilmiöt ovat niin moninaisia, tarvitaan suuri määrä erilaisia niitä mallintavia teoreettisia jakaumia. Asia ei kuitenkaan ole niin. Yleisesti tullaan toimeen varsin pienellä eri teoreettisten jakaumien määrällä, kuten jatkossa tulemme huomamaan. Oleellista on kiinnittää huomiota eri teoreettisten jakaumien yleisiin ominaisuuksiin ja käyttömahdollisuuksiin.

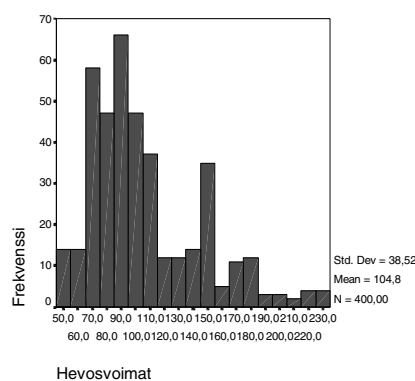
2.3 Empiirinen jakauma vs. teoreettinen jakauma

Empiirinen jakauma liittyy konkreettiseen tilastoyksiköiden joukkoon; se kuvaa tietyn empiirisesti mitattavan ominaisuuden (so. satunnaismuuttujan) arvon esiintymistiheydessä havaittavaa vaihtelua. Tarkastellaan seuraavaksi esimerkkien avulla empiiristä jakaumaa. Kuva 1 esittää eri hinta- ja teholuokkaan kuuluvien autojen (0:sta 100:aan) kiihtyvyydejakaumaa, joka on jokseenkin normaalijakautunut. Vastaavasti Kuva 2 esittää samojen autojen hevosvoimajakautumaa, joka on selvästi positiivisesti vino.

Kuva 1: Autojen kiihtyvyydejakauma (likimäärin normaalijakauma)



Kuva 2: Autojen hevosvoimajakautuma (positiivisesti vino jakauma)

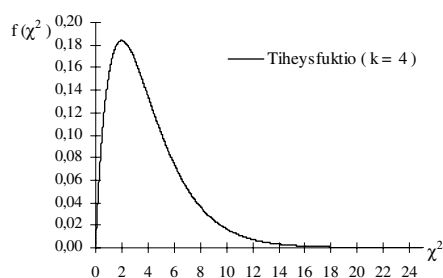


Ero keskenään melko voimakkaasti korreloivien jakaumien muodossa selittynee sillä, että kiihtyvyyteen vaikuttaa useampia tekijöitä kuin hevosvoimien suuruuteen. Yksi oleellinen tekijä on auton paino, joka ei vähennä moottorin tehoa, mutta hidastaa kiihtyvyyttä. Toinen tekijä on auton muodon aiheuttama ilmanvastus jne. Voidaan näin ollen todeta, että mitattavan empiirisen ominaisuuden jakauma on taipuvaisempi olemaan sitä todennäköisemmin normaalisti jakautu-

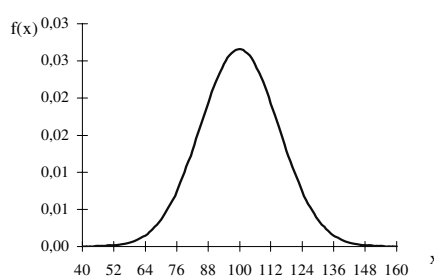
nut mitä useamman tekijän määräämästä ilmiöstä siinä on kysymys — asia, johon palataan normaalijakauman teoreettisen perustelun yhteydessä.

Teoreettinen jakauma on luonteeltaan yleisempi; se ei liity tiettyihin konkreettisiin tilastoyksiköihin, vaan kuvaa empiirisen tai matemaattisesti johdetun satunnaismuuttujan (esim. testisuureen) arvojen esiintymistiheyden todennäköisyyksiä. Tästä syystä teoreettista jakaumaa voidaan kutsua myös todennäköisyysjakaumaksi. Teoreettisia jakaumia käytetään testisuureiden ja empiiristen jakaumien matemaattisina malleina. Kun tunnetaan ko. empiirisen ilmiön tai tilastollisen testisuureen noudattama teoreettinen jakauma, mahdollistaa tämä niihin liittyvien arvojen esiintymistodennäköisyyksien määrittämisen.

Kuva 3: Esimerkki teoreettisesta jakaumasta — χ^2 -jakauma vapausasteilla 4



Kuva 4: Esimerkki teoreettisesta jakaumasta — Normaalijakauma (muuttujan X keskiarvo = 100 ja varianssi = 15²)



Kuva 3:ssa on kysymys teoreettisesta jakaumasta, χ^2 -jakaumasta, jota käytetään mm. empiirisiä ilmiöitä kuvaavien muuttujien otosvarianssin satunnaisvaihtelun mallintamiseen. Kuva 4 puolestaan esittää teoreettista normaalijakaumaa, joka tässä esimerkissä mallintaa Wechslerin aikuisille tarkoitetun älykkyystestin kokonaispistemäärän jakaumaa perusjoukon keskiarvon ollessa 100 ja varianssin 15².

Kuten jo yllä totesimmekin, tilastollinen testaus ja päätöksenteko perustuu teoreettisiin jakaumiin. Oleellinen ero empiiristen muuttujien mallintamiseen nähden on se, että tilastollisessa testauksessa ja päätöksenteossa tiettyä teoreettista jakaumaa (todennäköisyysjakaumaa) noudattavat matemaattisesti johdetut testisuureet empiirisiä ilmiöitä kuvaavien muuttujien sijasta. Tällöin mallinnetaan tietyn tilastollisen testisuureen satunnaiskäyttäytymistä, testisuureen, jonka arvoon perustuu tilastollinen päätöksenteko — asia, johon palataan tarkemmin hypoteesien testausta ja testisuureita käsittelevissä luvuissa 7, s. 66 ja 8, s. 73.

2.4 Teoreettisen jakauman peruskäsitteet

Määriteltäessä teoreettisiin jakaumiin liittyviä peruskäsitteitä, on hyödyllistä todeta samalla, että niillä on vastineensa empiiristen jakaumien puolella. Taulukko 1 esittää nämä käsitteet rinnasteisesti.

Taulukko 1: Teoreettisen ja empiirisen jakauman peruskäsitteiden välinen vastavuus

Empiirinen jakauma	Teoreettinen jakauma
• Suhteellinen frekvenssi tai %	• Todennäköisyys (p)
• (Prosenttinen) frekvenssijakauma	• Tiheysfunktio (f)
• (Prosenttinen) summajakauma	• Kertymäfunktio (F)

Kuten Taulukko 1 osoittaa, vastaa suhteellista frekvenssiä, frekvenssijakaumaa ja summa- eli kumulatiivista jakaumaa vastaavat teoreettiset käsitteet todennäköisyys, tiheysfunktio ja kertymäfunktio. Nämä teoreettiset peruskäsitteet ovat hyvin yleisiä — ts. ne ovat yhteisiä periaatteessa miltei kaikille teoreettisille jakaumille. Tarkastellaan seuraavaksi heuristisella tasolla hieman näiden käsitteiden luonnetta ja palataan niihin jatkossa sitten erityisemmin aina kunkin teoreettisen jakauman tarkastelun yhteydessä.

2.4.1 Todennäköisyys (p)

Tässä esityksessä ei tulla tarkastelemaan varsinaisesti alkeistodennäköisyyslaskennan teoreettisia perusteita, koska niillä on enimmäkseen merkitystä vain arpajais- ja peliteoriassa. Sen sijaan keskitytään todennäköisyyden arvojen laskeamiseen erilaisten todennäköisyysmallien avulla käytännön tutkimuksen vaatimissa tilanteissa. Koska tilastotieteelliset menetelmät perustuvat suuressa määrin todennäköisyyslaskentaan, on syytä kuitenkin luoda yleisluontoinen katsaus todennäköisyyden käsitteeseen.

2.4.1.1 Satunnaismalli vs. deterministinen malli

Todennäköisyyslaskennassa käsitellään *stokastisia* eli *satunnaismalleja* ja tutkitaan niiden säännöllisiä (hallittavia) ominaisuuksia. Satunnaismallien käyttöön täytyy turvautua lähinnä siksi, koska monien tieteenalojen tutkimat ilmiöt ja myös muutkin, kuten esim. elinkeinoelämän ilmiöt, eivät ole tarkasti ennustettavia eli kuvattavissa *determinististen mallien* avulla. Tunnettaessa tietyt ilmiöiden ja tapahtumien alkuehdot, determinististen mallien avulla voidaan ennustaa tarkasti näiden ilmiöiden ja tapahtumien lopputuloksia. Sen tähden niitä voidaan kuvata myös matemaattisten yhtälöiden avulla.

Silloin kun näitä tapahtumien ja ilmiöiden alkuehtoja ei tunneta, joudutaan käyttämään satunnaismalleja, jolloin pyritään määrittämään tapahtumien lopputuloksille tiettyjä esiintymistodennäköisyyksiä. Yksinkertaisena esimerkkinä satunnaistapahtumasta, jonka alkuehtoja ei tunneta, voidaan esittää nopanheitto. Monimutkaisempi satunnaismallin avulla kuvattava ilmiö on esim. vaikkapa sääilmiöiden määräytyminen.

2.4.1.2 Klassinen todennäköisyys

Kaikki tunnettuja ennusteita kuten esimerkiksi, että 20 % tietyn lääketieteellisen diagnoosin saaneista elää vielä kolmen vuoden kuluttua. Tällöin ajattelemme, että keskimäärin joka viiden potilas elää vielä kolmen vuoden kuluttua diagnoosin teon jälkeen. Tulkitsemme todennäköisyyden siis suhteellisen frekvenssin avulla, jossa suhteutamme tarkastelunalaisten esiintyneiden tapausten A lukumäärän f kaikkien tapauksen lukumäärään n eli $p = f(A)/n$. Tätä todennäköisyyden määritelmää kutsutaan ns. klassisen todennäköisyyden määritelmäksi. Näin ollen todennäköisyys, että diagnoosin saanut elää vielä kolmen vuoden kuluttua on $1/5 = 0,2$. Voimme siis todeta, että todennäköisyys voidaan esittää väliin $[0,1]$ kuuluvana suhdelukuna tai vastaavana prosenttilukuna.

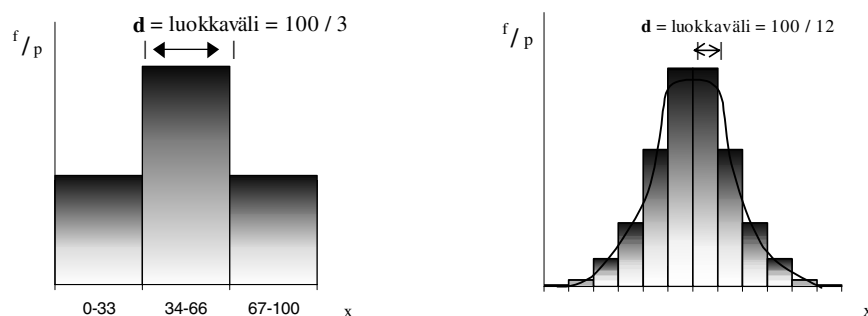
Liitettäessä todennäköisyyden käsite ylempänä esitettyyn stokastisten ts. satunnaismallien käsitteeseen, päädyimme lyhyesti toteamaan sen, miten erityistieteet soveltavat näitä käsitteitä etsiessään empiiristä tukea teoreettisille malleilleen. Erityistieteissä voidaan pyrkiä muotoilemaan tutkimusalaan koskevan kokemuksen kautta matemaattinen malli, tässä tapauksessa siis satunnaismalli, joka kuvaa tiettyihin ilmiöihin / tapahtumiin sisältyviä todennäköisyyksiä. Näiden matemaattisten mallien (=teoreettisten jakaumien) avulla pystymme sitten laskemaan haluttujen tapahtumien todennäköisyydet. Oleellista siis on, että näiden todennäköisyyksien tulkinta tapahtuu suhteellisen frekvenssin periaatteella ja että niiden laskemisen malli on muodostettu empiirisen kokemuksen ja teoreettisten oletusten avulla.

2.4.2 Tiheysfunktio $f(X)$

Kuvitelkaamme yleisesti, että voimme muuttujan frekvenssijakaumassa kasvat-
taa tilastoyksiköiden lukumäärää n rajatta ja antaa samalla luokkavälin d lä-
hestyä nollaa. Tällöin luokkavälin pienentyessä muuttujan frekvenssijakauman
muodossa tapahtuisi todennäköisesti huomattavaa muutosta.

Asian havainnollistamiseksi tarkastellaan Kuva 5:n erityistapausta, jossa muut-
tuja X vaihtelee välillä 0–100. Vasemmanpuoleisessa tilanteessa luokkavälin d
arvo on $100/3$ ja oikeanpuoleisessa se on pienentynyt neljänneksen verran eli
 $100/12$. Luokkavälin pienentyessä kohti nollaa huomaamme histogrammin
muodon lähestyvän jatkuvaa kuvauskäyrää.

Kuva 5: Luokkavälin d vaikutus muuttujan jakaumaan



Näin ollen voidaan ajatella, että on mahdollista johtaa käyrää kuvaava matemaattinen funktio asettamalla tiettyjä oletuksia muuttujan käyttäytymiselle. Näin muodostettua funktiota kutsutaan muuttujan X , jota merkitään $f(X)$:llä.

2.4.2.1 Diskreetti vs. jatkuva jakauma

Koska empiiriset muuttujat ja siten niiden jakaumat voivat olla joko *jatkuvia* tai *epäjatkuvia* eli *diskreettejä*, on tämä otettava myös teoreettisten jakaumien tiheysfunktioiden määrittelyssä huomioon. Jatkuvat muuttujat voivat tietyllä välillä saada, minkä tahansa reaalilukuarvon. Täten niiden tiheysfunktion kuvaajakin on jatkuva käyrä. Diskreetit muuttujat puolestaan voivat saada vain tiettyjä arvoja, esim. kokonaislukuarvoja. Jos teoreettinen jakauma ja sitä noudattava satunnaismuuttuja X on diskreetti, niin sen arvoihin x_1, x_2, \dots, x_n liittyy aina tietty esiintymistodennäköisyys p_1, p_2, \dots, p_n . Tällöin tiheysfunktion kuvaajaan ei ole jatkuva, vaan epäjatkuva, jolloin sitä kutsutaan *pistetodennäköisyysfunktioksi* ja merkitään p_i :llä (ks. Kuva 14, s. 15).

Yksinkertaisena esimerkkinä tapahtumasta, jota kuvaava muuttuja on diskreetti, on nopanheitto. Nopanheiton todennäköisyysjakaumaa siis kuvaa pistetodennäköisyysfunktio. Tarkemmin sanottuna kysymyksessä on tasajakaumaa noudattava diskreetti satunnaismuuttuja, joka vaihtelee välillä 1–6, jossa jokaiseen arvoon liittyy yhtä suuri esiintymistodennäköisyys $p_i = 1/6$.

2.4.2.2 Tasajakauma

Eräs teoreettisista jakaumista on jatkuva (tai diskreetti) tasajakauma. Tasajakauman merkitys on käytännössä yleensä melko vähäinen, mutta koska se on kuitenkin yksinkertaisuudessaan teoreettisen jakauman peruskäsitteitä hyvin havainnollistava, tarkastellaan sen avulla teoreettisten jakaumien perusominaisuuksia, joista tässä alaluvussa ensimmäisenä käsitellään tiheysfunktioita. Jos satunnaismuuttuja X noudattaa tasajakaumaa, merkitään: $X \sim \text{Tas}[a,b]$.

2.4.2.3 Jatkuvan tasajakauman tiheysfunktio

Jatkuvan tasajakauman tiheysfunktio määritellään seuraavasti:

$$f(X) = \begin{cases} \frac{1}{b-a}, & \text{kun } X \in (a \leq X \leq b) \\ 0 & \text{muualla} \end{cases}$$

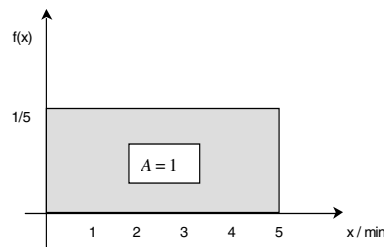
Näin ollen tasajakaumaa noudattavan satunnaismuuttujan X kaikilla arvoilla x , jotka kuuluvat välille $(a \leq X \leq b)$, on yhtä suuri esiintymistodennäköisyys.

Tarkastellaan jatkuvan tasajakauman tiheysfunktioita lähemmin pienen esimerkin valossa. Olkoon julkinen liikenne järjestetty tiettyyn esikaupunkiin siten, että bussit lähtevät pääte pysäkeiltä viiden minuutin väliajoin ilman erityisempiä aikatauluja. Matkustajat eivät tiedä bussien lähtöaikaa, vaan ainoastaan sen, että ne lähtevät viiden minuutin välein. Kuvatkoon muuttuja X sitä aikaa, jonka matkustaja joutuu odottamaan bussia. Jos havainnoisimme esim. 1000 matkustajan bussin odottamiseen käyttämän ajan ja muodostaisimme tästä havaintoaineistosta minuutin luokituksella prosenttisen frekvenssijakauman, olisi jokaisen luokan prosenttifrekvenssi suunnilleen sama. Tähän samaan tulokseen päädyimme ilman mitään havainnointia lähtemällä siitä, että jokainen X :n arvo välillä $[0,5]$ on yhtä todennäköinen, kuten voimme odottaa.

Asettamalla tiheysfunktiolle ehdon, että tiheysfunktion kuvaajan ja x -akselin väliin jäävän pinta-alan on oltava yksi (ks. Kuva 6), saamme liitettyä ylläesitettyyn tasajakauman tiheysfunktion määritelmään esimerkkiämme vastaavat arvot ja näin tiheysfunktiolle (alla vasemmalla) esitetyn lausekkeen:

Kuva 6: Tiheysfunktion kuvaaja (Odotusajan todennäköisyysjakauma)

$$f(X) = \begin{cases} \frac{1}{5-0}, & \text{kun } X \in [0,5] \\ 0 & \text{muualla} \end{cases}$$



Kuten Kuva 6 osoittaa, tiheysfunktio on vakio välillä $[0, 5]$ saaden arvon $1/5$. Täten jatkuvaa teoreettista tasajakaumaa noudattava satunnaismuuttuja X mallintaa hyvin bussin odotusaikaa kuvaavaa muuttujaa.

2.4.2.4 Todennäköisyyksien määrittäminen tiheysfunktion avulla

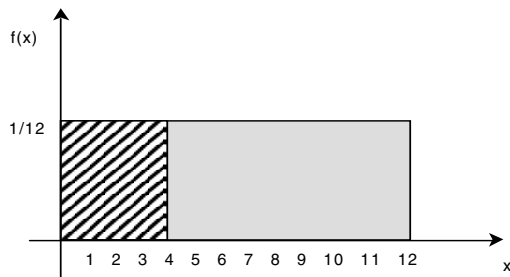
Tiheysfunktion avulla voidaan laskea erilaisten tapahtumien todennäköisyyksiä. Oletetaan, että syntymäaikaamuuttuja noudattaa tasajakaumaa, eli $X \sim \text{Tas}[0,12]$. Voimme esimerkiksi laskea, mikä on todennäköisyys, että satunnaisesti valittu henkilö on syntynyt huhtikuussa tai aikaisemmin. Merkitään

$$P(X \leq 4) = \frac{4}{12} = \frac{1}{3}$$

Todennäköisyys
Määrää tapahtuman

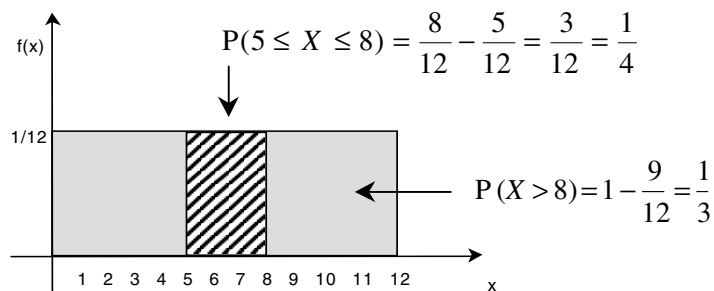
Kuva 7 havainnollistaa todennäköisyyden määräämisen; sen arvo on yhtä kuin x -akselin ja tiheysfunktion kuvaajan sekä suorien $x = 0$ ja $y = 4$ rajoittaman alueen pinta-ala, joka siis on $1/3$.

Kuva 7: Todennäköisyyden määrittäminen tiheysfunktion avulla



Voidaan myös laskea, mikä on todennäköisyys, että satunnaisesti valittu henkilö on syntynyt kesäkuun ja elokuun välisenä aikana — tai esimerkiksi elokuun jälkeen. Merkitään ja lasketaan, kuten Kuva 8:n yhteydessä.

Kuva 8: Todennäköisyyksien määrittäminen tiheysfunktion avulla



Kuten voidaan esimerkeistä todeta, todennäköisyydet ovat kolmea tyyppiä, joissa c ja d ovat annettuja vakioita (ks. myös 2.4.3.2):

$$\mathbf{P}(X \leq c), \quad \mathbf{P}(X > c) \quad \text{ja} \quad \mathbf{P}(c < X \leq d) .$$

Jatkuvan muuttujan tapauksessa todennäköisyydet $\mathbf{P}(X \leq c)$ ja $\mathbf{P}(X < c)$ ovat samoja, mutta myöhemmin, kun todennäköisyyksien määräämiseen käytetään kertymäfunktiota pitäydytään vain ensimmäisessä.

2.4.2.5 Tasajakauman odotusarvo ja varianssi

Jatkuvaa tasajakaumaa välillä $[a, b]$ noudattavan muuttujan X odotusarvolle ja varianssille on olemassa lausekkeet:

$$E(X) = \frac{a + b}{2} \quad \text{ja} \quad D^2(X) = \frac{(b - a)^2}{12} .$$

Täten esimerkiksi syntymäaikamuuttujan (mikäli yksikkönä käytetään kk) odotusarvo ja varianssi ovat

$$E(X) = \frac{0 + 12}{2} = 6 \quad \text{ja} \quad D^2(X) = \frac{(12 - 0)^2}{12} = 12 .$$

Samoihin arvoihin päädyttäisiin, mikäli käytettäisiin odotusarvon ja varianssin yleisiä — jakaumasta riippumattomia — integraalilauseita (ks. luku 2.5.1).

2.4.3 Kertymäfunktio $F(x)$

Empiiristen jakaumien puolella kumulatiivinen summakäyrä kertoo, kuinka monella prosentilla tilastoyksiköistä on pienempi tai yhtä suuri arvo kuin tietty luku x (ks. Kuva 10). Teoreettisten jakaumien puolella kertymäfunktio ilmaisee saman asian (ks. Kuva 9). Kertymäfunktion yleinen muoto on

$$F(x) = \mathbf{P}(X \leq x) .$$

Yhtälö siis ilmaisee, että muuttujan X kertymäfunktio pisteessä x on tietty todennäköisyys. Antamalla argumentille x eri arvoja saadaan kertymäfunktion kuvaaja (ks. Kuva 9).

Diskreetin muuttujan tapauksessa saadaan kertymäfunktiolle seuraava lauseke

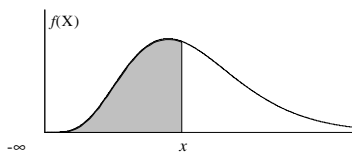
$$F(x) = \sum_{x_i \leq x} p_i .$$

Yhtälöä tarkastelemalla voidaan todeta, että kertymäfunktion arvo pisteessä x saadaan laskemalla yhteen kaikki ne \mathbf{p}_i :n arvot, joita vastaavat x_i -arvot eivät ylitä x :ää. Kertymäfunktion arvon määrittäminen diskreetin muuttujan tapauksessa siis perustuu siihen, että muodostetaan x :ää edeltäviin yksittäisiin x_i -arvoihin liittyvien todennäköisyyksien \mathbf{p}_i summa. Kuva 15 esittää diskreetin muuttujan kertymäfunktion kuvaajan (s. 15).

Jatkuvan muuttujan tapauksessa saadaan kertymäfunktiolle seuraava integraalilauseke

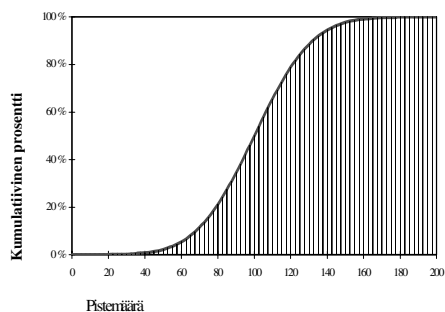
$$F(x) = \int_{-\infty}^x f(\mathbf{X})d\mathbf{X},$$

joka määrittellään tiheysfunktion $f(\mathbf{X})$ avulla, ja joka kuvaa pisteen x vasemmalle puolelle jäävän pinta-alan todennäköisyysmassan (ks. kuva alla).

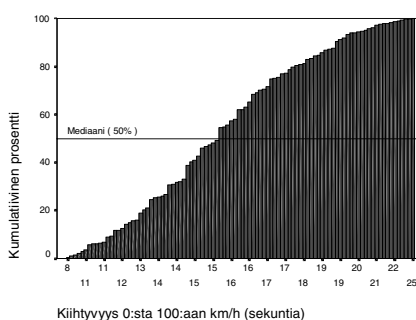


Pisteen x vasemmalle puolelle jäävä pinta-ala ja sitä vastaava kertynyt todennäköisyys saadaan lasketuksi juuri integraalilaskennan avulla.

Kuva 9: WAIS-R:n ÄÖ:n kertymäfunktio (teoreettinen)



Kuva 10: Autojen kiihtyvyyden kertymäfunktio (empiirinen)



Kuva 9 esittää teoreettisen normaalijakauman kertymäfunktion kuvaajaa, joka tässä tapauksessa mallintaa Wechslerin aikuisille tarkoitetun älykkyystestin (WAIS-R) kokonaispistemäärää (ks myös Kuva 4). Kuva 10 puolestaan esittää empiirisen muuttujan kumulatiivista jakaumaa (autojen kiihtyvyydjakauma, ks. myös Kuva 1).

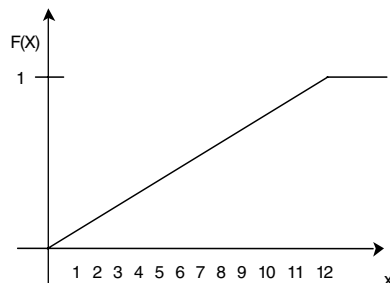
2.4.3.1 Tasajakauman kertymäfunktio

Jatkuvalle tasajakaumalle välillä $[a, b]$ voidaan johtaa kertymäfunktion lauseke:

$$F(x) = \begin{cases} 0, & \text{kun } x < a \\ \frac{x-a}{b-a}, & \text{kun } x \in [a, b] \\ 1, & \text{kun } x > b \end{cases}.$$

Kun tunnetaan tasajakaumaa noudattavan muuttujan vaihteluväli $[a, b]$, voidaan antaa argumentille x eri arvoja, jolloin saadaan kertymäfunktion kuvaaja. Tasajakaumaa oletetusti noudattavan syntymäaikamuuttujan vaihteluväli on $[0, 12]$. Täten voidaan muodostaa syntymäaikamuuttujan kertymäfunktion kuvaaja (ks. Kuva 11).

Kuva 11: Syntymäaikamuuttujan $X \sim \text{Tas}[0, 12]$ kertymäfunktion kuvaaja



Funktio tietyllä arvolla x antaa siis aina tähän pisteeseen mennessä kertyneen todennäköisyysmassan arvon. Kun $x < a$, on funktion arvo aina nolla; kun taas $x > b$, on funktion arvo aina yksi. Jo yllä esitetyt muotoa $\mathbf{P}(X \leq c)$, $\mathbf{P}(X > c)$ ja $\mathbf{P}(c < X \leq d)$ olevat todennäköisyydet voidaan määrätä kätevimmin juuri kertymäfunktion avulla. Yllä esitettyä bussiesimerkkiä käyttäen (ks. luku 2.4.2.2), voitaisiin laskea todennäköisyys, että matkustaja joutuu odottamaan bussia 3-5 minuuttia. Määrätään kertymäfunktion arvot pisteessä 5 ja 4 ja muodostetaan niiden erotus. Eli: $\mathbf{P}(3 < X \leq 5) = F(5) - F(3) = 5/5 - 3/5 = 2/5$.

2.4.3.2 Kertymäfunktion ominaisuuksia

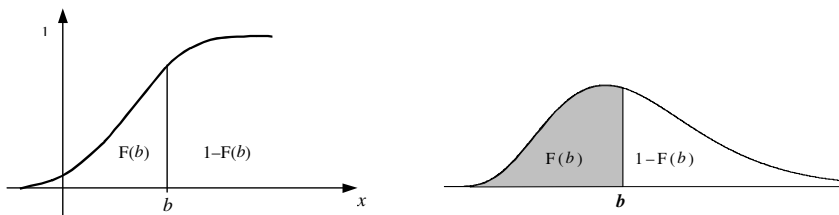
Riippumatta teoreettisesta jakaumasta on kaikilla kertymäfunktioilla 2 tärkeää ominaisuutta:

1. $\mathbf{P}(X > b) = 1 - F(b)$
2. $\mathbf{P}(a < X \leq b) = F(b) - F(a)$,

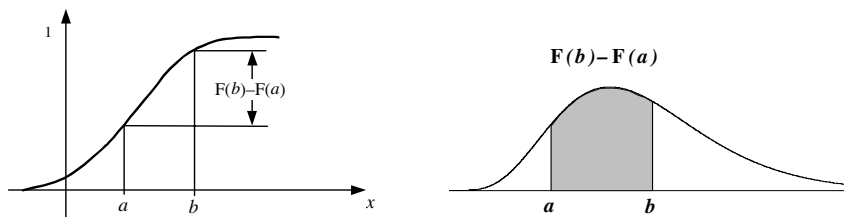
Jossa \mathbf{P} tarkoittaa todennäköisyyttä ja suluissa oleva lauseke määrää tapahtuman. Eli kertymäfunktion avulla voidaan määrätä muuttujan arvojen esiintymiseen liittyviä erityyppisiä todennäköisyyksiä.

Kuva 12 havainnollistaa kertymäfunktion ominaisuuteen 1 liittyvän todennäköisyyden määrittämistä. Tilanne esitetään myös tiheysfunktion kannalta, joka määrittelee kertymäfunktion, kuten yllä totesimme.

Kuva 12: Kertymäfunktion ja tiheysfunktion todennäköisyyksien $F(b)$ ja $1-F(b)$ määrittämisessä (Ominaisuus 1)



Kuva 13: Kertymäfunktion ja tiheysfunktion tietyn välin $[a,b]$ todennäköisyyden määrittämisessä (Ominaisuus 2)



Kuva 13:ssä esitetään puolestaan kertymäfunktion ominaisuuteen 2 liittyvän todennäköisyyden — eli tietyn välin $[a,b]$ todennäköisyyden — määrittäminen.

Lisäksi koskien diskreettiä jakaumaa, on olemassa vielä kolmas ominaisuus:

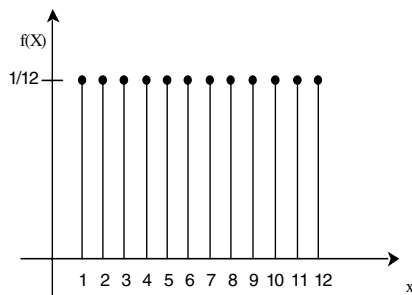
$$3. \quad \mathbf{P}(X = a_i) = \mathbf{F}(a_i) - \mathbf{F}(a_{i-1}) .$$

Tämä ominaisuus määrää diskreetin muuttujan tietyn yksittäisen arvon esiintymistodennäköisyyden eli ns. pistetodennäköisyyden, joka määrätään pisteen \mathbf{a}_i ja sitä välittömästi edeltävän pisteen \mathbf{a}_{i-1} kertymäfunktioiden arvojen erotuksena. Näin ei voida tietysti menetellä jatkuvan muuttujan tapauksessa, sillä jatkuvan muuttujan jatkuvuusominaisuudesta johtuen ei sille ole määriteltävissä mitään tarkkaa jakauman pistettä, jolle voitaisiin laskea pistetodennäköisyys. Jatkuvalle muuttujalle määrätään siksi vain tietyn välin todennäköisyyksiä.

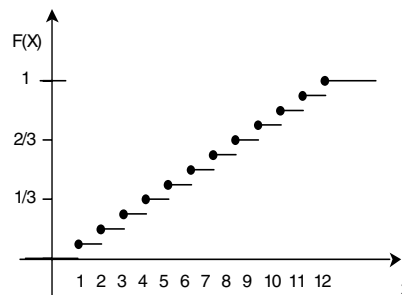
2.4.4 Diskreetti tasajakauma

Yllä olemme käsitäneet syntymäaikamuuttujan jatkuvana, sillä kuukausiluokituksesta huolimatta taustalla oleva aikamuuttuja on jatkuva luonteeltaan. On mahdollista ajatella syntymäaikamuuttuja myös diskreettinä muuttujana, eli syntymäkuukausimuuttujana, jolloin jokainen kuukausi muodostaisi epäjatkuvan kategorian. Tällöin ei olisi syntymäkuukausimuuttujan kohdalla mielekäästä puhua tiheysfunktioista, vaan pistetodennäköisyysfunktioista, jota Kuva 14 esittää. Kertymäfunktion kuvaajakin on silloin toisenlainen ja sitä kutsutaan *porrasfunktioiksi*, kuten Kuva 15 havainnollistaa.

Kuva 14: Diskreetin syntymäkuukausimuuttujan tiheysfunktion kuvaaja



Kuva 15: Diskreetin syntymäkuukausimuuttujan porrasfunktion kuvaaja



2.5 Teoreettiset tunnusluvut

Kuten empiirisille jakaumille, voidaan myös teoreettiselle jakaumille määritellä erilaisia tunnuslukuja, esimerkiksi keski- ja hajontalukuja. Teoreettisen jakauman yhteydessä näitä kutsutaan *teoreettisiksi tunnusluvuiksi*. Tarkastellaan näistä lähemmin keskiarvoa ja keskihajontaa sekä moodia, mediaania, fraktiileja että jakauman muotoa kuvaavia tunnuslukuja.

2.5.1 Odotusarvo ja varianssi

Taulukko 2 esittää rinnasteisesti satunnaismuuttujan X empiiriset ja teoreettiset perustunnusluvut — eli otoskeskiarvon / odotusarvon ja otosvariانسsin / variانسsin. Odotusarvoa $E(X)$ merkitään μ :llä (äännetään: *myy*) ja varianssia $D^2(X)$ σ^2 :lla (äännetään: *sigma toiseen*). Perusjoukon / populaation keskiarvosta puhuttaessa tarkoitetaan juuri μ :tä. Samoin σ^2 :lla tarkoitetaan perusjoukon / populaation varianssia.

Taulukko 2: Teoreettiset ja empiiriset tunnusluvut

Empiiriset (otos)tunnusluvut	Teoreettiset tunnusluvut
• Otoskeskiarvo \bar{X}	• Satunnaismuuttujan X odotusarvo $E(X) = \mu$
• Otosvariانسsi s^2	• Satunnaismuuttujan X variانسsi $D^2(X) = \sigma^2$

Jatkuvan satunnaismuuttujan X keskiarvolle ja variانسsille (= keskihajonnan neliö) on voimassa seuraavat kaavat:

$$E(X) = \int_{-\infty}^{\infty} X f(X) dX$$

$$D^2(X) = \sigma^2 = \int_{-\infty}^{\infty} (X - \mu)^2 f(X) dX = E(X^2) - \mu^2.$$

Odotusarvon E symboli tulee Expected value -termistä ja keskihajonnan D symboli vastaavasti Deviation -sanasta.

Muuttuja X noudattaa tasajakaumaa $[0,5]$ (ks. tasajakaumasta luku 2.4.2.2). Tämän jakauman keskiarvo ja variانسsi ovat

$$E(X) = \int_0^5 X \frac{1}{5} dX = \frac{1}{5} \cdot \frac{1}{2} X^2 \Big|_0^5 = \frac{1}{10} (5^2 - 0^2) = 2,5$$

$$D^2(X) = \int_0^5 X^2 \frac{1}{5} dX - 2,5^2 = \frac{1}{5} \cdot \frac{1}{3} X^3 \Big|_0^5 - 2,5^2 = \frac{25}{12} = 2,08.$$

Diskreetin satunnaismuuttujan muuttujan X odotusarvolle ja varianssille on olemassa seuraavat lauseet

$$E(X) = \mu = \sum_i p_i x_i$$

$$D^2(X) = \sigma^2 = \sum_i p_i \cdot (x_i - \mu)^2 .$$

Eräs diskreetti muuttuja X saa arvoja väliltä 1-4. Arvojen x_1, \dots, x_4 esiintymistodennäköisyydet ovat vastaavasti: $p_1 = 0,2$, $p_2 = 0,3$, $p_3 = 0,4$ ja $p_4 = 0,1$. Tämän muuttujan odotusarvo (keskiarvo) ja varianssi ovat seuraavat:

$$E(X) = \mu = 0,2 \cdot 1 + 0,3 \cdot 2 + 0,4 \cdot 3 + 0,1 \cdot 4 = \mathbf{2,4}$$

$$D^2(X) = \sigma^2 = 0,2 \cdot (1-2,4)^2 + 0,3 \cdot (2-2,4)^2 + 0,4 \cdot (3-2,4)^2 + 0,1 \cdot (4-2,4)^2 = \mathbf{0,84} .$$

Teoreettisten jakaumien odotusarvo ja varianssi voidaan määrätä yleisesti ylläesitettyillä lausekkeilla. Käytännössä niitä ei kuitenkaan juuri tarvita, koska yleensä eri teoreettisten jakaumien odotusarvolle ja varianssille on johdettu yksinkertaisemmat lausekkeet.

2.5.2 Muita tunnuslukuja

Moodi (Mo) on se muuttujan arvo, jolla tiheysfunktio saavuttaa maksimiarvonsa tai vastaavasti empiirisen muuttujan yleisin arvo/luokka.

Mediaani (Md) on jakauman suuruusjärjestykseen asetettujen lukujen keskimäinen arvo (= 50%:n fraktiili). Mikäli jakauman n on pariton saadaan mediaani seuraavasti

$$Md = x_{((n+1)/2)} ,$$

ja vastaavasti kahden keskimäisen keskiarvo, mikäli n on parillinen eli

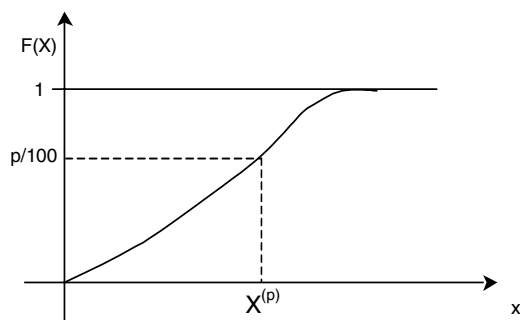
$$Md = (x_{(n/2)} + x_{(n/2+1)}) / 2 .$$

Muuttujan X p %:n (teoreettinen) **fraktiili** eli **prosenttipiste** $X^{(p)}$ toteuttaa yhtälön

$$F(X^{(p)}) = \frac{p}{100} .$$

Muuttujan fraktiili $X^{(p)}$ on siten arvo, johon asti todennäköisyysmassan suuruus on $p\%$. Kuva 16 esittää $p\%$:n fraktiilin $X^{(p)}$ määrittämistä kertymäfunktion tai kumulatiivisen summajakauman avulla.

Kuva 16: $p\%$:n fraktiilin määrittäminen



Eniten käytettyjä fraktiileita ovat 25% :n ja 75% :n sijaintikohdat jakaumassa. Näitä pisteitä kutsutaan usein jakauman *ala-* ja *yläkvartiileiksi* ja merkitään Q_1 :llä ja Q_3 :lla. Mediaania merkitään tällöin Q_2 :lla.

- **ESIM: Muuttuja $X \sim \text{Tas}[0,20]$. Määrää $F(X^{(65)})$.**

Yhtälöä soveltamalla saadaan:

$$F(X^{(p)}) = \frac{p}{100} = \frac{X^{(p)}}{20} = \frac{65}{100} \Leftrightarrow X^{(p)} = \frac{20 \times 65}{100} = 13.$$

Huom. Muiden tässä kirjassa esiintyvien teoreettisten jakaumien kohdalla fraktiilin määrittäminen ei ole aivan yhtä yksinkertaista kuin tasa-jakauman kohdalla. Muiden jakaumien yhteydessä fraktiilipisteen yhtälö ratkaistaan taulukosta saatavan kertymäfunktion arvon avulla, kuten tulemme jatkossa muiden jakaumien yhteydessä toteamaan.

Muuttujan **kriittinen arvo** on piste, jonka yläpuolelle jää p :n suuruinen alue. Siten esimerkiksi 90% :n fraktiili on $0,1$:n kriittinen arvo. Kriittisiä arvoja käytetään estimoinnissa ja hypoteesin testauksessa.

Vinousluvut osoittavat jakauman symmetrisyyttä / epäsymmetrisyyttä keskiarvonsa suhteen.

- **Skewness (vinous):** Kuvaa jakauman symmetrisyyttä / epäsymmetrisyyttä keskiarvonsa suhteen. Jakauma on likimäärin symmetrinen keskiarvonsa suh-

teen, mikäli arvo on lähellä nollaa (välillä -0.50 ; $+0.50$). Skewness on negatiivinen mikäli jakauma on oikealle vino ja vastaavasti positiivinen mikäli jakauma on vasemmalle vino. Skewness otosestimaatti lasketaan seuraavalla kaavalla

$$\text{Skewness} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^3}{s^3}.$$

Huipukkuusluvut osoittavat jakauman huipukkuutta eli muuttujan arvojen tiivistymistä keskiarvon ympärille tai tasasuuntautuneisuutta eli muuttujan arvojen kertymistä jakauman molempiin päihin.

- **Kurtosis (huipukkuus)**: Likimäärin normaalijakautuneen muuttujan kurtosis arvo on lähellä nollaa (tai kolmea mikäli alla esitettyssä kaavassa ei suoriteta ko. erotusta). Yleisesti huipukasta jakautumaa kuvaavat arvot ovat positiivisia ja tasaisempaa jakautumaa kuvaavat negatiivisia. Kurtosis otosestimaatti lasketaan seuraavalla kaavalla

$$\text{Kurtosis} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^4}{s^4} - 3.$$

Esimerkkinä skewness ja kurtosis arvoista voidaan ottaa Kuva 1 ja Kuva 2 sivulla 4. Kuva 1:ssä on kysymys jokseenkin normaalisti jakautuneesta autojen kiihtyvyyssjakaumasta, jolle saadaan skewness arvo $0,21$ ja kurtosis arvo $0,39$. Näin ollen autojen kiihtyvyyssjakauma on selvästi normaalijakautuneen suuntaainen, vaikkakin ehkä lievästi huipukas. Kuva 2:ssa on puolestaan kysymys selvästi positiivisesti eli oikealle vinosti jakautuneesta autojen hevosvoimajakau-
kaumasta. Nyt skewness arvo on $1,04$ ja kurtosis $0,59$. Näitä kahta jakauman muotoa kuvaavaa tunnuslukua on syytä käyttää aina yhdessä, sillä ne tukevat toistensa tulkintaa. Lisäksi skewness ja kurtosis soveltuvat vain suurehkojen otoskokojen yhteydessä käytettäväksi.

2.6 Normaalijakauma

Normaalijakaumaa voidaan pitää yleisimpänä teoreettisista jatkuvista jakaumista. Normaalijakaumalla onkin perinteisesti ollut hyvin keskeinen asema tilastotieteessä. Sen lisäksi, että monet empiiriset muuttujat ja testisuureet noudattavat likimäärin normaalijakaumaa, on monilla muillakin teoreettisilla jakaumilla siihen tietynlainen suhde, kuten tullaan myöhemmin toteamaan. Jos muuttuja X noudattaa normaalijakaumaa parametrein μ ja σ^2 , käytetään merkintää $X \sim N(\mu, \sigma^2)$. Toisena parametrinä käytetään varianssia σ^2 , ei keskihajontaa σ , koska useamman muuttujan normaalijakauman tapauksessa toinen parametri on kovarianssimatriisi, jonka yksiulotteinen vastine on varianssi. Muuttujan hajontalukuna on kuitenkin syytä käyttää keskihajontaa.

2.6.1 Keskeinen raja-arvolause (central limit theorem)

Keskeinen raja-arvolause antaa matemaattisesti esitettävissä olevan perustelun normaalijakaumalle. Olkoot X_i :t ($i=1, \dots, k$) erittäin yleisiä oletuksia täyttäviä, toisistaan riippumattomia muuttujia. Tällöin näiden muuttujien summan

$$X_{\text{sum}} = X_1 + X_2 + X_3 + \dots + X_k$$

jakauma lähestyy normaalijakaumaa, kun yhteenlaskettavien lukumäärä k kasvaa rajatta. Käytännössä jos yhteenlaskettavia on vähintään 30, on X_{sum} :n jakauma riittävän lähellä normaalijakaumaa. Jos X_i :t ovat esim. tasajakaumaa noudattavia muuttujia, riittää likimääräisen normaalijakauman saavuttamiseksi huomattavasti pienempi yhteenlaskettavien lukumäärä, esim. 10.

Mitä lähempänä X_i -muuttujien jakaumat ovat normaalijakaumaa, sitä pienempi saa yhteenlaskettavien lukumäärä k olla, jotta X :n normalisuus saavutetaan. Yleisesti siis, jos X_i -muuttujien jakaumat ovat tuntemattomia, niin $X_{\text{sum}} \sim N(\mu, \sigma^2)$, kun $k > 30$.

2.6.2 Muuttujatyypit jotka noudattavat likimäärin normaalijakaumaa

Tarkastellaan keskeisen raja-arvolauseen valossa eräitä keskeisiä muuttujatyyppejä, jotka noudattavat likimäärin normaalijakaumaa.

2.6.2.1 Virhemuuttujat

Erilaisten virheiden, kuten mittaus- ja havainnointivirheiden, voidaan ajatella muodostuvan useasta komponentista. Erilaisia virhelähteitä ovat esimerkiksi epätarkat mittalaitteet, väsymys, keskittymisen puute, huolimattomuus, huonot valaistusolosuhteet ja muut häiriön lähteet, kuten mm. melu. Erilaisten virhelähteiden ajatellaan vaikuttavan mittaus- ja havainnointivirheen muodostumiseen

toisistaan riippumattomasti. Näin ollen oleellista on siis huomata, että virheen määrän satunnaisvaihtelua kuvaava muuttuja on luonteeltaan summamuuttuja, jossa virheen kokonaismäärä muodostuu useista toisistaan *riippumattomista* virheen lähteistä. Tältä pohjalta ns. klassinen testi- / mittausteoria päättyy keskeiseen peruslauseeseensa:

$$\text{Mittaustulos} = \text{Todellinen arvo} + \text{Mittausvirheen arvo} ,$$

eli havaittu mittatulos käsitetään ylläselvitettyyn nojautuen todellisen arvon ja virheen summaksi.

2.6.2.2 Ominaisuudet

Monet fyysiset ja psyykkiset ominaisuudet noudattavat normaalijakaumaa. Etenkin ominaisuudet, joiden voidaan ajatella olevan monien tekijöiden määriä, ovat taipuvaisia jakautumaan normaalisti. Esimerkiksi älykkyyden ajatellaan jakautuvan normaalisti. Samoin pituus ja fyysinen suorituskyky ovat ominaisuuksia, jotka ovat moninaisten geneettisten ja ympäristötekijöiden määriä. Kuitenkaan kaikki tämän tyyppisen muuttujat eivät noudata normaalijakaumaa. Esimerkiksi ihmisen paino ei ole normaalisti jakautunut, vaikkakin se on monien tekijöiden määräämä ominaisuus. Tämä johtuu siitä, että ihmisen painoon vaikuttaa yksi tekijä, jolla on muihin tekijöiden nähden varsin voimakas vaikutus, nimittäin ravinto.

2.6.2.3 Otoskeskiarvon otantajakauma

Otoskeskiarvo \bar{X} voidaan käsitellä satunnaismuuttujaksi, jonka arvot vaihtelevat μ :n ympärillä otantaprosessia toistettaessa. Otoskeskiarvon otantajakauma kuvaa tätä satunnaisvaihtelua. Koska otoskeskiarvo on luonteeltaan summamuuttuja, on sen otantajakauma keskeisen raja-arvolauseen mukaan $\mathbf{N}(\mu, \sigma^2/n)$. Vas-

taavasti standardoitu otoskeskiarvo $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ noudattaa $\mathbf{N}(0, 1)$ -jakaumaa.

2.6.3 Normaalijakauman tiheysfunktio

Yleisessä muodossaan normaalijakauman tiheysfunktio on seuraavanlainen:

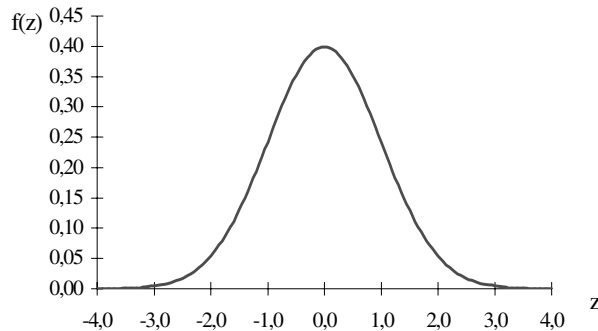
$$f(X|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{(X-\mu)^2}{2\sigma^2}} .$$

Tämä lauseke voidaan johtaa keskeisestä raja-arvolauseesta, tosin johto on varsin monimutkainen.

Tiheysfunktio sisältää kaksi parametriä: *populaation* keskiarvo μ ja *populaation* varianssi σ^2 . Symboli e on luonnollisen logaritmijärjestelmän vakioinen kantaluku (2,718).

Erityisesti jos $\mu = 0$ ja $\sigma^2 = 1$, on kysymys *standardoidun* normaalijakauman tiheysfunktion lausekkeesta, jolloin tiheysfunktiossa X :n tilalla käytetään Z :aa.

Kuva 17: Standardoidun normaalijakauman tiheysfunktion kuvaaja



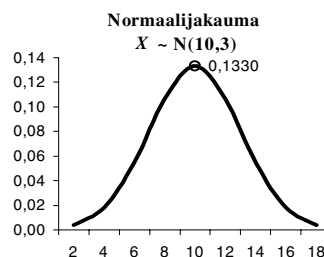
Kuva 17 sisältää standardoidun normaalijakauman tiheysfunktion kuvaajan. Tällöin Z -muuttuja noudattaa normaalijakaumaa parametrien arvoilla $\mathbf{N}(0, 1)$.

Kun $X = \mu$, saavuttaa tiheysfunktio maksimiarvonsa eli kuvaajan lakipisteen, joka on tällöin yhtä kuin tiheysfunktion osatermi

$$\frac{1}{\sigma\sqrt{2\pi}}. \text{ Tällöin puolestaan osatermi } e^{-\frac{(X-\mu)^2}{2\sigma^2}} = 1, \text{ koska } e^0 = 1.$$

Tarkastellaan normaalijakauman tiheysfunktiota esimerkin avulla. Kun muuttuja X noudattaa normaalijakaumaa $\mathbf{N}(10, 3^2)$, saadaan X :n arvoille [2, 4, 6, 8, 10, 12, 14, 16, 18] ko. parametrien arvoilla seuraavat tiheysfunktion arvot ja seuraava tiheysfunktion kuvaaja:

x	$f(X 10, 3^2)$	=NORMDIST(x;10;3;0)
2	0,0038	=NORMDIST(2;10;3;0)
4	0,0180	=NORMDIST(4;10;3;0)
6	0,0547	=NORMDIST(6;10;3;0)
8	0,1065	=NORMDIST(8;10;3;0)
10	0,1330	=NORMDIST(10;10;3;0)
12	0,1065	=NORMDIST(12;10;3;0)
14	0,0547	=NORMDIST(14;10;3;0)
16	0,0180	=NORMDIST(16;10;3;0)
18	0,0038	=NORMDIST(18;10;3;0)



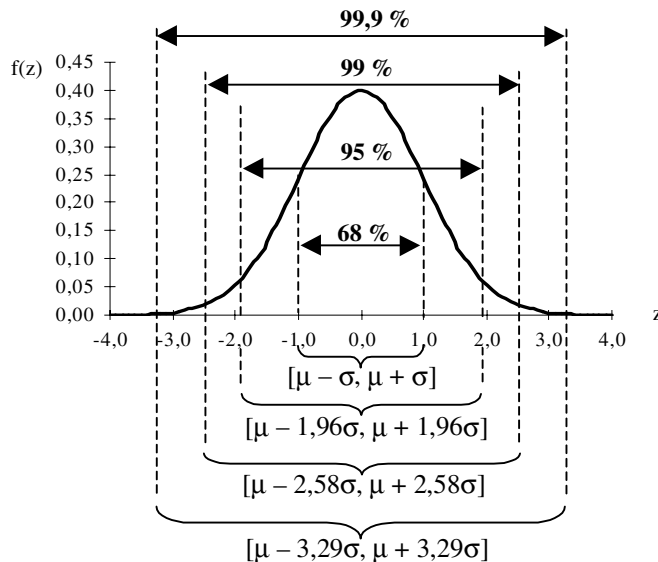
Tiheysfunktion kuvaajaa tarkasteltaessa, voidaan todeta, kun $X = \mu$, saa tiheysfunktio suurimman arvonsa 0,1330.

Taulukon kolmannessa sarakkeessa on esitettyä Excelin normaalijakauman tiheysfunktio NORMDIST, jonka avulla toisen sarakkeen arvot on laskettu. Funktion nimi on syötetty Excelin soluun, jota ennen on kirjoitettu yhtäsuuruusmerkki ja tämän jälkeen vastaavat argumenttien arvot suluissa puolipisteillä erotettuna. Viimeinen argumentin arvo '0' tarkoittaa, että funktiota käytetään nimenomaan tiheysfunktiona; arvo '1' tarkoittaisi, että NORMDIST –funktio käytettäisiin sen sijaan kertymäfunktiona.

2.6.4 Normaalijakauman ominaisuuksista

Normaalijakauman kokonaistodennäköisyys (=1) jaetaan sopimuksenvaraisesti seuraaviin odotusarvon suhteen symmetrisiin luottamusrajoihin, joita käytetään hyväksi tilastollisessa päätöksenteossa (Kuva 18).

Kuva 18: Standardoidun ja (hakasulkeissa) yleisen normaalijakauman luottamusrajoja



Kuva 18 esittää hakasuluissa myös yleistä normaalijakaumaa koskevat sovitut luottamusrajat, joita usein kutsutaan kahden ensimmäisen välin osalta yhden ja kahden *sigman* ($=\sigma$) väleiksi. Nimitys kuvaa *sigman* poikkeaman suuruutta *myyistä* ($=\mu$). Toisen välin tapauksessa *sigman* kertojana toimiva *z*:n arvo 1,96 pyöristetään usein kakkoseksi.

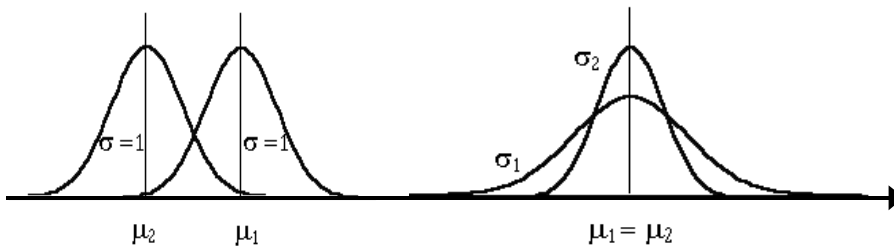
2.6.5 Parametrien μ ja σ^2 vaikutus tiheysfunktioon

Seuraavaksi tutkitaan, miten parametrit μ ja σ^2 vaikuttavat tiheysfunktion kuvaajaan. Oletetaan, että muuttujat X ja Y noudattavat normaalijakaumaa eli $X \sim N(\mu_1, \sigma_1^2)$ ja $Y \sim N(\mu_2, \sigma_2^2)$. Kuva 19 esittää kaksi sellaista perustilannetta, jotka kuvaavat ko. parametrien vaikutusta tiheysfunktion sijaintiin ja muotoon.

Kuva 19: Parametrien μ ja σ^2 vaikutus normaalijakauman tiheysfunktioon

Tilanne A: $\mu_1 > \mu_2; \sigma_1^2 = \sigma_2^2$

Tilanne B: $\mu_1 = \mu_2; \sigma_1^2 > \sigma_2^2$



Tilanteessa **A** $f(X | \mu_1, \sigma_1^2)$:n huippu on enemmän oikealla kuin $f(Y | \mu_2, \sigma_2^2)$:n huippu, koska μ_1 on suurempi kuin μ_2 . Jakaumien muoto on kuitenkin sama, koska varianssit ovat samat. Tilanteessa **B** $f(X | \mu_1, \sigma_1^2)$:n ja $f(Y | \mu_2, \sigma_2^2)$:n huiput ovat päällekkäin, koska μ_1 ja μ_2 ovat yhtä suuria. Mutta koska muuttujan X varianssi on suurempi kuin Y :n, on sen tiheysfunktion kuvaaja leveämpi.

2.6.6 Taulukon korvaaminen Excelin funktioilla

Excelissä on funktiota, joiden avulla voidaan määrätä erilaisiin teoreettisiin jakaumiin liittyviä todennäköisyyksiä ja fraktiilipisteitä — ts. niillä voidaan korvata taulukkokirjat. Normaalijakaumaan liittyvät Excelin funktiot ovat:

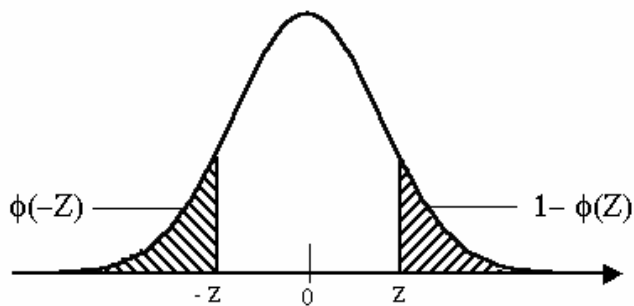
- **=NORMDIST(x; keskiarvo; keskihajonta; 0 tai 1)** antaa x :n arvoon liittyvän normaalijakauman tiheysfunktion arvon, mikäli 0 on valittuna viimeiseksi argumentin arvoksi. Jos viimeisen argumentin arvo on 1, antaa funktio x :ää vastaavan kertymäfunktion arvon. (suom. NORM. JAKAUMA).
- **=NORMINV(p; keskiarvo; keskihajonta)** antaa $100 \cdot p$ %:n fraktiilin, palauttaa normaalijakauman kertymäfunktion käänteisarvon. (NORM. JAKAUMA. KÄÄNT).
- **=NORMSDIST(z)** antaa standardoituun normaalijakaumaan liittyviä kertymäfunktion arvoja eri Z -muuttujan arvoille. (suom. NORM. JAKAUMA. NORMIT).

- $\text{=NORMSINV}(p)$ antaa Z -muuttujan $100 \cdot p$ %:n fraktiilin, eli palauttaa standardoidun normaalijakauman kertymäfunktion käänteisarvon. (suom. NORM. JAKAUMA. NORMIT. KÄÄNT).

2.6.7 Todennäköisyyksien määrittäminen (Z -jakaumassa)

Koska normaalijakauma on *symmetrinen* keskiarvonsa $\mu = 0$ suhteen, on kertymäfunktion arvo $\phi(-z)$ yhtäpitävä $1 - \phi(z)$ kanssa. Kuva 20 havainnollistaa tämän.

Kuva 20: Normaalijakauman symmetrisyyden hyväksikäytettävyys



- **ESIM:** Mikä on kertymäfunktion arvo, kun z :n arvo on $-1,5$?

Merkitään $\phi(-1,5) = P(Z \leq -1,5)$.

Yleensä normaalijakauman kertymäfunktion arvoja eri z :n arvoille saadaan taulukkokirjoista. Liitteessä 1 esitetään standardoidun normaalijakauman kertymäfunktion arvoja. Etsitään tämän taulukon ensimmäisestä sarakkeesta z :n arvo $-1,5$ ja katsotaan siihen liittyvä kertymäfunktion arvo vastaavalta riviltä, joka siis on 0,9332.

Mikäli käytettävissämme olisi sellainen taulukkokirja, jossa standardoidun normaalijakauman kertymäfunktioiden arvot olisi laskettu vain positiivisille Z :n arvoille, määrättäisiin $\phi(-1,5)$:n arvo normaalijakauman symmetrisyyttä hyväksi käyttäen seuraavasti

$$1 - \phi(1,5) = 1 - 0,9332 = \underline{0,0668} .$$

Excelillä $\phi(-1,5)$ voidaan määrätä kirjoittamalla soluun seuraava kaava

$\text{=NORMSDIST}(-1,5)$, josta saadaan tulos 0,0668.

- **ESIM: Millä todennäköisyydellä Z kuuluu välille $-2,2$ ja $-1,2$?**

Merkitään $P(-2,2 < Z \leq -1,2)$.

Liitteen 1 taulukkoa käyttäen määrätään kyseisten z -pisteiden kertymäfunktion arvot ja lasketaan todennäköisyys näiden erotuksena, eli

$$\phi(-1,2) - \phi(-2,2) = 0,1151 - 0,0139 = \underline{0,1012}.$$

Mikäli käytettävissämme olisi sellainen taulukko, jossa standardoidun normaalijakauman kertymäfunktioiden arvot olisi laskettu vain positiivisille Z :n arvoille, laskettaisiin ko. todennäköisyys normaalijakauman symmetrisyyttä hyväksi käyttäen seuraavasti

$$(1 - \phi(1,2)) - (1 - \phi(2,2)) = (1 - 0,8849) - (1 - 0,9861) = 0,1151 - 0,0139 = \underline{0,1012}.$$

Exceliä käyttäen kirjoitetaan soluun kaava:

=NORMSDIST(-1,2) – NORMSDIST(-2,2) ja saadaan tulos 0,1012.

- **ESIM: Määrää Z -jakauman 95 %:n fraktiilipiste — eli $F(Z^{(95)})$.**

Liitteen 1 taulukkoa käyttäen voidaan ratkaista kertymäfunktion yhtälö $\phi(Z^{(95)}) = p/100$ etsimällä taulukosta kertymäfunktion arvoa **0,95** lähimpänä olevan z :n arvon, joka on 1,64. Eli

$$\phi(1,64) = 95 / 100; (0,95 = 0,95).$$

(Symmetrisyydestä johtuen on 5 %:n fraktiili vastaavasti -1,64).

Exceliä käyttäen kirjoitetaan soluun kaava: =NORMSINV(0,95), joka antaa tuloksen 1,64.

(ja vastaavasti 5 %:n fraktiilille =NORMSINV(0,05) = -1,64).

3 Lineaarinen muunnos

3.1 Yleistä

Lineaarinen muunnos tarkoittaa, muuttujan arvojen kertomista tai jakamista jollain vakiolla — tai vakion lisäämistä (tai vähentämistä) muuttujan arvoihin (arvoista). Eli

$$Y = a + bX, \text{ jossa } a \text{ ja } b \text{ ovat vakioita,}$$

jolloin muunnettavan jakauman X arvojen suhde muunnetun jakaumaan Y arvoihin nähden säilyy vakioisena koko X :n vaihteluvälillä. Ts. lineaarisessa muunnoksessa muuttujan arvoja ei koroteta esim. tiettyyn potenssiin, niistä ei oteta neliöjuurta tai käytetä muutokseen muita epälineaarisia funktiota, kuten \cos tai \sin jne., jotka muuttavat X :n ja Y :n arvojen väliset suhteet epävakioiseksi. Huomion arvoinen seikka on, että lineaarinen muunnos on pätevä ja riippumaton lähtöjakauman muodosta.

Tyypillisenä esimerkkinä lineaarisesta muunnoksesta voidaan ottaa lämpötilan mittaaminen, jossa celsius- ja fahrenheit –asteikoiden välisen muunnoksen ilmaisee yhtälö

$$\text{Fahrenheit} = 32 + 1,8 * \text{Celcius} .$$

Liitteessä 8 esitetään lineaarisen muunnoksen yksi käytännöllinen sovellus, nimittäin asteikkovälin muunnos sekä ohjeet sen suorittamiseen Excelissä.

3.1.1 Lineaarisen muunnoksen vaikutus muuttujan keskiarvoon ja varianssiin

Vakion lisääminen tai vähentäminen muuttujaan (tai muuttujasta) X korottaa (tai vähentää) sen keskiarvoa vakion suuruuden verran — varianssin säilyessä samana. Muuttujan X kertominen tai jakaminen suurentaa tai pienentää sen keskiarvoa ko. kertoimen tai jakajan suuruuden verran — varianssin kasvaessa (tai pienentyessä) kertoimen (tai jakajan) neliön verran.

Voidaan osoittaa, että teoreettiset tunnusluvut käyttäytyvät yleisesti (jakaumasta riippumatta) lineaarisessa muunnoksessa $Y = a + bX$ samalla tavalla kuin vastaavat empiiriset tunnusluvut eli

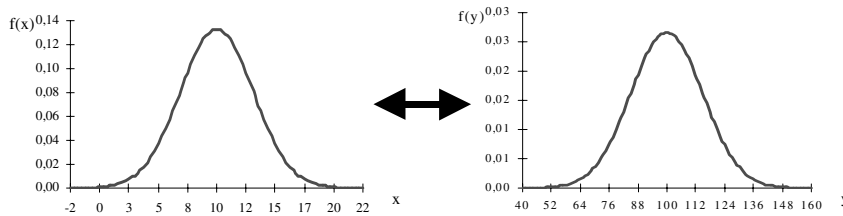
$$E(a + bX) = a + b E(X)$$

$$D^2(a + bX) = b^2 D^2(X) .$$

Erityisesti normaalijakauman tapauksessa myös X :n jakauman muoto säilyy lineaarisessa muunnoksessa samana. Kun muuttuja X noudattaa normaalijakaumaa odotusarvolla μ ja varianssilla σ^2 , tällöin myös lineaarisesti muunnet-

tu uusi muuttuja $Y = a + bX \sim N(a + b\mu, b^2\sigma^2)$ keskiarvolla $a + b\mu$ ja varianssilla $b^2\sigma^2$, jossa a ja b ovat mielivaltaisesti valittavissa olevia vakioita. Kuva 21 havainnollistaa vielä muunnosta.

Kuva 21: Jakauman muoto ja tunnusluvut normaalijakauman lineaarisessa muunnoksessa



Kuten kuvaa tarkastelemalla voi todeta, muuntuu X :n jakauma, joka siis $\sim N(10, 3^2)$ lineaarisesti Y :n jakaumaksi, joka puolestaan $\sim N(100, 15^2)$ valituilla vakion arvoilla $a = 50$ ja $b = 5$.

Voidaan vielä lisäksi todeta, että mikäli $\mu = 0$ ja $\sigma^2 = 1$, silloin lauseke $N(a + b\mu, b^2\sigma^2)$ yksinkertaistuu $N(a, b^2)$:ksi.

3.2 Muuttujan standardoiminen

Lineaarista muunnosta tarvitaan (mm.) todennäköisyyksien ja fraktiilien laskemiseksi silloin, kun muuttajat noudattavat normaalijakaumaa yleisesti, eli kun $X \sim N(\mu, \sigma^2)$, ts. silloin kun ei ole kysymys Z -muuttujasta. Normaalijakaumaa noudattava muuttuja X voidaan palauttaa Z -muuttujaksi eli *standardoida* ja laskea sitten siihen liittyviä todennäköisyyksiä standardoidussa normaalijakaumassa. Muuttuja X standardoidaan seuraavasti:

$$Z = \frac{(X - \mu)}{\sigma}.$$

Eli jokainen muuttujan X arvo vähennetään keskiarvostaan ja jaetaan keskihajonnallaan, jolloin saadaan X :n arvoja vastaavat Z -muuttujan arvot. Ts. standardoinnissa muuttujan X arvot suhteutetaan keskihajontaansa, jolloin voidaan tarkastella niiden sijaintia jakaumassa. Standardoidun muuttujan Z keskiarvoksi tulee 0 ja varianssiksi 1 .

- **ESIM:** Z :n arvo $0,76$ tarkoittaa sen olevan $0,76$ keskihajonnan mittaa (pituutta) keskiarvon yläpuolella. Kun taas Z :n arvo $-2,87$ tarkoittaa

taa sen olevan 2,87 keskihajonnan mittaa keskiarvon alapuolella.

3.3 Lineaarisen muunnoksen vaikutus muuttujan fraktiilin määräämiseen

Muuttujan X fraktiilin lauseke voidaan johtaa standardoidun muuttujan Z lausekkeesta:

$$Z = \frac{X - \mu}{\sigma} \Leftrightarrow X^{(p)} = \mu + \sigma Z^{(p)}.$$

Eli kerrotaan standardointikaavan molemmat puolet σ :lla ja asetetaan X :n ja Z :n yläindekseiksi toisiaan vastaavat fraktiilin (prosenttipisteiden) arvot.

• **ESIM:** Muuttuja $X \sim N(100, 15^2)$. Määritä 38 %:n fraktiilipiste.

Käsin laskemalla määritetään liitteen 1 taulukon avulla $X^{(38)}$:aa vastaava Z :n arvo eli $Z^{(38)}$. Tämä tapahtuu siten, että etsitään ensin taulukosta 38%:a likimäärin vastaava kertymäfunktion arvo. Tämän jälkeen katsotaan sitä vastaava Z :n arvo, eli $Z^{(38)}$, joka on $-0,31$. Lopuksi sijoitetaan arvot kaavaan ja lasketaan — eli: $100 + 15 \times (-0,31) = \underline{95,35}$.

Excelissä syötetään soluun kaava `=NORMINV(0,38; 100; 15)` ja saadaan tulokseksi 95,42.

3.3.1 Lineaarisen muunnoksen vaikutus muuttujan todennäköisvyksien määräämiseen

Kuten on jo todettukin, normaalijakaumaa noudattavan muuttujan X arvoihin liittyviä todennäköisyyksiä voidaan laskea standardoimalla ne Z :n arvoiksi ja käyttämällä standardoidun normaalijakauman kertymäfunktiota. Eli yleisesti esitettynä

$$P(X \leq x_c) = P\left(Z \leq \frac{x_c - \mu}{\sigma}\right) = \Phi\left(\frac{x_c - \mu}{\sigma}\right),$$

jossa x_c tarkoittaa standardoitavaa arvoa. Tarkastellaan asiaa esimerkin valossa.

- **ESIM:** Mikä on todennäköisyys, että satunnaisesti poimittu yksilö saa pistemäärään, joka on suurempi kuin 14, eräässä hyvinvointia tutkivassa mittarissa, jonka tiedetään jakautuvan perusjoukon tasolla normaalisti parametrein $\mu = 10$ ja $\sigma = 3$?

Käsin laskettaessa merkitään ensin — ja lasketaan

$$P(X > 14) = P\left(Z > \frac{14 - 10}{3}\right) = P(Z > 1,333) .$$

Eli suoritetaan lineaarinen muunnos käänteisesti jakaumasta $X \sim (10, 3^2)$ jakaumaan $Z \sim (0, 1)$. Sitten katsotaan liitteen 1 taulukosta Z :n arvoa vastaava todennäköisyys ja vähennetään se kokonaistodennäköisyydestä (=1) eli

$$P(X > 14) = 1 - P\left(Z > \frac{14 - 10}{3}\right) = 1 - \Phi\left(\frac{14 - 10}{3}\right) =$$

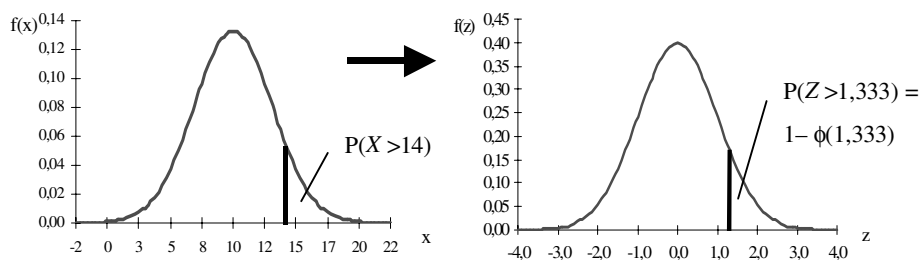
$$1 - \Phi(1,333) = 1 - 0,9082 = \underline{\underline{0,0918}} .$$

Excelissä laskettaessa syötetään soluun kaava:

$$= 1 - \text{NORMDIST}(14; 10; 3; 1) = \underline{\underline{0,0912}} .$$

Kuva 22 havainnollistaa sen, mitä esimerkin laskuoperaatioissa tapahtuu.

Kuva 22: Muuttujan $X \sim N(10, 3^2)$ standardoiminen



4 Otantajakauma

4.1 Yleistä

Ennen kuin voimme varsinaisesti käsitellä (luvussa 5) estimointia ja tunnuslukujen *estimaattoreita*, on määriteltävä otantajakauman käsite ja muutama keskeinen otantajakaumatyyppi sekä tutkittava jonkin verran niiden ominaisuuksia. Koska otoksen perusteella emme voi saada tarkasti selville perusjoukon (eli populaation) tunnuslukujen arvoja, on tyydyttävä otoksesta laskettuihin tunnuslukuihin, eli *estimaatteihin*, jotka mahdollisimman luotettavasti kuvaisivat perusjoukon tunnuslukujen todellisia arvoja. Perusjoukon tunnusluvun arvoa estimaattaessa, tapahtuu kuitenkin estimaatin eli otostunnusluvun arvossa vaihtelua otoksesta toiseen, jolloin ko. tunnusluvun eri estimaatit muodostavat jakauman. Näin muodostunutta jakaumaa kutsutaan estimaattorin otantajakaumaksi. Estimaatin ja populaation tunnusluvun erotusta kutsutaan

$$\text{Otantavirhe } e = \text{Estimaatti} - \text{Populaation tunnusluku.}$$

4.2 Otoskeskiarvon otantajakauma

Tarkastellaan otantajakauman ja otantavirheen käsitteitä tarkemmin esimerkin avulla, jossa on simuloitu otoskeskiarvon satunnaiskäyttäytymistä (Taulukko 3). Poimitaan eräästä populaatiosta, joka $\sim N(5, 3^2)$, kymmenen seitsemän tilastoyksikön satunnaisotosta ja tutkitaan otoskeskiarvon \bar{X} kykyä estimoida populaation keskiarvoa μ ja siinä havaittavaa satunnaisvirhettä, otantavirhettä e .

Taulukko 3: Otoskeskiarvon satunnaiskäyttäytymisen simulointi

Otoksen järj. nro	Satunnaisotokset populaatiosta (n = 7)							\bar{X}	s	$e = (\bar{X} - \mu)$
1	6	2	8	1	3	5	8	4,71	2,81	0,29
2	10	4	10	4	7	3	8	6,57	2,94	-1,57
3	0	0	4	9	8	6	3	4,29	3,59	0,71
4	8	9	6	1	1	2	4	4,43	3,31	0,57
5	4	5	1	1	9	5	1	3,71	2,98	1,29
6	8	8	3	5	10	2	3	5,57	3,10	-0,57
7	1	6	2	8	7	4	9	5,29	3,04	-0,29
8	3	2	10	7	9	8	10	7,00	3,27	-2,00
9	7	9	10	10	5	1	5	6,71	3,30	-1,71
10	5	0	7	2	9	0	8	4,43	3,78	0,57
Otantajakauman ja -virheen otoskeskiarvot								5,27		0,27
Otantajakauman ja -virheen keskihajonta								1,15		1,15

Taulukkoa 3 tutkimalla voidaan havaita, että *otantajakauman otoskeskiarvo 5,27* tulee lähelle *populaation keskiarvoa 5,0* jo seitsemän tilastoyksikön otosten pohjalta ja että otantajakauman otoskeskihajonta **1,15** lähestyy sen teoreettista arvoa **1,134**, joka saadaan laskettua seuraavasti:

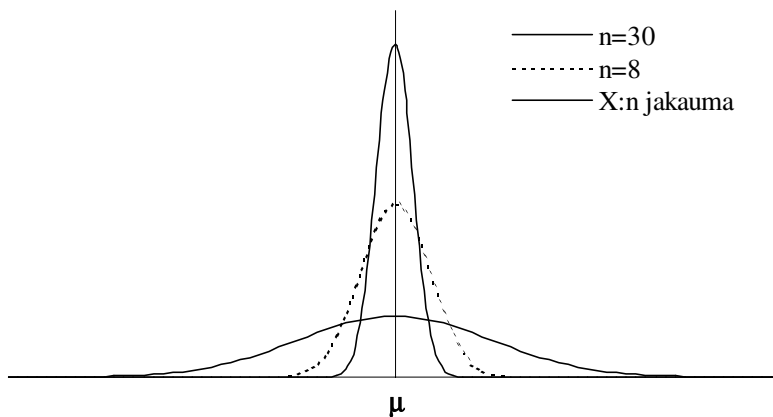
$$D(\bar{X}) = \sqrt{\frac{\sigma^2}{n}} = \sqrt{\frac{9}{7}} = 1,134.$$

Taulukosta voidaan edelleen havaita, että virhemuuttujan **e** otoskeskiarvo tulee lähelle nollaa (**0,27**) ja sen havaittu otoskeskihajonta on yhtä suuri kuin otantajakauman (**1,15**).

Otantajakauman varianssin lauseke pitää sisällään estimoinnin kannalta hyvin tärkeän asian: nimittäin, kun otoskoko kasvaa, pienenee vastaavasti otantajakauman varianssi $D^2(\bar{X})$ — ja täten pienenee myös \bar{X} :n estimaattien vaihtelu μ :n ympärillä, mikä puolestaan johtaa tarkempiin μ :n estimaatteihin. Oleellista on oivaltaa, että otantavirhe pienenee, kun otoskoko **n** kasvaa, kuten Kuva 23 havainnollistaa. Kuva 23:ssä on piirrettynä X :n jakauma sekä siitä muodostettujen otoskeskiarvojen otantajakaumat otoskoon arvoilla 8 ja 30.

Otoskeskiarvon otantajakauma on normaalisti jakautunut $N(\mu, \sigma^2/n)$, mikäli X jakautuu populaatiossa normaalisti. Kuitenkin, vaikkei X :n jakauma olisikaan perusjoukossa normaalisti jakautunut, noudattaa \bar{X} :n otantajakauma likimäärin normaalijakaumaa, mikäli otoskoko on kyllin suuri. Yleisesti pidetään otoskoko 30 sinä rajana, jonka jälkeen \bar{X} :n otantajakauma $\sim N(\mu, \sigma^2/n)$.

Kuva 23: \bar{X} :n otantajakauma otoskoon eri arvoilla



4.3 Prosenttiluvun otantajakauma

Otosprosenttiluku \mathbf{p} estimoii populaation prosenttiluvun $\boldsymbol{\theta}$ arvoa. Symbolit $\boldsymbol{\theta}$ ja \mathbf{p} ilmaisevat tietyn tutkittavan ominaisuuden \mathbf{A} suhteellisen osuuden populaatiossa ja vastaavasti sen estimaatin otoksessa. $\boldsymbol{\theta}$ ja \mathbf{p} ovat ajateltavissa dikotomisen (2-luokkaisen) tai 0–100 välillä vaihtelevan muuttujan keskiarvoiksi. Tällöin kun tilastoyksiköllä on ominaisuus \mathbf{A} , hän saa arvon 100 (tai 1) ja kun tilastoyksiköllä ei ole ominaisuutta \mathbf{A} , hän saa arvon 0.

Seuraavassa 22 tilastoyksikön populaatiossa on $\boldsymbol{\theta}$:n arvo 33%, eli ominaisuuden \mathbf{A} (=1) suhteellinen osuus on 8/22. Erään yksinkertaisen satunnaisotoksen (n=7) pohjalta saatu otosprosenttiluku $\mathbf{p} = 2/7 \times 100 = 28,6\%$.

1	0	0	1
0	0	1	0
1	0	0	0
0	1	0	1
0	0	1	0

Toistettaessa yksinkertaista satunnaisotantaa, esiintyy otosprosenttiluvun \mathbf{p} arvoissa $\boldsymbol{\theta}$:n ympärille keskittyvää *satunnaisvaihtelua* otoksesta toiseen. Näin muodostuu otosprosenttiluvun otantajakauma. Otosprosenttiluvun otantajakauma on summamuuttujan kaltainen keskeinen raja-arvolauseen mukaan, sillä mitattavan ominaisuuden \mathbf{A} satunnaisesti vaihteleva osuus (lukumäärä) suhteessa otoskokoan muodostuu kussakin otoksessa summautuvasti — eli yleisesti:

$$\mathbf{p}_i = \frac{\text{Ominaisuus } \mathbf{A}_{\text{havaittu frekvenssi}}}{\mathbf{n}_{\text{otos}}} \cdot 100 .$$

Muodostuva \mathbf{p} :n otantajakauma ei ole ilman muuta normaalisti jakautunut, vaan noudattaa pienillä otoskoon arvoilla binomijakaumaa, kuten seuraavassa luvussa 0 tulemme toteamaan. \mathbf{p} :n otantajakauman likimääräinen normalisuus riippuu populaation prosenttiluvun $\boldsymbol{\theta}$ arvosta — ja *keskeisen raja-arvolauseen* mukaan otoskoosta. Täten, jotta \mathbf{p} :n noudattaisi normaalijakaumaa, täytyy populaation prosenttiluvun $\boldsymbol{\theta}$ ja otoskoon välillä olla Taulukko 4:n esittämä suhde. Lisäksi oletetaan ominaisuuden \mathbf{A} omaavien tilastoyksiköiden populaation olevan suhteellisen suuri (jopa ääretön) tai pienen populaation tapauksessa (kuten ylläesitetyn 22 tilastoyksikön tapauksessa), että jokainen poimittu yksikkö palautetaan takaisin populaatioon ennen seuraavan yksikön poimintaa. Ilman tätä jo poimitun yksikön palautusta (pienen populaation kohdalla) \mathbf{p} :n otantaja-

kauma ei noudata binomijakaumaa vaan sen sijaan hypergeometristä jakaumaa, jota emme käsittele tässä kirjassa.

Taulukko 4: Prosenttiluvun otantajakauman normaalijakaumaan perustuvan approksimaation sallittavuus eri θ :n ja otoskoon arvoilla

Pienempi luvuista θ , $100-\theta$	Otoskoon n alaraja
50	30
40	50
30	80
20	200
10	600
5	1400

Taulukko 4:een liittyen perusjoukon prosenttiluvuista katsotaan aina sitä kumpi on pienempi. Esimerkiksi, jos $\theta = 70$, otetaan sen komplementtiosa eli 30 ($100-70 = 30$).

Riippumatta n :n arvosta ovat otosprosenttiluvun p odotusarvo (teoreettinen keskiarvo) ja varianssi:

$$E(p) = \theta \quad \text{ja} \quad D^2(p) = \frac{\theta(100 - \theta)}{n} .$$

Näin ollen yksinkertaisella satunnaisotoksella määrätyn prosenttiluvun p otantajakauma

$$\sim N(\theta, \theta(100-\theta)/n) ,$$

mikäli n ja θ toteuttaa Taulukko 4:n ehdot. Voidaan todeta lisäksi, että estimaattorin p (teoreettinen) varianssi $D^2(p)$ riippuu θ :sta lausekkeen $\theta(100-\theta)$ kautta, joka θ :n eri arvoilla saa Taulukko 5:n esittämiä arvoja.

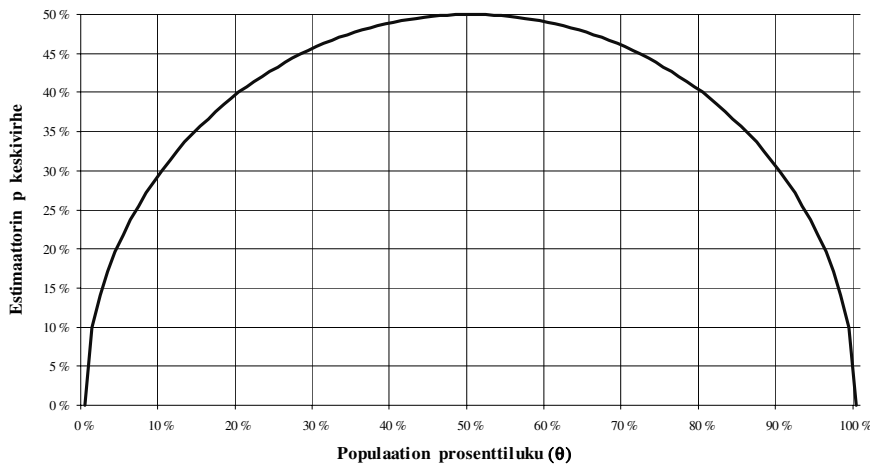
Taulukko 5: Estimaattorin p varianssin riippuvuus θ :n arvosta, kun $n = 1$

%-luku					
θ	50	40, 60	30, 70	20, 80	10, 90
$\theta(100-\theta)$	2500	2400	2100	1600	900
$\sqrt{\theta(100-\theta)}$	50	49	45,8	40	30

Tarkastelemalla taulukkoa huomaamme, että mitä lähempänä θ on 50:tä sitä suurempaa on p :n vaihtelu otoksesta toiseen.

Kuva 24 esittää, kuinka populaation prosenttiluvun estimaattorin p keskivirhe (ts. p :n keskihajonta) on kuvattavissa populaation prosenttiluvun epälineaarisena funktiona.

Kuva 24: Populaation prosenttiluvun θ ja sen estimaattorin p keskivirheen välinen yhteys, kun $n = 1$



4.4 Binomijakauma

Binomijakauma on luonteeltaan *diskreetti* todennäköisyysjakauma. Se liittyy sellaiseen satunnaisilmiöön, jolla on vain kaksi tulosmahdollisuutta, kuten esimerkiksi oikein-väärin –tyyppiset vastaukset. Usein binomijakauman yhteydessä mainitaan tästä syystä *Bernoullin koe*, joka määritellään toistokokeeksi, jossa on vain kaksi tulosmahdollisuutta — joko tapahtuma A esiintyy ($=1$) tai ei esiinny ($=0$). Tapahtumaan A voidaan sitten liittää tietty esiintymistodennäköisyys. Kun tunnetaan tapahtuman A esiintymistodennäköisyys, muodostuu Bernoullin koetta toistettaessa satunnaismuuttujan X jakauma, jonka arvot ovat esiintyneiden tapahtumien A lukumääriä. Tätä jakaumaa kutsutaan binomijakaumaksi. Kun satunnaismuuttuja X noudattaa binomijakaumaa, merkitsemme: $X \sim \text{Bin}(n,p)$, jossa parametri p tarkoittaa tapahtuman A todennäköisyyttä ja parametri n Bernoullin kokeen toistojen lukumäärää.

4.4.1 Binomijakauma ja otosprosenttiluvun otantajakauma

Binomijakauman parametrin p estimaattori on otoksesta laskettu suhteellinen frekvenssi eli $\hat{p} = X/n$, joka (sadalla kerrottuna) on siis sama kuin populaation prosenttiluvun estimaattori, kuten edellisessä luvussa 4.3 jo todettiin. Näin päädytään binomijakauman ja otosprosenttiluvun p otantajakauman väliseen yhteyteen. Binomijakauma voidaan tästä syystä käsittää otosprosenttiluvun matemaattiseksi malliksi. Myöhemmin tässä luvussa (ks. 4.4.3, 4.4.6 ja Kuva 26) tullaan toteamaan, että binomijakauman kuvaaja lähestyy parametrin n kasvaessa normaalijakaumaa. Milloin otosprosenttiluvun otantajakauma ei noudata riittävän tarkasti normaalijakaumaa, ts. kun Taulukko 4:n ehdot otoskoon suhteen eivät toteudu, voidaan otosprosenttiluvun ei-normaalista satunnaiskäyttämistä kuvata binomijakauman avulla.

4.4.2 Binomijakauman pistetodennäköisyysfunktio

Binomijakaumaa noudattavan satunnaismuuttujan X arvojen esiintymistodennäköisyydet määrätään seuraavan *pistetodennäköisyysfunktion* avulla

$$b(k; n, p) = \binom{n}{k} p^k q^{n-k},$$

missä $q = 1 - p$ ja $k = 0, 1, \dots, n$.

Binomijakauman pistetodennäköisyysfunktion avulla voidaan määrätä todennäköisyys, että ilmiön tapahtuessa n kertaa, esiintyy tapahtuma A , jonka todennäköisyys on p , k kertaa. Lukuja n ja p kutsutaan siis binomijakauman parametreiksi.

- **ESIM:** Erään ATK-laiteyrityksen valmistamista tuotteista $n. 10\%$ palautetaan takuuajana valmistevian vuoksi. Poimitaan myyntiin menevistä laitteista seitsemän kappaleen satunnaisotos. Mikä on todennäköisyys, että saadaan vähintään 2 viallista tuotetta?

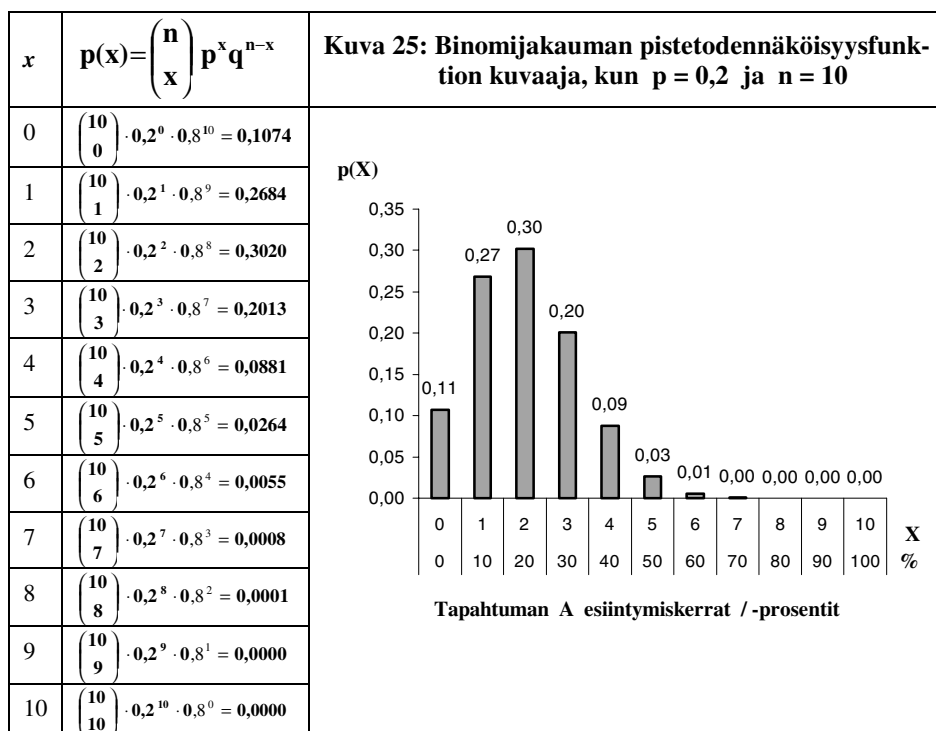
Merkitään $P(X \geq 2) = 1 - P(X < 2) = 1 - P(X \leq 1)$ — ja lasketaan:

$$1 - p(x) = \binom{7}{1} \cdot 0,1^1 \cdot 0,9^{7-1} = 1 - \frac{7!}{1!(7-1)!} \cdot 0,1 \cdot 0,9^6 = 1 - 0,3720 = 0,6280$$

Kysytty todennäköisyys siis on 0,6280.

Excelissä kirjoitetaan funktio: $=1 - \text{BINOMDIST}(1;7;0,1;FALSE) =$
0,6280.

Muodostetaan satunnaismuuttujan $X \sim \text{Bin}(10, 0,2)$ teoreettinen binomitodennäköisyysjakauma ja esitetään se graafisesti (Kuva 25).



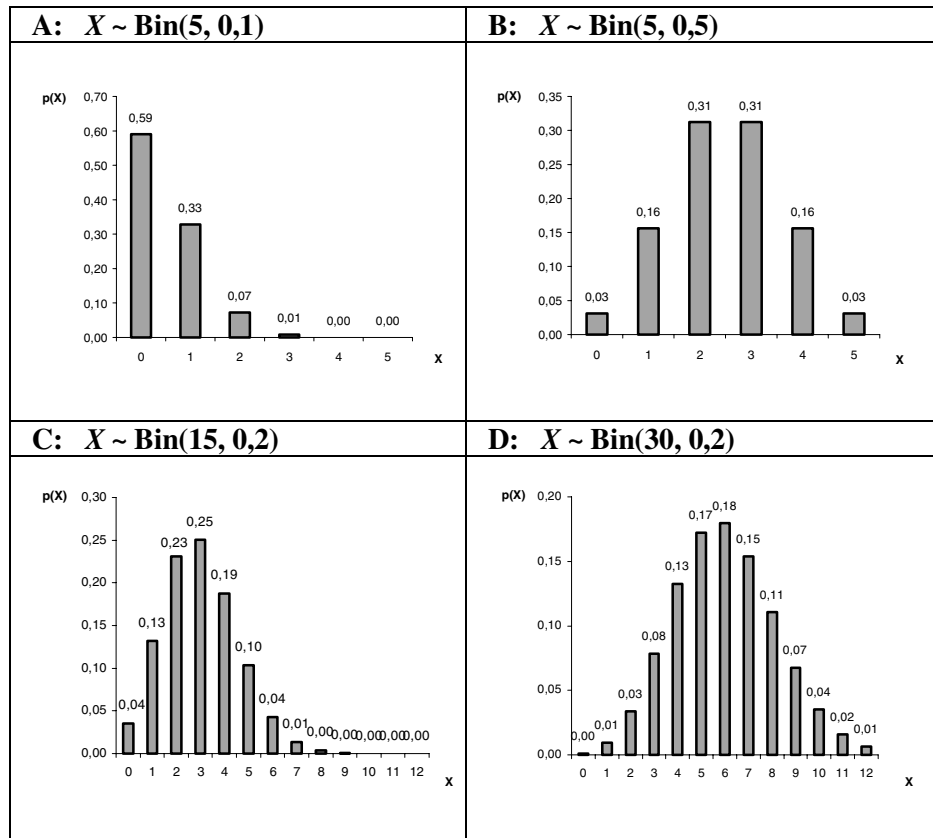
Pistetodennäköisyysfunktion kuvaajassa (Kuva 25) x -akselilla on satunnaismuuttujan X toistokokeessa mahdollisesti saamat arvot (eli tapahtuman A esiintymiskerrat yhdessä toistokokeessa) ja niiden alla vastaavat prosentiosuudet, joiden esiintymistiheyden todennäköisyys esitetään puolestaan y -akselilla.

Tarkasteltaessa pistetodennäköisyysfunktion arvoja suhteessa eri X :n arvoihin, voidaan todeta (otosprosenttiluvun mallintamisen kannalta), että todennäköisyys perusjoukon prosenttiluvun ali- tai yliestimoitumiseen (eli että saadaan ominaisuuden A estimaatiksi jokin muu kuin 2) on kohtalaisen suuri. Tällöin otantavirhe on helposti 10 %:n luokkaa. Tähän vaikuttaa oleellisesti otoskoko, jonka kasvaessa otantavirhe binomijakauma:otantavirhe pienentyy. Eli esimerkiksi jos $p=0,2$ ja $n=10$, kuten yllä, vaihtelee ominaisuuden A lukumäärän ja otoskoon osamäärä melko paljon, mikäli A :n lukumäärä otoksessa (tai toistokokeessa) poikkeaa odotusarvosta 2 vain yhdenkin A :n esiintymiskerran verran (esim. $1/10=0,1$). Mutta jos n onkin esimerkiksi 100, voi A :n suhteellinen lukumäärä otoksessa (tai toistokokeessa) poiketa odotusarvosta 20 useammankin A :n esiintymiskerran verran ilman, että osamäärä poikkeaa odotusarvostaan niin paljon (esim. $16/100=0,16$). Toisin sanoen otantavirhe pienenee otoskoon kasvaessa.

4.4.3 Binomijakauman muodon riippuvuus parametreistä n ja p

Binomijakauman vinous vähenee ja symmetrisyys lisääntyy sitä enemmän mitä lähempänä parametrin p arvo on lukua 0,5. Myös parametrin n arvon suurentaminen vaikuttaa binomijakaumaan symmetrisyyttä lisäävästi. Kuva 26 havainnollistaa tätä neljän eri erityistapauksen avulla.

Kuva 26: Esimerkkejä binomijakaumasta parametrien n ja p eri arvoilla.



4.4.4 Binomijakauman kertymäfunktio

Binomijakauman kertymäfunktio satunnaismuuttujalle X , joka $\sim \text{Bin}(n, p)$, on seuraava:

$$F(k) = P(X \leq k) = \sum_{i=0}^k \binom{n}{i} p^i q^{n-i}.$$

Binomijakauman kertymäfunktiota koskevat samat yleiset ominaisuudet, kuin muitakin kertymäfunktiota: sen avulla voidaan laskea jo yllä (ks. luku 2.4.3.2) esitetyn tyyppisiä todennäköisyyksiä. Tarkastellaan seuraavaksi esimerkin avulla binomijakauman kertymäfunktion käyttöä.

- **ESIM:** Määrää kertymäfunktion arvo satunnaismuuttujalle $X \sim \text{Bin}(15, 0,3)$ pisteessä 6.

Koska binomijakauman kertymäfunktion arvo lasketaan pisteeseen k mennessä kertyneiden pistetodennäköisyysfunktion arvojen summana, olisi kovin työlästä laskea yksittäiset arvot ja muodostaa sitten niiden summa. Sen sijaan voidaan katsoa liitteen 3 taulukosta suoraan ko. arvo, kun tiedetään n , p ja k . Tässä tapauksessa kertymäfunktion arvoksi saadaan 0,8689.

Laskettaessa Excelin funktiolla kirjoitetaan:
=BINOMDIST(6;15;0,3;TRUE) = 0,8689.

- **ESIM:** Määrää satunnaismuuttujan $X \sim \text{Bin}(12, 0,4)$ se arvo (kriittinen arvo), jossa binomijakauman kertymäfunktion arvo on suurempi tai yhtä suuri kuin 0,95.

Kertymäfunktion arvoa 0,95 vastaava satunnaismuuttujan X arvo on se jakauman yläalueen luku, jonka voimme odottaa esiintyvän toistokokeessa 5 %:n todennäköisyydellä. Toistojen lukumäärän ollessa 12 ja tapahtuman A esiintymistodennäköisyyden 0,4 voimme liitteen Taulukkoa 4 hyväksikäyttäen määrätä ko. satunnaismuuttujan (kriittisen) arvon), joka on 8.

Käytettäessä Excelin funktiota kirjoitetaan:
=CRITBINOM(12;0,4;0,95) = 8.

4.4.5 Binomijakauman odotusarvo ja varianssi

Binomijakauman odotusarvolle ja varianssille on olemassa seuraavat lausekkeet:

$$E(X) = np \quad \text{ja} \quad D^2(X) = npq.$$

Muuttuja X noudattaa binomijakaumaa parametreilla $n=40$ ja $p=0,25$. Tällöin muuttujan X odotusarvoksi saadaan $40 \cdot 0,25 = 10$ ja keskihajonnaksi vastaavasti $\sqrt{40 \cdot 0,25 \cdot 0,75} = 2,739$.

4.4.6 Binomi- ja normaalijakauman välinen yhteys

Kun satunnaismuuttuja $X \sim \text{Bin}(n,p)$, tällöin standardoitu satunnaismuuttuja

$$Z = \frac{X - np}{\sqrt{npq}},$$

jossa $q = 1-p$, noudattaa otoskoon n kasvaessa yhä tarkemmin standardoitua normaalijakaumaa $N(0,1)$. Kuva 26 havainnollistaa, kuinka otoskoon kasvaessa riittävästi binomijakauma alkaa lähestymään normaalijakaumaa. (ks. myös Taulukko 4).

4.4.7 Excelin funktiot, jotka liittyvät binomijakaumaan

=BINOMDIST(x; n; p ; TRUE / FALSE) antaa binomijakauman pistetodennäköisyysfunktion arvoja mikäli funktion viimeinen argumentti on FALSE ja kertymäfunktion arvoja mikäli funktion viimeinen argumentti on TRUE. x on tietyn tapahtuman A esiintymiskerrat toistokokeessa, jossa toistojen lukumäärän ilmaisee parametri n . Parametri p ilmaisee puolestaan ko. tapahtuman esiintymistodennäköisyyden. (suom. BINOMIJAKAUMA)

=CRITBINOM(n; p; α) antaa binomijakauman kriittisiä arvoja. Funktion ensimmäinen parametri n ilmaisee kokeen toistojen lukumäärän ja p tapahtuman A esiintymistodennäköisyyden. Viimeinen argumentti α asettaa kriittistä arvontodennäköisyyden. (suom. BINOMIJAKAUMA.KRIT)

4.5 *t*-jakauma

4.5.1 *t*-jakauman yhteydet muihin jakaumiin

Kuten jo tiedetään, standardoitu otoskeskiarvo $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ noudattaa $N(0,1)$ -

jakaumaa. Käytännössä on usein niin, että populaation keskihajonta σ on tuntematon, jolloin se joudutaan estimoimaan otoksesta, eli laskemaan otoskeski-

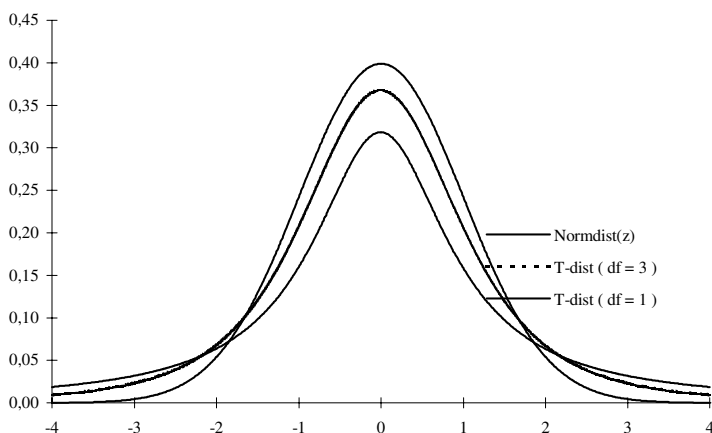
hajonta $s = \sqrt{\frac{\sum(x_i - \bar{X})^2}{n-1}}$, joka sijoitetaan standardoidun otoskeskiarvon kaavaan

σ :n tilalle ja saadaan *t*-satunnaismuuttuja

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}},$$

joka ei noudata $N(0,1)$ -jakaumaa, vaan *t*-jakaumaa. Vaikka *t*-jakauma muistuttaa suuresti normaalijakaumaa, on sen kuvaaja kuitenkin riippuvainen otoskoosta (vapausasteista $df = n-1$). Kuva 27 osoittaa, kuinka *t*-jakauma on pienillä otoskoon arvoilla matalampi ja leveämpi kuin normaalijakauma, ja kuinka se lähestyy normaalijakaumaa otoskoon kasvaessa.

Kuva 27: $N(0,1)$ -jakauma ja *t*-jakauma, kun $df = 1$ ja 3



Etenkin pienellä otoskoolla *t*-jakauman varianssi on suurempi kuin *Z*-jakauman, mistä syystä *t*-jakauma on origon kohdalla matalampi. Tämä johtuu otosvarianssin s^2 noudattaman otantajakauman, χ^2 -jakauman, suurehkosta

epäsymmetrisyydestä pienten vapausteiden **df** (otoskokojen) yhteydessä, jolloin otosvarianssin s^2 arvojen todennäköisyys olla pienempiä kuin vastaavan populaation varianssin σ^2 on lisääntynyt. Näin olen **t**-jakauman normaalijakaumaan nähden leventyneet hännät johtuvat juuri siitä, että populaation varianssin estimaatit ovat pienten otoskokojen yhteydessä taipuvaisia olemaan liian pieniä. Tällöin σ :n paikalle kaavassa $Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$ asetettuna ne olisivat jakajana liian pieniä ja tuottaisivat normaalijakauman esiintymistodennäköisyyksiin nähden liian usein itseisarvoltaan suuria **Z**:n arvoja. Tästä syystä, kun σ :n arvoa ei tunneta, täytyy käyttää $Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$:n sijasta $t = \frac{\bar{X} - \mu}{s / \sqrt{n}}$; ($df = n - 1$):tä, ja korjata **t**-jakauman muoto otoskokoa vastaavaksi.

4.5.2 **t**-jakauman fraktiilit

Fraktiilipisteiden arvotkaan eivät ole **t**-jakauman kohdalla vakioita, vaan riippuvaisia niin ikään vapausasteista — fraktiilien itseisarvot pienenevät vapausasteiden kasvaessa ja siten lähestyvät vastaavia normaalijakauman fraktiilien arvoja.

- **ESIM: Määrää t-jakauman 95 %:n fraktiilipisteet, kun otoskoot ovat 3, 10 ja 30.**

Taulukkoa (liite 2) käyttäen, etsitään **t**-jakauman fraktiilit $t^{(95)}$ vastaten jakauman 95 %:n pistettä ja annettuja otoskokoja (vapausasteita $df = n - 1$).

Saadaan $t_{df=2}^{(95)} = 4,303$; $t_{df=9}^{(95)} = 2,262$; $t_{df=29}^{(95)} = 2,045$.

Excelin funktiota käyttäen kirjoitetaan soluun: $=TINV(0,05; 2) = 4,303$.
 $=TINV(0,05; 9) = 2,262$.
 $=TINV(0,05; 29) = 2,045$.

4.5.3 t-jakauman tiheysfunktio

On huomattava, että **t**-jakauma ei ole minkään tietyn estimaattorin otantajakauma, vaan estimaattoreista \bar{X} ja **s** johdetun lausekkeen otantajakauma, joka siis muuttuu otoskoon funktiona. Tiheysfunktio **t**-jakaumalle on seuraava

$$s_k(\mathbf{t}) = c_k \left(1 + \frac{\mathbf{t}^2}{k}\right)^{-\frac{k+1}{2}},$$

jossa c_k on **t**:stä riippumaton normeerausvakio. Sen arvo määrätään siten, että $\int_{-\infty}^{\infty} s_k(\mathbf{t}) d\mathbf{t} = 1$. Kokonaislukua **k** kutsutaan vapausasteluvuksi (degree of freedom), jonka arvo vaihtelee sen mukaan, mistä **t**-muuttujatyypistä (testisuureesta) on kysymys.

4.5.4 Excelin t-jakaumaan liittyvät funktiot

=TDIST(x; k; 1 tai 2) antaa **x**:ään liittyvän todennäköisyyden arvon ko. vapausasteilla **k**. Viimeisen argumentin arvo voi olla joko 1 tai 2. Arvo 1 on yksisuuntaisen testin todennäköisyys ja arvo 2 vastaavasti kaksisuuntaisen testin todennäköisyys.
(suom. TJAKAUMA)

=TINV(p; k) antaa **t**-jakauman kriittisiä arvoja 2-suuntaisen testin tapauksessa todennäköisyydelle **p** vapausasteilla **k**.
(suom. TJAKAUMA.KÄÄNT)

=TTEST(Alue1; Alue2; suunta; 1, 2 tai 3) suorittaa kahden riippumattoman otoksen tai riippuvien otosten **t**-testit riippuen viimeisestä argumentin arvosta. Kun viimeinen argumentin arvo on 2 saadaan homogeneisten ryhmävarianssien **t**-testin riskitason todennäköisyys; kun arvo 3 saadaan heterogeenisten ryhmävarianssien **t**-testin riskitason todennäköisyys; kun arvo 1 saadaan riippuvien otosten **t**-testin riskitason todennäköisyys. Toinen argumentti eli ”suunta” määrää onko testi 1- vai 2-suuntainen (arvo 1 =yksisuuntainen; arvo 2 =kaksisuuntainen). Alue1 ja Alue2 viittaavat solualueisiin, joissa testattava data sijaitsee.
(suom. TTESTI)

4.6 χ^2 -jakauma — eli otosvarianssin s^2 otantajakauma

t -jakauman yhteydessä viitattiinkin jo alustavasti χ^2 -jakaumaan (lue: khiin neliö), kun todettiin, että etenkin pienillä otoskoon arvoilla otosvarianssin s^2 otantajakauma on epäsymmetrinen, ja että näin ollen lisääntyy todennäköisyys saada liian pieniä σ^2 -estimaatteja. Vaikka χ^2 -jakaumaa kutsutaankin otosvarianssin s^2 otantajakaumaksi, ei sen sovellettavuusalue kuitenkaan rajoitu yksinomaan populaation varianssista σ^2 tehtäviin johtopäätöksiin, vaan sitä voidaan soveltaa monenlaisiin tilastollisiin testitilanteisiin.

Edellä käsitellyt estimaattorit \bar{X} ja \mathbf{p} olivat havaintojen x_1, \dots, x_n lineaarisia lausekkeita. Varianssi s^2 puolestaan on havaintojen suhteen kvadraattinen eli neliömuotoa oleva. Jotta otosvarianssin satunnaiskäyttäytyminen otoksesta toiseen voitaisiin hallita ja ymmärtää, tarkastellaan sen matemaattista mallia — ”standardoidun” neliösumman jakaumaa.

Olkoot Z_1, Z_2, \dots, Z_k riippumattomia, $\mathbf{N}(\mathbf{0}, \mathbf{1})$ -jakaumaa noudattavia satunnaismuuttujia. Korotetaan \mathbf{k} muuttujaa neliöön ja muodostetaan niiden summa eli $Z_1^2 + Z_2^2, \dots, + Z_k^2$.

Näin muodostuvaa ”standardoitua” neliösummaa kutsumme $\chi^2(\mathbf{k})$ -satunnaismuuttujaksi vapausasteluvulla \mathbf{k} . $\chi^2(\mathbf{k})$ -satunnaismuuttujan jakaumaa kutsumme juuri χ^2 -jakaumaksi (vapausasteluvulla \mathbf{k}).

4.6.1 χ^2 -jakauman tiheysfunktio

χ^2 -jakauman tiheysfunktion lauseke on seuraavanlainen

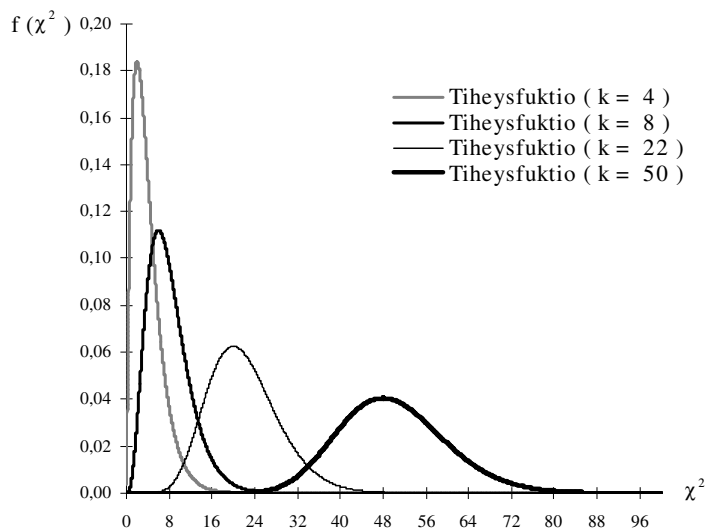
$$f(\chi^2) = \mathbf{h}_{\mathbf{k}} e^{-\frac{\chi^2}{2}} (\chi^2)^{\frac{\mathbf{k}}{2}-1},$$

jossa vakio $\mathbf{h}_{\mathbf{k}}$ määrätään siten, että $\int_0^{\infty} f(\chi^2) d\chi^2 = 1$. χ^2 -jakauma on määritel-

ty ainoastaan positiivisille arvoille, koska neliöiden summana χ^2 :n arvo on aina vähintään nolla. Vapausasteluku \mathbf{k} vaikuttaa voimakkaasti tiheysfunktion kuvaajan muotoon. Kun \mathbf{k} on 1 tai 2, on kuvaaja jatkuvasti oikealle laskeva käyrä. Kahta suuremmilla arvoilla tiheysfunktion kuvaaja ensin nousee **moodin** ollessa pisteessä $\chi^2 = \mathbf{k}-2$, minkä jälkeen kuvaaja on laskeva.

Kuva 28 havainnollistaa kuinka eri vapausasteen arvot vaikuttavat χ^2 -jakauman tiheysfunktion kuvaajaan.

Kuva 28: χ^2 -jakauman tiheysfunktion kuvaajia k :n eri arvoilla



Kuvaa tarkastelemalla voi todeta, että χ^2 -jakauman tiheysfunktion epäsymmetrisyys korostuu erityisesti pienillä vapausasteiden k arvoilla, jolloin sitä noudattavan otosvarianssin s^2 arvojen todennäköisyys olla pienempiä kuin vastaavan populaation varianssin σ^2 on korostunut.

Toisaalta kuvasta voi havaita myös sen yleisen ilmiön, mikä tuli jo t -jakaumankin kohdalla ilmi, että vapausasteiden kasvaessa, χ^2 -jakauman tiheysfunktion epäsymmetrisyys alkaa vähenemään ja tiheysfunktion kuvaaja alkaa lähestymään normaalijakaumaa. Tämä siis perustuu siihen, että $\chi^2(k)$ on summamuuuttuja. Normaalijakaumaan perustuva approksimaatio on kyllin tarkka 30:tä suuremmilla k :n arvoilla.

4.6.2 χ^2 -jakauman odotusarvo ja varianssi

Voidaan osoittaa, että muuttujan $\chi^2(k)$ odotusarvo ja varianssi ovat

$$E(\chi^2(k)) = k \quad \text{ja} \quad D^2(\chi^2(k)) = 2k,$$

joista on todettavissa, että kun vapausasteluku kasvaa, kasvaa $E(\chi^2(\mathbf{k}))$:n ja $D^2(\chi^2(\mathbf{k}))$:n arvot vastaavasti, koska \mathbf{k} ilmaisee $\chi^2(\mathbf{k})$ lausekkeessa olevien ei-negatiivisten yhteenlaskettavien lukumäärän.

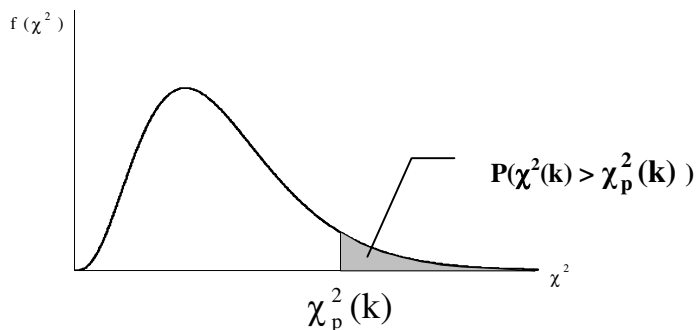
4.6.3 χ^2 -jakauman kriittiset arvot

$\chi^2(\mathbf{k})$ -jakauman teoreettinen kriittinen arvo $\chi_p^2(\mathbf{k})$ toteuttaa yhtälön

$$P(\chi^2(\mathbf{k}) > \chi_p^2(\mathbf{k})) = p.$$

Kuva 29 havainnollistaa vielä määritelmän. Huomattakoon tässä vielä, että esim. 90%:n fraktiili on 0,1:n kriittinen arvo.

Kuva 29: Kriittisen arvon $\chi_p^2(\mathbf{k})$ määritelmä



On jälleen huomattava, kuten oli asia t -jakaumankin kohdalla, että fraktiilit ja kriittisten arvot muuttuvat vapausasteiden mukaan.

- **ESIM: Määrää χ^2 -jakauman 65 %:n fraktiilipisteet vapausasteille 3 ja 25.**

Koska taulukkoarvoja ei ole yleisesti saatavilla kaikille fraktiiliarvoille, määrätään ko. pisteet Excelin funktiolla. Kirjoitetaan soluun kaavat

$$\begin{aligned} &= \text{CHIINV}(0,65; 3) = \underline{1,642}. \\ &= \text{CHIINV}(0,65; 25) = \underline{21,752}. \end{aligned}$$

Esimerkin tulokset voidaan myös ymmärtää kertymäfunktion käsitettä hyväksikäyttäen. Fraktiilipiste vastaa siis sitä pistettä, jonka vasemmalla puolella on

65% todennäköisyysmassaa. Ero vastausten suuruudessa — huolimatta siitä, että ne ovat samoja fraktilipisteitä — kertoo juuri, kuinka merkittävästi vapausasteet vaikuttavat χ^2 :n arvoihin.

- **ESIM: Määrää seuraavat todennäköisyydet, kun $k=13$:**
 $P(\chi^2(k) > 15)$, $P(3 < \chi^2(k) \leq 7)$ ja $P(\chi^2(k) < 15)$

Hyödynnetään jälleen Excelin funktiota ja kirjoitetaan soluun seuraavat kaavat:

$$=1 - \text{CHIDIST}(15; 13) = \underline{0,307}$$

$$=\text{CHIDIST}(3; 13) - \text{CHIDIST}(7; 13) = \underline{0,096}$$

$$=\text{CHIDIST}(15; 13) = \underline{0,693}$$

Tässä tehtävässä siis hyödynnetään niitä kertymäfunktion yleisiä ominaisuuksia, jotka koskevat periaatteessa kaikkia teoreettisia jakaumia sekä kerrataan ne todennäköisyystyypit, jotka kertymäfunktion avulla voidaan määrätä. Excelin CHIDIST –funktio antaa juuri χ^2 -jakauman kertymäfunktion arvoja.

4.6.4 Otosvarianssin s^2 ja χ^2 -jakauman välinen yhteys

Tähän mennessä ei ole vielä tarkasteltu otosvarianssin yhteyttä χ^2 -jakaumaan. Nämä yhteydet tulevat esille, kun kirjoitetaan s^2 seuraavissa muodoissa

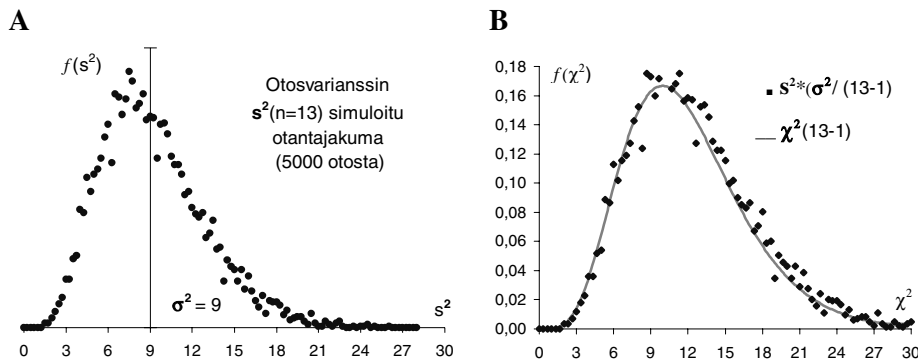
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{(n-1)} = \left[\frac{\sigma^2}{n-1} \cdot \sum_{i=1}^n \left(\frac{x_i - \bar{X}}{\sigma} \right)^2 \right] = \left[\frac{\sigma^2}{n-1} \cdot \sum_{i=1}^n Z_i^2 \right] = \left[\frac{\sigma^2}{n-1} \cdot \chi^2(n-1) \right],$$

jossa x_1, \dots, x_n ovat yksinkertainen satunnaisotos jakaumasta $N(\mu, \sigma^2)$. On osoitettavissa, että neliösumma $\sum_{i=1}^n ((x_i - \bar{X})/\sigma)^2$ eli $\sum_{i=1}^n Z_i^2$ noudattaa $\chi^2(n-1)$ -jakaumaa. Näin ollen otosvarianssi s^2 on esitettävissä vakion $\sigma^2/(n-1)$ ja teoreettisen muuttujan $\chi^2(n-1)$ tulona. Tämän lauseen kohdalla ei voida luopua populaation normaalisuusoletuksesta suurillakaan otoskoon arvoilla.

Havainnollistetaan lausekkeen yhteyksiä vielä graafisesti. Kuva 30 sisältää esitykset A ja B, joista ensimmäisessä on simuloitu otosvarianssin satunnaiskäyt-

täytymistä generoimalla tietokoneen avulla normaalijakaumasta $N(10,3^2)$ 5000 otosta ($n=13$). Kustakin otoksesta on laskettu varianssi ja tämän jälkeen muodostettu pienellä luokkavälillä s^2 -jakauma, joka siis demonstroi otosvarianssin satunnaisvaihtelua otoksesta toiseen (otoskoolla 13).

Kuva 30: Tietokoneella generoitu otosvarianssin otantajakauma (A) ja sitä mallintava χ^2 -jakauma (B)



Populaation varianssi σ^2 on 9, joka on merkitty kuvaan A pystyviivalla. Kuvas-
sa B esitetään sama generoitu s^2 -jakauma, mutta nyt kerrottuna vakiolla
 $\sigma^2/(n-1)$, jolloin s^2 -jakauma yhtyy teoreettiseen $\chi^2(n-1)$ -jakaumaan, kuten ku-
va B osoittaa. Teoreettinen $\chi^2(12)$ -jakauma siis mallintaa otosvarianssin empii-
ristä jakaumaa ja määrää siihen liittyvät todennäköisyydet. Esitettyjä tuloksia
hyväksi käyttäen on mahdollista mm. määrätä estimoitavalle populaation va-
rianssille σ^2 luottamusväli (ks. Luku 6.4).

4.6.5 Excelin funktiot, jotka liittyvät χ^2 -jakaumaan

=CHIDIST(χ^2 -arvo; vapausasteet) antaa kriittistä arvoa (χ^2 -arvo)
vastaavan todennäköisyyden ko. vapausasteilla.
(suom. CHIJAKAUMA)

=CHIINV(todennäköisyys; vapausasteet) antaa χ^2 -jakauman kriitti-
sen arvon, kun todennäköisyys ja vapausasteet tunnetaan.
(suom. CHIJAKAUMA.KÄÄNT)

=CHITEST(Alue1; Alue2) suorittaa χ^2 -riippumattomuustestin ja an-
taa siihen liittyvän riskitodennäköisyyden eli nk.p-arvon. **Alue 1:**lle va-
litaan havaittujen frekvenssien alue ja **Alue2:**lle odotettujen frekvenssi-
en alue.
(suom. CHITESTI)

4.7 F-jakauma — eli varianssitestien otantajakauma

Myöhemmin tässä kirjassa käsiteltävää varianssien yhtäsuuruustestiä sekä yleisemminkin toisaalla käsiteltäviä tilastollisia menetelmiä, kuten varianssianalyysiä varten, on määriteltävä **F-jakauma**, joka on varianssitestien testisuureiden otantajakauma. Kuten myöhemmin tulemme keskiarvotestien kohdalla huomaamaan, varianssitestejä tarvitaan mm. selvittämään, onko keskiarvotestien oletus ryhmävarianssien yhtäsuuruutta koskien voimassa. Tässä yhteydessä käsittelemme F-jakaumaa kuitenkin vain otantajakaumana.

Tarkastellaan kahta eri tilastoyksikköryhmien A ja B muodostamaa normaali-jakaumaa $\mathbf{N}(\boldsymbol{\mu}_1, \boldsymbol{\sigma}^2)$ ja $\mathbf{N}(\boldsymbol{\mu}_2, \boldsymbol{\sigma}^2)$, jotka kuvaavat tietyn muuttujan X arvojen jakautumaa näissä ryhmissä eli osapopulaatioissa. Varianssit ovat molempien osapopulaatioiden kohdalla samat, mutta keskiarvot *voivat* poiketa toisistaan. Poimitaan yksinkertaisella satunnaisotannalla näistä osapopulaatiosta toisistaan riippumattomat otokset $x_{11}, x_{12}, \dots, x_{1n}$ ja $x_{21}, x_{22}, \dots, x_{2n}$, joista laskettavat otosvarianssit voidaan esittää jo yllä χ^2 -jakauman yhteydessä esitetyn lauseen pohjalta seuraavasti

$$s_1^2 = \frac{\sum_j (x_{1j} - \bar{X}_1)^2}{(n_1 - 1)} = \frac{\sigma^2 \chi_1^2(\mathbf{k})}{(n_1 - 1)}$$

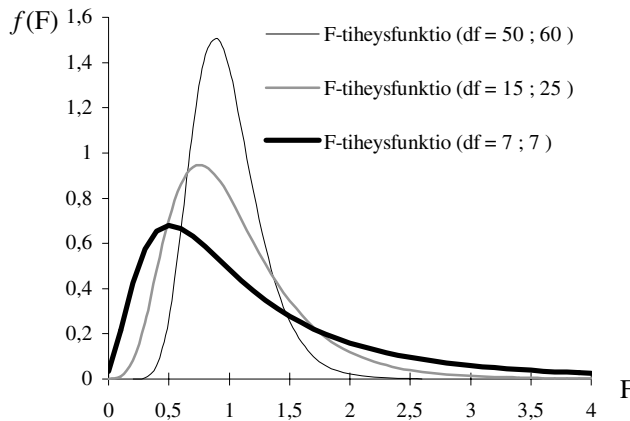
$$s_2^2 = \frac{\sum_j (x_{2j} - \bar{X}_2)^2}{(n_2 - 1)} = \frac{\sigma^2 \chi_2^2(\mathbf{k})}{(n_2 - 1)}.$$

Muodostettaessa otosvarianssien s_1^2 ja s_2^2 osamäärä

$$F = \frac{s_1^2}{s_2^2} = \frac{\left(\frac{\sigma^2 \chi_1^2(\mathbf{k})}{(n_1 - 1)} \right)}{\left(\frac{\sigma^2 \chi_2^2(\mathbf{k})}{(n_2 - 1)} \right)},$$

saadaan **F-satunnaismuuttuja**, joka noudattaa **F-jakaumaa** vapausasteluvuilla $\mathbf{k}_1 = n_1 - 1$ ja $\mathbf{k}_2 = n_2 - 1$, jotka puolestaan ovat **F:ssä** olevien otosvarianssien vapausasteluvut. Tästä jakaumasta käytetään merkintää **F(k₁ k₂)-jakauma**, jonka tiheysfunktion lausekkeen ohitamme. **F-satunnaismuuttujan** lauseketta tarkastellessa voimme myös todeta **F-jakauman** ja χ^2 -jakauman välisen yhteyden, mikä perustuu juuri siihen, että s^2 noudattaa χ^2 -jakaumaa. Kuva 31 esittää muutamia **F(k₁ k₂)-jakauman** tiheysfunktion kuvaajia eri vapausasteilla.

Kuva 31: $F(k_1, k_2)$ -jakauman tiheysfunktion kuvaajia k_1 :n ja k_2 :n eri arvoilla

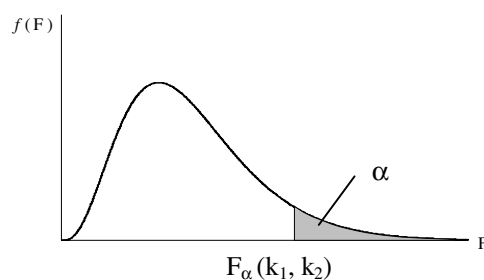


Estimaattoreiden s_1^2 ja s_2^2 osamäärä $F = s_1^2/s_2^2$ vaihtelee ykkösen molemmin puolin. Mitä suuremmat vapausasteluvut k_1 ja k_2 ovat sitä epätodennäköisemmin F :n arvo poikkeaa huomattavasti ykkösestä. Kuva 31 tuo tämän asian myös selvästi esiin, jossa vapausesteiden arvoilla $k_1 = 50$ ja $k_2 = 60$ tiheysfunktio keskittyy melko huipukkaasti ykkösen ympärille.

4.7.1 F-jakauman kriittiset arvot

F-jakauman sovelluksissa tarvitsemme jakauman kriittisiä arvoja eri vapausasteille k_1 ja k_2 . Kriittiset arvot saadaan hypoteesin testauksen yhteydessä liitteen 6 taulukosta tietyille riskitasoille tai voidaan määrätä Excelin FINV –funktion avulla yleisemminkin. Merkitään F-jakauman kriittistä arvoa $F_\alpha(k_1, k_2)$:lla joka toteuttaa yhtälön $P(F > F_\alpha(k_1, k_2)) = \alpha$, jossa α tarkoittaa haluttua riskitasoa.

Kuva 32: $F(k_1, k_2)$ -jakauman kriittinen arvo $F_\alpha(k_1, k_2)$



Kuva 32 esittää F-jakauman kriittisen arvon määrittelyn, jossa on nyt siis huomattava kahden vapausasteluvun yhtäaikainen vaikutus kriittisen arvon määräytymiseen. Yleisimmin F-jakaumaan liittyvien tilastomenetelmien koh-

dalla hypoteesin testaukset ovat yksisuuntaisia niin päin, että F-suhteen osamäärän oletetaan saavan vain ykköistä suurempia arvoja. Tällöin kriittinen arvo $F_{\alpha}(k_1, k_2)$ määrätään vain jakauman oikeasta 'hännästä', kuten Kuva 32:ssa. Kuitenkin, mikäli testitilanne niin vaatii, voidaan kriittinen arvo määrätä myös jakauman vasemmanpuoleisesta 'hännästä' seuraavalla kaavalla:

$$F_{1-\alpha}(k_1, k_2) = \frac{1}{F_{\alpha}(k_2, k_1)} .$$

Siten esimerkiksi

$$F_{0,95}(10, 4) = \frac{1}{F_{0,05}(4, 10)} = \frac{1}{3,48} = 0,28 .$$

On huomattava, että muodostettaessa $F_{\alpha}(k_1, k_2)$:n käänteisluku $1 / F_{\alpha}(k_2, k_1)$, täytyy vapausasteiden järjestys nimittäjässä kääntää, koska myös F-osamäärässä osoittaja ja nimittäjä (eli varianssit s_1^2 ja s_2^2) vaihtavat paikkaa.

Jos testitilanne vaatii kaksisuuntaista testausta, eli että F-suhteen osamäärä voi vaihdella ykkösen molemmin puolin, määrätään kriittiset arvot jakauman molemmista päistä juuri määritellyillä kriittisten arvojen lausekkeilla. Tällöin riskitaso α puolitetaan eli molemmissa päissä se on $\alpha/2$.

- **ESIM: Määrää F-jakauman kriittinen arvo $F_{\alpha}(k_1, k_2)$, kun $\alpha=0,05$ (95%:n fraktiilipiste), kun vapausasteet ovat $k_1=10$ ja $k_2=18$.**

Koska ko. kriittinen arvo löytyy yleensä taulukkokirjoista, etsitään k_1 :n arvo liitteen 6 taulukon sarakkeesta ja k_2 :n vastaavasti riviltä, joiden yhtymäkohdasta löytyy tarvittava kriittinen arvo eli $F_{0,05}(10, 18)$, joka on 2,412.

Excelin funktiolla laskettaessa kirjoitetaan soluun =FINV(0,05;10;18) = 2,412.

- **ESIM: Määrää $P(0.5 < F(k_1, k_2) \leq 4.2)$, kun k_1 ja k_2 ovat (4;5), (10;15) ja (100;100).**

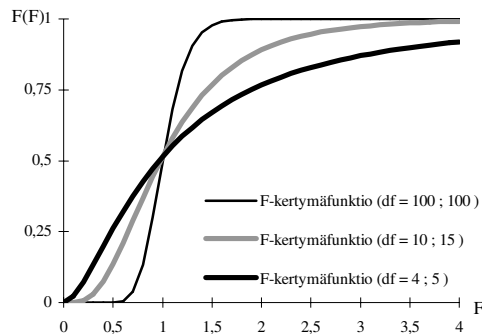
Koska Excelin funktioilla voidaan määrätä F-jakauman kertymäfunktion arvoja, niin voidaan laskea näin ollen F-jakaumaan liittyviä todennä-

köisyyksiä. Eli todennäköisyydet, että F :n arvo kuuluu välille 0,5–4,2 ko. eri vapausasteilla saadaan kirjoittamalla Excelin soluun seuraavasti

$$\begin{aligned} &= \text{FDIST}(0,5;4;5) - \text{FDIST}(4,2;4;5) &&= \underline{0,6655} \\ &= \text{FDIST}(0,5;10;15) - \text{FDIST}(4,2; 10;15) &&= \underline{0,8582} \\ &= \text{FDIST}(0,5;100;100) - \text{FDIST}(4,2;100;100) &&= \underline{0,9997} \end{aligned}$$

Voidaan todeta, että vaikka väli 0,5–4,2 säilyy vakiona sisältyy F -satunnaismuuttujan arvo siihen sitä suuremmalla todennäköisyydellä mitä suuremmat vapausasteiden arvot ovat. Ts. väli 0,5–4,2 pitää sisällään sitä suuremman todennäköisyysmassan mitä suuremmat vapausasteet ovat. Kuva 33 esittää kertymäfunktion kuvaajat esimerkkiin liittyvillä vapausasteilla.

Kuva 33: F-jakauman kertymäfunktion kuvaajia eri vapausasteilla



Kuvasta voi helposti havaita, kuinka vapausasteiden kasvattaminen lisää tietyn välin todennäköisyyttä pitää sisällään F -satunnaismuuttujan arvo.

4.7.2 Excelin funktiot, jotka liittyvät F-jakaumaan

=FDIST(x; k₁; k₂) antaa F -jakauman kriittiseen arvoon x liittyvän todennäköisyyden vapausasteilla k_1 ja k_2 (oikeanpuoleisen ”hännän” pinta-ala).

=FINV(α; k₁; k₂) antaa F -jakauman kriittisiä arvoja todennäköisyydelle α vapausasteille k_1 ja k_2 .

=FTEST(Alueen1; Alueen2) testaa Alueen1 ja Alueen2 arvojen varianssien yhtäsuuruuden yksisuuntaisella testillä antaen riskitodennäköisyyden eli nk. p-arvon.

5 Estimointi

5.1 Johdanto

Eräs tilastotieteen tärkeimmistä tehtävistä on kehittää menetelmiä, joiden avulla kyetään tekemään otoksen perusteella mahdollisimman luotettavia johtopäätöksiä perusjoukon tunnusluvuista eli parametreistä. Tätä tilastotieteen osa-aluetta kutsutaan *estimointiteoriaksi*. Kuten edellisen luvun alussa jo totesimme, emme otoksen perusteella voi saada tarkasti selville perusjoukon tunnuslukujen arvoja, vaan meidän on tyydyttävä otoksesta laskettuihin tunnuslukuihin eli *estimaatteihin*, jotka mahdollisimman luotettavasti kuvaisivat perusjoukon tunnuslukujen todellisia arvoja. Kutsumme tätä siis *estimoinniksi*. On kaksi erilaista tapaa estimoida perusjoukon tunnuslukuja: *piste-estimointi* ja *väliestimointi*.

5.2 Piste-estimointi

Piste-estimoinnissa pyritään etsimään *estimaattori*, jonka arvot ovat mahdollisimman lähellä populaation parametrin arvoa. Tietyllä parametrillä voi olla useita estimaattoreita. Esimerkiksi populaation keskiarvon μ estimaattorina voidaan käyttää otosmediaania M_d ja otoskeskiarvoa \bar{X} , joiden otantajakaumat kuvaavat niiden satunnaiskäyttäytymistä.

5.2.1 Estimaattorin ”hyvyys”

Tarkastellaan yleisesti parametriä T ja sen estimaattoria \hat{T} , joka on siis satunnaismuuttuja. \hat{T} :n tärkeimmät tunnusluvut ovat odotusarvo (= teoreettinen keskiarvo) ja varianssi, jotka voidaan laskea otantajakauman avulla. Näiden kahden tunnusluvun avulla voidaan arvioida estimaattorin ”hyvyyttä”. Estimaattorin odotusarvon avulla määritellään harhattomuus ja estimaattorin varianssin avulla *tarkkuus*.

Harhattomuus — Estimaattori \hat{T} on parametrin T harhaton estimaattori, jos \hat{T} :n odotusarvo on yhtä kuin T eli kaavana

$$E(\hat{T}) = T .$$

On huomattava, että harhattomuus on perusjoukosta riippumaton ominaisuus eli se on voimassa kaikille perusjoukoille. Sanallisesti harhattomuutta voidaan kuvailla vielä esim. siten, että harhaton estimaattori antaa keskimäärin oikeita tu-

loksia. Eli positiiviset ja negatiiviset otantavirheet kumoavat toisensa ”pitkässä juoksussa”.

Jos estimaattorin \hat{T} harha $B(\hat{T})$ määritellään seuraavasti (harha = bias)

$$B(\hat{T}) = E(\hat{T}) - T.$$

Jos estimaattori on harhainen eli $B(\hat{T}) \neq 0$, niin estimaattorissa on systemaattista virhettä (joko positiivista tai negatiivista). Lähes kaikki yleisimmät estimaattorit ovat harhattomia. Tunnetuin esimerkki *harhaisesta* estimaattorista on otoskeskihajonta s , joka *ei ole* σ :n harhaton estimaattori.

Taulukko 6 esittää muutamia keskeisiä harhattomia estimaattoreita.

Taulukko 6: Esimerkkejä estimaattoreista, jotka ovat harhattomia

\bar{X} on μ :n harhaton estimaattori
s^2 on σ^2 :n harhaton estimaattori
\mathbf{p} on $\boldsymbol{\theta}$:n harhaton estimaattori
\mathbf{r} on $\boldsymbol{\rho}$:n harhaton estimaattori

Tarkkuus — Toinen tapa määrittää eri estimaattorien keskinäistä paremmuutta on vertaamalla niiden otantajakaumien varianssien suuruutta toisiinsa.

Tarkastellaan asiaa vertaamalla kahta populaation keskiluvun estimaattoria \hat{M}_d :tä ja \bar{X} :ää. Olkoon populaation jakauma $\sim N(\mu, \sigma^2)$ tällöin ko. estimaattorien otantajakaumat ja niiden varianssit ovat seuraavat:

$$\hat{M}_d \sim N\left(\mu, \frac{\pi\sigma^2}{2n}\right) \quad \text{ja} \quad \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right).$$

Vertaamalla \hat{M}_d :n ja \bar{X} :n keskivirheitä, voimme todeta, että

$$\frac{\pi\sigma^2}{2n} > \frac{\sigma^2}{n},$$

eli \hat{M}_d :n varianssi on suurempi. Näin ollen johtopäätös on se, että \bar{X} estimoi tarkemmin populaation parametrin keskiarvoa.

5.3 Väliestimointi

Väliestimoinnissa määrätään satunnaisotannalla poimitusta otoksesta X_1, X_2, \dots, X_n riippuvat luottamusvälin ala- ja ylärajat, joiden väliin populaation parametri halutulla todennäköisyydellä kuuluu. Yleisesti esitettynä populaation parametrin T luottamusvälin ylä- ja alarajat \hat{T}_{alaraja} ja $\hat{T}_{\text{yläraja}}$ toteuttavat yhtälön

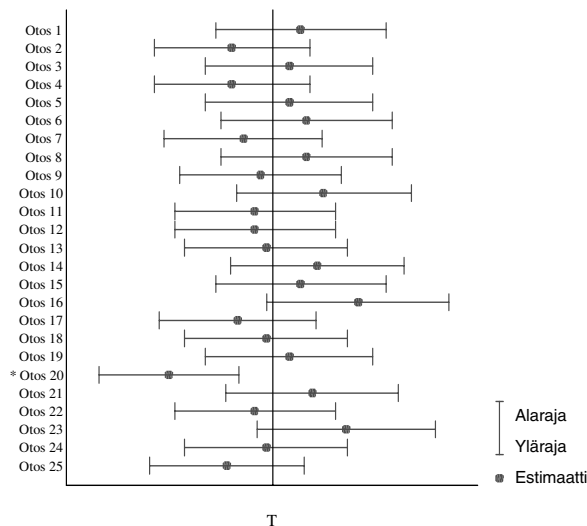
$$P(\hat{T}_{\text{alaraja}} \leq T \leq \hat{T}_{\text{yläraja}}) = 1 - \alpha$$

jossa α ilmaisee halutun *riskitason*, esim. 0,05 eli 5 %. Lukua $(1-\alpha)100\%$, esim. 95%, kutsutaan *luottamustasoksi*. Siten puhutaan esimerkiksi keskiarvon 95 %:n luottamusvälistä.

Koska siis populaation parametrin estimaattorin arvossa esiintyy vaihtelua otoksesta toiseen (jota estimaattorin keskivirhe kuvaa), on väli $(\hat{T}_{\text{alaraja}}, \hat{T}_{\text{yläraja}})$ luonteeltaan satunnaisväli. Tästä syystä eri otoksista saadaan parametrin T suhteen sijainniltaan vaihtelevia luottamusvälejä. (ks. Kuva 34, s. 56). Silloin, kun populaation varianssikin täytyy estimoida, ts. kun estimaattorin keskivirhettä ei tunneta, esiintyy lisäksi luottamusvälin pituudessa satunnaisvaihtelua sijaintivaihtelun lisäksi.

Kuva 34 esittää, kuinka 25:stä eri otoksesta määrätyt 95 %:n luottamusvälit sisältävät 24 tapauksessa populaation parametrin T arvon. Tämä tulos osuu hyvin yhteen teorian kanssa, sillä voimme odottaa $0,95 \times 25$ tapauksessa luottamusvälin sisältävän parametrin T . Kuten kuvasta näkyy otos 20 ei sisällä populaation parametrin arvoa.

Kuva 34: Toisistaan riippumattomista otoksista laskettuja luottamusvälejä populaation parametrille T



6 Luottamusvälien laskeminen

6.1 Luottamusvälin määrittäminen populaation keskiarvolle μ , kun σ^2 tunnetaan

Z -jakaumaa noudattavalle satunnaismuuttujalle on voimassa yhtälö

$$P(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1 - \alpha,$$

joka siis määrää todennäköisyyden $(1 - \alpha)$, että Z kuuluu välille $-Z_{\alpha/2}, Z_{\alpha/2}$.

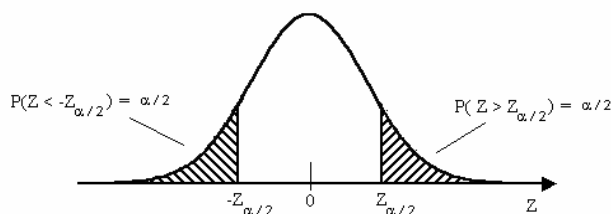
Näin ollen esimerkiksi standardoitu otoskeskiarvo $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ ($= Z$) kuuluu 95

%:n varmuudella arvojen $-Z_{0,05/2}$ ($= -1,96$) ja $Z_{0,95/2}$ ($= 1,96$) välille. Tällöin standardoitu otoskeskiarvo toteuttaa yhtälön

$$P(-1,96 \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq 1,96) = 0,95.$$

Kuva 35 havainnollistaa yhtälön suhdetta teoreettiseen normaalijakaumaan.

Kuva 35: Standardoidun normaalijakauman luottamusvälin päätepisteiden (kriittisten arvojen) määritelmät



Kuten on jo todettu, normaalijakaumaa noudattavan satunnaismuuttujan X otoskeskiarvo \bar{X} noudattaa jakaumaa $N(\mu, \sigma^2 / n)$. Mikäli tunnetaan myös X :n varianssi σ^2 , voidaan otoskoosta riippumatta² laskea sellaiset luottamusrajat,

² Jos σ^2 tunnetaan, pieni otoskoko ei vaikuta siihen, noudattaako \bar{X} :n otantajakauma normaalijakaumaa $N(\mu, \sigma^2 / n)$.

joiden sisään \bar{X} :n estimoima μ halutun suurella todennäköisyydellä kuuluu. Yllä selvitettyyn nojautuen voidaan johtaa yleisesti voimassa oleva yhtälö

$$P\left[\bar{X} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right] = 1 - \alpha ,$$

jossa $-Z_{\alpha/2}$ ja $Z_{\alpha/2}$ ovat sellaiset standardoidun normaalijakauman kriittiset arvot (halutulla luottamustasolla $(1-\alpha)100\%$), jotka toteuttavat jo Kuva 35:ssä esitetyn yhtälön

$$P(|Z| > Z_{\alpha/2}) = \alpha .$$

- **ESIM:** Määrää populaation keskiarvolle μ 95 %:n luottamusväli, kun otoskeskiarvon ($n=50$) arvoksi saatiin 464 ja varianssin σ^2 tiedetään olevan 100^2 .

Liitteen 1 taulukosta saadaan kriittiset arvot $-Z_{0,05/2}$ ja $Z_{0,05/2}$, jotka ovat $-1,96$ ja $1,96$. Ylläesitetyn pohjalta sijoitetaan kriittiset arvot, otostunnusluvut sekä otoskoko yhtälöön ja lasketaan

$$P\left[464 - 1,96 \cdot \frac{100}{\sqrt{50}} \leq \mu \leq 464 + 1,96 \cdot \frac{100}{\sqrt{50}}\right] = 1 - 0,05,$$

saaden

$$P[436,3 \leq \mu \leq 491,7] = 0,95 .$$

Eli μ sijaitsee 95 %:n varmuudella välillä 436,3 ja 491,7.

Tutkitaan seuraavaksi, miten otoskoko vaikuttaa luottamusvälin pituuteen. Oletetaan, että muuttujalla X on populaatiossa keskiarvo 10 ja varianssi 3^2 . Lasketaan erisuuria otoskokoja (n) käyttäen useita luottamusvälejä osoittamaan, miltä väliltä populaation keskiarvo kunkin otoskoon tapauksessa voidaan olettaa löytyvän. Luottamustasoksi valitaan 5 %. Taulukko 7 esittää eri otoskoon arvoille laskettuja luottamusvälin ala- ja ylärajan arvoja sekä vastaavat luottamusvälin pituudet.

Taulukko 7: Luottamusvälin päätepisteitä ja pituus eri otoskoon arvoilla, kun $X \sim N(10, 3^2)$ ja $\alpha = 0,05$.

n	Alaraja 95 %	Yläraja 95 %	Välin pituus (d)
1	4,12	15,88	11,76
10	8,14	11,86	3,72
30	8,93	11,07	2,15
70	9,30	10,70	1,41
300	9,66	10,34	0,68
800	9,93	10,07	0,13

Luottamusvälin pituus **d** saadaan määrättyä seuraavasti

$$d = 2 \cdot Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} .$$

Taulukko 7 tuo esiin, että otoskoon kasvaessa luottamusvälin ala- ja yläraja lähestyvät toisiaan eli luottamusvälin pituus **d** lyhentyy. Tämä luonnollisesti tarkoittaa, että otoskoko kasvatamalla estimaattorin \bar{X} keskivirhe σ / \sqrt{n} pienentyy ja näin ollen estimointi tarkentuu ja muuttuu luotettavammaksi.

6.2 Populaation keskiarvon μ luottamusvälin määrittäminen, kun σ^2 on tuntematon

Kun populaation varianssia σ^2 ei tunneta (mikä on hyvin tavallista), käytetään sen paikalla otoksesta estimoitua otosvariانسsia $s^2 = \sum_{i=1}^n (x_i - \bar{X})^2 / (n-1)$. Tällöin populaation keskiarvon μ luottamusvälin kaavassa

$$\mu: n \text{ lv} : \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

asetetaan σ :n paikalle **s** ja lisäksi **Z**-jakauman kriittisen arvon ($Z_{\alpha/2}$) paikalle vastaava **t**-jakauman kriittinen arvo $t_{\alpha/2}$ vapausasteilla **df = n-1** ja saadaan

$$\mu: n \text{ lv} : \bar{X} \pm t_{\alpha/2} \cdot \frac{s}{\sqrt{n}} .$$

Jos otoskoko on pienempi kuin 30, on syytä varmistaa, että X :n jakauma on populaatiossa normaalin, sillä muutoin, varsinkin paljon 30:tä pienemmillä otoksilla, on s^2 taipuvainen aliestimoimaan σ^2 :n ja sen vaihtelu on yleisesti suurempaa, kuten yllä luvussa 4.6 todettiin.

Suuremmilla otoksilla kuin 30 populaation normalisuutta ei välttämättä tarvitse olettaa, koska s^2 :n satunnaisvaihtelu otoksesta toiseen on huomattavasti pienempää ja toisaalta koska σ^2 :n aliestimoinnin vaara vähentyy s^2 :n otantajakauman lähestyessä normaalijakaumaa (ks. Kuva 28, s. 46 ja Kuva 30, s. 49).

- **ESIM:** Tutkittaessa erään oppilaitoksen opiskelijoiden ($n = 45$) alkoholiin käyttämää rahamäärää viikossa, saatiin otoskeskiarvoksi 86 mk ja otoskeskihajonnaksi 76 mk. Laske 95 %:n luottamusväli oppilaitoksen opiskelijoiden muodostaman populaation keskiarvolle.

Liitteen 2 taulukkoa käyttämällä etsitään t -jakauman fraktiilipiste $t^{(97,5)}$ vapausasteilla $df = 45 - 1$, joka on **2,015**. Asetetaan arvot kaavaan, jolloin saadaan

$$\text{Alaraja} = 86 - 2,02 \cdot \frac{76}{\sqrt{45}} = \underline{\underline{63,2 \text{ mk}}}$$

$$\text{Yläaraja} = 86 + 2,02 \cdot \frac{76}{\sqrt{45}} = \underline{\underline{108,8 \text{ mk}}},$$

jossa luottamusvälin pituus $d = \underline{\underline{45,63 \text{ mk}}}$.

Excelin funktioilla laskettaessa kirjoitetaan soluihin kaavat:

$$= 86 - \text{TINV}(0,05; 44) * 76 / \text{SQRT}(45) = \underline{\underline{63,2}}$$

$$= 86 + \text{TINV}(0,05; 44) * 76 / \text{SQRT}(45) = \underline{\underline{108,8}}$$

$$= 108,8 - 63,2 = \underline{\underline{45,63}}.$$

6.3 Luottamusvälin määrittäminen populaation prosenttiluvulle θ

Kuten yllä (ks. Luku 4.3) todettiin estimaattorin \mathbf{p} otantajakauma noudattaa likimäärin $\mathbf{N}(\theta, \theta(100-\theta)/n)$, mikäli otoskoko on riittävän suuri suhteessa $\theta:n$ eli vastaavan populaation prosenttiluvun arvoon (ks. Taulukko 4, s. 35)³. Koska estimaattorin \mathbf{p} varianssi $\theta(100-\theta)/n$ riippuu tuntemattomasta parametrimestä θ , emme kykene suoraan määrittämään $\mathbf{p}:n$ tarkkaa varianssia. Mutta jos käytämme varianssin lausekkeessa $\theta:n$ paikalla otoksesta estimoitua $\mathbf{p}:tä$ saamme populaation varianssille silti riittävän approksimaation. Otoskoon ollessa riittävä standardoitu muuttuja

$$Z = \frac{\mathbf{p} - \theta}{\sqrt{\frac{\mathbf{p}(100 - \mathbf{p})}{n}}}$$

noudattaa likimäärin $\mathbf{N}(\mathbf{0},1)$ -jakaumaa. Jos asetamme tämän lausekkeen $Z:n$ paikalle jo tutussa yhtälössä

$$P(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1 - \alpha ,$$

saadaan tulokseksi, kun ensin ratkaistaan molemmat epäyhtälöt $\theta:n$ suhteen, seuraava yhtälö

$$P \left[\mathbf{p} - Z_{\alpha/2} \cdot \sqrt{\frac{\mathbf{p}(100 - \mathbf{p})}{n}} \leq \theta \leq \mathbf{p} + Z_{\alpha/2} \cdot \sqrt{\frac{\mathbf{p}(100 - \mathbf{p})}{n}} \right] = 1 - \alpha .$$

Täten prosenttiluvun θ luottamusvälin päätepisteet ovat

$$\mathbf{p} \pm Z_{\alpha/2} \cdot \sqrt{\frac{\mathbf{p}(100 - \mathbf{p})}{n}} .$$

³ Mikäli otoskoko ei täytä Taulukko 4:n ehtoja, noudattaa $\mathbf{p}:n$ otantajakauma binomijakaumaa, kuten jo todettiin luvussa 0. Koska \mathbf{p} on otosprosenttiluvun kohdalla ominaisuuden \mathbf{A} esiintymiskertojen lukumäärä suhteessa otoskokoon eli $\mathbf{p} = f(\mathbf{A})/n$, lähestyy sen otantajakauma keskeisen raja-arvolauseen mukaan otoskoon kasvaessa normaalijakaumaa. Tätä tulosta hyödynnetään luottamusvälin määrittämisessä populaation prosenttiluvun arvolle.

- **ESIM:** Tietyn puolueen kannatuksen arvioimiseksi poimittiin 900 hengen satunnaisotos äänioikeutettujen muodostamasta populaatiosta. Kannatusprosentiksi saatiin 16,7 % Määrää 95 %:n luottamusväli populaation prosenttiluvun θ arvolle.

Liitteen 1 taulukosta saadaan jo tutut $N(0,1)$ -jakauman kriittiset arvot. Asetetaan saatu piste-estimaatti 16,7 % ja kriittiset arvot $(-1,96, 1,96)$ kaavaan ja lasketaan

$$\text{Alaraja} = 16,7 - 1,96 \cdot \sqrt{\frac{16,7(100 - 16,7)}{900}} = \underline{14,3 \%}$$

$$\text{Yläraja} = 16,7 + 1,96 \cdot \sqrt{\frac{16,7(100 - 16,7)}{900}} = \underline{19,1 \%} .$$

Voimme todeta 95 %:n varmuudella, että ko. puoleen kannatusprosentti äänestysikäisten muodostamassa populaatiossa sijaitsee välillä (14,3 %, 19,1 %).

6.4 Luottamusvälin määrittäminen populaation varianssille σ^2

Käsitellessämme $\chi^2(k)$ -jakauman yhteyttä otosvarianssin s^2 otantajakaumaan luvussa 4.6.4, totesimme, että esitettyjä tuloksia hyväksi käyttäen on mahdollista määrätä estimoitavalle populaation varianssille σ^2 luottamusväli. Vaikka käytännön tutkimuksessa tarvitaan perin harvoin luottamusväliä nimenomaan populaation varianssille, on kuitenkin erityistilanteita, joissa halutaan varmuutta populaation varianssin satunnaiskäyttäytymisen alueesta. Seuraava yhtälö ilmaisee halutulla luottamustasolla α populaation varianssin luottamusvälin päätepisteiden arvot

$$\mathbf{P} \left[\frac{s^2}{\left(\frac{\chi_{\alpha/2}^2}{n-1} \right)} \leq \sigma^2 \leq \frac{s^2}{\left(\frac{\chi_{\alpha/2}^2}{n-1} \right)} \right] = 1 - \alpha .$$

- **ESIM:** Eräässä tutkimuksessa ($n = 41$), jossa mitattiin lasten syntymäpituutta, saatiin otosvarianssille s^2 arvo **9,61**. Määrää **95%:n** luottamusväli syntymäpituuden populaation varianssin arvolle.

Määrätään χ^2 -jakauman ala- että ylärajan kriittiset arvot vapausasteilla 40 joko liitteen 5 taulukosta tai esim. Excelin CHIINV -funktiolla. Saadaan $\chi_{0,025}^2 = 59,34$ ja $\chi_{0,975}^2 = 24,42$. Tämän jälkeen asetetaan vaaditut arvot yhtälöön eli

$$P \left[\frac{9,61}{\left(\frac{59,34}{41-1}\right)} \leq \sigma^2 \leq \frac{9,61}{\left(\frac{24,42}{41-1}\right)} \right] = 0,95 ,$$

ja lasketaan saaden

$$P [6,48 \leq \sigma^2 \leq 15,74] = 0,95 .$$

Luottamusvälin päätepisteet ovat **6,48** ja **15,74**, joiden väliltä populaation varianssi **9,61** löytyy 95 %:n varmuudella. On huomattava, että tämän luottamusvälin pituus **d** on melko suuri. Ellei otoskoko ole melko suuri tai s^2 vastaavasti melko pieni, ei ole mahdollista muodostaa kovin tarkkoja väliestimaatteja populaation varianssille.

6.5 Luottamusvälin määrittäminen populaation korrelaatiokertoimelle ρ

Voimme muodostaa luottamusvälin populaation korrelaatiokertoimelle ρ (lue: *rho*) ratkaisemalla myöhemmin luvussa 8.6.3 tarkemmin esiteltävän kaavan

$$Z = \frac{r^* - \rho^*}{\sqrt{\frac{1}{n-3}}} ,$$

ρ^* :n suhteen, jolloin saadaan luottamusvälin päätepisteiden kaavaksi

$$\rho^* = r^* \pm Z_{\alpha/2} \cdot \frac{1}{\sqrt{n-3}} .$$

Kaavan otoskorrelaatiokertoimen r sekä populaation korrelaatiokertoimen ρ yläindeksinä olevat tähtisymbolit tarkoittavat, että kertoimien arvoille on suoritettu ns. Fisher-transformaatio, joka on seuraavaa muotoa

$$r^* = (0,5) \log_e \left| \frac{1+r}{1-r} \right|.$$

Tämä muunnos on välttämätöntä silloin, kun ρ :n arvo poikkeaa nolasta, jolloin sen estimaattorin r otantajakauma ei noudata symmetrisesti t -jakaumaa. Fisher-transformaatio muuntaa r :n arvon siten, että sen otantajakauma noudattaa tarkemmin normaalijakaumaa mahdollistaen täten normaalijakauma approksimaatiot.

- **ESIM:** Eräessä tutkimuksessa ($n=203$) todettiin, että alkoholiin kuukaudessa käytetyn rahamäärän ja kuukausiansion välinen korrelaatio on 0,43. Määrää 99,9%:n luottamusväli perusjoukon korrelaatiokertoimen ρ arvolle.

Haetaan taulukosta 1 normaalijakauman kriittinen arvo, joka vastaa 99,9%:n luottamustasoa. Tämä on $Z_{0,001/2} = 3,29$.

Seuraavaksi suoritetaan otoskorrelaatiokertoimelle Fisher-transformaatio joko käsin eli

$$r^* = (0,5) \log_e \left| \frac{1+0,43}{1-0,43} \right| = 0,46,$$

tai Excelin FISHER –funktiolla, jolloin kirjoitetaan

$$=FISHER(0,43) = 0,46.$$

Seuraavaksi sijoitetaan arvot kaavaan, lasketaan ja saadaan

$$\begin{aligned} \rho^* &= 0,46 \pm 3,29 \cdot \frac{1}{\sqrt{203-3}} \\ &= 0,46 \pm 0,23 \\ &= 0,23 \leq \rho^* \leq 0,69. \end{aligned}$$

Saadut ala- ja ylärajan arvot täytyy vielä FISHER-transformoida käänteisesti takaisin ρ^* :stä ρ :hen. Tämä tehdään joko käsin kaavalla

$$\rho_{\text{alaraja}} = \frac{e^{2\rho^*_{\text{alaraja}}} - 1}{e^{2\rho^*_{\text{alaraja}}} + 1} \rightarrow \frac{2,718^{2 \cdot 0,23} - 1}{2,718^{2 \cdot 0,23} + 1} = 0,226$$

$$\rho_{\text{yläraja}} = \frac{e^{2\rho^*_{\text{yläraja}}} - 1}{e^{2\rho^*_{\text{yläraja}}} + 1} \rightarrow \frac{2,718^{2 \cdot 0,69} - 1}{2,718^{2 \cdot 0,69} + 1} = 0,598$$

tai Excelin FISHERINV –funktiolla jolloin kirjoitetaan soluun kaavat

$$\begin{aligned} &= \text{FISHERINV}(0,23) = \underline{0,226} \\ &= \text{FISHERINV}(0,69) = \underline{0,598} . \end{aligned}$$

Täten saadaan estimoiduksi luottamusvälin päätepisteet, joiden sisään populaation korrelaatiokertoimen arvo ρ sisältyy

$$P[0,226 \leq \rho \leq 0,598] = 0,999$$

99,9 %:n todennäköisyydellä. Voidaan myös todeta, että $\rho = 0$ ei sisälly lainkaan tälle välille, mikä antaa samalla varmuuden siitä, että ρ :n arvo on nollaa selvästi suurempi.

7 Tilastollinen hypoteesin testaus

7.1 Tilastollisen testauksen lähtökohta

Tilastollinen hypoteesi tarkoittaa yhden tai useamman perusjoukon tuntemattomia tunnuslukuja eli parametrejä koskevia matemaattisia tai loogisia väittämiä. Tarkastellaan kahta tällaista väittämää.

- 1) Presidentinvaaleissa kahden toiselle kierrokselle selviytyneen ehdokkaan kannatusprosentit koko äänestäjäkunnassa ovat viikkoa ennen äänestystä yhtä suuret.
- 2) Tyttöjen ja poikien vasemman puoleisen aivolohkon keskimääräinen aktiiviteetti on sama.

Perusjoukkojen suuruuden vuoksi näiden väittämien totuusarvoa ei voida suoraan todentaa. Jos väittämien totuusarvoista halutaan jotain lausua, on sen perustuttava otannalla saatuun informaatioon. Otoksista lasketut otostunnusluvut ovat kuitenkin vain parametrien estimaatteja, joihin sisältyy otantavirhettä. Niinpä otostunnuslukujen perusteella ei voida sanoa mitään ehdottoman tarkkaa näistä väittämistä vaan ainoastaan todennäköisyyksiä. Eli lyhyesti: Tilastollinen hypoteesin testaus käsittelee perusjoukkojen parametrejä koskevien väittämien paikkansapitävyyden tutkimista otosinformaation perusteella (ks. luku 4, s. 32 ja luku 5, s. 54).

Perusjoukon parametrejä koskevia väittämiä kutsutaan *tilastollisiksi hypoteeseiksi*, joita ovat *nollahypoteesi* H_0 ja *vaihtoehtoinen hypoteesi* H_1 . Tilastolliset hypoteesit johdetaan yleensä *tutkimushypoteesista*, joka tyypillisesti on muotoa ”käsittelyissä on eroa” tai ”tietty tekijä vaikuttaa johonkin toiseen tekijään”. Nollahypoteesi *ei ole* tällöin tutkimushypoteesi, vaan se on parametrien matemaattisesti yksinkertaisin tilanne, esimerkiksi että perusjoukon kaksi prosenttilukua ovat yhtä suuret. Vaihtoehtoinen hypoteesi puolestaan ilmoittaa parametrien muut kyseeseen tulevat arvot.

Nollahypoteesin määrittämisen taustalla on matemaattis-tilastollisen testiteorian asettamat vaatimukset. Kun nollahypoteesi on tosi, voidaan sen totuusarvoa mitaavalle testisuurelle johtaa otantajakauma ja sen perusteella tehdä johtopäätöksiä nollahypoteesin ja vaihtoehtoisen hypoteesin paikkansapitävyydestä.

7.2 Hypoteesit

Hypoteesit ovat tuntemattomia parametrejä koskevia väittämiä. Niitä on kahdenlaisia: *nollahypoteesi* H_0 ja *vaihtoehtoinen hypoteesi* H_1 . Nollahypoteesi ilmaisee yleensä parametrien yksinkertaisimman asetelman. Esimerkiksi olkoon

presidentinvaalien toisen kierroksen kahden ehdokkaan kannatusprosentit kaikkien äänioikeutettujen muodostamassa perusjoukossa θ_1 ja θ_2 ($=100-\theta_1$). Tällöin nollahypoteesi on, että ehdokkaiden kannatusprosentit perusjoukossa ovat yhtä suuret eli

$$H_0: \theta_1 = \theta_2 = 50 \%$$

Vaihtoehtoinen hypoteesi H_1 voi olla joko *yksisuuntainen* tai *kaksisuuntainen*. Kaksisuuntaista vaihtoehtoista hypoteesia käytetään, kun parametreista ei ole käytettävissä mitään vankkaa etukäteisinformaatiota. Tässä esimerkissä on paikallaan käyttää kaksisuuntaista vaihtoehtoista hypoteesia eli

$$H_1: \theta_1 \neq 50\% \text{ (jolloin myös } \theta_2 \neq 50\% \text{)} .$$

Tilastollinen testaus osoittaa, kumpaa hypoteesia päädytään otosinformaation perusteella kannattamaan. Johtopäätös perustuu luonnollisesti otoksesta lasketun ensimmäisen (tai toisen) ehdokkaan kannatusprosentin p poikkeamaan nollahypoteesin ilmoittamasta arvosta 50 %. Mutta kuinka suuri pitää poikkeaman olla, jotta nollahypoteesista voidaan luopua? Vastauksen antaa *testisuureen* arvo ja siihen liittyvä *riskitodennäköisyys* α (alpha).

7.2.1 Esimerkkejä tilastollisista hypoteeseista

Esimerkki 1: Henkilön väitettyjä ESP-kykyjä (ESP=extra sensory perception) mitattiin siten, että hänelle annettiin pelkän tuntoaistin perusteella tunnistettavaksi kolmenvärisiä kortteja. Merkitään θ :lla henkilön todennäköisyyttä tietää (tai arvata) kortin oikea väri. Tällöin tilastolliset hypoteesit ovat

$$H_0: \theta = 33,3\% \quad , \quad H_1: \theta > 33,3\% .$$

Nollahypoteesin mukaiseen arvoon päästään pelkästään arvaamalla. Vaihtoehtoinen hypoteesi on yksisuuntainen, koska 33,3 prosenttia pienemmät arvot eivät oikeastaan voi tulla kysymykseen.

Esimerkki 2: Kahden eri mittarin on tarkoitus mitata fyysistä suorituskyykyä. Halutaan testata tilastollisesti, onko mittareiden välillä systemaattista eroa. Merkitään näiden mittareiden keskiarvoja μ_1 :llä ja μ_2 :lla siinä perusjoukossa, johon mittareita soveltaa. Tilastolliset hypoteesit ovat nyt

$$H_0: \mu_1 = \mu_2 \quad , \quad H_1: \mu_1 \neq \mu_2 .$$

Huom. Testiä ei pidä muotoilla korrelaatiokertoimen avulla.

Esimerkki 3: Lääketehdas ilmoittaa, että sen tuottama lääke parantaa vähintään 80 prosentissa tapauksia. Tilastolliset hypoteesit ovat

$$\mathbf{H}_0: \theta = 80 \% \quad , \quad \mathbf{H}_1: \theta < 80 \% .$$

Vaihtoehtoinen hypoteesi on yksisuuntainen, koska ollaan kiinnostuneita nimenomaan siitä mahdollisuudesta, että lääketehdas liioittelee.

7.3 Tilastollisen testin mahdolliset virheet

Tilastolliseen päätöksentekoon otosinformaation perusteella liittyy aina virhemahdollisuus. Asettamalla vastakkain todellisuus (\mathbf{H}_0 tai \mathbf{H}_1) ja johtopäätös (\mathbf{H}_0 tai \mathbf{H}_1) saadaan seuraava nelikenttä.

Taulukko 8: Tilastollisen testauksen vaihtoehdot ja niiden todennäköisyydet

		Todellisuus	
		\mathbf{H}_0	\mathbf{H}_1
Johtopäätös otosinformaation perusteella	\mathbf{H}_0	Oikea johtopäätös, $tn = 1 - \alpha$	2. lajin virhe, $tn = \beta$
	\mathbf{H}_1	1. lajin virhe, $tn = \alpha$	Oikea johtopäätös, $tn = 1 - \beta$

(tn = todennäköisyys)

Taulukko 8:n nelikenttä esittää neljää eri vaihtoehtoista johtopäätöstä toisaalta sen mukaan, vallitseeko todellisuudessa \mathbf{H}_0 vai \mathbf{H}_1 ja toisaalta sen mukaan, mikä on testin antama tulos otosinformaation perusteella. Neljästä vaihtoehdosta kaksi antaa oikean johtopäätöksen ja kaksi virheellisen tuloksen. Jos oikea \mathbf{H}_0 hylätään eli päädytään \mathbf{H}_1 :n kannalle, on kyseessä *1. lajin virhe*. Jos todellisuudessa \mathbf{H}_1 -hypoteesi on oikea ja testin perusteella päädytään \mathbf{H}_0 -hypoteesin kannalle, on kyseessä *2. lajin virhe*.

Tilastollisessa testauksessa 1. lajin virheelle annetaan määrätty todennäköisyys, joka voidaan sietää päätöksenteossa. Tätä todennäköisyyttä kutsutaan riskitodennäköisyydeksi ja merkitään α :lla (alpha) Riskitodennäköisyydet voidaan luokitella neljään luokkaan, joiden rajoina ovat α :n arvot: 10 %, 5 %, 1 % ja

0,1 %. Näitä luokkia kutsutaan *riskitasoiksi* tai *merkitsevyystasoiksi* ja ne esitetään kootusti Taulukko 9:ssä.

Taulukko 9: Testin riskitasot eli merkitsevyystasot

Riskitodennäköisyys α	Riskitaso eli merkitsevyystaso	Tähti-merkintä
$\alpha > 0,10$	Tilastollisesti ei-merkitsevä	
$0,05 < \alpha \leq 0,10$	Tilastollisesti oireellinen	
$0,01 < \alpha \leq 0,05$	Tilastollisesti melkein merkitsevä	*
$0,001 < \alpha \leq 0,01$	Tilastollisesti merkitsevä	**
$\alpha \leq 0,001$	Tilastollisesti erittäin merkitsevä	***

Riskitodennäköisyyttä α kutsutaan myös testin **p**-arvoksi (p-value). Usein ”tilastollisesti oireellinen”-taso unohdetaan varsinkin, jos otoskoko ei ole kovin pieni ja tilastollisesta merkitsevyydestä puhutaan vasta, kun riskitodennäköisyys on pienempi kuin 0,05. Tähtimerkintöjä on kätevä käyttää, jos testejä on paljon ja halutaan säästää tilaa.

Toisen lajin virheen todennäköisyys, jota merkitään β :lla, riippuu parametrin todellisesta arvosta eikä se ole yhtä helposti hallittavissa kuin 1. lajin virhe. Otosinformaation perusteella voidaan laskea arvio 2. lajin virheen todennäköisyydelle. Testin *voimakkuus* (Power) liittyy läheisesti 2. lajin virheeseen ja sitä käsitellään luvussa 9.

7.4 Testisuure

Testisuure on otostunnusluvun ja nollahypoteesin määräämän parametrin arvon matemaattinen lauseke, jonka jakauma (nk. otantajakauma) tunnetaan silloin, kun nollahypoteesi on tosi. Edellä olleen kannatusprosenttia koskevan esimerkin tapauksessa, jossa nollahypoteesi on $\theta_1 = 50\%$, testisuure Z on

$$Z = \frac{p - 50}{\sqrt{\frac{50(100 - 50)}{n}}} = \sqrt{n}(p - 50)/50 ,$$

jossa p on otosprosenttiluku ja n on otoskoko. Tämän testisuureen ensimmäinen lauseke on jonkin verran koukeroinen, mutta sillä on haluttu tuoda esille se seikka, että Z on poikkeama $p - 50$ jaettuna p :n keskivirheellä.

Mitä enemmän otosprosenttiluku p poikkeaa 50 %:sta, sitä suurempi testisuureen Z itseisarvo on. Tehtävänä on määrätä sellainen raja, joka testisuureen Z itseisarvon tulee ylittää, jotta nollahypoteesi voidaan hylätä. Testisuureita ei täs-

sä kirjassa matemaattisesti johdeta. Kuitenkin kaikille prosenttiluku- ja keskiarvotestien testisuureille on voimassa seuraava lauseke:

$$\text{Testisuure} = \frac{\text{poikkeama } H_0\text{:n mukaisesta tilanteesta}}{\text{estimaattorin keskivirhe}}$$

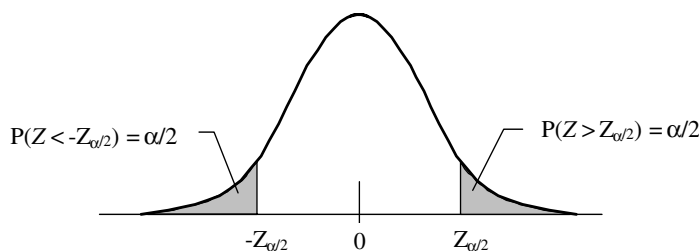
Tässä lausekkeessa osoittaja on yleensä itsestään selvä. Keskivirheen lausekkeet ovat sen sijaan usein komplisoidumpia.

7.5 Kriittinen alue

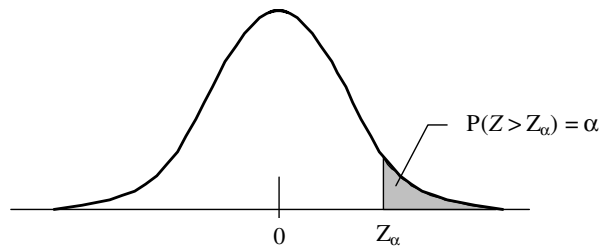
Mikä on itseisarvoltaan riittävän suuri testisuureen arvo, jotta nollahypoteesi voidaan hylätä? Tähän saadaan vastaus otantajakauman avulla. Jos nollahypoteesi on tosi, noudattaa testisuure Z normaalijakaumaa $N(0,1)$. Tämän jakauman perusteella voidaan päätellä, mitkä Z :n arvot ovat itseisarvoltaan niin suuria, että ne ovat suhteellisen harvinaisia Z :n arvoiksi nollahypoteesin vallitessa. Valitaan riskitodennäköisyys α sen mukaan, kuinka suuri 1. lajin virheen todennäköisyys siedetään. Tämän jälkeen katsotaan $N(0,1)$ -normaalijakauman taulukosta (liite 1) tai Excel-funktion avulla *kriittinen arvo*. Kaksisuuntaisen testin tapauksessa kriittiset arvot ovat $\pm Z_{\alpha/2}$ (ks. väliestimointi 5.3, s. 56).

Kriittisen alueen muodostavat ne testisuureen arvot, jotka johtavat nollahypoteesin hylkäämiseen. Kuva 36 esittää kriittisen alueen kaksisuuntaisen testin tapauksessa.

Kuva 36: Kriittisten arvojen ja kriittisen alueen määrittäminen kaksisuuntaisen vaihtoehdoisen hypoteesin tapauksessa



Kuva 37: Kriittisen arvon ja kriittisen alueen määrittäminen yksisuuntaisen vaihtoehdoisen hypoteesin tapauksessa



Jos kyseessä on yksisuuntainen vaihtoehtoinen hypoteesi, määrätään kriittinen alue $N(0,1)$ -normaalijakauman oikeasta tai vasemmasta hännästä riippuen siitä onko vaihtoehtoinen hypoteesi muotoa ”suurempi kuin” tai ”pienempi kuin”. Kuva 37:ssä esitetään yksisuuntainen vaihtoehtoinen hypoteesi, joka on muotoa ”suurempi kuin” eli kriittinen alue on määrätty oikeasta hännästä.

Seuraavassa taulukossa on standardoidun normaalijakauman kriittiset arvot kolmella riskitasolla sekä yksi- ja kaksisuuntaisen vaihtoehtoisen hypoteesin tapauksissa.

Taulukko 10. Standardoidun normaalijakauman kriittisiä arvoja

	Yksisuuntainen testi			Kaksisuuntainen testi		
	5 %	1 %	0,1 %	5 %	1 %	0,1 %
Kriittinen arvo	1,64	2,33	3,09	$\pm 1,96$	$\pm 2,58$	$\pm 3,29$

Huomataan että yksisuuntaisen testin kriittiset arvot ovat pienempiä kuin vastaavat kaksisuuntaisen testin kriittisten arvojen itseisarvot. Siten yksisuuntaisessa testissä päästään nollahypoteesi hylkäämään helpommin kuin kaksisuuntaisessa testissä.

7.6 Esimerkki tilastollisesta testauksesta

Tarkastellaan edellä ollutta, kahden presidenttiehdokkaan kannatusprosentteihin liittyvää tilastollista testausta, jonka testisuure $Z = \sqrt{n} (p - 50)/50$ esitettiin kappaleessa 7.4. Annetaan otosprosenttiluvulle p ja otoskoolle n eräitä arvoja ja katsotaan millaisiin johtopäätöksiin Taulukko 10:n kriittiset arvot antavat aihetta. Kyseessä on kaksisuuntainen vaihtoehtoinen hypoteesi.

Taulukko 11: Z-testisuureen arvot ja johtopäätökset

Otoskoko n	Otosprosenttiluku p				
	51	52	53	54	55
100	0,20 em	0,40 em	0,60 em	0,80 em	1,00 em
500	0,45 em	0,89 em	1,34 em	1,79 em	2,24 *

1000	0,63 em	1,26 em	1,90 em	2,53 *	3,16 **
2000	0,89 em	1,79 em	2,68 **	3,58 ***	4,47 ***
4000	1,26 em	2,53 *	3,79 ***	5,06 ***	6,32 ***

Taulukko 11:een on merkitty Z -testisuureen arvot ja testin tulos eri n :n ja p :n arvoilla. Merkintä ”em” tarkoittaa ei-merkitsevää tulosta ja tilastollisesti merkitsevät Z :n arvot on merkitty tähtimerkinnällä. Esimerkiksi 54 %:n kannatus tuhannen hengen otoksessa antaa Z :n arvoksi 2,53, mikä johtaa nollihypoteesin hylkäämiseen viiden prosentin riskitasolla (yhden tähden tapaus). Samaa Z :n arvoon päästään, jos p :n arvo on 52 % ja otoskoko neljä tuhat henkeä. Tilastollisen testin johtopäätös edellisessä esimerkissä voidaan ilmaista myös seuraavasti: 54 % poikkeaa 50 %:sta *tilastollisesti melkein merkitsevästi* (ks. Taulukko 9).

Testisuureen Z arvo 2,53 on varsin lähellä yhden prosentin kriittistä arvoa 2,58. Z :n arvoon 2,53 liittyvä tarkka riskitodennäköisyys α saadaan Excelin NORMSDIST –funktion avulla, joka antaa standardoidun normaalijakauman kertymäfunktionarvon halutussa pisteessä. Helpoimmin haluttu tulos saadaan seuraavasti:

$$= 2 - \text{NORMSDIST}(-2,53) = \underline{0,0114} \text{ eli } \underline{1,14 \%} .$$

Huom. Käytä tilastollisen testauksen yhteydessä aina sanaa *merkitsevä* — ei esim. sanaa *merkittävä*. Tulosten käytännön merkittävyys on aivan eri asia kuin tilastollinen merkitsevyys.

7.7 Tilastollinen testaus käytännössä

On useita tapoja, joilla tilastollinen päätöksenteko käytännössä tehdään joko tietokoneen avustuksella tai ilman sitä. Seuraavassa esitellään nämä vaihtoehdot.

- 1) Tietokoneen tilasto-ohjelma laskee havaintomatriisista lähtien testisuureen arvon, p -arvon (=riskitodennäköisyys α) ja mahdollisesti kertoo sanallisesti testin tuloksen. Esimerkkinä alla on sukupuolen ja terveyskeskukseen tulon syyn ristiintaulukon tilastollisen testin osa (Tixel-ohjelma).

Kontingenssikerroin = 0,194
 Khiin neliö = 28,01 Vap. ast. = 9
 P-arvo = 0,001 Tilastollisesti erittäin merkitsevä

- 2) Mittatulokset ovat Excel-tilulukossa ja käytetään jotain Excel-funktiota, joka antaa suoraan p -arvon. Esimerkkinä varianssien yhtäsuuruuden testaus (ks. luku 8.8).

- 3) Käytettävissä ovat valmiiksi lasketut tunnusluvut, jotka sijoitetaan testisuureen kaavaan ja saadaan testisuureen arvo. Testin tulos ja **p**-arvo katsotaan joko Excel-funktion avulla tai esim. tämän kirjan liitteenä olevista taulukoista.
- 4) Käytetään Tixel-ohjelman testipohjaa, johon syötetään lasketut tunnusluvut. Tulokseksi saadaan paitsi testin tulos selkokielellä myös luottamusväli (ks. 8.2).

7.8 Testin oletukset eli vaatimukset

Jotta testisuureen otantajakauma voidaan määrätä, on tiettyjen *oletusten* eli *vaatimusten* oltava voimassa. Nämä oletukset vaihtelevat testistä toiseen, mutta eräät oletukset ovat yhteisiä kaikille testeille. Tällainen oletus on esim. otoksen *satunnainen poiminta*. Ilman sitä ei tilastollista testausta voida suorittaa.

Eräät oletukset ovat ehdollisia. Esimerkiksi keskiarvo-testeissä muuttujan jakauman normaalisuutta koskeva oletus on välttämätön vain, jos otoskoko **n** on pieni (esim. **n** pienempi kuin 30).

8 Tilastolliset testit

8.1 Testien ryhmittely

Tilastolliset testit voidaan luokitella neljän kriteerin mukaan.

1. **Mistä tunnusluvusta on kyse** — Tämän jaottelun mukaan testien tärkeimmät pääryhmät ovat
 - Keskiarvotestit
 - Prosenttilukutestit
 - Riippuvuuslukutestit
 - Varianssitestit
 2. **Mikä on vertailtavien ryhmien lukumäärä** — Tämän jaottelun mukaan puhutaan yhden, kahden ja useamman ryhmän (otoksen) testeistä.
 3. **Ovatko ryhmät (otokset) toisistaan riippumattomia vai riippuvia** — Otokset ovat toisistaan riippumattomia, jos otoksilla ei ole mitään ”tekemistä keskenään”. Esimerkiksi, jos naisista ja miehistä poimitaan erilliset otokset, niin otokset ovat keskenään riippumattomia. Näin on asia myös silloin,
-

kun otos jaetaan analyysivaiheessa naisten ja miesten ryhmiin (otoksiin). Riippuvia otoksia syntyy seuraavista syistä:

- a) **Ennen-jälkeen asetelma** — Siinä tutkimusyksiköt mitataan ennen kuin altistava tekijä on vaikuttanut. Mittaus toistetaan kun tekijä on vaikuttanut. Esimerkiksi koehenkilöiden asenteet mitataan ennen koulutustapahtumaa ja samat asenteet mitataan koulutustapahtuman jälkeen. Tällöin halutaan tutkia, onko asenteissa tapahtunut tilastollisesti merkitsevää muutosta
 - b) **Vastinparimenettely** — Vastinparimenettelyllä on tarkoitus vähentää satunnaisvaihtelua tilastollisessa testauksessa ja saada siten eri käsitteilyjen vaikutukset paremmin esille. Vastinparimenettelyssä yksiköt jaetaan relevanttien tekijöiden suhteen homogeenisiin pareihin ja parien sisällä arvotaan yksiköt arvotaan eri käsittelyille. Vertailtavien käsittelyjen lukumäärä voi olla myös kahta suurempi.
 - c) **Eri muuttujien vertailu** — Esimerkiksi kaksikymmentä maistajaa arvioi kolmea eri tuotetta (esim. jugurttia) siten, että kukin maistajaa antaa arvosanan jokaisesta tuotteesta. Tällöin mittaukset ovat keskenään riippuvia. Samoin jos koehenkilöiden vasemmasta ja oikeasta aivolohkosta tehdään mittauksia, niin mittauksien tulokset ovat keskenään riippuvia.
4. **Parametriset ja non-parametriset testit** — Parametrisissä testeissä on enemmän jakaumiin liittyviä oletuksia, minkä vuoksi non-parametrisia testejä kutsutaan myös jakaumista vapaiksi testeiksi. Parametristen testien voimakkuus on yleensä suurempi kuin non-parametristen testien.

Seuraavissa luvuissa eri testien käsittely on ryhmitelty ensin sen mukaan mistä tunnusluvusta on kysy. Tämän ryhmittelyn sisällä testejä on käsitelty kohdan otosten lukumäärän ja riippuvuuden mukaisissa alaluvuissa.

Kunkin testin yhteydessä on esitetty seuraavat seikat:

- Nollahypoteesi
- Vaihtoehtoiset hypoteesit
- Oletukset
- Testisuureen kaava
- Tixel-ohjelman testipohja ja esimerkkejä.

Huom. Vaikka tilastollinen testaus voidaan osittain korvata luottamusvälitarkastelulla, ei tätä pidä ottaa yleiseksi käytännöksi. On nimittäin tilanteita, joissa nämä kaksi menettelyä johtavat eri tilastolliseen johtopäätökseen. Tixel-

ohjelman testipohjissa on kuitenkin myös 95 %:n luottamusvälit, joten tilastollisen testin tulosta voidaan helposti verrata luottamusväliin.

8.2 Tixelin testipohja

Tixel-ohjelmassa on eri testejä varten testipohjia, joihin syötetään tarvittavat tunnusluvut ja otoskoon arvot. Kun viimeinen luku syötetty, näkyy testipohjassa seuraavat tulosteet:

- Testisuureen arvo
- ja 2-suuntaisen testin riskitodennäköisyydet
- testin tuloksen sanallinen luonnehdinta
- luottamusväli

Esimerkkinä testipohjasta on kahden riippumattoman otoksen prosenttilukutesti.

KAHDEN RIIPPUMATTOMAN OTOKSEN PROSENTTILUKUTESTI				
AINEISTO:				
MUUTTUJA:				
	1. O T O S	2. O T O S	TESTISUURE	
%-luku:			#DIV/0!	
Otoskoko:			RISKITASO	RISKITASO
			1-suunt.	2-suunt.
		p-arvo:	#DIV/0!	#DIV/0!
			#DIV/0!	#DIV/0!
Perusjoukkojen prosenttilukujen erotuksen				
		piste-estimaatti:	0	
		95%:n luottamusväli:	#DIV/0!	#DIV/0!
<p>Ohje: Täytä harmaalla merkityt solut (prosenttiluvut ilman %-merkkiä). Jatko: A vaa Tixel-valikko ja valitse vaihtoehto.</p>				

Neljään rasteroituun soluun syötetään arvot, jonka jälkeen #DIV/0! – merkittyihin soluihin ilmestyvät oikeat tulokset.

Yksisuuntaiselle vaihtoehdoiselle hypoteesille riskitaso lasketaan aina testisuureen itseisarvosta. Tämä ei aiheuta virheitä, jos otostestisuure on vaihtoehdoisen hypoteesin osoittamalla alueella. Jos sen sijaan näin ei ole, tulee nollahypoteesi hyväksyä itsestään selvyytensä riippumatta siitä mitä riskitodennäköisyys näyttää.

Syöttöpohjan avulla voidaan myös havainnollistaa, miten muutokset syöttötiedoissa, esim. keskiarvoissa, keskihajonnoissa ja otoskoissa vaikuttavat testin lopputulokseen.

8.3 Prosenttilukutestit

Seuraavassa käsitellään yhden otoksen prosenttilukutestiä sekä kahden riippumattoman otoksen ja kahden riippuvan otoksen prosenttilukutestejä. Useamman kuin kahden korreloimattoman otoksen testi esitetään yhteensopivuustestien yhteydessä (ks. luku 8.7).

8.3.1 Yhden otoksen prosenttilukutesti

Testi testaa poikkeako otoksesta laskettu prosenttiluku p tilastollisesti merkittävästi nollassa nollahypoteesin ilmoittamasta prosenttiluvusta θ_0 . Näiden kahden prosenttiluvun lisäksi on tunnettava otoskoko n .

Oletukset. Seuraavien oletusten on oltava voimassa:

- 1) Otos on poimittu satunnaisesti perusjoukosta.
- 2) Jotta testisuure noudattaisi likimäärin normaalijakaumaa, on otoskoon n oltava riittävän suuri. Tämä raja riippuu prosenttiluvun todellisesta arvosta, joka on kuitenkin tuntematon. Sen sijaan käytetään nollassa nollahypoteesin mukaista arvoa θ_0 . Saadaan seuraavaa taulukko.

Taulukko 12: Otoskoon n alaraja normaalijakauman käytölle.

Pienempi luvuista θ_0 , $100-\theta_0$	Otoskoon n alaraja
50	30
40	50
30	80
20	200
10	600
5	1400

Jos otoskoko on liian pieni tämän taulukon mukaan, on testauksessa käytettävä binomijakaumaa. Tätä käsitellään luvussa 8.3.1.1. (ks. binomijakaumasta luku 4.4, s. 36; otosprosenttiluvun otantajakaumasta ks. luku 4.3, s. 34).

Hypoteesit. Nollahypoteesi ilmaisee, että perusjoukon prosenttiluvun arvo annettu luku θ_0 . Vaihtoehtoinen hypoteesi puolestaan tuo esiin ne vaihtoehdot, jotka tulevat kysymykseen, kun nollassa nollahypoteesi ei ole tosi.

Testisuure. Testisuure on suhde, jonka osoittajana on otosprosenttiluvun poikkeama nollassa nollahypoteesin ilmoittamasta arvosta ja jakajana otosprosenttiluvun keskiluvun silloin, kun nollassa nollahypoteesi on tosi. Testisuureen kaava on

$$Z = \frac{p - \theta_0}{\sqrt{\theta_0(100 - \theta_0) / n}}$$

Jos nollihypoteesi on tosi, testisuure Z noudattaa normaalijakaumaa $N(0,1)$.

- **ESIM:** Vaaleissa oli kaksi ehdokasta A ja B, joiden kannatusosuudet tuhannen hengen otoksessa olivat 47 % ja 53 %. Poikkeavatko prosenttiluvut tilastollisesti merkitsevästi toisistaan?

Kyseessä on yhden otoksen prosenttilukutesti, jossa otoksen prosenttiluku on 53 % (tai 47 %) ja nollihypoteesin mukainen prosenttiluku 50 %. Sijoitetaan luvut testisuureen kaavan, joka antaa seuraavan tuloksen

$$Z = \frac{53 - 50}{\sqrt{50(100 - 50) / 1000}} = 1,9$$

Tixelin syöttöpohja antaa seuraavan tuloksen. Vaihtoehtoinen hypoteesi on kaksisuuntainen, koska etukäteen ei voitu mitään vaihtoehtoa sulkea pois.

YHDEN OTOKSEN PROSENTTILUKUTESTI			
	OTOS	NOLLAHYPOTEESI	TESTISUURE
%-luku:	53	50	1,90
Otoskoko:	1000		
			RISKITASO
			1-suunt. 2-suunt.
		p-arvo:	2,890 % 5,781 %
			Melkein merkitsev. Oireellinen
Perusjoukon prosenttiluvun 95%:n luottamusväli:			49,91 56,09

Kaksisuuntaisen testin riskitaso 5,78 % osoittaa, että prosenttiluvut eivät eroa tilastollisesti merkitsevästi toisistaan. Ero on ainoastaan tilastollisesti oireellinen. Ehdokkaan B kannatuksen 95 %:n luottamusväli on (49,9 % – 56,1 %). Siten 50 % on mahdollinen arvo B:n kannatukselle, mikä on sopusoinnussa testin tuloksen kanssa.

Jos otoskoko kasvaa 1500 henkeen, kasvaa testisuure arvoon 2,32 ja tulos on melkein merkitsevä. Siten voidaan kohtuullisella varmuudella sanoa, että perusjoukossa ei vallitse tasatilanne, vaan toisen ehdokkaan kannatus on suurempi. Taulukosta nähdään, että 50 % ei kuulu teoreettisen prosenttiluvun luottamusväliin.

YHDEN OTOKSEN PROSENTTILUKUTESTI				
	OTOS	NOLLAHYPOTEESI	TESTISUURE	
%-luku:	53	50		2,32
Otoskoko:	1500		RISKITASO	RISKITASO
			1-suunt.	2-suunt.
		p-arvo:	1,008 %	2,016 %
			Melkein	Melkein
			merkitsevä	merkitsevä
Perusjoukon prosenttiluvun 95%:n luottamusväli:			50,47	55,53

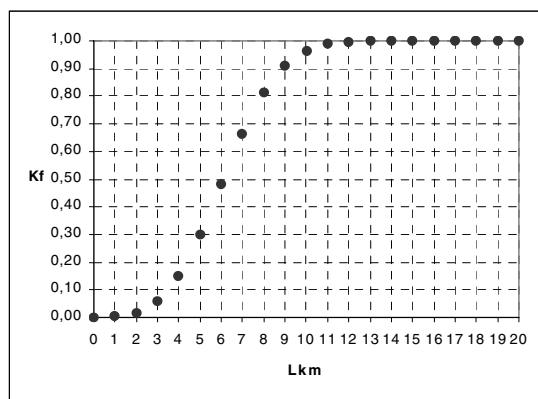
8.3.1.1 Pieni otoskoko.

Jos otoskoko ei täytä Taulukko 12:n vaatimuksia, on testi rakennettava binomijakauman pohjalle (ks. myös luku 4.3, s. 34). Esimerkiksi tutkitaan henkilön ESP-kykyjä (ESP=extra sensory perception). Testataan häntä korteilla, joissa hän tietää olevan kolme väri vaihtoehtoa. Korteja on käytössä kaksikymmentä ja satunnaisuuttujana X on oikein tiedettyjen värien lukumäärä. Nollahypoteesi on pelkästään arvaamalla saatu oikeiden vastausten prosenttiosuus eli $\theta = 33,3\%$. Vaihtoehtoisen hypoteesin mukaan θ on suurempi kuin $33,3\%$.

Taulukko 12 mukaan toistokokeiden lukumäärä pitäisi olla vähintään 80, jotta normaalijakaumaan perustuvaa testiä voitaisiin käyttää. Tässä tapauksessa se on vain kaksikymmentä, joten on käytettävä binomijakaumaa.

Testisuurena käytetään nyt suoraan X :n arvoa eli oikeiden vastausten lukumäärää. Jos nollahypoteesi on tosi, noudattaa X binomijakaumaa **Bin(20; 0,33)**. Koska vaihtoehtoinen hypoteesi on oikealle yksisuuntainen, määrätään kriittinen alue X :n vaihtelualueen 0-20 oikeasta reunasta. Apuna käytetään binomijakauman todennäköisyyksiä (tiheys- tai kertymäfunktioita). Kuva 38:stä nähdään likimäärin kertymäfunktion arvot.

Kuva 38: Binomijakauman kertymäfunktion kuvaaja



Eniten kiinnostaa kertymäfunktion arvo pisteessä $x = 10$. Tämä tarkka arvo on 0,9626, joten todennäköisyys, että X :n arvo on suurempi kuin kymmenen, on

$$1 - 0,9626 = 0,0374 \text{ eli } 3,74 \%$$

Tämä sillä ehdolla, että nollahypoteesi on tosi. Jos X :n havaittu arvo on suurempi kuin kymmenen, hylkäämme nollahypoteesin 3,74 %:n riskitasolla. Excelissä voidaan määrätä kriittinen arvo halutulle riskitasolle ko. parametrein. Kirjoitetaan seuraava funktio Excelin soluun:

$$= \text{CRITBINOM}(20; 0,33; 0,95) = 10 ,$$

jolloin todetaan, että myös 5 %:n riskitasolla haluttu kriittinen arvo on 10.

8.3.2 Kahden riippumattoman otoksen prosenttilukutesti

Kaksi otosta ovat riippumattomia, jos niillä ei ole mitään yhteyttä keskenään. Seuraavassa alaluvussa 8.3.3 käsitellään kahden riippuvan otoksen tapausta. Tarkastellaan kahta perusjoukkoa (esim. naiset ja miehet) ja tietyn ominaisuuden omaavien prosentuaalisia osuuksia näissä perusjoukoissa (tupakoivien osuudet naisten ja miesten perusjoukoissa). Olkoot nämä prosentit θ_1 ja θ_2 . Nollahypoteesin mukaan prosenttiluvut ovat yhtä suuret eli

$$H_0: \theta_1 = \theta_2$$

Näiden hypoteesien testaamiseksi poimitaan perusjoukoista n_1 ja n_2 suuruiset riippumattomat otokset. Testisuureen kaava on

$$Z = \frac{p_1 - p_2}{\sqrt{p(100 - p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

jossa p on prosenttiosuus, kun otokset yhdistetään eli $p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$.

- **ESIM:** Erään ominaisuuden prosenttiluvut miesten ja naisten otoksissa olivat vastaavasti 30% ja 40%. Poikkeavatko prosenttiluvut tilastollisesti merkitsevästi toisistaan, kun otoskoot olivat 200 ja 180?

Yhteisen prosenttiluvun p arvoksi saadaan $p = (200 \cdot 30 + 180 \cdot 40) / (200 + 180) = 34,7$. Asetetaan arvot kaavaan ja saadaan

$$Z = \frac{30 - 40}{\sqrt{34,7(100 - 34,7)\left(\frac{1}{200} + \frac{1}{180}\right)}} = -2,04 .$$

Tixelin testipohja antaa seuraavan tuloksen:

KAHDEN KORRELOIMATTOMAN OTOKSEN PROSENTTILUKUTESTI				
	1. OTOS	2. OTOS	TESTISUURE	
%-luku:	30	40		-2,04
Otoskoko:	200	180	RISKITASO	RISKITASO
			1-suunt.	2-suunt.
			p-arvo: 2,048 %	4,096 %
			Melkein merkitsevä	Melkein merkitsevä
Perusjoukkojen prosenttilukujen erotuksen				
			piste-estimaatti:	-10
			95%:n luottamusväli:	-19,59 -0,41

Riskitasot osoittavat, että prosenttilukujen ero on tilastollisesti melkein merkitsevä. Luottamusväli haarukoi erotuksen välille (-19.6% , $-0,4\%$), joten erotuksen arvo nolla ei kuulu luottamusväliin. Tämä on sopusoinnussa testin antaman tuloksen kanssa.

8.3.3 Kahden riippuvan otoksen prosenttilukutesti

Kahden riippuvan otoksen asetelma syntyy esimerkiksi vastinparimenettelyn kautta tai ennen-jälkeen –tutkimusasetelmasta. Yleisesti näihin mittauspareihin viitataan nimillä Mittaus_1 ja Mittaus_2. Mitattava ominaisuus on kaksiluokkainen, esimerkiksi kannattaako henkilö tiettyä esitystä vai ei. Prosenttiluku kertoo toisen vaihtoehdon kannatuksen suhteellisen osuuden.

Vaikka tilastollinen nollahypoteesi on sama kuin edellisen alaluvun tapauksessa, eli että kahteen mittaukseen liittyvien perusjoukkojen prosenttiluvut ovat yhtäsuuret, on testisuure kuitenkin eri. Sitä varten on laskettava seuraava 2*2-ristiintaulukko. Olkoot 0 ja 1 ensimmäisen ja toisen mittauksen tulosvaihtoehdot. Saadaan seuraava ristiintaulukko, jonka lukumääriä merkitään **a**, **b**, **c** ja **d**:

Taulukko 13: Riippuvien mittausten ristiintaulukko

		Mittaus 2		
		0	1	yht.
Mittaus 1	0	a	b	a+b
	1	c	d	c+d
yht.		a+c	b+d	n

Nollahypoteesi voidaan nyt määritellä uudelleen siten, että muutostodennäköisyys nolhasta ykköseen on yhtä suuri kuin todennäköisyys muutokselle ykkösestä nolnaan. Testiä varten keskitytään siten ainoastaan frekvensseihin **b** ja **c**. Nollahypoteesin mukaan näiden frekvenssien ero johtuu ainoastaan sattumasta.

Testisuurena käytetään yhden otoksen testisuureta, jolla testataan poikkeako otosprosenttiluku $p = b/(b+c)*100$ tilastollisesti merkitsevästi 50 prosentista. Testisuurelle voidaan johtaa vielä yksinkertaisempi kaava frekvenssien funktiona

$$Z = \frac{b - c}{\sqrt{b + c}},$$

joka noudattaa normaalijakaumaa $N(0,1)$, jos nollahypoteesi on tosi.

Huom. Tätä testiä kutsutaan myös McNemarin testiksi. Seuraavassa esimerkissä havainnollistamme McNemarin testin käyttöä.

- **ESIM:** 148 hengen otos mitattiin ennen koulutustilaisuutta ja sen jälkeen. Koulutuksen sisällön kiinnostavuutta koskevien Kyllä vastausten (=1) osuus oli ennen koulutustilaisuutta 47 % ja koulutustilaisuuden jälkeen 39 %. Voidaanko sanoa, että koulutusohjelman sisältöä koskeva kiinnostus on lisääntynyt sattumaa varmemmin koulutustilaisuuden vaikutuksesta ?

Näillä tiedoilla ei testiä vielä voida tehdä, vaan tarvitaan ennen- ja jälkeen – muuttujien ristiintaulukko, joka on esitettyssä testipohjassa. Tästä ristiintaulukosta poimitaan (viitaten Taulukko 13:een) frekvenssien **b** ja **c** arvot, jotka asetetaan McNemarin kaavaan, joka antaa seuraavan tuloksen

$$Z = \frac{28 - 16}{\sqrt{28 + 16}} = 1,81 .$$

Tixelin testipohja päättyy samaan tulokseen:

	1	2	%
1	42	28	47%
2	16	62	53%
%	39%	61%	
Testisuure	Riskitaso	Riskitaso	
1,80906807	1-suunt.	2-suunt.	
	0,03523522	0,070470433	

Kaksisuuntainen testi ei anna prosenttiluvuille tilastollisesti merkitsevää eroa. Sen sijaan yksisuuntainen testin mukaan prosenttiluvut eroavat tilastollisesti melkein merkitsevästi. Käyttämällä taulukkoa (Liite1) tai Excelin NORMSDIST –funktiota on myös mahdollista määrätä testisuureen arvoon liittyvä riskitodennäköisyys. Eli esim. Excelissä kirjoitetaan soluun seuraava kaava

$$=1-NORMSDIST(1,81) = \underline{0,035} .$$

Eli voidaan tehdä johtopäätös, että koulutustilaisuus on lisännyt kiinnostusta sattumaa varmemmin 3,5 %:n riskitasolla.

Tämä on siis yksisuuntaisen testin riskitodennäköisyys (joita NORMSDIST –funktio antaa). Mikäli testi on kaksisuuntainen saadaan riskitodennäköisyys kertomalla saatu todennäköisyys kahdella, jolloin riskitodennäköisyys on 0,070. Tällöin tulos ei olisi tilastollisesti merkitsevä.

8.4 Keskiarvotestit

Keskiarvotesteillä testataan keskiarvoihin liittyviä hypoteeseja. Otosten lukumäärän ja riippumattomuuden tai riippuvuuden mukaan keskiarvotestit ryhmitellään seuraavasti:

- Yhden otoksen testi
- Kahden riippumattoman otoksen testi
- Kahden riippuvan otoksen testi
- Useamman kuin kahden riippumattoman otoksen testi eli yksisuuntainen varianssianalyysi

Näistä käsitellään kolmea ensimmäistä tarkemmin – yksisuuntaisesta varianssianalyysistä esitetään ainoastaan syöttöpohja. Lukija voi kysyä, mihin on jäänyt useamman kuin kahden riippuvan otoksen testi. Tämä testin käsittely edellyttää varianssianalyysin syvällisempää tuntemista ja sen vuoksi sitä ei tässä käsitellä.

8.4.1 Yhden otoksen keskiarvotesti

Muuttujan X keskiarvoa perusjoukossa merkitään μ :llä, joka on tuntematon. Halutaan tutkia, olisiko μ :n arvo tietty tunnettu nollahypoteesin mukainen arvo μ_0 . Tähän etsitään vastausta perusjoukosta poimitun otoksen ja siitä lasketun otoskeskiarvon \bar{X} avulla. Seuraavassa on esimerkkejä tällaisesta tutkimusasetelmasta:

- Tiedetään, että edellisenä vuonna työstä poissaolopäivien keskiarvo kuka kulta kohden tietyllä teollisuuden sektorilla oli μ_0 . Seuraavan vuoden alkupuolella tuli ”kentältä” viestejä, että työstä poissaolojen määrä on rajussa nousussa. Koska mitään virallisia tilastoja ei ollut nopeasti käytettävissä, päätetään tehdä asiasta otantatutkimus ja testata tilastollisesti, onko työstä poissaolopäivien keskiarvo kasvanut μ_0 :sta.
- Usein yhden otoksen keskiarvotestiä sovelletaan *erotusmuuttujaan*. Esimerkiksi otokseen poimitujen henkilöiden fyysinen suorituskyky mitataan ennen valmennusjaksoa ja sen jälkeen. Fyysisen suorituskyvyn muutoksen mittari muodostetaan jälkeen ja ennen mittausten erotuksena, jolloin eräs nollahypoteesi on, että systemaattista muutosta suorituskyvyssä ei ole ollut. Jos valmennusmenetelmä lupaa kymmenen prosentin kasvun fyysisessä suorituskyvyssä, on nollahypoteesi luonnollisesti muotoiltava eri tavalla. Ennen- ja jälkeen –tutkimusasetelmaa käsitellään tarkemmin kahden riippuvan otoksen keskiarvotestin yhteydessä.

Hypoteesit. Nollahypoteesi ilmaisee, että μ :n arvo on annettu vakio μ_0 . Vaihtoehtoinen hypoteesi voi olla joko yksi- tai kaksisuuntainen.

Oletukset. Oletukset koskevat otoksen poimintaa ja muuttujan jakaumaa.

- Otos on perusjoukosta satunnaisesti.
- Muuttuja noudattaa perusjoukossa normaalijakaumaa $N(\mu, \sigma^2)$

Näistä kahdesta oletuksesta satunnaisotantaa koskeva on ehdottomasti tärkeämpi. Perusjoukon jakauman normalisuus on välttämätön vain jos otoskoko n hyvin pieni (esim. $n < 30$).

Testiä varten on otoskeskiarvon lisäksi tunnettava keskihajonta. Yleensä se estimoidaan (lasketaan) otoksesta, mutta joskus keskihajonta on etukäteen tunnettu. Tämä seikka vaikuttaa testisuureeseen ja otosjakaumaan. Testisuure t on samaa muotoa edelläkin: otoskeskiarvon poikkeama nollahypoteesista jaetaan keskivirheellä eli

$$t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}} = \sqrt{n} (\bar{X} - \mu_0)/s ,$$

jossa s on otoskeskihajonta. Testisuure t noudattaa t -jakaumaa vapausasteilla $df = n-1$, jos nollahypoteesi on tosi.

On muistettava, että t -jakauma lähestyy standardoitua normaalijakaumaa, kun vapausasteet kasvavat. Jo sadan havainnon otoskoolla ero on merkityksetön.

Huom. Jos muuttujan perusjoukon keskihajonta σ on tunnettu, sijoitetaan se testisuureen lausekkeessa otoskeskihajonnan s paikalle, jolloin otantajakauma on normaalin $N(0,1)$ — eli tällöin testisuure noudattaa Z -jakaumaa.

- ESIM:** Testataan poikkeako otoskeskiarvo 98 tilastollisesti merkitsevästi nollahypoteesin mukaisesta arvosta 100. Kahden sadan hengen otoksesta laskettu keskihajonta on 8,5.

Voidaan laskea käsin käyttäen ko. t -testisuureen kaavaa — tai luvut voidaan asettaa esimerkiksi Tixelin syöttöpohjaan, kuten alla on menetelty.

YHDEN OTOKSEN KESKIARVOTESTI		
OTOS	NOLLAHYPOTEESI	TESTISUURE
Keskiarvo: 98	100	-3,32756132
Keskihajonta: 8,5		VAP. ASTE
Otoskoko: 200		199
	RISKITASO	RISKITASO
	1-suunt.	2-suunt.
p-arvo:	0,052 %	0,104 %
	Erittäin merkitsevä	Merkitsevä
Perusjoukon keskiarvon 95%:n luottamusväli:	96,81	99,19

Voidaan todeta, että otoskeskiarvo 98 poikkeaa H_0 :n mukaisesta tilanteesta eli sadasta tilastollisesti merkitsevästi $p = \underline{0,104\%}$.

8.4.2 Kahden riippumattoman otoksen keskiarvotesti

Merkitään muuttujan X keskiarvoa ja keskihajontaa kahdessa perusjoukossa μ_i :llä ja σ_i :llä ($i=1,2$). Poimitaan perusjoukoista n_1 ja n_2 suuruiset riippumattomat otokset, joista lasketaan otoskeskiarvot \bar{X}_1 ja \bar{X}_2 sekä vastaavat keskihajonnat s_1 ja s_2 . Tehtävänä on testata tämän otosinformaation perusteella, onko perusjoukkojen keskiarvojen erotus jokin annettu vakio δ , joka yleensä on nolla. Eli lähes poikkeuksetta testataan, ovatko perusjoukkojen keskiarvot yhtä suuret.

Hypoteesit. Nollahypoteesina on, että perusjoukkojen keskiarvojen erotus on annettu vakio δ eli

$$H_0 : \mu_1 - \mu_2 = \delta \text{ (on yleensä nolla).}$$

Vaihtoehtoinen hypoteesi H_1 voi olla joko yksi- tai kaksisuuntainen.

Oletukset. Testin oletukset ovat seuraavat

- Perusjoukoista on poimittu satunnaisotokset, jotka ovat toisistaan riippumattomat. Riippumattomuus tarkoittaa, että otoksilla ei ole mitään yhteyttä toisiinsa.
- Muuttujan jakaumat näissä kahdessa perusjoukossa ovat normaaliset
- Muuttujan varianssit perusjoukossa ovat yhtä suuret eli $\sigma_1^2 = \sigma_2^2$. Tällöin tietenkin myös perusjoukkojen keskihajonnat ovat yhtä suuret. Syy miksi tässä korostetaan varianssien yhtäsuuruutta on se, että varianssien yhtäsuuruudelle on olemassa tilastollinen testi (ks. kappale 8.8). Jos tämä oletus ei ole voimassa, on testisuuretta hieman muokattava. Tätä vaihtoehtoa ei tässä tarkemmin käsitellä, mutta luvussa 8.5.2 on esitelty Excelin komento, jolla tämä tapaus voidaan ratkaista.

Testisuure on

$$t = \frac{\bar{X}_1 - \bar{X}_2 - \delta}{\sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

joka noudattaa t -jakaumaa vapausasteilla $df = n_1 + n_2 - 2$, jos nollahypoteesi on tosi.

- **ESIM:** Testataan eroavatko kahden sektorin otostutkimuksista lasketut keskipalkat 10200 mk ja 10500 mk tilastollisesti merkitsevästi toisistaan. Vastaavat keskihajonnat ovat 1200 mk ja 1300 mk sekä otoskoot 200 ja 150 henkeä.

Asetetaan annetut arvot kaavaan ja saadaan testisuureen arvoksi

$$t = \frac{10200 - 10500}{\sqrt{\frac{199 \cdot 1200^2 + 149 \cdot 1300^2}{200 + 150 - 2} \left(\frac{1}{200} + \frac{1}{150} \right)}} = -2,23$$

Tixelin testipohjaan syötetään keskiarvojen lisäksi keskihajonnat ja otoskoot.

	1. OTOS	2. OTOS	TESTISUURE	
Keskiarvo:	10200	10500	-2,23304285	
Keskihajonta:	1200	1300	VAP. ASTE	
Otoskoko:	200	150	348	
			RISKITASO	RISKITASO
			1-suunt.	2-suunt.
			p-arvo: 1,309 %	2,618 %
			Melkein merkitsevä Melkein merkitsevä	
Perusjoukkojen keskiarvojen erotuksen				
	piste-esimaatti		-300	
	95%:n luottamusväli:		-564,2	-35,8

Kaksisuuntaisen testin riskitaso osoittaa, että keskipalkoissa on tilastollisesti melkein merkitsevä ero riskitasolla $p = 2,618 \%$.

8.4.3 Perusjoukon varianssit tunnetaan

Jos perusjoukon varianssit tunnetaan, voidaan otoskeskiarvojen erotuksen keskivirhe laskea tarkasti ilman satunnaisvirhettä. Testisuureta merkitään tällöin Z :lla, koska otantajakauma on standardoitu normaalijakauma. Testisuureeksi saadaan

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - \delta}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

joka siis noudattaa $N(0,1)$ -jakaumaa, jos nollahypoteesi on tosi. Mikäli siis perusjoukon varianssit tunnetaan, mikä on suhteellisen harvinaista, niin Z - ja t -testisuureet eroavat merkittävämmiin vain silloin, kun otoskoot ovat hyvin pienet (esim. alle kahdenkymmenen).

8.5 Kahden riippuvan otoksen keskiarvotesti

Luvussa 8.1 on käsitelty eri tapauksia, joissa otosten (mittausten) välille syntyy riippuvuutta. Vaikka edellisen luvun testisuure voidaan silloinkin teknisesti laskea, on testin antama riskitodennäköisyys kuitenkin väärä.

- **ESIM:** Tarkastellaan esimerkkiä, jossa viideltä satunnaisesti valitulta koehenkilöltä mitattiin reaktioaika selvänä ($=X$) ja yhden promillen humalatilassa ($=Y$). Halutaan testata, onko promillen humalatilassa reaktioaikoja hidastava vaikutus tilastollisesti merkitsevä. Mittaukset ovat tässä tapauksessa riippuvia, koska koehenkilöiltä tehtiin kaksi mittausta. Tulokset ovat seuraavassa taulukossa.

Taulukko 14. Reaktioaikoja koskevat mittaukset

Henkilö	Selvänä X	Promillen humala Y	Erotus d
1	0,62	0,71	0,09
2	0,81	0,80	-0,01
3	0,74	0,79	0,05
4	0,53	0,55	0,02
5	0,69	0,73	0,04
Keskiarvo	0,678	0,716	0,038
Keskihajonta	0,108	0,100	0,037
Korrelaatio(X,Y) = 0,9395			

Taulukosta käy ilmi myös erotusmuuttujan $d = Y - X$ arvot sekä näiden kolmen muuttujan keskiarvot ja keskihajonnat. Jatkossa tarvitaan erityisesti erotusmuuttujan keskiarvoa \bar{d} ja keskihajontaa s_d .

Huom. Vaikka erotusmuuttujan keskiarvo on alkuperäisten muuttujien keskiarvojen erotus, niin vastaava ei päde keskihajontoihin. Tämän luvun lopussa palataan vielä siihen, miten erotusmuuttujan keskihajonta voidaan johtaa alkuperäisten muuttujien tunnuslukujen funktiona.

Tilastollinen testi perustuu erotusmuuttujan arvoihin, joihin sovelletaan yhden otoksen testiä. Testisuure on

$$t = \frac{\bar{d}}{s_d/\sqrt{n}},$$

joka noudattaa t -jakaumaa vapausasteilla $df = n-1$, jos nollahypoteesi on tosi.

Taulukko 14:n tiedoista saadaan t -testisuureen⁴ arvoksi

$$t = \frac{0,038}{0,037/\sqrt{5}} = 2,296,$$

jonka yksisuuntaisen testin riskitaso on 0,0417 (eli siis 4,17 %), ja joka voidaan määrätä esim. liitteen 2 taulukosta (vapausasteilla 4) tai Excelin TDIST – funktiolla seuraavasti.

$$=TDIST(2,296; 4; 1) = \underline{0,0417}$$

Riskitason pohjalta voidaan päätellä, että keskiarvojen ero on tilastollisesti *melkein* merkitsevä.

8.5.1 Erotusmuuttujan keskihajonta tunnuslukujen avulla

Jos erotusmuuttujan d tilastoyksikkökohtaiset arvot eivät ole käytettävissä, voidaan d :n keskihajonta laskea myös muuttujien X ja Y tunnuslukujen avulla. Voimassa on seuraava kaava:

$$s_d^2 = s_x^2 + s_y^2 - 2r_{xy} s_x s_y.$$

Sijoittamalla Taulukko 14:n tunnusluvut tähän kaavaan saadaan seuraava tulos

$$s_d^2 = 0,108^2 + 0,1^2 - 2 * 0,9395 * 0,108 * 0,1 = 0,00137,$$

jonka neliöjuurena saadaan erotusmuuttujan d keskihajonnaksi **0,037** eli sama arvo kuin Taulukko 14:ssä.

⁴ Excelissä olisi mahdollisuus käyttää TTEST-funktiota testin suorittamiseen; ks. funktion käytöstä luku 4.5.4, s. 44. Myös Excelin analyysityökalut (ks. 8.5.2, s. 90) mahdollistavat tämän ja muun tyyppisten keskiarvotestien suorittamisen.

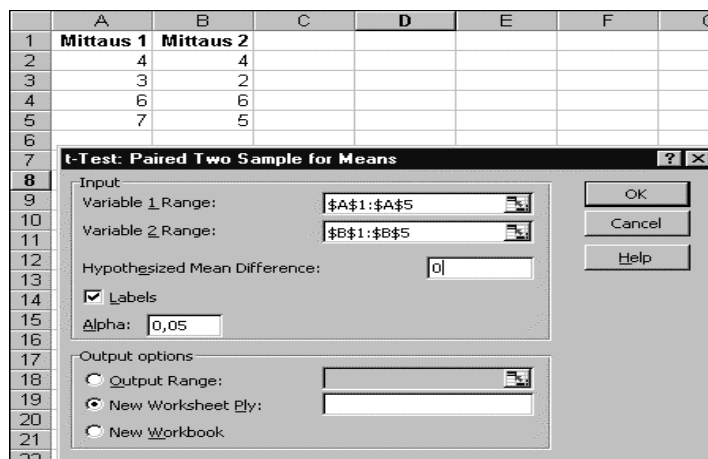
8.5.2 Excelin keskiarvotestit

Excelin Tools/Data Analysis –valikossa on kahden otoksen keskiarvotestejä neljälle eri tilanteelle. Kaikki testit edellyttävät havaintoarvojen tuntemista. Testejä ei voi siis tehdä valmiiksi laskettujen tunnuslukujen avulla.

1) t-Test: Paired Two Sample for Means

Testi on *riippuvien otosten* keskiarvotesti (ks. luku 8.5). Seuraavassa kuvassa on testin määrittelyikkuna.

Kuva 39: Riippuvien otosten keskiarvotesti (Excel)



Määrittelyikkunan *Range* –kohtiin merkitään solualueet, joissa lukusarjat sijaitsevat. *Hypothesized Mean Difference* –kohtaan merkitään perusjoukkojen keskiarvojen erotus eli yleensä nolla. *Labels* –kohta valitaan, jos lukusarjojen tunnukset ovat mukana *Range* –kohdassa. Alpha on riskitaso.

2) t-Test: Two-Sample Assuming Equal Variances

Kahden riippumattomien otoksen t-testi, jossa oletetaan perusjoukon varianssien olevan yhtä suuret. (ks. luku 8.4.2).

3) t-Test: Two-Sample Assuming Unequal Variances

Kahden riippumattomien otoksen t-testi, jossa perusjoukon varianssit voivat olla erisuuret (ks. luku 8.4.2).

4) z-test: Two Sample for Means

Kahden riippumattoman otoksen keskiarvotesti, kun muuttujan varianssit molemmissa perusjoukoissa oletetaan tunnetuiksi. Tällöin **Z**-testisuure noudattaa standardoitua normaalijakaumaa, ei **t**-jakaumaa (ks. luku 8.4.3).

8.5.3 Yksisuuntainen varianssianalyysi

Tässä on tarkoitus hyvin lyhyesti käsitellä useamman kuin kahden riippumattoman otoksen keskiarvotestiä, joka tunnetaan myös nimellä yksisuuntainen varianssianalyysi. Yksisuuntainen –termi johtuu siitä, että tilastoyksiköt on ryhmitelty yhden ryhmittelymuuttujan, esim. ikäryhmän mukaan. Varianssianalyysin nimi tulee siitä, että testisuure määritellään tiettyjen varianssien avulla, jolloin otantajakauma on F-jakauma (ks. luku 4.7, s. 50).

Tarkastellaan esimerkkiä, jossa on mitattu erään pankin asiakkaiden asiakastytyväisyyttä kolmessa konttorissa. Asteikko neljästä kymmeneen ja halutaan tietää, ovatko keskiarvoerot tilastollisesti merkitseviä. Testiä varten tarvitaan, kuten kahdenkin otoksen tapauksessa, ryhmien keskiarvot, keskihajonnat ja lukumäärät. Tixel-ohjelmassa on testipohja, johon nämä luvut voidaan syöttää.

Taulukko 15: Yksisuuntaisen varianssianalyysin esimerkki

Ryhmä	Lkm	Keskiarvo	Keskihajonta
Konttori 1	70	7,7	1,2
Konttori 2	60	7,9	1,3
Konttori 3	80	8,2	1,1
F = 3,35 vap. asteet: 2 ja 207			
p = 0,0372 Tilastollisesti melkein merkitsevä			

F-testisuureen arvo 3,35 on tilastollisesti melkein merkitsevä ($\alpha = 3,72\%$), joten kolmen konttorin asiakastytyväisyyden välillä on systemaattista eroa.

8.6 *Riippuvuuslukutestit*

Yleisimmät tavat mitata kahden muuttujan välistä riippuvuutta ovat khiin neliöön (χ^2) perustuvat *kontingenssikerroin C* tai *Cramerin phi ϕ_c* ja lineaarista riippuvuutta mittaava *korrelaatiokerroin (r)*.

8.6.1 Korrelaatiokerroin

Kahden muuttujan X ja Y korrelaatiokerroin mittaa muuttujien välistä lineaarista eli suoraviivaista riippuvuutta. Otokorrelaatiokerroin r on yhtä kuin muuttujien kovarianssi (s_{XY}) jaettuna keskihajontojen tulolla eli

$$r = \frac{s_{XY}}{s_X s_Y},$$

jossa kovarianssille s_{XY} on voimassa lauseke

$$S_{XY} = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{n - 1} = \frac{\sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n}}{n - 1}.$$

Huom. Tarkemmin sanottuna tässä on kyseessä *Pearsonin tulomomenttikorrelaatiokerroin*. On muitakin korrelaatiokerroimia kuten esim. Spearmanin järjestyskorrelaatiokerroin.

Merkitään perusjoukon korrelaatiokerrointa ρ :lla (rho). Nollahypoteesi yleensä asettaa ρ :lle arvon nolla eli nollahypoteesin mukaan muuttujien välillä ei ole lineaarista riippuvuutta. Luvussa 8.6.3 tarkastellaan yleisempää tapausta, jossa nollahypoteesin korrelaatiokerroimelle asettama arvo on nolosta poikkeava.

Korrelaatiokerroimen testi olettaa, että otos on satunnaisesti poimittu kahden muuttujan normaalijakaumasta. Testisuure on

$$t = \frac{r \sqrt{n - 2}}{\sqrt{1 - r^2}},$$

joka noudattaa t -jakaumaa vapausasteilla $df = n - 2$, jos nollahypoteesi on tosi.

- **ESIM:** Viidenkymmenen kahden havainnon otoksesta laskettu korrelaatiokerroin on $-0,25$. Poikkeako se tilastollisesti merkitsevästi nolosta? Oletetaan että vaihtoehtoinen hypoteesi on kaksisuuntainen.

Sijoitetaan annetut arvot kaavaan ja saadaan t -testisuureen arvoksi

$$t = \frac{-0,25 \sqrt{50}}{\sqrt{1 - (-0,25)^2}} = -1,826.$$

Liitteen 2 taulukosta (vapausasteilla 50) nähdään, että t -jakauman kriittinen arvo viiden prosentin riskitasolla kaksisuuntaisen testin tapauksessa on 2,009, joten nollahypoteesia ei hylätä. Sama seikka nähdään myös Excelin TDIST-funktiolla, joka antaa testisuureen arvoon liittyvän todennäköisyyden. Saadaan

$$=TDIST(1,826;50;2)=0,0738 \text{ eli } \underline{7,38\%}.$$

TDIST –funktion käytössä on huomattava, että ensimmäinen argumentti on testisuureen arvon *itseisarvo*.

8.6.2 Korrelaatiokertoimen kriittiset arvot

Usein korrelaatiokertoimia on testattavana suuri määrä, esimerkiksi korrelaatiomatriisin muodossa. Tällöin on järkevää ratkaista testisuureen lauseke otoskorrelaatiokertoimen r funktiona ja laskea testin kriittiset arvot suoraan otoskorrelaatiokertoimelle. Saadaan seuraava lauseke r :n kriittiselle arvolle.

$$r = \frac{t_{\alpha/2}}{\sqrt{n - 2 + t_{\alpha/2}^2}}.$$

Sijoittamalla tähän kaavaan t -arvot kolmella tyypillisellä riskitasolla ja 1- ja 2-suuntaisen testin tapauksessa sekä otoskoon n eräillä arvoilla, saadaan seuraava taulukko. Kaksisuuntaisessa tapauksessa on kyse kriittisen arvon itseisarvosta.

Taulukko 16: Otoskorrelaatiokertoimen kriittiset arvot eräillä otoskoon n arvoilla

n	5 %		1 %		0,10 %	
	1-suunt.	2-suunt.	1-suunt.	2-suunt.	1-suunt.	2-suunt.
10	0,549	0,632	0,715	0,765	0,847	0,872
20	0,378	0,444	0,516	0,561	0,648	0,679
30	0,306	0,361	0,423	0,463	0,541	0,570
50	0,235	0,279	0,328	0,361	0,427	0,451
100	0,165	0,197	0,232	0,256	0,305	0,324
150	0,135	0,160	0,190	0,210	0,250	0,266
200	0,117	0,139	0,164	0,182	0,217	0,231
500	0,074	0,088	0,104	0,115	0,138	0,147
1000	0,052	0,062	0,074	0,081	0,098	0,104
2000	0,037	0,044	0,052	0,058	0,069	0,074

Taulukosta huomataan, että kriittiset arvot pienenevät, kun otoskoko kasvaa. Yksisuuntaisen testin kriittiset arvot ovat pienempiä kuin kaksisuuntaisen testin arvot. Liitteessä 7 esitetään laajempi korrelaatiokertoimien kriittisten arvojen taulukko.

Tilastollisen testin todennäköisyyyslaskelmat pätevät vain, kun kyseessä on *yksi* tilastollinen testaus (nollahypoteesi ja vaihtoehtoinen hypoteesi). Jos testejä on samanaikaisesti useita, kuten esimerkiksi korrelaatiomatriisin tapauksessa, ei valittu riskitodennäköisyys enää pidä paikkansa. Tätä seikkaa voidaan havainnollistaa generoimalla 100 havaintoa kymmenestä riippumattomasta normaali-jakaumaa noudattavasta muuttujasta ja laskemalla tästä aineistosta 10×10 korrelaatiomatriisi. Saatu matriisi on seuraavassa taulukossa.

Taulukko 17: Kymmenen riippumattoman muuttujan otoskorrelaatiomatriisi

	x1	x2	x3	x4	x5	x6	x7	x8	x9	x10
x1	1,000	-0,127	-0,055	0,114	-0,088	-0,046	-0,040	-0,115	0,117	0,098
x2	-0,127	1,000	-0,019	0,197	-0,003	0,130	0,072	0,066	-0,144	0,003
x3	-0,055	-0,019	1,000	0,029	0,010	0,027	0,067	0,060	-0,182	-0,151
x4	0,114	0,197	0,029	1,000	0,007	0,097	-0,214	0,126	0,038	0,202
x5	-0,088	-0,003	0,010	0,007	1,000	0,059	0,074	0,064	-0,003	0,050
x6	-0,046	0,130	0,027	0,097	0,059	1,000	-0,250	0,047	0,020	-0,124
x7	-0,040	0,072	0,067	-0,214	0,074	-0,250	1,000	0,251	0,058	-0,099
x8	-0,115	0,066	0,060	0,126	0,064	0,047	0,251	1,000	0,006	0,054
x9	0,117	-0,144	-0,182	0,038	-0,003	0,020	0,058	0,006	1,000	0,051
x10	0,098	0,003	-0,151	0,202	0,050	-0,124	-0,099	0,054	0,051	1,000

Taulukko 16:n mukaan korrelaatiokertoimen kriittinen arvo 5 %:n riskitasolla (2-suunt. testi) on 0,197. Itseisarvoltaan tätä suuremmat korrelaatiot on merkitty korrelaatiomatriisiin lihavoituina. Kun otetaan korrelaatiomatriisin symmetrisyys huomioon, niin siinä on 45 korrelaatiota, joista neljä poikkeaa tilastollisesti melkein merkitsevästi nolasta. Todellisuudessa kaikki muuttujien korrelaatiot ovat perusjoukon tasolla tasan nollia.

8.6.3 Korrelaatiokerroin nollahypoteesin mukaan eri suuri kuin nolla

Harvoin tulee vastaan tilanteita, jossa nollahypoteesin mukainen korrelaatiokerroin on nolasta poikkeava. Tähän liittyvä tekniikka on kuitenkin hyvä ottaa tässä esille, koska sitä tarvitaan korrelaatiokertoimien yhtä suuruuden testauksessa ja korrelaatiokertoimen luottamusvälin määrittämisessä.

Merkitään nollahypoteesin mukaista korrelaatiokertoimen arvoa ρ_0 :lla. Edellisen luvun testisuureta tai samantapaista suureta ei nyt voida käyttää, koska tällaisen muuttujan otantajakauma on vino eikä siten noudata t-jakaumaa. Korrelaatiokertoimeen nyt sovellettava nk. Fisherin transformaatiota, jolla saadaan likimäärin normaalijakaumaa noudattava suure. Määritellään r^* -muuttuja seuraavalla Fisherin transformaatiolla

$$r^* = (0,5) \log_e \left| \frac{1+r}{1-r} \right|,$$

jossa \log_e on luonnollinen logaritmi (ln). Jos nollahypoteesi on tosi, noudattaa r^* -muuttuja normaalijakaumaa odotusarvon ja varianssin ollessa vastaavasti

$$E(r^*) = \rho_0^* = 0,5 \ln \left(\frac{1+\rho_0}{1-\rho_0} \right) \quad \text{ja} \quad D^2(r^*) = \sigma^{*2} = \frac{1}{n-3}.$$

Testisuure on

$$Z = \frac{r^* - \rho_0^*}{\sigma^*} = (r^* - \rho_0^*) \sqrt{n - 3} ,$$

joka noudattaa likimäärin normaalijakaumaa $N(0,1)$, jos nollahypoteesi on tosi.

Huom. Excelissä on **FISHER** -niminen funktio, jolla saadaan Fisherin transformaatio suoritetuksi.

- **ESIM:** Lukuisat tutkimukset ovat osoittaneet, että kahden muuttujan välinen korrelaatio on 0,5. Eräästä viidenkymmenen kahden havainnon otoksesta saadaan otoskorrelaatiokertoimen arvoksi 0,3. Onko poikkeama tilastollisesti merkitsevä?

Vaihtoehtoinen hypoteesi on yksisuuntainen eli $H_1: \rho < 0,5$.

Tehdään ensin Fisher-transformaatio molemmille korrelaatiokertoimille:

$$r^* = 0,5 \ln \left(\frac{1 + 0,3}{1 - 0,3} \right) = 0,310 \quad \text{ja} \quad \rho_0^* = 0,5 \ln \left(\frac{1 + 0,5}{1 - 0,5} \right) = 0,549 .$$

Testisuureen arvoksi saadaan

$$Z = (0,310 - 0,549) \sqrt{52 - 3} = -1,673 .$$

Koska yksisuuntaisen testin kriittinen arvo 5 %:n riskitasolla on $-1,645$, hylätään nollahypoteesi viiden prosentin riskitasolla.

8.6.4 Kahden riippumattoman korrelaatiokertoimen yhtäsuuruuden testaus

Olko ρ_1 ja ρ_2 muuttujien X ja Y korrelaatiot kahdessa perusjoukossa. Halutaan testata nollahypoteesi, että nämä korrelaatiot ovat yhtä suuret eli

$$H_0: \rho_1 = \rho_2$$

Tämän nollahypoteesin testaamiseksi poimitaan molemmista perusjoukoista riippumattomat otokset suuruudeltaan n_1 ja n_2 yksikköä. Olko otoskorrelaatiot r_1 ja r_2 . Testisuureta varten tehdään Fisherin transformaatiot molemmille korrelaatiokertoimille

$$r_1^* = 0,5 \ln \left(\frac{1 + r_1}{1 - r_1} \right) \quad \text{ja} \quad r_2^* = 0,5 \ln \left(\frac{1 + r_2}{1 - r_2} \right) .$$

Testisuure on nyt

$$Z = \frac{r_1^* - r_2^*}{\sqrt{\frac{1}{n_1 - 3} + \frac{1}{n_2 - 3}}},$$

joka noudattaa standardoitua normaalijakaumaa, jos nollahypoteesi on tosi.

- **ESIM: Äidin pituuden ja lapsen syntymäpituuden korrelaatiot tyttöjen ja poikien otoksissa olivat vastaavasti 0,403 (n=55) ja 0,148 (n=65). Onko korrelaatiokertoimien ero tilastollisesti merkitsevä?**

Fisher-transformoidaan tyttöjen ja poikien otosten korrelaatiokertoimet

$$r_1^* = 0,5 \ln \left(\frac{1 + 0,403}{1 - 0,403} \right) = 0,427 \quad \text{ja} \quad r_2^* = 0,5 \ln \left(\frac{1 + 0,148}{1 - 0,148} \right) = 0,149$$

ja lasketaan korrelaatiokertoimen erotuksen keskivirhe, joka on 0,188. Asetetaan arvot kaavaan ja saadaan Z -testisuureen arvo

$$Z = \frac{0,427 - 0,149}{0,188} = 1,479,$$

joka ei ole ylitä edes viiden prosentin riskitasolla kaksisuuntaisen testin kriittistä arvoa **1,96**. Näin ollen tyttöjen ja poikien korrelaatiokertoimien ero ei ole tilastollisesti merkitsevä.

8.6.5 χ^2 -riippumattomuustesti

χ^2 -riippumattomuustestillä testataan, onko kahden muuttujan välillä tilastollista riippuvuutta. Testi perustuu muuttujien ristiintaulukkoon eli frekvensseihin, joiden muuttujat ovat mitta-asteikoiltaan joko laadullisia tai luokiteltuja kvantitatiivisia muuttujia (esim. ikäluokka).

Riippuvuuslukuna laadullisten muuttujien tapauksessa käytetään yleisimmin joko kontingenssikerrointa C tai Cramerin phi-kerrointa ϕ_c , jotka molemmat perustuvat χ^2 -testisuureeseen. Näitä riippuvuuslukuja käsitellään alla luvussa 8.6.5.1.

Tarkastellaan yleistä ristiintaulukkoa, jossa on g riviä ja h saraketta. Muuttujan X luokkaan x_i ja muuttujan Y luokkaan y_j kuuluvien tilastoyksiköiden lukumäärä eli frekvenssi olkoon f_{ij} . Frekvenssit voidaan esittää seuraavassa taulukkomuodossa.

Taulukko 18: Ristiintaulukko

		X				$yht.$
		x_1	x_2	\dots	x_h	
Y	y_1	f_{11}	f_{12}	\dots	f_{1h}	f_{1*}
	y_2	f_{21}	f_{22}	\dots	f_{2h}	f_{2*}
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
	y_g	f_{g1}	f_{g2}	\dots	f_{gh}	f_{g*}
	$yht.$	f_{*1}	f_{*2}	\dots	f_{*h}	n

Taulukkoon on merkitty myös *reuna- eli marginaalifrekvenssit*, jotka ovat frekvenssien summia rivi- ja sarakesuunnassa.

Havaittuja frekvenssejä f_{ij} verrataan *odotettuihin frekvensseihin* e_{ij} , jotka ovat odotettuja siinä tilanteessa, että riippuvuutta muuttujien X ja Y välillä ei esiinny. Odotetuille frekvensseille saadaan lauseke, kun oletetaan kunkin sarakkeen odotettujen frekvenssien noudattavan prosentuaalista reunajakaumaa. Odotetun frekvenssin lauseke on

$$e_{ij} = \frac{f_{i*} f_{*j}}{n}$$

Seuraavassa taulukossa on sukupuolen ja viikossa alkoholiin käytettyjen markkojen ristiintaulukko ja odotetut frekvenssit (opiskelija-aineisto). Esimerkiksi e_{11} :n arvoksi saadaan $750 \times 1203 / 3963 = 227,7$.

Taulukko 19: Sukupuolen ja alkoholimarkkojen ristiintaulukko ja odotetut frekvenssit

Havaitut frekvenssit				Odotetut frekvenssit			
	<i>Mies</i>	<i>Nainen</i>	<i>Yht.</i>		<i>Mies</i>	<i>Nainen</i>	<i>Yht.</i>
0	134	616	750	0	227,7	522,3	750
1-20	141	617	758	1-20	230,1	527,9	758
21-50	218	750	968	21-50	293,8	674,2	968
51-100	320	502	822	51-100	249,5	572,5	822
101-200	248	216	464	101-200	140,9	323,1	464
201-	142	59	201	201-	61	140	201
Yht.	1203	2760	3963	Yht.	1203	2760	3963

Taulukosta huomataan, että odotettujen frekvenssien reunafrekvenssit ovat samat kuin alkuperäisen taulukon.

Mitä enemmän havaitut frekvenssit poikkeavat odotetuista frekvensseistä, sitä voimakkaampaa on muuttujien välinen riippuvuus. Esimerkiksi miespuolisia absolutisteja aineistossa on 134, kun heidän lukumääränsä tulisi riippumattomuuden vallitessa olla 228. Vastaavasti naispuolisia suurkäyttäjiä on 59, kun odotettu frekvenssi on 140. Nämä suurehkot erot havaittujen ja odotettujen frekvenssien välillä viittaavat siihen, että sukupuoli ja alkoholiin käytettyjen markkojen välillä on selvä tilastollinen riippuvuus.

Riippuvuuden kokonaismäärää mitataan χ^2 -lausekkeella

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(f_{ij} - e_{ij})^2}{e_{ij}} .$$

Edellä olleen esimerkin tapauksessa saadaan χ^2 :n arvoksi

$$\chi^2 = \frac{(134 - 227,7)^2}{227,7} + \frac{(616 - 522,3)^2}{522,3} + \dots + \frac{(59 - 140)^2}{140} = 432,9$$

On selvää, että mitä suurempi χ^2 :n arvo on sitä voimakkaampi on kahden muuttujan välinen riippuvuus. Kuinka suuri pitää χ^2 :n arvo olla, jotta riippuvuus olisi tilastollisesti merkitsevä? Tähän antaa vastauksen χ^2 -riippumattomuustesti.

Nollahypoteesi: Perusjoukossa muuttujien välillä ei ole riippuvuutta.

Vaihtoehtoinen hypoteesi: Perusjoukossa muuttujien välillä on riippuvuutta. On huomattava, että vaihtoehtoinen hypoteesi on yksisuuntainen.

Oletukset:

- a) Otos on poimittu perusjoukosta satunnaisesti.
- b) Korkeintaan kaksikymmentä prosenttia odotetuista frekvensseistä on viittä pienempiä.

Testisuure: on edellä määritelty χ^2 :n lauseke. Jos nollahypoteesi on tosi, noudattaa tämä lauseke likimäärin χ^2 -jakaumaa vapausasteilla $df = (\mathbf{g}-1) \times (\mathbf{h}-1)$, jossa \mathbf{g} ja \mathbf{h} ovat vastaavasti ristiintaulukon rivien ja sarakkeiden lukumäärät.

Kriittinen alue: Kriittinen alue erotetaan χ^2 -jakauman oikeasta hännästä. Kriittinen arvo voidaan määrätä Excelin CHIINV-funktiolla tai katsoa liitteestä 5 vapausasteilla $df = (\mathbf{g}-1) \times (\mathbf{h}-1)$.

Edellä olleessa esimerkissä vapausasteet ovat $df = (6-1)(2-1) = 5$. Käyttäen Exceliä saadaan kriittiseksi arvoksi 0,1 %:n riskitodennäköisyydellä

$$=CHINV(0,001; 5) = 20,51.$$

Havaittu χ^2 arvo **432,9** kuuluu selvästi kriittiselle alueelle, joten sukupuolen ja alkoholiin käytettyjen markkojen välillä vallitsee tilastollisesti erittäin merkitsevä riippuvuus.

8.6.5.1 Ristiintaulukon riippuvuusluvut.

χ^2 :n perustuvat riippuvuusluvut ovat *kontingenssikerroin C* ja *Cramerin phi ϕ_c* . Näiden kaavat ovat seuraavat:

$$C = \sqrt{\frac{\chi^2}{n + \chi^2}} \quad \text{ja} \quad \phi_c = \sqrt{\frac{\chi^2}{n(k-1)}} ,$$

jossa k on pienempi ristiintaulukon rivien tai sarakkeiden lukumääristä.

Kontingenssikertoimen huono puoli on, että sen suurin teoreettinen arvo ei ole yksi, vaan luku, joka riippuu ristiintaulukon rivien ja sarakkeiden lukumääristä seuraavan kaavan mukaan:

$$C_{\max} = \sqrt{\frac{k-1}{k}} , \text{ jossa } k = \min(g, h).$$

Seuraavassa taulukossa on kontingenssikertoimen maksimiarvoja eräillä k :n arvoilla.

Taulukko 20: Kontingenssikertoimen maksimiarvoja ($k =$ pienempi rivien ja sarakkeiden lukumääristä)

k	2	3	4	5	6
C_{\max}	0,71	0,82	0,87	0,89	0,91

Edellä olleessa esimerkissä (sukupuolen ja alkoholimarkat) on χ^2 :n arvo 432,94, joten kontingenssikerroin on

$$C = \sqrt{\frac{432,94}{3963 + 432,94}} = 0,314.$$

Cramerin phiin ϕ_c vaihteluväli on nolosta ykköseen ja se on siksi suositeltavampi riippuvuusluku ristiintaulukon tapauksessa. Esimerkin tapauksessa saadaan Cramerin phiin arvoksi

$$\phi_c = \sqrt{\frac{432,94}{(2-1) 3963}} = 0,331.$$

8.7 χ^2 -yhteensopivuustesti

χ^2 -yhteensopivuustestillä testataan, noudattaako muuttuja tiettyä nollahypoteesin määräämää jakaumaa. Jos kyseessä on kvalitatiivinen muuttuja, ovat muuttujan luokat valmiina. Kvantitatiivisen muuttujan tapauksessa on luokat määrittävä luokkarajojen avulla. χ^2 -yhteensopivuustesti perustuu otoksesta havaittujen frekvenssien ja nollahypoteesista johdettujen teoreettisten (odotettujen) frekvenssien vertailuun. Olkoot jakauman luokat $\mathbf{E}_1, \dots, \mathbf{E}_k$. Nollahypoteesin mukaan luokkien esiintymistodennäköisyydet ovat tietyt arvot θ_i ($i=1, \dots, k$), joiden summa on yksi.

Merkitään nollahypoteesin mukaisia frekvenssejä \mathbf{n} yksikön otoksessa \mathbf{e}_i llä, jolloin $\mathbf{e}_i = \mathbf{n}\theta_i$. Otoksesta lasketut frekvenssit ovat \mathbf{f}_i ($i=1, \dots, k$). Testisuure on

$$\chi^2 = \sum_{i=1}^k \frac{(\mathbf{f}_i - \mathbf{e}_i)^2}{\mathbf{e}_i},$$

noudattaa χ^2 -jakaumaa vapausasteilla $\mathbf{k} - 1$, jos nollahypoteesi on tosi. Jos testisuure ylittää kriittisen arvon, nollahypoteesi hylätään.

- **ESIM:** Olkoon muuttujana syntymäkuukausi ja nollahypoteesina se, että jokaisen kuukauden todennäköisyys on sama eli $\theta_i = 1/12$. Sadan kahdenkymmenen lapsen satunnaisotoksesta lasketut frekvenssit \mathbf{f}_i ovat seuraavassa taulukossa samoin kuin odotetut frekvenssit \mathbf{e}_i .

Taulukko 21. Havaitut ja odotetut frekvenssit

kk	1	2	3	4	5	6	7	8	9	10	11	12	Summa
\mathbf{f}_i	7	11	10	10	10	15	11	7	13	7	11	8	120
\mathbf{e}_i	10	10	10	10	10	10	10	10	10	10	10	10	120

χ^2 -testisuureen arvoksi saadaan

$$\chi^2 = \frac{(7-10)^2}{10} + \frac{(11-10)^2}{10} + \dots + \frac{(8-10)^2}{10} = 6,8.$$

Excelissä määräten χ^2 :n kriittinen arvo viiden prosentin riskitasolla ja vapausasteilla 11 on

$$=CHIINV(0,05; 11) = 19,68 ,$$

joten nollahypoteesia ei hylätä. Siten tämän aineiston perusteella ei voida hylätä nollahypoteesia, jonka mukaan syntymäkuukausien todennäköisyydet ovat yhtä suuret. Suuremmalla aineistolla päädyttäisiin varmaan eri johtopäätökseen.

- **ESIM: Presidentinvaaleissa on kolme tasavahvaa ehdokasta ja lisäksi neljä muuta ehdokasta, joiden kannatus yhteensä jää alle kahdenkymmenen prosentin. Kahden tuhannen hengen mielipidetutkimus antaa kolmelle kärkiehdokkaalle seuraavat kannatusprosentit: 30 %, 27 % ja 26 %. Onko näiden prosenttien ero tilastollisesti merkitsevä?**

Lasketaan kolmea kärkiehdokasta kannattaneiden lukumäärät otoksessa. Ne ovat vastaavasti 600, 540 ja 520. Nollahypoteesin mukaan kannatusprosentit perusjoukossa ovat yhtä suuret, jolloin kolmella odotetulla frekvenssillä on sama arvo eli $(600+540+520)/3=553,3$. χ^2 -testisuureen arvoksi saadaan

$$\chi^2 = \frac{(600 - 553,3)^2}{553,3} + \frac{(540 - 553,3)^2}{553,3} + \frac{(520 - 553,3)^2}{553,3} = 6,27.$$

Excelin **CHIDIST**-funktioilla saadaan todennäköisyys, että nollahypoteesin ollessa tosi, khiin neliön arvo on vähintään 6,27. Tämä todennäköisyys on

$$=CHIDIST(6,27; 2)=0,0436 \text{ eli } \underline{4,36\%}.$$

Siten voidaan päätellä, että kolmen kärkiehdokkaan kannatusprosenttien ero on tilastollisesti melkein merkitsevä.

8.8 Kahden varianssin yhtäsuuruuden testi

Varianssi on (toista astetta oleva) muuttujan vaihtelua mittaava tunnusluku. Yleisimmin käytetty hajontalukuhan on keskihajonta, joka on varianssin neliöjuuri ja jonka mittayksikkö on sama kuin muuttujan mittayksikkö.

Jos halutaan testata, vaihtelevatko muuttujan arvot kahdessa tilastoyksiköiden joukossa tilastollisesti merkitsevästi, käytetään testaamisessa variansseja – ei

keskihajontoja. Näin sen vuoksi koska nimenomaan varianssien osamäärän otantajakauma tunnetaan.

Kahden riippumattoman otoksen keskiarvotestin eräänä oletuksena on perusjoukkojen varianssien yhtäsuuruus. Tämän oletuksen paikkansapitävyys voidaan testata varianssien yhtäsuuruuden testillä.

Toisena esimerkkinä varianssien yhtä suuruuden testaamisesta olkoon tilanne, jossa tutkitaan, miten lisääntynyt kilpailu oppilaiden välillä vaikuttaa oppimistulosten vaihteluun. Tutkimushypoteesi on, että kilpailu lisää vaihtelua. Tämän todentamiseksi tehdään tutkimus, jossa toista oppilasryhmää opetetaan tavallisella opetusmenetelmällä ja toista ryhmää kilpailua lisäävällä opetusmenetelmällä. Tällöin nollahypoteesiksi asetetaan, että oppimistulosta mittaavan muuttujan varianssit ovat samat molemmissa ryhmissä.

Hypoteesit: Olkoon muuttujan X jakaumat kahdessa perusjoukossa normaaliset $N(\mu_1, \sigma_1^2)$ ja $N(\mu_2, \sigma_2^2)$. Nollahypoteesin mukaan varianssit ja siten myös keskihajonnat ovat yhtä suuret.

$$\mathbf{H}_0: \sigma_1 = \sigma_2.$$

Vaihtoehtoinen hypoteesi voi olla yksi- tai kaksisuuntainen.

$$\mathbf{H}_1: \sigma_1 > \sigma_2 \quad \text{tai} \quad \mathbf{H}_1: \sigma_1 \neq \sigma_2$$

Huom. Hypoteesit eivät liity mitenkään keskiarvoihin, joten perusjoukkojen keskiarvot voivat olla joko yhtäsuuret tai erisuuret.

Oletukset:

- a) Otokset poimittu satunnaisesti ja toisistaan riippumattomasti.
- b) Muuttujan jakaumat perusjoukoissa ovat normaaliset.

Testisuure. Testisuure on otosvarienssien suhde

$$F = \frac{s_1^2}{s_2^2},$$

joka noudattaa F -jakaumaa vapausasteilla $df_1 = n_1 - 1$ ja $df_2 = n_2 - 1$, jos nollahypoteesi on tosi (ks. 4.7, s. 50).

Huom. Otosvarienssien laskennassa jakajana on käytettävä havaintojen lukumäärää vähennettynä yhdellä, koska täten saadaan perusjoukon varianssien harhattomat estimaatit. Excelissä varianssin laskemiseen käytetään **VAR** –funktioita, ei **VARP** –funktioita, joka on perusjoukon varianssin funktio.

On huomattava, että yksisuuntaisen vaihtohtoisen hypoteesin tapauksessa testisuureen F osoittajassa on nimenomaan sen otoksen varianssi, jonka perusjoukon vastineen vaihtohtoinen hypoteesi asettaa suuremmaksi. Testisuureta *ei* saa määrätä siten, että suurempi otosvarianssi jaetaan pienemmällä otosvarianssilla (ks. myös 4.7.1, s. 51).

Excelissä on kaksi funktiota, joilla voidaan varianssien yhtä suuruus. **FTEST** –funktion argumentteina ovat havaintoarvot ja se antaa kaksisuuntaisen testin riskitason. **FDIST**–funktion argumentteina ovat varianssien suhde ja vapausasteiden lukumäärät. Funktion tulos on todennäköisyys, että nollahypoteesin vallitessa varianssien suhde on suurempi kuin annettu argumentin arvo (ks. lisää F-jakauman funktioista 4.7.2 s. 53).

Seuraavassa taulukossa on havainnollistettu näiden molempien funktioiden käyttöä. Siinä on kaksi viiden havainnon otosta, joiden alapuolelle on laskettu varianssit. **FTEST** –funktio antaa kaksisuuntaisen testin todennäköisyydeksi 0,1084 ja **FDIST** –funktio yksisuuntaisen testin todennäköisyydeksi 0,0542, joka on puolet kaksisuuntaisen testin todennäköisyydestä.

Taulukko 22. Varianssien yhtäsuuruuden testaus

	A	B	C	D	E
1	Otos 1	Otos 2		Excel-funktio	Funktion arvo
2	11	11		=FTEST(A2:A6;B2:B6)	0,1084
3	21	13			
4	18	15			
5	15	12			
6	19	14			
7	Varianssit:				
8	15,2	2,5		=FDIST(A9/B9;4;4)	0,0542

Käsin laskettaessa määrätään ensin varianssien osamäärä ja tämän jälkeen katsotaan ko. osoittaja sekä nimittäjä vapausasteilla (n_1-1 ja n_2-1) liitteen 6 taulukosta, ylittääkö osamäärä, ts. F -testisuure, kriittisen arvon vähintään 5 %:n riskitasolla.

- **ESIM: Onko naisten (n=40) varianssi eräässä asennemittarissa suurempi kuin miesten (n=30)? Naisten $s^2 = 7,56$ ja miesten $s^2 = 4,27$.**

Asetetaan hypoteesit $H_0: \sigma_N = \sigma_M$, $H_1: \sigma_N > \sigma_M$ ja lasketaan testisuureen arvo $F = \frac{s_1^2}{s_2^2} = \frac{7,56}{4,27} = 1,77$. Taulukosta 6 todetaan, että kriittinen

arvo 5%:n riskitasolle on 1,79. Koska testisuureen arvo ei ylitä tätä, todetaan

taan, että naiset eivät ole *merkittävästi* heterogeenisempiä — eli H_0 jää voimaan.

9 Testin voimakkuusfunktio (Power)

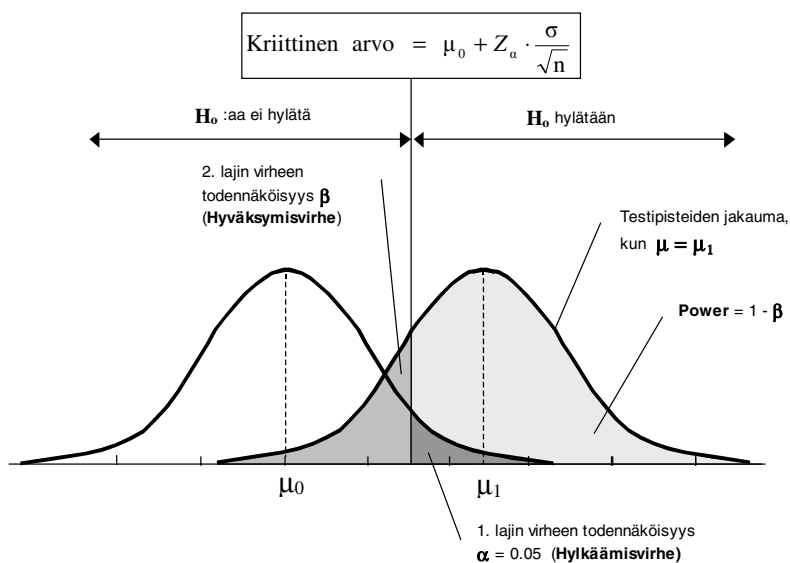
9.1 Hyväksymisvirheen todennäköisyys β

Tavanomaisesti tilastollisen päättelyn käytännössä ollaan kiinnostuneita vain minimoimaan tai kontrolloimaan nollahypoteesin hylkäämisvirheen todennäköisyys eli ns. tilastollinen riskitaso. Kokeita ja muita tilastollisia testitilanteita suunniteltaessa ja analysoitaessa ollaan taipuvaisia ohittamaan se tärkeä tosiasia, että on myös olemassa 2. lajin virhe ja siihen liittyvä todennäköisyys. Koska testin voimakkuus tilastollisena käsitteenä liittyy oleellisesti juuri 2. lajin virheen todennäköisyyteen (vaikkakaan se ei ole riippumaton 1. lajin virheen todennäköisyydestä), on sekin jäänyt tästä syystä vähemmälle huomiolle.

Testin voimakkuuden ymmärtämiseksi on hyvä aluksi tarkastella 2. lajin virheen eli hyväksymisvirheen luonnetta. Hyväksymisvirhe tapahtuu, kun vaihtoehtoinen hypoteesi H_1 on perusjoukon tasolla tosi, mutta otoksen perusteella nollahypoteesia H_0 ei päädytä hylkäämään. Kuva 40:n vasemmanpuoleinen jakauma edustaa tilannetta, jossa H_0 on tosi; oikean puoleinen jakauma puolestaan kuvaa tilannetta, jossa H_1 on tosi. Toisen lajin virheen todennäköisyys, β , muodostuu siitä H_1 :n jakauman osasta, joka *ei ylitä* 1. lajin virheen todennäköisyyden pohjalta määrättyä kriittistä arvoa (päätöskriteerin arvoa).

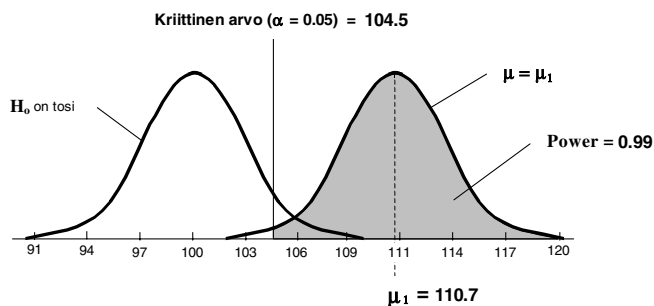
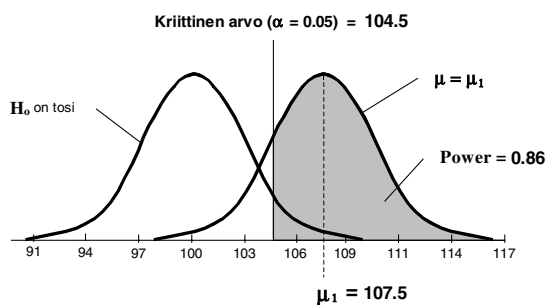
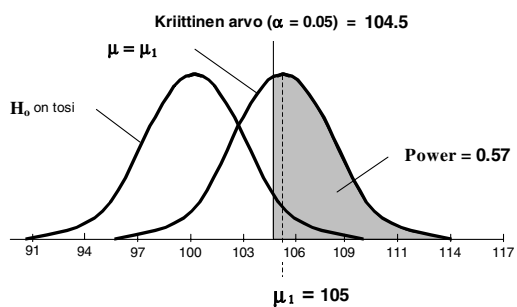
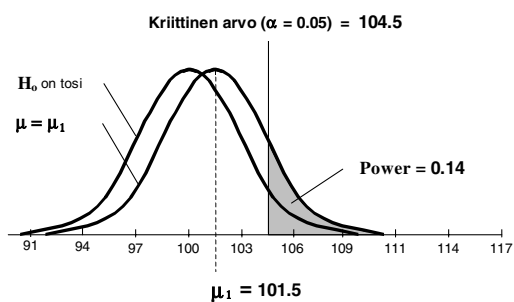
Kuten Kuva 40:sta voi havaita, testin voimakkuusfunktion arvo muodostuu siitä H_1 :n jakauman osasta, joka ylittää kriittisen arvon. Näin ollen testin voimakkuuden arvo, tietyssä pisteessä μ , on $1-\beta$.

Kuva 40: Ensimmäisen lajin virheen (α), toisen lajin virheen (β) ja testin voimakkuuden suhde toisiinsa



Tarkastellaan seuraavaksi Kuva 41:n esittämien eri tilanteiden havainnollistamana, kuinka testin voimakkuusfunktion arvo vaihtelee μ :n eri arvojen mukaan α :n (=0,05), otoskoon (=30) ja σ^2 :n (=15) pysyessä vakiona.

Kuva 41: Testin voimakkuus (Power) neljän eri μ :n arvon tapauksessa



Tilastollisen testauksen yhteydessä ei voida samalla tavalla määrätä kiinteää arvoa hyväksymisvirheen todennäköisyydelle (β) kuin voidaan tehdä hylkäämisvirheen todennäköisyydelle (α). Tämä johtuu siitä, että β :n arvo riippuu tuntemattomasta perusjoukon parametrusta (esim. μ), joka joudutaan estimoimaan otoksesta. Näin ollen empiirisessä tutkimuksessa testin voimakkuuden arvoa ei

voida määrätä etukäteen, vaan sille lasketaan estimaatti havaitun otostunnusluvun pohjalta.

9.2 Testin voimakkuuteen vaikuttavia tekijöitä

Testin voimakkuuden arvo (eli todeksi oletetun H_1 :n vallitsevuuden todennäköisyys), riippuu useista tekijöistä. Näitä tekijöitä, joiden funktiona testin voimakkuus voidaan kuvata, ovat:

- α eli 1. lajin virheen todennäköisyys
- μ :n arvo eli populaation keskiarvo
- n eli otoskoko
- σ^2 :n arvo eli populaation varianssi
- käytetty testityyppi.

Tarkastellaan seuraavaksi hieman spesifimmin näitä testin voimakkuuteen vaikuttavia tekijöitä.

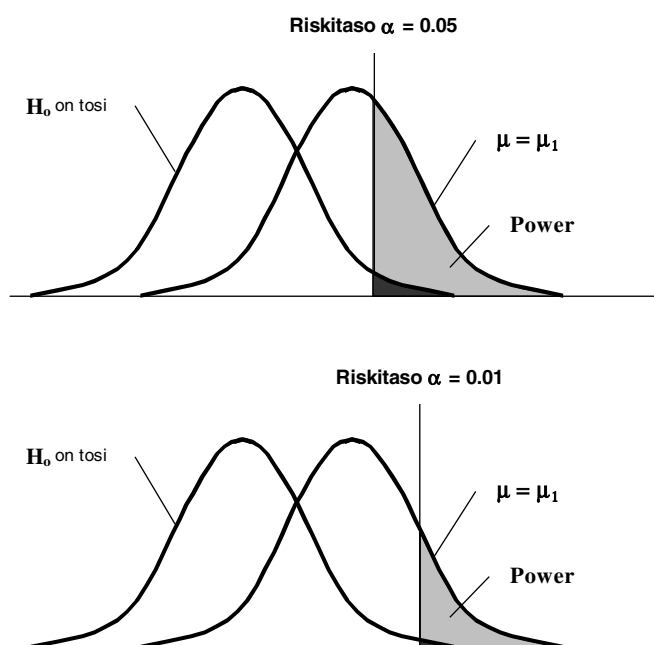
9.2.1 α :n vaikutus testin voimakkuuteen

Pienennettäessä alfan eli 1. lajin virheen todennäköisyyden arvoa, pienenee testin voimakkuuden arvokin. Kuva 42 havainnollistaa tätä. Ylemmässä kuvassa, jossa tilastollisen päätöksen teon riskitasoksi on asetettu 5 %, testin voimakkuus on huomattavasti suurempi kuin alemmaman kuvan tilanteessa, jossa riskitasoksi on asetettu 1 %.

Tämä tosiasia muodostaa testin voimakkuusfunktion käytön kannalta sen ongelman, että ei ole mitään selvää kriteeriä sille, mihin α :n arvoon sen laskeminen pitäisi perustua. Eräs sääntö voisi olla se, että lasketaan testin voimakkuus sille riskitasolle, jota käytetään H_0 :n hyväksymisen / hylkäämisen perustana.

Testin voimakkuutta ei ole mielekäästä laskea tarkalle (empiirisestä aineistosta määrätylle) tilastolliselle riskitasolle, koska tällöin testin voimakkuudeksi tulee aina 0,5.

Kuva 42: Tilastollisen riskitason α yhteys testin voimakkuuteen



Kuva 42 tuo esiin konkreettisesti sen tosiasian, että haluttaessa lisää varmuutta H_0 :n hylkäämisen perustaksi, joudutaan tyytymään pienempään testivoimakkuuteen. Toisaalta mikäli vaikutuksen suuruus on riittävän suurta (jolloin Power $\rightarrow 1$), ei testin voimakkuuden arvo välttämättä laske käytännöllisesti katsottuna merkittävässä määrin, vaikka se laskettaisiin huomattavasti pienemmällekin riskitasolle. Yleisesti ilmaistuna α :n ja β :n suhde on seuraava:

Kun α (eli 1. lajin virheen todennäköisyys) pienentyy, niin β (eli 2. lajin virheen todennäköisyys) kasvaa — ja täten testin voimakkuus $1-\beta$ pienentyy.

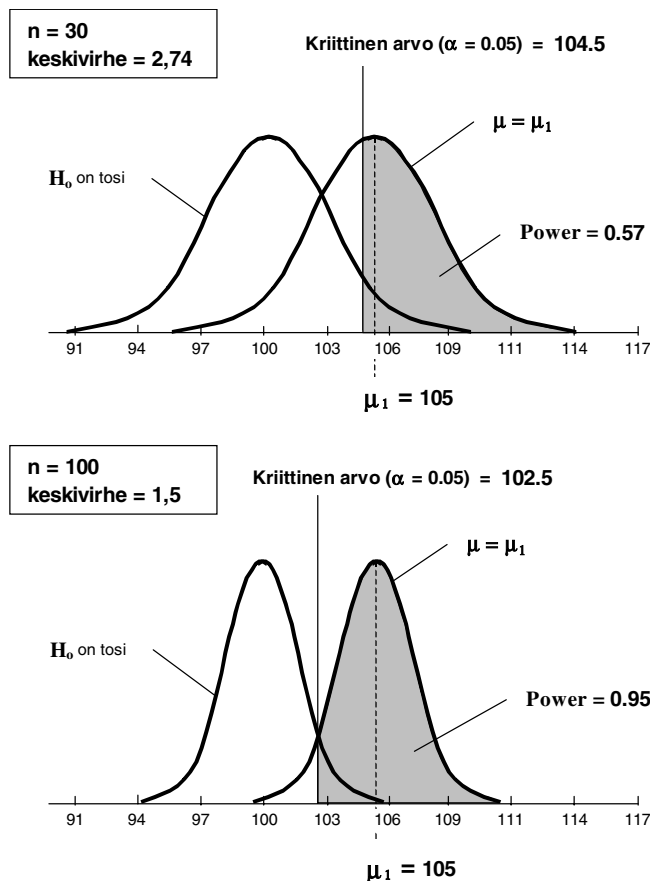
9.2.2 μ :n arvo

Testin voimakkuusfunktiolle on voimassa lauseke $P(\mu) = 1 - \beta(\mu)$, joka kuvaa oikean johtopäätöksen todennäköisyyden vaihtelua μ :n funktiona. Näin ollen, kuten Kuva 41:stä voi todeta, H_1 :n jakauman sijainti suhteessa H_0 :n jakaumaan riippuu yksinomaan siitä, mikä μ :n arvo on. Erisuuria μ :n arvoja vastaavat erisuuret testivoimakkuusfunktion arvot. Kuva 45 ja Kuva 46 havainnollistavat sekä yksi- että kaksisuuntaisen testauksen tilanteissa kuinka testin voimakkuus vaihtelee μ :n funktiona.

9.2.3 Otoskoko (n)

Otoskoon suurentaminen pienentää estimaattorin (esim. \bar{X} :n) keskivirhettä ja siten mahdollistaa testin voimakkuuden kasvamisen. Kuva 44 havainnollistaa tätä. Keskivirheen pienentyessä pienentyy myös kriittinen arvo, mikä kasvattaa täten sitä H_1 :n jakaumaan sisältyvää aluetta, joka sijaitsee tämän kriittisen arvon yläpuolella. Tämä alue siis muodostaa testin voimakkuuden eli oikean johtopäätöksen todennäköisyyden.

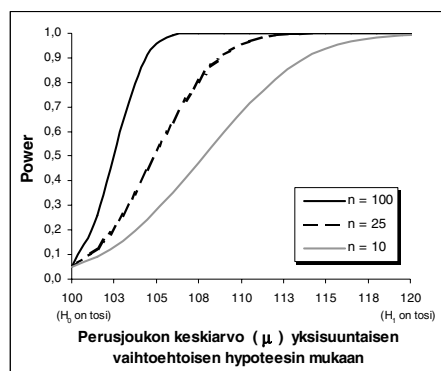
Kuva 44: Otoskoon yhteys testin voimakkuuteen



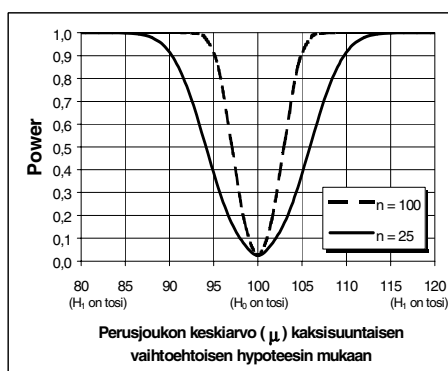
Toisaalta, kun μ :n keskivirhe pienentyy, pienentyy myös H_0 :n ja H_1 :n jakaumien limittyvä pinta-ala, jolloin suurempi osa H_1 :n jakaumasta siirtyy kriittisen arvon yläpuolelle. Näin ollen — muiden testin voimakkuuteen vaikuttavien tekijöiden säilyessä muuttumattomina — otoskoon suurentaminen 25:stä 100:aan nostaa testin voimakkuuden 0,57:stä 0,95:ään, kuten Kuva 44 esittää.

Kuva 45 esittää yksisuuntaisen H_1 :n tapauksessa, kuinka testin voimakkuus on toisaalta riippuvaista μ :n arvosta ja toisaalta otoskoosta. Voidaan siis todeta, kuinka samalla μ :n arvolla päädytään huomattavasti suurempaan testin voimakkuuteen, mikäli kasvatetaan reilusti otoskokoa.

Kuva 45: Testin voimakkuus ja μ :n arvo 1-suuntaisessa testissä



Kuva 46: Testin voimakkuus ja μ :n arvo 2-suuntaisessa testissä



Kuva 46 puolestaan tuo esiin vastaavan tilanteen kuin Kuva 45, mutta nyt kaksisuuntaisen H_1 :n tapauksessa. Ero yksi- ja kaksisuuntaisen testivoimakkuuden määrittämisen välillä on periaatteessa aivan sama kuin yksi- ja kaksisuuntaisen hylkäämisvirheen todennäköisyydenkin määrittämisessäkin: jaetaan α :n arvo kahdella ja määrätään nyt symmetrisesti kriittiset arvot jakauman molemmista päistä, jolloin ne itseisarvoisesti suurenevat. Tällöin, H_0 :aa vastaavaa μ :n arvoa suuremmilla arvoilla (>100 , Kuva 45:ssä ja Kuva 46:ssa) yksisuuntainen H_1 antaa kaksisuuntaista H_1 :stä suuremman voimakkuuden testille. Kun taas H_0 :aa vastaavaa μ :n arvoa pienemmillä μ :n arvoilla yksisuuntainen H_1 antaa hyvin nopeasti nollaa lähestyvän testivoimakkuuden, kuten Kuva 45 osoittaa.

9.2.4 σ^2 :n arvo

Kaikki tekijät, jotka pienentävät tutkittavan muuttujan varianssia, lisäävät testin voimakkuutta. Vaikka käytännössä voi olla melko vaikeaa pienentää tai ylipäättään muuttaa perusjoukon varianssia, on se kuitenkin usein kokeellisissa tutkimusasetelmissä mahdollista. Kokeellisissa tutkimuksissa usein pyritään kontrolloimaan (vakioimaan) joidenkin tekijöiden vaikutus tutkittavaan ilmiöön samalla, kun joidenkin toisten tekijöiden varioinnin vaikutusta tarkastellaan. Tällöin, vaikka tutkittavalla ilmiöllä olisikin tietty varianssi perusjoukon tasolla, on mahdollista pienentää tätä varianssia siten, että eliminoidaan sellaisten tekijöiden vaikutus tutkittavaan ilmiön varianssiin, jotka eivät ole tutkimuskysymyksen kannalta relevantteja tai jotka tuottavat siihen nähden virhevaihtelua.

Esimerkiksi ajokokemuksen määrän vaikutus voitaisiin haluta eliminoida tutkittaessa, miten sääolosuhteet ja kelloaika vaikuttavat autoilijoiden liikennetarkkaavuuteen. Tällöin se osa liikennetarkkaavuusmittarin varianssista, jonka ajokokemus selittää, eliminoidaisiin tarkastelusta. Näin saataisiin mielekkäällä tavalla pienennettyä liikennetarkkaavuusmittarin varianssia ja voitaisiin tutkia erityisemmin juuri sääolosuhteiden ja kelloajan vaikutusta ilman, että jäisi merkittävässä määrin epäselväksi, mikä osa liikennetarkkaavuusmittarin varianssia on luettavissa sääolosuhteiden ja kellon ajan vaikuttamaksi ja mikä osa ajokokemuksen.

Yllä kuvatulla menettelyllä on siis yhteytensä myös testin voimakkuuteen: mikäli kyetään kontrolloimalla joidenkin tekijöiden vaikutus pienentämään testimuuttujan varianssia, kasvaa testin voimakkuus vastaavasti.

9.2.5 Käytetty testityyppi

Joskus tilastollista tutkimusta tehtäessä voi herätä kysymys siitä, mitä testityyppejä pitäisi käyttää. Yleisesti ajatellen pitäisi suosia sitä testiä, kumpi antaa samasta datasta laskettuna suuremman testivoimakkuuden. Tarkastelematta asiaa sen seikkaperäisemmin, voidaan yleisesti todeta, että kahdesta testistä voimakkaampi on se, kumpi hyödyntää datan perusteellisemmin. Yleensä ns. parametrisillä testeillä on parempi testivoimakkuus verrattuna non-parametrisiin testeihin. Tämä pätee erityisesti silloin, kun parametristen testien käytön edellytykset, kuten testeihin liittyvät muuttujien jakautumaa koskevat oletukset, ovat voimassa.

9.3 *Testin voimakkuusfunktion arvon laskeminen tietyille vaihtoehdoiselle hypoteesille*

Testivoimakkuuden arvoja voidaan laskea käytännössä erilaisille testeille. Täten testivoimakkuus on laskettavissa eri teoreettisia jakaumia noudattavien testisuureiden kuvaamille testitilanteille. Tässä esityksessä pitäydytään kuitenkin yksinkertaisuuden vuoksi vain yhden otoksen keskiarvotestissä (z-testi).

- **ESIM:** Lasketaan testivoimakkuusfunktion arvo toiselle Kuva 41:ssä esitetylle testitilanteelle, jossa on kysymys juuri yhden otoksen keskiarvotestistä. Testitilanteeseen liittyvät arvot ovat $\mu_0 = 100$; $\sigma = 15$; $n = 30$ ja $\alpha = 0,05$. Ajatellaan lisäksi, että ko. μ_1 :n arvo 105 onkin nyt otoksesta laskettu μ :n estimaatti eli \bar{X} .
-

Ensiksi tarvitaan nollahypoteesin hylkäämiseen liittyvä kriittinen arvo eli sellainen μ :n arvo, joka on riittävän suuri, jotta voimme varmuudella hylätä H_0 :n. Tämä lasketaan seuraavasti:

$$\text{Kriittinen arvo} = \mu_0 + Z_\alpha \cdot \frac{\sigma}{\sqrt{n}} .$$

Kuva 41:n esimerkkiä noudattaen, jossa siis $\mu_0 = 100$ ja $\sigma = 15$ ja otoskoko $n = 30$, saamme 5 %:n riskitasolla kriittiseksi arvoksi seuraavan arvon:

$$\text{Kriittinen arvo} = 100 + 1,64 \cdot \frac{15}{\sqrt{30}} \approx 104,5 .$$

Täten, jos otoskeskiarvo on suurempi kuin 104,5 hylkäämme 95 %:n varmuudella nollahypoteesin. Mikäli otoskeskiarvo on puolestaan pienempi kuin kriittinen arvo, jää nollahypoteesi voimaan. Tällöin on kuitenkin 2. lajin virheen todennäköisyyden suuruinen mahdollisuus, että hyväksymme H_0 :n virheellisesti — olettaessa, että H_1 on tosi perusjoukon tasolla.

Lasketaan seuraavaksi tämä 2. lajin virheen todennäköisyys. Ensin vähennetään kriittisestä arvosta 104,5 μ :tä estimoiva otoskeskiarvo 105 ja jaetaan erotus sitten keskivirheellä — Eli

$$Z_1 = \frac{\text{Kriittinen arvo} - \mu_1}{\frac{\sigma}{\sqrt{n}}} = \frac{104,5 - 105}{\frac{15}{\sqrt{30}}} \approx -0,183 .$$

Tämän jälkeen määrätään normaalijakauman kertymäfunktion arvo Z :n arvolle $-0,183$ liitteen 1 taulukosta ja saadaan näin 2. lajin virheen todennäköisyys

$$F(Z_1) = F(-0,183) \approx \underline{0,429} .$$

Testin voimakkuuden arvo saadaan 2. lajin todennäköisyydestä seuraavasti:

$$\text{Testin voimakkuus } P(\mu) = 1 - \beta(\mu) = 1 - 0,429 \approx \underline{0,571} .$$

Kuten estimoidusta otoskeskiarvosta 105 päädytään toteamaan, voidaan H_0 hylätä (5 %:n riskitasolla). Tuloksen yleistettävyyttä ja varmuutta ajatellen on kuitenkin huomattava, että testin voimakkuuden arvo on vain n. 0,57, mikä täten osoittaa (olettaessa H_1 todeksi eli \bar{X} :n arvo 105 todella perusjoukossa vallitsevaksi), että oikean johtopäätöksen todennäköisyys on koh-

talaisen matala — eli n. 43 %:ssa tapauksista testi antaisi väärän johtopäätöksen (ks. Kuva 41, s. 106).

Näin ollen siis, vaikka H_0 hylätään ja H_1 astuu voimaan, ei H_1 :n vallitsevuudesta voida olla erityisen varmoja. Lisäksi on vielä huomattava, että testin voimakkuuden arvo laskettiin 5 %:n riskitasolle. Pienemmälle riskitasolle — kuten esim. 1 %:n riskitasolle — laskettuna sen arvo olisi vieläkin pienempi (n. 0,30). Tällöin lisäksi H_0 olisi jäänyt ko. tapauksessa voimaan., koska kriittiseksi arvo olisi määräytynyt tällöin 106,4 (ks. Kuva 42, s. 109).

On myös syytä vielä todeta, että koska käytimme tässä esimerkissä otoskeskiarvoa \bar{X} kiinteän μ :n arvon sijasta, on näin laskettu testivoimakkuuden arvokin nyt vain estimaatti, jossa on vaihtelua otoksesta toiseen. Yllä olemme kuitenkin käsitelleet aina tilanteita, jossa tiettyyn H_1 :een liittyvä μ_1 :n arvo on oletettu tunnetuksi.

10 Otoksoon määrääminen

Otoksoon määrääminen voidaan suorittaa kahta eri tekniikkaa ja periaatetta hyväksi käyttäen. Otokskoko voidaan määrätä testivoimakkuusfunktion avulla tai luottamusvälin avulla. Tarkastellaan ensiksi testivoimakkuusfunktion tapausta.

10.1 Otoksoon määrääminen testin voimakkuusfunktiota ($1-\beta$) hyväksi käyttäen

Käytännön elämässä tulee usein vastaan tilanne, jossa haluttaisiin tietää esim., että miten paljon jokin asia vaikuttaa tiettyyn toiseen asiaan tai ilmiöön. Ei siis olla kiinnostuneita vain siitä, vaikuttaako jokin johonkin, vaan halutaan tietää tarkemmin, miten suurta tuo vaikutus on. Tällaista tietoa kaivattaisiin usein päätöksenteon perustaksi. Esimerkiksi kenkiä myyvässä yrityksessä haluttaisiin tietää, miten paljon uudella mainosohjelmalla on vaikutusta heidän kenkiensä viikoittaiseen kokonaisynttiin. Vastaavasti kasvatus- ja koulutustyössä kaivataan usein luotettavaa arviota siitä, miten tehokkaita erilaiset opetus- ja kasvatusmenetelmät ovat tehtävässään.

Myös erityistieteellisessä tutkimuksessa tulee usein vastaan tilanne, jossa ollaan kiinnostuneita ei vain tilastollisesti merkitsevästä vaikutuksesta, vaan myös sen aktuaalisesta suuruudesta ja täten myös vaihtoehtoisen hypoteesin vallitsevuuden todennäköisyydestä.

Kun tiedetään kuinka suuri vaikutus jollain menettelyllä johonkin asiaan haluttaisiin saavuttaa ja kun tiedetään, että otoskoko vaikuttaa testivoimakkuuteen, voidaan otoskoko valita sellaiseksi, että testivoimakkuusfunktion arvo (eli oikean johtopäätöksen todennäköisyys) saadaan asetetuksi halutun suuruiseksi. Toi-

sin sanoen järjestetään tilastollinen testitilanne, jossa testin voimakkuus voidaan asettaa otoskoon valinnan avulla niin suureksi, että mikäli haluttu vaikutuksen suuruus todella on olemassa, on vaihtoehdoisen hypoteesin hyväksymisen todennäköisyys riittävän suuri.

Tarkastellaan otoskoon määrittämistä erikseen keskiarvon ja prosenttiluvun tapauksessa. Yksinkertaisuuden vuoksi käsitellään vain yksisuuntaista vaihtoehdoista hypoteesia.

10.1.1 Keskiarvo

Tarkastellaan perusjoukon keskiarvoon μ liittyvää yhden otoksen testiasetelmaa, kun riskitaso eli ensimmäisen lajin virheen todennäköisyys on α :

$$H_0: \mu = \mu_0 \text{ (= annettu arvo) , } H_1: \mu > \mu_0 \text{ (tai } H_1: \mu < \mu_0 \text{).}$$

Oletetaan lisäksi että keskihajonta σ on tunnettu.

Otoskoko n nyt voidaan määrätä seuraavan ehdon perusteella. Olkoon $\mu_1 (>\mu_0)$ sellainen perusjoukon keskiarvo, jonka kohdalla halutaan, että testi kohtuullisen suurella todennäköisyydellä eli voimakkuusfunktion arvolla $1-\beta$ hylkää nollahypoteesin H_0 , toisin sanoen hyväksyy oikean vaihtoehdoisen hypoteesin H_1 . Kuten yhdeksännessä luvussa todettiin, on β hyväksymisvirheen todennäköisyys. Voimakkuusfunktion ”kohtuullisen suuri” arvo ja sitä kautta myös toisen lajin virheen todennäköisyys voidaan määrätä vapaasti. Pitää kuitenkin muistaa, että mitä suurempi testin voimakkuus on pisteessä μ_1 , sitä suurempi otoskoko saadaan.

Näistä lähtökohdista käsin voidaan otoskoolle johtaa seuraava kaava:

$$n = \frac{(Z_\alpha + Z_\beta)^2 \sigma^2}{(\mu_1 - \mu_0)^2} ,$$

jossa Z_α ja Z_β ovat ensimmäisen ja toisen lajin virheisiin liittyvät standardoidun normaalijakauman kriittiset arvot. Taulukko 23 esittää tavanomaisesti käytetyt normaalijakauman kriittiset arvot yksisuuntaisen vaihtoehdoisen hypoteesin tilanteelle.

Taulukko 23: Eräitä standardoidun normaalijakauman kriittisiä arvoja (yksisuuntainen vaihtoehdoisen hypoteesi)

Todennäköisyys (α tai β)	0,10	0,05	0,01	0,001
Kriittinen arvo	1,28	1,64	2,33	3,09

- **ESIM:** Kunta on valmis hankkimaan matematiikan ATK-ohjelman, jos se parantaa matematiikan oppimistasoa. Sitä varten on tarkoitus järjestää koe, jossa satunnaisesti valitut opiskelijat käyttävät kyseistä ohjelmaa ja sen jälkeen mitataan heidän matemaattinen tietotaso tietyllä vakioasteikolla. Tämän asteikon jakauma normaalilla opetuksella on $N(50;10^2)$. Kuinka suuri otoskoon tulee olla, jos oppimistason viiden yksikön suuruinen parannus halutaan todentuvan 90%:n varmuudella? Tilastollisen testiin hylkäämisvirheeseen, eli että kunta ostaa nolla-vaikutteisen ATK-ohjelman varaudutaan viiden prosentin todennäköisyydellä.

Otoskoon määrääminen suoritetaan sijoittamalla esimerkin ja Taulukko 23:n antamat luvut otoskoon kaavaan:

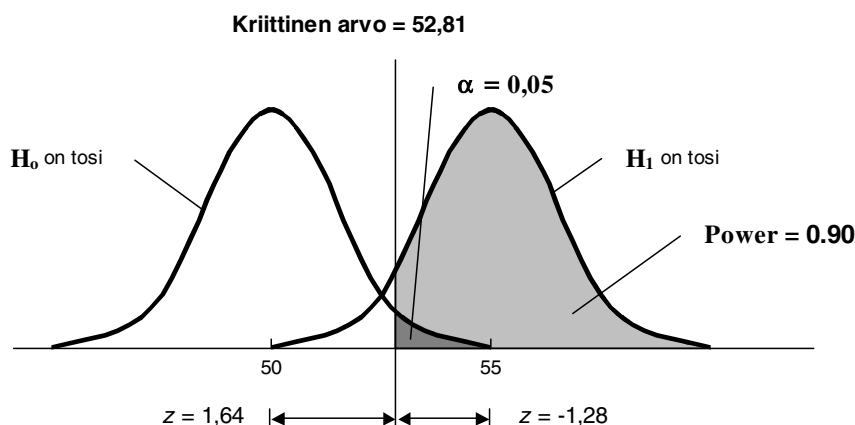
$$n = \frac{(1,64 + 1,28)^2 10^2}{(55 - 50)^2} = 34 \cdot$$

Siten kokeen oppilasmääräksi valitaan 34.

Kuva 47 havainnollistaa vielä yrityksen tarvitseman koeasetelman (ks. myös 9.2.3., s. 110 ja Kuva 44). Kuvassa esitettävä kriittinen arvo, joka siis toimii nollahypoteesin hylkäämisen tai hyväksymisen kriteerinä, saadaan laskettua otoskoon määräämisen jälkeen, kuten luvussa 9.3 on esitetty, seuraavasti:

$$\text{Kriittinen arvo} = \mu_0 + Z_\alpha \cdot \frac{\sigma}{\sqrt{n}} = 50 + 1,64 \cdot \frac{10}{\sqrt{34}} \approx 52,81$$

Kuva 47: Yrityksen tarvitsema koeasetelma tilastollisen päätöksen teon kannalta



Eräs ongelma otoskoon kaavan soveltamisessa on, että se edellyttää muuttujan keskihajonnan σ tuntemista. Tämä voidaan kuitenkin välttää sillä, että ilmaistaan poikkeama $\mu_1 - \mu_0$ keskihajontaan suhteutettuna. Otoskoon kaava voidaan nimittäin kirjoittaa seuraavaan muotoon

$$n = \frac{(Z_\alpha + Z_\beta)^2 \sigma^2}{(\mu_1 - \mu_0)^2} = \frac{(Z_\alpha + Z_\beta)^2}{\left(\frac{\mu_1 - \mu_0}{\sigma}\right)^2}$$

Siten kaavan käytössä riittää, kun standardoidulle poikkeamalle $(\mu_1 - \mu_0) / \sigma$ annetaan haluttu arvo. Edellisessä esimerkissä tämä arvo oli 0,5 eli 50%.

10.1.2 Prosenttiluku

Prosenttiluvun θ tapauksessa yksisuuntaisen testin asetelma riskitasolla α seuraava:

$$H_0: \theta = \theta_0 \text{ (= annettu arvo)}, \quad H_1: \theta > \theta_0 \text{ (tai } H_1: \theta < \theta_0)$$

Halutaan, että testin voimakkuus pisteessä θ_1 on $1 - \beta$. Otoskoon lauseke prosenttiluvun tapauksessa saadaan edellisen alaluvun kaavasta pienellä muunnoksella. Varianssin σ^2 paikalle sijoitetaan dikotomisen muuttujan (jonka arvot ovat 0 ja 100) varianssi $\theta_{01}(100 - \theta_{01})$, jossa θ_{01} on θ_0 :n ja θ_1 :n keskiarvo. Otoskoon kaavaksi saadaan siten

$$n = \frac{(Z_\alpha + Z_\beta)^2 \theta_{01}(100 - \theta_{01})}{(\theta_1 - \theta_0)^2}.$$

- **ESIM:** Lääketehdas väitti erään lääkkeen parantavan 80%:ssa tapauksista. Väitteen testaamiseksi suoritetaan potilaskoe, jossa hylkäämisvirheen todennäköisyydeksi α valitaan 0,05. Määrää otoskoko siten, että testin voimakkuus pisteessä 70% on 0,95

Eli jos lääke todellisuudessa tehoaa vain 70%:ssa tapauksista, niin tilastollinen testi hylkää nollahypoteesin ($\theta=80%$) todennäköisyydellä 0,95. Otoskoon arvoksi saadaan

$$n = \frac{(1,64 + 1,64)^2 75(100 - 75)}{(70 - 80)^2} = 202.$$

10.2 Otoskoon määrääminen luottamusvälin avulla

Yllä tarkastelimme, otoskoon määräämistä testivoimakkuusfunktion avulla määrättyinä. On myös mahdollista määrätä otoskoko populaation parametrin luottamusvälin avulla. Tarkastelemme tässä alaluvussa otoskoon määräämistä keskiarvon, prosenttiluvun ja korrelaatiokertoimen tapauksissa.

10.2.1 Keskiarvo

Haluttaessa tietää kuinka suuri otoskoon on oltava, jotta populaation keskiarvo μ voitaisiin estimoida halutulla varmuudella, voidaan luottamusvälin pituuden d kaava ratkaista otoskoon n suhteen ja saada kaava otoskoon määräämiselle eli

$$d = 2 \cdot Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \quad \Leftrightarrow \quad n = \frac{4(Z_{\alpha/2})^2 \cdot \sigma^2}{d^2}.$$

Tämän kaavan käyttö edellyttää tutkittavan muuttujan X varianssin σ^2 tuntemista. Koska populaation varianssia σ^2 ei käytännössä useinkaan tunneta, voidaan yrittää käyttää joidenkin aikaisempien tutkimusten perustalta saatua arviota. Tämä tuottaa tietysti jonkin verran epävarmuutta otoskoon määräämiseen,

mutta ei kuitenkaan enää vaikuta estimaattorin tarkkuuteen estimoitaessa otoksen perusteella populaation parametriä. Otskokoa määrättäessä ei siis voida käyttää otoksesta laskettavaa varianssia s^2 estimaattina σ^2 :lle, koska otosta ei luonnollisesti ole vielä kyseisessä tutkimuksen vaiheessa poimittu.

- **ESIM:** Terveysviranomaiset haluavat selvittää, mikä on 12-15 – vuotiaiden päihderiippuvaisten nuorten huumausaineisiin käyttämä viikoittainen rahamäärä. Aikaisemmasta tutkimuksesta on tietona, että viikoittaisen markkamäärän keskihajonta on n. 100 mk:aa. Mikä pitäisi olla otoskoon, kun terveysviranomaiset tahtovat estimoida 99 %:n varmuudella viikoittaisen rahamäärän 12,5 mk:n tarkkuudella?

Laskentaa varten tarvitaan seuraavat arvot: luottamusvälin pituus $d=25$, kriittinen arvo 1 %:n riskitasolle eli $Z_{0,01/2} = 2,58$ (liitteen 1 taulukosta) ja varianssi $\sigma^2 = 100^2$, jotka asetetaan kaavaan.

$$n = \frac{4(2,58)^2 \cdot 100^2}{25^2} = \frac{26656}{625} \approx 426.$$

Vaadituksi otoskooksi saadaan 426 (päihderiippuvaista nuorta).

10.2.2 Prosenttiluku

Kuten keskiarvonkin kohdalla, populaation prosenttiluvun θ luottamusvälin pituuden d kaava voidaan ratkaista otoskoon suhteen eli

$$d = 2 \cdot Z_{\alpha/2} \cdot \sqrt{\frac{\theta(100-\theta)}{n}} \Leftrightarrow n = \frac{4(Z_{\alpha/2})^2 \theta(100-\theta)}{d^2}.$$

Näin saatu otoskoon lauseke riippuu estimoitavasta parametristä θ , joka on tuntematon. Kuten keskiarvon tapauksessa (σ^2), on nytkin käytettävä enemmän tai vähemmän karkeaa arvioita. Jos tiedämme θ :n sijaitsevan tietyllä välillä $[\theta_1, \theta_2]$, voimme käyttää θ :n arviona sitä välin arvoa, joka on lähinnä 50:tä (jolloin θ :n varianssiksi $\theta(100-\theta)$ saadaan suurempi arvo). Tällöin otoskoko ei tule ainakaan aliestimoiduksi

- **ESIM:** Erään poliittisen puolueen kannatusprosentin θ tiedetään olevan äänestysikäisten muodostamassa perusjoukossa 30–40 %:n

välillä. Mikä pitäisi olla otokseen, jotta perusjoukon prosenttiluku voitaisiin estimoida 95 %:n todennäköisyydellä puolet tarkemmin?

Laskun suorittamista varten tarvitsemme seuraavat arvot: $d=(40-30)/2=5$; kriittinen arvo 5 %:n riskitasolle eli $Z_{0,05/2}=1,96$ (liitteen 1 taulukosta) ja se θ :n arvo, kumpi on lähempänä 50:tä eli 40. Asetetaan arvot kaavaan ja lasketaan

$$n = \frac{4(1,96)^2 40(100 - 40)}{5^2} = \frac{36879}{25} \approx 1475 .$$

Vaadituksi otokseksi saadaan $n = 1475$ (äänestysikäistä henkilöä), mikä voidaan pyöristää 1400:ksi.

10.2.3 Korrelaatiokerroin

Samaa periaatetta noudattaen, ratkaisemme populaation korrelaatiokertoimen ρ luottamusvälin pituuden \mathbf{d}^* kaavan otoskoon \mathbf{n} suhteen eli

$$\mathbf{d}^* = 2 \cdot Z_{\alpha/2} \cdot \frac{1}{\sqrt{\mathbf{n} - 3}} \Leftrightarrow \mathbf{n} = \frac{4 \cdot (Z_{\alpha/2})^2 + 3}{(\mathbf{d}^*)^2}.$$

Saatu otoskoon lauseke riippuu valitusta varmuustasosta $Z_{\alpha/2}$ ja siitä tarkkuudesta, eli luottamusvälin pituudesta \mathbf{d}^* , jolla tahdomme populaation korrelaatiokertoimen arvon ρ estimoida. On huomattava, että \mathbf{d}^* edustaa nyt Fisher-transformoinnin

kautta muodostunutta pidentynyttä luottamusvälin pituutta ja se pitää siksi lyhentää eli Fisher-transformoida käänteisesti vastaamaan varsinaista luottamusvälin pituutta (ks. 6.5, s. 63).

Koska \mathbf{d}^* :n arvo (eli siis estimointitarkkuus) asetetaan ρ :n yksikköinä (esim. 0,2), voidaan se käänteistransformoida vastaavaksi \mathbf{d} :n arvoksi samaan tapaan, kuten transformoimme ρ :n arvoja luottamusvälin päätepisteitä määrättäessä. Tämä voidaan suorittaa Excelin FISHERINV –funktiolla, jolloin kirjoitetaan soluun =FISHERINV(0,2) ja saadaan $\mathbf{d} = 0,1974$. Usein kuitenkin käytännössä haluttaisiin saada \mathbf{d} :n arvoksi tietty arvo (esim. juuri 0,2), joka ei sitten enää muuntuisi käänteistransformaatioissa. Tällöin täytyy määrätä sellainen \mathbf{d}^* :n arvo, jonka muunnoksena saadaan haluttu \mathbf{d} :n arvo. Tämä tietysti saadaan määrättyä päinvastaisesti menetellen eli suorana FISHER-transformaationa

$$\mathbf{d}^* = \text{FISHER}(\mathbf{d}),$$

jolloin asetetaan haluttua varmuusväliä \mathbf{d} vastaava pidennetty väli \mathbf{d}^* laskukaavaan. Eli esim. 0,2 yksikön tapauksessa saataisiin \mathbf{d}^* :n arvoksi 0,2027. voidaan huomata, että pienillä \mathbf{d} :n arvoilla muutos on melko vähäistä, mutta ei kuitenkaan estimoinnin kannalta merkityksetöntä.

- **ESIM:** Tutkija haluaa selvittää 0,1 yksikön tarkkuudella, miten paljon kuukausiansioiden ja koetun onnellisuuden välillä on korrelaatiota työikäisten muodostamassa perusjoukossa. Tutkijalla on ongelmana, miten suuri otoskoon pitäisi olla, jotta ko. tarkkuus populaation korrelaatiokertoimen arvon estimoinnissa saavutettaisiin. Hän lisäksi haluaisi estimointiin sellaisen varmuustason, että vain yhdessä otoksessa sadasta tulos voisi sattumalta poiketa asetetusta tarkkuusvälistä 0,1. Mikä on tutkijan tarvitsema otoskoko ?

Laskemiseen tarvittavat kaksi arvoa ovat: $d^*=0,1$ ja liitteen 1 taulukosta standardoidun normaalijakauman 99%:n kriittinen arvo ($Z_{0,01/2}$) 2,58. Ensiksi muunnetaan d :n arvo d^* :ksi eli = **FISHER(0,1) = 0,1003** Tämän jälkeen asetetaan arvot kaavaan ja lasketaan

$$n = \frac{4 \cdot (2,58)^2 + 3}{0,1003^2} \approx 2944 .$$

Saatu otoskoko 2944 (työkäistä) on melko suuri, johon ratkaisevasti vaikuttaa juuri 0,1 yksikön estimointitarkkuus. Varmuustaso vaikuttaa vähemmän sillä, jos riskitaso olisi 5% (eli $Z_{0,05/2} = 1,96$), olisi vaadittu otoskoko n. 1826. Mutta jos estimointitarkkuudeksi valittaisiin 0,2 yksikköä, riittäisi otoskooksi (1%:n riskitasollakin) n. 721.

Taulukko 24: Otoskokoja populaation korrelaatiokerroimen ρ estimoimiseksi eri d :n ja α :n arvoille (d on estimoinnin tarkkuus ρ :n yksikköinä).

d	Riskitaso α ($\times 100\%$)			
	10 %	5 %	1 %	0,10 %
0,05	5520	7334	11796	18493
0,1	1373	1824	2934	4600
0,15	605	804	1293	2027
0,2	336	447	719	1127
0,25	212	282	453	710
0,3	144	192	308	483
0,35	103	138	221	347
0,4	77	102	165	258
0,45	59	78	126	197
0,5	46	61	98	153

11 Otantamenetelmät

11.1 Otantamenetelmien käytön perusta

Otantamenetelmiä käytetään, kun halutaan arvioida (estimoida) mahdollisimman tarkkaan yhden tai yleensä useamman muuttujan tunnuslukuja (lähinnä keskiarvoja tai prosenttilukuja) *äärellisessä* perusjoukossa eli populaatiossa. Luvuissa 4-6 käsiteltiin osittain samaa asiaa, mutta tällöin muuttujan arvot äärellisessä perusjoukossa on korvattu muuttujan jakaumamallilla (esim. normaalijakaumalla), jossa havaintoja on ääretön määrä. Edellisissä luvuissa esitetyt tulokset, jotka koskivat otantajakaumia, estimoinnin tarkkuutta, luottamusvälejä ja hypoteesin testausta, ovat pienin muutoksin voimassa myös äärellisen perusjoukon tapauksessa. Tärkein muutos on, että estimaattoreiden keski-*virheisiin* tulee mukaan *äärellisen perusjoukon korjaustermi*.

Tässä luvussa käsitellään lyhyesti otantamenetelmiä keskittyen lähinnä yksinkertaiseen satunnaisotantaan ja ositettuun otantaan.

Tarkastellaan ensin estimointia äärellisessä perusjoukossa hieman yleisemmästä näkökulmasta. Periaatteessa on kolme tapaa selvittää tuntematon tunnusluku: 1) *kokonaistutkimus*, 2) *otantatutkimus* ja 3) *muut menetelmät*. Kokonaistutkimuksessa mitataan tai ainakin pyritään mittaamaan perusjoukon kaikki yksiköt. Otantatutkimuksessa poimitaan perusjoukosta satunnaisesti osajoukko eli *otos*, jonka pohjalta lasketaan arvio eli *estimaatti* perusjoukon tunnusluvulle. Muista menetelmistä mainittakoon yhden tai useamman asiantuntijan näkemykseen perustuvat arviot ja nk. katu-gallupit. Näitä menetelmiä ei käsitellä yhteydessä.

On selvää että kokonaistutkimus on kalliimpi ja tulokset saadaan hitaammin kuin otantatutkimuksessa. Usein otantatutkimus on käytännössä ainoa vaihtoehto. On ehkä yllättävää, mutta otantatutkimuksella saadaan jopa luotettavampi tulos kuin kokonaistutkimuksella.

Kokonais- ja otantatutkimuksen vertailemiseksi tarkastellaan ensin tunnusluvun arviointiin liittyviä virhelähteitä, joista tärkeimmät ovat 1) *otantavirhe*, 2) *mittausvirhe* ja 3) *kato*.

Otantavirhe on yhtä kuin otoksesta lasketun tunnusluvun poikkeama perusjoukon todellisesta arvosta. Esimerkiksi jos otoksesta laskettu puolueen kannatusprosentti on 21 % ja todellinen arvo on 22 %, on otantavirhe -1 %-yksikköä. Yksittäisen otoksen kohdalla otantavirheen suuruutta ei tiedetä, mutta otantavirheen satunnaiskäyttäytyminen yleensä hallitaan. Voidaan esimerkiksi sanoa, kuinka suuret otantavirheet ovat harvinaisia. Kokonaistutkimuksessa ei esiinny otantavirhettä.

Mittausvirhettä syntyy esimerkiksi, kun vastaaja salaa todellisen puoluekantansa ja ilmoittaa toisen puolueen. Mittausvirhettä voidaan vähentää käyttämällä luotettavampia mittausmenetelmiä. Esimerkiksi jos haastattelutilanteessa tataan vastausten anonyymisyys, mittausvirheen todennäköisyys pienenee. Koska

otantatutkimuksessa mittausten kokonaismäärä on pienempi kuin kokonaistutkimuksessa, voidaan otantatutkimuksessa käyttää tarkempia mittausten menetelmiä. Siten otantatutkimuksessa mittausten virhettä esiintyy yleensä vähemmän kuin kokonaistutkimuksessa.

Katoa esiintyy, jos otokseen tai kokonaistutkimuksessa perusjoukkoon kuuluvaa henkilöä ei tavoiteta tai hän kieltäytyy vastaamasta. Kokonaistutkimuksessa vastausten suuri määrä tuudittaa helposti vääriin luotettavuuden tunteeseen ja katoon liittyvät ongelmat helposti unohtetaan. Otantatutkimuksessa on tapana tehdä vähintään yksi uusintakysely (karhuaminen) kadon piiriin kuuluvista ja siten kadon vaikutukset voidaan arvioida ja ottaa huomioon estimaateissa. Kolmesta edellä käsitellystä virhelähteestä kahdessa (mittausvirhe ja kato) otantatutkimus johtaa yleensä pienempää virheeseen kuin kokonaistutkimus.

11.2 Otantamenetelmät

Otantamenetelmä sisältää kaksi seikkaa: tapa jolla otos poimitaan ja kaava jolla arvio (estimaatti) tuntemattomalle tunnusluvulle lasketaan. Tärkeimmät otantamenetelmät ovat

- 1) **Yksinkertainen satunnaisotanta** (YSO, simple random sample)
- 2) **Systemaattinen eli tasavälinen otanta** (systematic sampling)
- 3) **Ositettu otanta** (stratified sampling)
- 4) **Ryväsotanta** (cluster sampling)

Seuraavassa tullaan käyttämään samoja merkintöjä kuin edelläkin tässä kirjassa. Uutena merkintänä on perusjoukon suuruus **N**. Muistettakoon että otoskokoa merkitään vastaavalla pienellä kirjaimella eli **n**:llä.

11.3 Yksinkertainen satunnaisotanta

Yksinkertaisessa satunnaisotannassa tilastoyksiköt numeroidaan 1:stä N:ään tai käytetään jo olemassa olevaa numerointia (esim. asiakasnumerot). Seuraavassa esimerkissä otoksen poiminta tehdään asiakasnumeron avulla (Anro). Asiakasnumeron tulee olla yksikäsitteinen eli kahdella asiakkaalla ei saa olla samaa asiakasnumeroa.

Taulukko 25: Asiakasrekisterin alkua Excel-ympäristössä

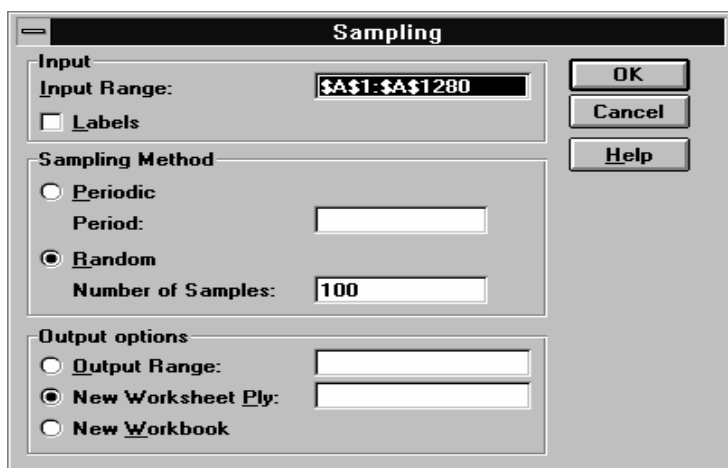
	A	B	C
1		Asiakasrekisteri	
2	Anro	Asiakas	Osoite
3	30123	Asiakas 1	Osoite 1
4	30345	Asiakas 2	Osoite 2
5	20678	Asiakas 3	Osoite 3
6	10123	Asiakas 4	Osoite 4
7	40234	Asiakas 5	Osoite 5
8	60678	Asiakas 6	Osoite 6
9	20123	Asiakas 7	Osoite 7

Huom1. Excelin poiminta edellyttää, että poiminnassa käytettävän luvun on oltava numeromuodossa - siinä ei saa olla esim. kirjaimia.

Yksinkertaisessa satunnaisotannassa jokaisella yksiköllä on yhtä suuri todennäköisyys tulla poimituksi otokseen. Käytännössä otoksen poiminta tehdään satunnaislukujen avulla. Aikaisemmin käytettiin satunnaislukutaulukoita, mutta nykyään esim. Excel-taulukkolaskennassa on Sampling –menetelmä, jolla poiminta voidaan tehdä.

Valitaan **Tools/Data Analysis** –valikosta (Työkalut/Tietojen analysointi) **Sampling** –vaihtoehto (Otanta), jolloin saadaan seuraava valintaikkuna.

Kuva 48: Otannan määrittelyikkuna Excelissä



Määrittele **Input Range** –alueeksi solut, joissa on perusjoukon numerointi (esimerkissä A1:A1280). **Sampling Method** on Random ja **Number of Samples** (termi on muuten väärä - pitäisi olla Sample Size) on otoskoko **n**.

Poimitettujen lukujen sijoittamiselle on kolme vaihtoehtoa:

- Sijoitetaan tulostusalueelle (Output Range), jonka määrittelemiseksi riittää ensimmäinen solu.
- Sijoitetaan uudelle työarkille (New Worksheet Ply).
- Sijoitetaan uuteen työkirjaan (New Workbook).

Huom1. Excelin poiminnassa sama luku voi esiintyä otoksessa useamman kerran. Siksi on syytä määritellä poiminnassa otoskooksi hieman suurempi luku kuin varsinainen otoskoko. Poimituista luvuista voidaan poistaa toistot **Data/Filter/Advanced Filter** komennolla (Tiedot/Suodata/Erikoissuodatus) ja merkitsemällä rasti **Unique records only** –kohtaan (Kuva 49). Toisaalta toisto-

jen poisto voidaan tehdä myös seuraavassa vaiheessa, kun otokseen poimituihin tilastoyksiköihin liitetään rekisteristä lisätietoja.

Otokseen poimituista yksiköistä (esim. asiakkaista) tarvitaan yleensä muitakin tietokannassa olevia tietoja. Nämä saadaan poimittua **Data/Filter/Advanced Filter** –komennolla, jonka valintaikkuna on esitetään Kuva 49:ssä.

Kuva 49: Tietojen erikoissuodatus Excelissä



Copy to Another Location –vaihtoehto valitaan, koska otokseen poimituista yksiköistä halutaan oma rekisteri.

List Range on alkuperäinen rekisteri siten, että siihen sisältyy myös kenttien nimet.

Criteria Range –alueeseen sijoitetaan otokseen poimitut luvut. Luvuissa voi toistoja. Alueen ensimmäiseen soluun sijoitetaan *kentän nimi* (Taulukko 25:n esimerkissä s. 124 kentän nimi on **Anro**).

Copy to –alueeseen merkitään tulosalueen vasen ylänurkka.

Unique Records Only –ruutuun merkitään rasti, koska toistoja ei haluta mukaan.

11.4 Systemaattinen otanta

Systemaattisessa eli tasavälisessä otannassa ei perusjoukon tilastoyksiköitä tarvitse numeroida. Riittää kun yksiköt ovat peräkkäisessä järjestyksessä, esim. asiakasrekisteri, jossa ei ole asiakasnumeroa — tai laatikossa olevat osoitekortit.

Oletetaan, että *poimintaväli* $k = N/n$ on kokonaisluku. Esimerkiksi jos $N = 10000$ ja $n = 200$, on poimintaväli $k = 50$. Perusjoukon k ensimmäistä yksikköä numeroidaan ja näistä poimitaan otokseen satunnaisesti yksi, jonka järjestysnumero olkoon u . Seuraavat yksiköt poimitaan otokseen u :sta lukien poimintavälin k välein. Siten otokseen poimitujen yksiköiden järjestysnumerot ovat u , $u+k$, $u+2k$, $u+3k$ jne.

Jos perusjoukossa on esim. systemaattista vaihtelua, saattaa systemaattinen otantaa perusjoukosta vääristyneen kuvan. Esimerkiksi jos kioskin päivittäisestä myynnistä valitaan joka seitsemäs päivä, voi virhe olla varsiin suuri puoleen tai toiseen.

11.5 Estimointi yksinkertaisessa satunnaisotannassa

Otoskeskiarvon \bar{X} keskivirheelle $D(\bar{X})$ voidaan matemaattisesti johtaa seuraava kaava:

$$D(\bar{X}) = \frac{\sigma}{\sqrt{n}} \sqrt{1 - \frac{n}{N}},$$

jossa σ on muuttujan keskihajonta perusjoukossa. Suhdetta n/N kutsutaan **otantasuhteeksi** (voidaan ilmaista myös prosentuaalisena) ja termiä $1 - n/N$ vastaavasti **äärellisen perusjoukon korjaustermiksi**.

Keskiarvon keskivirheen kaavasta voidaan tehdä seuraavia johtopäätöksiä:

1. Otoskoko n on jakajassa neliöjuuren alla. Jos keskivirhe halutaan otoskoko kasvatamalla esim. puolittaa, on otoskoko nelinkertaistettava.
2. Muuttujan keskihajonta vaikuttaa suorassa suhteessa keskivirheeseen. Käytännössä muuttujan keskihajontaa ei voida pienentää. Sen sijaan perusjoukko voidaan jakaa osiin eli **ositteisiin** siten, että ositteiden sisällä muuttujan vaihtelu on vähäisempää. Tähän perustuu ositettu otanta.
3. Perusjoukon koko N ei yleensä vaikuta keskivirheen suuruuteen, koska otantasuhde n/N on yleensä lähellä nollaa, jolloin vastaavasti termi $1 - n/N$ on lähellä ykköstä. Eli olipa perusjoukon suuruus 30000 henkeä tai 3 milj., niin keskivirhe on käytännössä sama.
4. Pelkästään otantasuhteen arvosta ei voida päätellä estimoinnin tarkkuutta. Esimerkiksi jos kahdesta eri otoksesta tiedetään ainoastaan otantasuhteet 1% ja 2%, niin ei voida päätellä, että jälkimmäinen otanta antaisi tarkemman tuloksen.

Yleensä otantasuhde on niin pieni (alle 10 %), että äärellisen perusjoukon korjaustermi voidaan unohtaa. Silloin keskiarvon keskivirhe yksinkertaistuu muotoon

$$D(\bar{X}) = \frac{\sigma}{\sqrt{n}},$$

joka on sama kuin keskivirheen lauseke jo luvussa 4.2, s. 32.

Prosenttiluvun p keskivirheelle on vastaavasti voimassa kaava

$$D(p) = \frac{\sqrt{\theta(100-\theta)}}{\sqrt{n}} \sqrt{1-\frac{n}{N}},$$

jossa θ on perusjoukon prosenttiluku (ks. myös 4.3 s. 34).

11.5.1 Luottamusväli

Keskivirheen avulla saadaan perusjoukon tuntemattomalle tunnusluvulle luottamusväli (ks. luku 5.3, s. 56 ja luku 6, s. 57). Yleisimmin käytetään 95 %:n luottamusväliä, jolloin on 95 %:n varmuus, että perusjoukon tunnusluku todella kuuluu otoksesta laskettuun väliin. Keskiarvon ja prosenttiluvun 95 %:n luottamusväli ovat likimäärin

$$\bar{X} \pm 1,96 * \text{keskivirhe} \quad \text{ja} \quad p \pm 1,96 * \text{keskivirhe}.$$

11.6 *Ositettu otanta*

Ositetussa otannassa perusjoukko jaetaan osiin eli ositteisiin, joissa kussakin tehdään oma otanta ja näin saadut tulokset yhdistetään koko perusjoukon tunnusluvun estimaatiksi. Seuraavassa esimerkkejä osituksesta:

- Yritykset jaetaan suuriin, keskisuuriin ja pieniin yrityksiin.
- Sairaanhoidopiiri muodostuu 25 kunnasta. Kukin kunta muodostaa oman ositteensa.
- Väestö jakautuu läänijaon mukaisiin ositteisiin

Huom. Ositusta ei pidä sotkea ryvästykseseen (ryväсотantaan). Ryväсотannassa rypäistä poimitaan otos (esim. kaikista Suomen kunnista poimitaan 30 kuntaa, joista poimitaan edelleen henkilötason otokset).

Ositetun otannan käytölle voidaan löytää useita syitä:

- Halutaan tarkempia estimaatteja. Ositetun otannan estimaatit verrattuna yksinkertaisella satunnaisotannalla saatuihin estimaatteihin ovat sitä tarkempia mitä enemmän ositteiden keskiarvot tai prosenttiluvut eroavat toisistaan.
 - Jokaisessa ositteessa halutaan tunnusluku estimoida halutulla tarkkuudella, minkä lisäksi halutaan saada estimaatti koko perusjoukon tunnusluvulle.
 - Tietojen keräämiseen liittyvät syyt: joissakin ositteissa käytetään kirjekselyä ja toisissa ositteissa henkilökohtaista käyntiä.
-

11.6.1 Estimointi ositetussa otannassa

Jokaisesta ositteesta poimitaan yksinkertainen satunnaisotos. Merkitään **i.** ositteen otoskeskiarvoa \bar{X}_i :lla. Nämä ovat tietenkin ositteiden keskiarvojen estimaatteja. Koko perusjoukon harhaton estimaatti saadaan ositteiden keskiarvojen painotettuna keskiarvona, kun painoina käytetään *ositepainoja* $W_i = N_i/N$ (N on perusjoukon koko eli N_i :den summa). Kaavana saadaan

$$\bar{X}_0 = \sum_i W_i \bar{X}_i = \frac{1}{N} \sum_i \sum_j \frac{N_i}{n_i} x_{ij},$$

jossa x_{ij} on **j.** tilastoyksikön arvo **i.** ositteessa.

Tästä kaavasta huomataan, että on kaksi tapaa laskea keskiarvon estimaatti ositetussa otannassa.

- 1) Painotetaan ositekeskiarvoja ositepainoilla.
- 2) Painotetaan havaintoarvoa x_{ij} osamäärällä N_i/n_i . Useimmissa tilasto-ohjelmissa on painomuuttujan käyttömahdollisuus, jolloin painomuuttujan arvo määrätään tämän osamäärän mukaisesti.

On huomattava, että tämä estimaatti ei yleensä ole kaikkien otokseen poimittujen lukujen keskiarvo. Ainoastaan *suhteellisen kiintiöinnin* tapauksessa ositetun otannan estimaatti on yhtä kuin otoksen lukujen keskiarvo.

Huom. Koska prosenttiluku voidaan tulkita keskiarvoksi, pätee edellä ollut kaava myös prosenttiluvulle.

Ositetun otannan estimaatin keskivirheelle on voimassa lauseke

$$D(\bar{X}_0) = \sqrt{\sum W_i^2 \frac{\sigma_i^2}{n_i} (1 - n_i/N_i)}.$$

Esimerkki. Vuonna 1993 tehtiin kyselytutkimus kansalaisten energia-asenteista. Koska ydinvoimapaikkakuntien Eurajoen ja Loviisan asukkaiden mielipiteet halettiin erityisesti selvittää, poimittiin otos ositettuna siten, että ositteita oli kolme: Eurajoki, Loviisa ja muu Suomi. Useamman sadan kysymyksen joukosta tarkastellaan tässä lähemmin kysymystä ”Pitääkö ydinvoiman käyttöä lisätä?” Sen vastausvaihtoehdot ovat viereisessä taulukossa.

- | | |
|---|-------------------------------------|
| 1 | Käyttöä pitäisi tuntuvasti lisätä |
| 2 | Käyttöä pitäisi hieman lisätä |
| 3 | Käytön nykyinen taso on sopiva |
| 4 | Käyttöä pitäisi hieman vähentää |
| 5 | Käyttöä pitäisi tuntuvasti vähentää |
| 6 | Käytöstä pitäisi luopua kokonaan |
| 7 | En osaa sanoa |

Jotta tästä kysymyksestä olisi järkevä laskea keskiarvoja, jätettiin vaihtoehto seitsemän ”En osaa sanoa” pois. Ositekohtaiset keskiarvot ja muut tiedot ovat seuraavassa taulukossa.

Taulukko 26: Pitääkö ydinvoiman käyttöä lisätä? Ositekohtaiset tiedot.

Osite	Ositteen koko N_i	Otoskoko n_i	Ositepaino W_i	N_i/n_i	Keskiarvo
Eurajoki	4 800	123	0,00120	39,0	2,642
Loviisa	7 700	122	0,00193	63,1	2,643
Muu Suomi	3 983 000	1 339	0,99687	2 974,6	3,445
Yhteensä	3 995 500	1 584	1,00000	2 522,4	3,315

Huomataan että ydinvoimapaikkakunnilla suhtautuminen ydinvoiman lisärakentamiseen on huomattavasti myönteisempää kuin muulla Suomessa. Taulukossa oleva yhteensä –keskiarvo 3,315 on laskettu tavallisena keskiarvona ja on tässä tapauksessa perusjoukon keskiarvon harhainen estimaatti. Oikea harhaton keskiarvo saadaan painottamalla ositekohtaisia keskiarvoja ositepainoilla ja tulokseksi saadaan 3,442.

11.6.2 Otoksen kiintiöinti

Kun otoskoko määrätty, on vielä päätettävä, miten otoskoko jaetaan ositteiden kesken. Tätä kutsutaan otoksen *kiintiöinniksi* eli *allokoinniksi*. On huomattava että ositettu otanta sinänsä sallii millaisen kiintiöinnin tahansa. Ositetun otannan estimaattori antaa aina harhattoman tuloksen. Otoksen kiintiöinti vaikuttaa sen sijaan estimaattorin keskivirheeseen. Yleensä halutaan estimoida perusjoukon tunnusluku mahdollisimman tarkasti. Joskus voi olla tavoitteena arvioida ositekohtaiset tunnusluvut samalla tarkkuudella.

Seuraavassa taulukossa on kolme yleisimmin käytettyä kiintiöintiä.

Taulukko 27: Otoksen kiintiöintimenetelmät

Tasainen kiintiöinti	$n_i = n/(\text{ositteiden lkm})$
Suhteellinen kiintiöinti	$n_i = n \cdot N_i/N$
Optimaalinen kiintiöinti	n_i on suhteellinen tuloon $N_i \cdot \sigma_i$

Tarkastellaan esimerkkiä, jossa tarkoituksena on arvioida USA:n sahojen tuotannon määrä v. 1943, kun edellisen vuoden tiedot ovat käytettävissä koko perusjoukon osalta. Sahat on jaettu kolmeen ositteeseen: suuret, keskisuuret ja pienet sahat. Taulukko 28 esittää esimerkkiin liittyvät ositekohtaiset tiedot:

Taulukko 28: Ositekohtaiset tiedot

Osite	Sahan tuotanto v. 1942	Sahojen lkm	Keskim. tuotanto v. 1942	Arvioitu kes- kihajonta
Suuret	5000-	538	11030	9000
Keskisuuret	1000-4999	4756	1780	1200
Pienet	-999	30964	204	300
Yhteensä		36258	571	1684

Verrataan keskiarvon estimoinnin tarkkuutta yksinkertaisen satunnaisotannan ja ositetun otannan kolmen kiintiöinnin tapauksissa, kun otoskoko on sata sahaa. Ositetun otannan kolme kiintiöintivaihtoehtoa antaa Taulukko 29:ssä esitetyt ositteiden otoskoot.

Taulukko 29: Sadan sahan otos allokoitu kolmella eri tavalla

Osite	Tasainen kiintiöinti	Suhteellinen Kiintiöinti	Optimaalinen kiintiöinti
Suuret	34	2	24
Keskisuuret	33	13	29
Pienet	33	85	47

Suhteellinen kiintiöinti ottaa huomioon vain ositteen koon, minkä vuoksi suurten sahojen ositteesta valitaan vain kaksi sahaa.

Vertailtavien menetelmien keskivirheet esitetään Taulukko 30:ssä.

Taulukko 30: Otanta- ja kiintiöintimenetelmien keskivirheet

Menetelmä	Keskivirhe
Yksinkertainen satunnaisotanta	168,0
Ositettu otanta, tasainen kiintiöinti	57,1
Ositettu otanta, suhteellinen kiintiöinti	107,7
Ositettu otanta, optimaalinen kiintiöinti	54,7

Huomataan, että kaikilla ositetun otannan kiintiöinneillä saadaan huomattavasti tarkemmat estimaatit kuin yksinkertaisella satunnaisotannalla. Optimaalisen kiintiöinnin keskivirhe on n. puolet suhteellisen kiintiöinnin keskivirheestä. Tämä ero johtuu ositteiden keskihajontojen suurista eroista.

11.7 Ryväсотanta

Ryväsotannassa perusjoukko jaetaan rypäisiin (esim. kuntajako). Rypäistä poimitaan otos (esim. 30 kuntaa) ja näistä otokseen poimituista rypäistä poimitaan henkilö- tai yritystason otokset. Tällöin on kyse *kaksiasteisesta ryväсотannasta*, sillä otantaa käytettiin kahdessa vaiheessa.

Ryväsotannalla voidaan vähentää kustannuksia varsinkin, jos otokseen poimitujen henkilöiden luona tai yrityksissä joudutaan käymään.

Koska rypäiden koko yleensä vaihtelee suuresti (esim. kunnat), on syytä poimia rypäät niiden suuruuteen suhteutetuilla todennäköisyyksillä. Rypäiden sisällä käytetään tilastoyksiköiden poiminnassa tasaista kiintiöintiä eli jokaisesta poimitusta rypästä poimitaan sama määrä tilastoyksiköitä riippumatta rypään koosta. Tällöin keskiarvon ja prosenttiluvun estimaattori on *itsepainottuva*, jolloin estimaatit voidaan laskea tavalliseen tapaan ilman mitään painotusta.

Ositettu otanta ja ryväсотanta voidaan yhdistää siten, että rypäät jaetaan ositteisiin (esim. kuntien läänijako). Jokaisessa ositteessa käytetään ryväсотantaa ja ositekohtaiset estimaatit yhdistetään lopuksi koko perusjoukon tunnusluvun estimaatiksi.

11.8 Otoksoon määrääminen äärellisen perusjoukon tapauksessa

Luvussa 10 (s. 114) käsiteltiin otoksoon määräämistä äärettömän perusjoukon tapauksessa. Vaikka tällöinkin käytännössä perusjoukko on äärellinen, ei perusjoukon suuruutta otettu näissä tarkasteluissa huomioon. Tämän luvun tulokset saadaan muunnettua vastaamaan äärellistä perusjoukkoa yksinkertaisella kaavalla. Määrätään ensin otoskoko luvussa 10 esitetyllä tavalla, jolloin siis perusjoukon suuruus jätetään huomioonottamatta. Merkitään näin saatua otoskokoa n_0 :lla. Lopullinen otoskoko n saadaan kaavalla

$$n = \frac{n_0}{1 + \frac{n_0}{N}}.$$

Seuraavassa taulukossa (Taulukko 31) on havainnollistettu perusjoukon suuruuden vaikutusta otoskookoon. Siinä otoksoon arvo äärettömän perusjoukon tapauksessa on $n_0=100$ ja todellinen otoskoko on laskettu viidellä perusjoukon suuruuden arvolla.

Taulukko 31: Perusjoukon suuruuden vaikutus otoskookoon

Perusjoukon suuruus N	1 000 000	100 000	10 000	1 000	500
Otoskoko n	100	100	99	91	83

Tixel-ohjelmassa on laskentapohjat (ks. Taulukko 32 ja Taulukko 33) otoskoon määrittämiseksi prosenttiluvun ja keskiarvon tapauksessa. Ne perustuvat lukujen 10.2.1 ja 10.2.2 kaavoihin ja perusjoukon suuruuden huomioonottamiseen. Laskentapohjilla voidaan paitsi määrätä otoskoko myös arvioida estimoinnin tarkkuus, kun otoskoko tunnetaan.

Kun laskentapohjia käytetään otoskoon määrittämiseen, täytetään neljästä rasteroidusta solusta kaksi ensimmäistä ja alimmainen solu, jolla määritellään haluttu estimoinnin tarkkuus (95 %:n luottamusvälin puolikas).

Prosenttiluvun tapauksessa on annettava arvio estimoitavasta prosenttiluvusta. Suurin otoskoko saadaan arvolla 50 ja otoskoko pienenee, mitä enemmän loitonnutaan viidestäkymmenestä.

Taulukko 32: Otoskoon määrittäminen prosenttiluvun tapauksessa (Tixel)

Otoskoon määrittäminen ja estimoinnin tarkkuus %-luvut		
Perusjoukon suuruus: Arvio estimoitavasta %-luvusta:	<input type="text"/>	Laskennan tulos:
Otoskoko: tai	<input type="text"/>	#DIV/0!
95 %:n luottamusvälin puolikas:	<input type="text"/>	#DIV/0!
<p>O hje: Täytä rasteroidut solut. Otoskoko ja 95 %:n luottamusvälin pituus ovat vaihtoehtoisia tapoja määrittellä otannon tarkkuus. Jos syötät otoskoon arvon, antaa laskenta lv:n puolikkaan pituuden. Jos syötät lv:n puolikkaan arvon, antaa laskenta otoskoon arvon.</p>		

Taulukko 33: Otoskoon määrittäminen keskiarvon tapauksessa (Tixel)

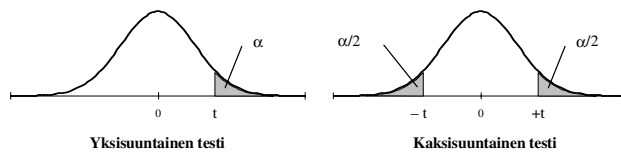
Otoskoon määrittäminen ja estimoinnin tarkkuus Keskiarvo		
Perusjoukon suuruus: Keskihajonta:	<input type="text"/>	Laskennan tulos:
Otoskoko: tai	<input type="text"/>	#DIV/0!
95 %:n luottamusvälin puolikas:	<input type="text"/>	#DIV/0!
<p>O hje: Täytä rasteroidut solut. Otoskoko ja 95 %:n luottamusvälin pituus ovat vaihtoehtoisia tapoja määrittellä otannon tarkkuus. Jos syötät otoskoon arvon, antaa laskenta lv:n puolikkaan pituuden. Jos syötät lv:n puolikkaan arvon, antaa laskenta otoskoon arvon.</p>		

Alla sijaitsevassa esimerkissä sadan tuhannen hengen perusjoukon prosenttiluku halutaan estimoida 3 %-yksikön tarkkuudella, kun karkea arvio prosenttiluvulle on 30 %. Otoskoon arvoksi saadaan 888 henkilöä.

Perusjoukon suuruus Arvio estimoitavasta %-luvusta	100 000 30	Laskennan tulos:
Otoskoko: tai	<input type="text"/>	888
95 %:n luottamusvälin puolikas	3	#DIV/0!

3,5	0,9998	0,9998	0,9998	0,9998	0,9998	0,9998	0,9998	0,9998	0,9998	0,9998
-----	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------

LIITE 2: t-jakauman kriittisiä arvoja



df	Merkitsevyystaso yksisuuntaisessa testissä									
	0,25	0,2	0,15	0,1	0,05	0,025	0,01	0,005	0,0005	0,00005
	Merkitsevyystaso kaksisuuntaisessa testissä									
	0,50	0,40	0,30	0,20	0,10	0,05	0,02	0,01	0,001	0,0001
1	1,000	1,376	1,963	3,078	6,314	12,706	31,821	63,656	636,578	6370,544
2	0,816	1,061	1,386	1,886	2,920	4,303	6,965	9,925	31,600	100,136
3	0,765	0,978	1,250	1,638	2,353	3,182	4,541	5,841	12,924	28,014
4	0,741	0,941	1,190	1,533	2,132	2,776	3,747	4,604	8,610	15,534
5	0,727	0,920	1,156	1,476	2,015	2,571	3,365	4,032	6,869	11,176
6	0,718	0,906	1,134	1,440	1,943	2,447	3,143	3,707	5,959	9,080
7	0,711	0,896	1,119	1,415	1,895	2,365	2,998	3,499	5,408	7,888
8	0,706	0,889	1,108	1,397	1,860	2,306	2,896	3,355	5,041	7,120
9	0,703	0,883	1,100	1,383	1,833	2,262	2,821	3,250	4,781	6,594
10	0,700	0,879	1,093	1,372	1,812	2,228	2,764	3,169	4,587	6,212
11	0,697	0,876	1,088	1,363	1,796	2,201	2,718	3,106	4,437	5,923
12	0,695	0,873	1,083	1,356	1,782	2,179	2,681	3,055	4,318	5,695
13	0,694	0,870	1,079	1,350	1,771	2,160	2,650	3,012	4,221	5,513
14	0,692	0,868	1,076	1,345	1,761	2,145	2,624	2,977	4,140	5,364
15	0,691	0,866	1,074	1,341	1,753	2,131	2,602	2,947	4,073	5,239
16	0,690	0,865	1,071	1,337	1,746	2,120	2,583	2,921	4,015	5,134
17	0,689	0,863	1,069	1,333	1,740	2,110	2,567	2,898	3,965	5,043
18	0,688	0,862	1,067	1,330	1,734	2,101	2,552	2,878	3,922	4,966
19	0,688	0,861	1,066	1,328	1,729	2,093	2,539	2,861	3,883	4,899
20	0,687	0,860	1,064	1,325	1,725	2,086	2,528	2,845	3,850	4,838
21	0,686	0,859	1,063	1,323	1,721	2,080	2,518	2,831	3,819	4,785
22	0,686	0,858	1,061	1,321	1,717	2,074	2,508	2,819	3,792	4,736
23	0,685	0,858	1,060	1,319	1,714	2,069	2,500	2,807	3,768	4,694
24	0,685	0,857	1,059	1,318	1,711	2,064	2,492	2,797	3,745	4,654
25	0,684	0,856	1,058	1,316	1,708	2,060	2,485	2,787	3,725	4,619
26	0,684	0,856	1,058	1,315	1,706	2,056	2,479	2,779	3,707	4,587
27	0,684	0,855	1,057	1,314	1,703	2,052	2,473	2,771	3,689	4,556
28	0,683	0,855	1,056	1,313	1,701	2,048	2,467	2,763	3,674	4,531
29	0,683	0,854	1,055	1,311	1,699	2,045	2,462	2,756	3,660	4,505
30	0,683	0,854	1,055	1,310	1,697	2,042	2,457	2,750	3,646	4,482
40	0,681	0,851	1,050	1,303	1,684	2,021	2,423	2,704	3,551	4,321
50	0,679	0,849	1,047	1,299	1,676	2,009	2,403	2,678	3,496	4,228
60	0,679	0,848	1,045	1,296	1,671	2,000	2,390	2,660	3,460	4,169
70	0,678	0,847	1,044	1,294	1,667	1,994	2,381	2,648	3,435	4,127
80	0,678	0,846	1,043	1,292	1,664	1,990	2,374	2,639	3,416	4,095
90	0,677	0,846	1,042	1,291	1,662	1,987	2,368	2,632	3,402	4,072
100	0,677	0,845	1,042	1,290	1,660	1,984	2,364	2,626	3,390	4,054
200	0,676	0,843	1,039	1,286	1,653	1,972	2,345	2,601	3,340	3,970
300	0,675	0,843	1,038	1,284	1,650	1,968	2,339	2,592	3,323	3,944
400	0,675	0,843	1,038	1,284	1,649	1,966	2,336	2,588	3,315	3,930
500	0,675	0,842	1,038	1,283	1,648	1,965	2,334	2,586	3,310	3,922
∞	0,674	0,842	1,036	1,282	1,645	1,960	2,326	2,576	3,290	3,891

Kun $df \rightarrow \infty$, lähestyvät t-jakauman kriittiset arvot normaalijakaumaa
LIITE 3: Binomijakauman kertymäfunktion arvoja

$$F(k) = P(X \leq k) = \sum_{i=0}^k \binom{n}{i} p^i q^{n-i}$$

		p-arvo												
n	k	0,05	0,1	0,15	0,2	0,25	0,3	0,4	0,5	0,6	0,7	0,8	0,9	
2	0	0,9025	0,8100	0,7225	0,6400	0,5625	0,4900	0,3600	0,2500	0,1600	0,0900	0,0400	0,0100	
	1	0,9975	0,9900	0,9775	0,9600	0,9375	0,9100	0,8400	0,7500	0,6400	0,5100	0,3600	0,1900	
	2	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
3	0	0,8574	0,7290	0,6141	0,5120	0,4219	0,3430	0,2160	0,1250	0,0640	0,0270	0,0080	0,0010	
	1	0,9928	0,9720	0,9393	0,8960	0,8438	0,7840	0,6480	0,5000	0,3520	0,2160	0,1040	0,0280	
	2	0,9999	0,9990	0,9966	0,9920	0,9844	0,9730	0,9360	0,8750	0,7840	0,6570	0,4880	0,2710	
	3	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
4	0	0,8145	0,6561	0,5220	0,4096	0,3164	0,2401	0,1296	0,0625	0,0256	0,0081	0,0016	0,0001	
	1	0,9860	0,9477	0,8905	0,8192	0,7383	0,6517	0,4752	0,3125	0,1792	0,0837	0,0272	0,0037	
	2	0,9988	0,9963	0,9880	0,9728	0,9492	0,9163	0,8208	0,6875	0,5248	0,3483	0,1808	0,0523	
	3	1,0000	0,9999	0,9995	0,9984	0,9961	0,9919	0,9744	0,9375	0,8704	0,7599	0,5904	0,3439	
	4	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
5	0	0,7738	0,5905	0,4437	0,3277	0,2373	0,1681	0,0778	0,0313	0,0102	0,0024	0,0003	0,0000	
	1	0,9774	0,9185	0,8352	0,7373	0,6328	0,5282	0,3370	0,1875	0,0870	0,0308	0,0067	0,0005	
	2	0,9988	0,9914	0,9734	0,9421	0,8965	0,8369	0,6826	0,5000	0,3174	0,1631	0,0579	0,0086	
	3	1,0000	0,9995	0,9978	0,9933	0,9844	0,9692	0,9130	0,8125	0,6630	0,4718	0,2627	0,0815	
	4	1,0000	1,0000	0,9999	0,9997	0,9990	0,9976	0,9898	0,9688	0,9222	0,8319	0,6723	0,4095	
	5	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
6	0	0,7351	0,5314	0,3771	0,2621	0,1780	0,1176	0,0467	0,0156	0,0041	0,0007	0,0001	0,0000	
	1	0,9672	0,8857	0,7765	0,6554	0,5339	0,4202	0,2333	0,1094	0,0410	0,0109	0,0016	0,0001	
	2	0,9978	0,9842	0,9527	0,9011	0,8306	0,7443	0,5443	0,3438	0,1792	0,0705	0,0170	0,0013	
	3	0,9999	0,9987	0,9941	0,9830	0,9624	0,9295	0,8208	0,6563	0,4557	0,2557	0,0989	0,0159	
	4	1,0000	0,9999	0,9996	0,9984	0,9954	0,9891	0,9590	0,8906	0,7667	0,5798	0,3446	0,1143	
	5	1,0000	1,0000	1,0000	0,9999	0,9998	0,9998	0,9993	0,9959	0,9844	0,9533	0,8824	0,7379	0,4686
	6	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000
7	0	0,6983	0,4783	0,3206	0,2097	0,1335	0,0824	0,0280	0,0078	0,0016	0,0002	0,0000	0,0000	
	1	0,9556	0,8503	0,7166	0,5767	0,4449	0,3294	0,1586	0,0625	0,0188	0,0038	0,0004	0,0000	
	2	0,9962	0,9743	0,9262	0,8520	0,7564	0,6471	0,4199	0,2266	0,0963	0,0288	0,0047	0,0002	
	3	0,9998	0,9973	0,9879	0,9667	0,9294	0,8740	0,7102	0,5000	0,2898	0,1260	0,0333	0,0027	
	4	1,0000	0,9998	0,9988	0,9953	0,9871	0,9712	0,9037	0,7734	0,5801	0,3529	0,1480	0,0257	
	5	1,0000	1,0000	0,9999	0,9996	0,9987	0,9962	0,9812	0,9375	0,8414	0,6706	0,4233	0,1497	
	6	1,0000	1,0000	1,0000	1,0000	0,9999	0,9998	0,9984	0,9922	0,9720	0,9176	0,7903	0,5217	
	7	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000
8	0	0,6634	0,4305	0,2725	0,1678	0,1001	0,0576	0,0168	0,0039	0,0007	0,0001	0,0000	0,0000	
	1	0,9428	0,8131	0,6572	0,5033	0,3671	0,2553	0,1064	0,0352	0,0085	0,0013	0,0001	0,0000	
	2	0,9942	0,9619	0,8948	0,7969	0,6785	0,5518	0,3154	0,1445	0,0498	0,0113	0,0012	0,0000	
	3	0,9996	0,9950	0,9786	0,9437	0,8862	0,8059	0,5941	0,3633	0,1737	0,0580	0,0104	0,0004	
	4	1,0000	0,9996	0,9971	0,9896	0,9727	0,9420	0,8263	0,6367	0,4059	0,1941	0,0563	0,0050	
	5	1,0000	1,0000	0,9998	0,9988	0,9958	0,9887	0,9502	0,8555	0,6846	0,4482	0,2031	0,0381	
	6	1,0000	1,0000	1,0000	0,9999	0,9996	0,9987	0,9915	0,9648	0,8936	0,7447	0,4967	0,1869	
	7	1,0000	1,0000	1,0000	1,0000	1,0000	0,9999	0,9993	0,9961	0,9832	0,9424	0,8322	0,5695	
	8	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000
9	0	0,6302	0,3874	0,2316	0,1342	0,0751	0,0404	0,0101	0,0020	0,0003	0,0000	0,0000	0,0000	
	1	0,9288	0,7748	0,5995	0,4362	0,3003	0,1960	0,0705	0,0195	0,0038	0,0004	0,0000	0,0000	
	2	0,9916	0,9470	0,8591	0,7382	0,6007	0,4628	0,2318	0,0898	0,0250	0,0043	0,0003	0,0000	
	3	0,9994	0,9917	0,9661	0,9144	0,8343	0,7297	0,4826	0,2539	0,0994	0,0253	0,0031	0,0001	
	4	1,0000	0,9991	0,9944	0,9804	0,9511	0,9012	0,7334	0,5000	0,2666	0,0988	0,0196	0,0009	
	5	1,0000	0,9999	0,9994	0,9969	0,9900	0,9747	0,9006	0,7461	0,5174	0,2703	0,0856	0,0083	
	6	1,0000	1,0000	1,0000	0,9997	0,9997	0,9987	0,9957	0,9750	0,9102	0,7682	0,5372	0,2618	0,0530
	7	1,0000	1,0000	1,0000	1,0000	0,9999	0,9996	0,9962	0,9805	0,9295	0,8040	0,5638	0,2252	0,0530

8	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9997	0,9980	0,9899	0,9596	0,8658	0,6126
9	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000
p-arvo													
n	k	0,05	0,1	0,15	0,2	0,25	0,3	0,4	0,5	0,6	0,7	0,8	0,9
10	0	0,5987	0,3487	0,1969	0,1074	0,0563	0,0282	0,0060	0,0010	0,0001	0,0000	0,0000	0,0000
	1	0,9139	0,7361	0,5443	0,3758	0,2440	0,1493	0,0464	0,0107	0,0017	0,0001	0,0000	0,0000
	2	0,9885	0,9298	0,8202	0,6778	0,5256	0,3828	0,1673	0,0547	0,0123	0,0016	0,0001	0,0000
	3	0,9990	0,9872	0,9500	0,8791	0,7759	0,6496	0,3823	0,1719	0,0548	0,0106	0,0009	0,0000
	4	0,9999	0,9984	0,9901	0,9672	0,9219	0,8497	0,6331	0,3770	0,1662	0,0473	0,0064	0,0001
	5	1,0000	0,9999	0,9986	0,9936	0,9803	0,9527	0,8338	0,6230	0,3669	0,1503	0,0328	0,0016
	6	1,0000	1,0000	0,9999	0,9991	0,9965	0,9894	0,9452	0,8281	0,6177	0,3504	0,1209	0,0128
	7	1,0000	1,0000	1,0000	0,9999	0,9996	0,9984	0,9877	0,9453	0,8327	0,6172	0,3222	0,0702
	8	1,0000	1,0000	1,0000	1,0000	1,0000	0,9999	0,9983	0,9893	0,9536	0,8507	0,6242	0,2639
	9	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9999	0,9990	0,9940	0,9718	0,8926	0,6513
10	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
12	0	0,5404	0,2824	0,1422	0,0687	0,0317	0,0138	0,0022	0,0002	0,0000	0,0000	0,0000	0,0000
	1	0,8816	0,6590	0,4435	0,2749	0,1584	0,0850	0,0196	0,0032	0,0003	0,0000	0,0000	0,0000
	2	0,9804	0,8891	0,7358	0,5583	0,3907	0,2528	0,0834	0,0193	0,0028	0,0002	0,0000	0,0000
	3	0,9978	0,9744	0,9078	0,7946	0,6488	0,4925	0,2253	0,0730	0,0153	0,0017	0,0001	0,0000
	4	0,9998	0,9957	0,9761	0,9274	0,8424	0,7237	0,4382	0,1938	0,0573	0,0095	0,0006	0,0000
	5	1,0000	0,9995	0,9954	0,9806	0,9456	0,8822	0,6652	0,3872	0,1582	0,0386	0,0039	0,0001
	6	1,0000	0,9999	0,9993	0,9961	0,9857	0,9614	0,8418	0,6128	0,3348	0,1178	0,0194	0,0005
	7	1,0000	1,0000	0,9999	0,9994	0,9972	0,9905	0,9427	0,8062	0,5618	0,2763	0,0726	0,0043
	8	1,0000	1,0000	1,0000	0,9999	0,9996	0,9983	0,9847	0,9270	0,7747	0,5075	0,2054	0,0256
	9	1,0000	1,0000	1,0000	1,0000	1,0000	0,9998	0,9972	0,9807	0,9166	0,7472	0,4417	0,1109
	10	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9997	0,9968	0,9804	0,9150	0,7251	0,3410
	11	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9998	0,9978	0,9862	0,9313	0,7176
12	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
15	0	0,4633	0,2059	0,0874	0,0352	0,0134	0,0047	0,0005	0,0000	0,0000	0,0000	0,0000	0,0000
	1	0,8290	0,5490	0,3186	0,1671	0,0802	0,0353	0,0052	0,0005	0,0000	0,0000	0,0000	0,0000
	2	0,9638	0,8159	0,6042	0,3980	0,2361	0,1268	0,0271	0,0037	0,0003	0,0000	0,0000	0,0000
	3	0,9945	0,9444	0,8227	0,6482	0,4613	0,2969	0,0905	0,0176	0,0019	0,0001	0,0000	0,0000
	4	0,9994	0,9873	0,9383	0,8358	0,6865	0,5155	0,2173	0,0592	0,0093	0,0007	0,0000	0,0000
	5	0,9999	0,9978	0,9832	0,9389	0,8516	0,7216	0,4032	0,1509	0,0338	0,0037	0,0001	0,0000
	6	1,0000	0,9997	0,9964	0,9819	0,9434	0,8689	0,6098	0,3036	0,0950	0,0152	0,0008	0,0000
	7	1,0000	1,0000	0,9994	0,9958	0,9827	0,9500	0,7869	0,5000	0,2131	0,0500	0,0042	0,0000
	8	1,0000	1,0000	0,9999	0,9992	0,9958	0,9848	0,9050	0,6964	0,3902	0,1311	0,0181	0,0003
	9	1,0000	1,0000	1,0000	0,9999	0,9992	0,9963	0,9662	0,8491	0,5968	0,2784	0,0611	0,0022
	10	1,0000	1,0000	1,0000	1,0000	0,9999	0,9993	0,9907	0,9408	0,7827	0,4845	0,1642	0,0127
	11	1,0000	1,0000	1,0000	1,0000	1,0000	0,9999	0,9981	0,9824	0,9095	0,7031	0,3518	0,0556
	12	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9997	0,9963	0,9729	0,8732	0,6020	0,1841
	13	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9995	0,9948	0,9647	0,8329	0,4510
	14	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9995	0,9953	0,9648	0,7941
15	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	
20	0	0,3585	0,1216	0,0388	0,0115	0,0032	0,0008	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000
	1	0,7358	0,3917	0,1756	0,0692	0,0243	0,0076	0,0005	0,0000	0,0000	0,0000	0,0000	0,0000
	2	0,9245	0,6769	0,4049	0,2061	0,0913	0,0355	0,0036	0,0002	0,0000	0,0000	0,0000	0,0000
	3	0,9841	0,8670	0,6477	0,4114	0,2252	0,1071	0,0160	0,0013	0,0000	0,0000	0,0000	0,0000
	4	0,9974	0,9568	0,8298	0,6296	0,4148	0,2375	0,0510	0,0059	0,0003	0,0000	0,0000	0,0000
	5	0,9997	0,9887	0,9327	0,8042	0,6172	0,4164	0,1256	0,0207	0,0016	0,0000	0,0000	0,0000
	6	1,0000	0,9976	0,9781	0,9133	0,7858	0,6080	0,2500	0,0577	0,0065	0,0003	0,0000	0,0000
	7	1,0000	0,9996	0,9941	0,9679	0,8982	0,7723	0,4159	0,1316	0,0210	0,0013	0,0000	0,0000
	8	1,0000	0,9999	0,9987	0,9900	0,9591	0,8867	0,5956	0,2517	0,0565	0,0051	0,0001	0,0000
	9	1,0000	1,0000	0,9998	0,9974	0,9861	0,9520	0,7553	0,4119	0,1275	0,0171	0,0006	0,0000
	10	1,0000	1,0000	1,0000	0,9994	0,9961	0,9829	0,8725	0,5881	0,2447	0,0480	0,0026	0,0000
	11	1,0000	1,0000	1,0000	0,9999	0,9991	0,9949	0,9435	0,7483	0,4044	0,1133	0,0100	0,0001
	12	1,0000	1,0000	1,0000	1,0000	0,9998	0,9987	0,9790	0,8684	0,5841	0,2277	0,0321	0,0004
	13	1,0000	1,0000	1,0000	1,0000	1,0000	0,9997	0,9935	0,9423	0,7500	0,3920	0,0867	0,0024
	14	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9984	0,9793	0,8744	0,5836	0,1958	0,0113
	15	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9997	0,9941	0,9490	0,7625	0,3704	0,0432
	16	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9987	0,9840	0,8929	0,5886	0,1330
	17	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9998	0,9964	0,9645	0,7939	0,3231
	18	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9995	0,9924	0,9308	0,6083
19	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	0,9992	0,9885	0,8784	

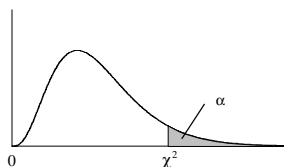
20	1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000
----	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

LIITE 4: Binomijakauman kriittisiä arvoja $n:n$, $\alpha:n$ ja $p:n$ eri arvoille

n	α	p-arvo											
		0,05	0,1	0,15	0,2	0,25	0,3	0,4	0,5	0,6	0,7	0,8	0,9
5	0,9	1	1	2	2	3	3	3	4	4	5	5	5
	0,95	1	2	2	3	3	3	4	4	5	5	5	5
	0,99	2	2	3	3	4	4	5	5	5	5	5	5
	0,999	3	3	4	4	4	5	5	5	5	5	5	5
6	0,9	1	2	2	2	3	3	4	5	5	6	6	6
	0,95	1	2	2	3	3	4	4	5	5	6	6	6
	0,99	2	3	3	4	4	5	5	6	6	6	6	6
	0,999	3	4	4	5	5	5	6	6	6	6	6	6
7	0,9	1	2	2	3	3	4	4	5	6	6	7	7
	0,95	1	2	3	3	4	4	5	6	6	7	7	7
	0,99	2	3	4	4	5	5	6	6	7	7	7	7
	0,999	3	4	5	5	6	6	7	7	7	7	7	7
8	0,9	1	2	3	3	4	4	5	6	7	7	8	8
	0,95	2	2	3	4	4	5	5	6	7	8	8	8
	0,99	2	3	4	5	5	6	6	7	8	8	8	8
	0,999	3	4	5	6	6	7	7	8	8	8	8	8
9	0,9	1	2	3	3	4	4	5	6	7	8	9	9
	0,95	2	3	3	4	4	5	6	7	8	8	9	9
	0,99	2	3	4	5	5	6	7	8	9	9	9	9
	0,999	3	4	5	6	7	7	8	9	9	9	9	9
10	0,9	1	2	3	4	4	5	6	7	8	9	10	10
	0,95	2	3	3	4	5	5	7	8	8	9	10	10
	0,99	3	4	4	5	6	7	8	9	9	10	10	10
	0,999	4	5	6	6	7	8	9	9	10	10	10	10
11	0,9	2	2	3	4	5	5	6	8	9	10	10	11
	0,95	2	3	4	5	5	6	7	8	9	10	11	11
	0,99	3	4	5	6	6	7	8	9	10	11	11	11
	0,999	4	5	6	7	8	8	9	10	11	11	11	11
12	0,9	2	3	3	4	5	6	7	8	9	10	11	12
	0,95	2	3	4	5	6	6	8	9	10	11	12	12
	0,99	3	4	5	6	7	7	9	10	11	12	12	12
	0,999	4	5	6	7	8	9	10	11	12	12	12	12
13	0,9	2	3	4	4	5	6	7	9	10	11	12	13
	0,95	2	3	4	5	6	7	8	9	11	12	13	13
	0,99	3	4	5	6	7	8	9	11	12	12	13	13
	0,999	4	5	7	8	8	9	11	12	13	13	13	13
14	0,9	2	3	4	5	6	6	8	9	11	12	13	14
	0,95	2	3	4	5	6	7	9	10	11	12	13	14
	0,99	3	4	6	7	8	8	10	11	12	13	14	14
	0,999	4	6	7	8	9	10	11	12	13	14	14	14
15	0,9	2	3	4	5	6	7	8	10	11	13	14	15
	0,95	2	4	5	6	7	8	9	11	12	13	14	15
	0,99	3	5	6	7	8	9	10	12	13	14	15	15
	0,999	4	6	7	8	9	10	12	13	14	15	15	15
16	0,9	2	3	4	5	6	7	9	11	12	13	15	16
	0,95	2	4	5	6	7	8	10	11	13	14	15	16
	0,99	3	5	6	7	8	9	11	13	14	15	16	16
	0,999	4	6	8	9	10	11	12	14	15	16	16	16
17	0,9	2	3	4	6	7	8	9	11	13	14	16	17
	0,95	3	4	5	6	7	8	10	12	13	15	16	17
	0,99	3	5	6	8	9	10	12	13	15	16	17	17
	0,999	5	6	8	9	10	11	13	15	16	17	17	17

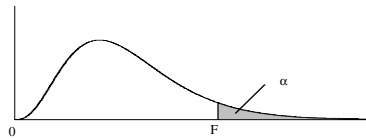
n	α	p-arvo											
		0,05	0,1	0,15	0,2	0,25	0,3	0,4	0,5	0,6	0,7	0,8	0,9
18	0,9	2	3	5	6	7	8	10	12	13	15	16	18
	0,95	3	4	5	7	8	9	11	12	14	16	17	18
	0,99	4	5	7	8	9	10	12	14	15	17	18	18
	0,999	5	7	8	9	11	12	14	15	17	18	18	18
19	0,9	2	4	5	6	7	8	10	12	14	16	17	19
	0,95	3	4	6	7	8	9	11	13	15	16	18	19
	0,99	4	5	7	8	9	11	13	14	16	18	19	19
	0,999	5	7	8	10	11	12	14	16	17	19	19	19
20	0,9	2	4	5	6	8	9	11	13	15	17	18	20
	0,95	3	4	6	7	8	9	12	14	16	17	19	20
	0,99	4	6	7	8	10	11	13	15	17	18	20	20
	0,999	5	7	9	10	11	13	15	17	18	19	20	20
21	0,9	2	4	5	7	8	9	11	13	15	17	19	21
	0,95	3	5	6	7	9	10	12	14	16	18	20	21
	0,99	4	6	7	9	10	11	14	16	18	19	20	21
	0,999	5	7	9	10	12	13	15	17	19	20	21	21
22	0,9	2	4	5	7	8	9	12	14	16	18	20	21
	0,95	3	5	6	8	9	10	13	15	17	19	20	22
	0,99	4	6	8	9	10	12	14	16	18	20	21	22
	0,999	5	7	9	11	12	14	16	18	20	21	22	22
23	0,9	3	4	6	7	8	10	12	15	17	19	21	22
	0,95	3	5	6	8	9	11	13	15	18	20	21	23
	0,99	4	6	8	9	11	12	15	17	19	21	22	23
	0,999	5	8	9	11	13	14	17	19	21	22	23	23
24	0,9	3	4	6	7	9	10	13	15	17	20	22	23
	0,95	3	5	7	8	10	11	14	16	18	20	22	24
	0,99	4	6	8	10	11	13	15	18	20	22	23	24
	0,999	5	8	10	11	13	14	17	19	21	23	24	24
25	0,9	3	4	6	8	9	10	13	16	18	20	22	24
	0,95	3	5	7	8	10	11	14	17	19	21	23	25
	0,99	4	6	8	10	12	13	16	18	20	22	24	25
	0,999	6	8	10	12	13	15	18	20	22	24	25	25
30	0,9	3	5	7	9	11	12	15	19	21	24	27	29
	0,95	4	6	8	10	12	13	16	19	22	25	27	29
	0,99	5	7	9	11	13	15	18	21	24	26	29	30
	0,999	6	9	11	13	15	17	20	23	26	28	30	30
40	0,9	4	6	9	11	14	16	20	24	28	32	35	38
	0,95	4	7	10	12	15	17	21	25	29	33	36	39
	0,99	6	9	12	14	17	19	23	27	31	34	37	40
	0,999	7	11	14	16	19	21	26	30	33	36	39	40
50	0,9	5	8	11	14	16	19	24	30	34	39	44	48
	0,95	5	9	12	15	18	20	26	31	36	40	44	48
	0,99	7	10	14	17	20	23	28	33	38	42	46	49
	0,999	8	13	16	19	23	25	31	36	40	44	48	50
70	0,9	6	10	14	18	22	26	33	40	47	54	60	66
	0,95	7	11	16	20	24	27	35	42	49	55	61	67
	0,99	8	13	18	22	26	30	38	45	51	58	63	68
	0,999	10	16	21	25	29	33	41	48	54	60	65	69
100	0,9	8	14	20	25	31	36	46	56	66	76	85	94
	0,95	9	15	21	27	32	38	48	58	68	77	86	95
	0,99	11	18	24	30	35	41	52	62	71	80	89	96
	0,999	13	20	27	33	39	45	55	65	75	83	91	98

LIITE 5: χ^2 -jakauman kriittisiä arvoja



df	$\alpha \rightarrow$												
	0,995	0,99	0,975	0,95	0,90	0,75	0,50	0,25	0,10	0,05	0,025	0,01	0,005
1	0,00	0,00	0,00	0,00	0,02	0,10	0,45	1,32	2,71	3,84	5,02	6,63	7,88
2	0,01	0,02	0,05	0,10	0,21	0,58	1,39	2,77	4,61	5,99	7,38	9,21	10,60
3	0,07	0,11	0,22	0,35	0,58	1,21	2,37	4,11	6,25	7,81	9,35	11,34	12,84
4	0,21	0,30	0,48	0,71	1,06	1,92	3,36	5,39	7,78	9,49	11,14	13,28	14,86
5	0,41	0,55	0,83	1,15	1,61	2,67	4,35	6,63	9,24	11,07	12,83	15,09	16,75
6	0,68	0,87	1,24	1,64	2,20	3,45	5,35	7,84	10,64	12,59	14,45	16,81	18,55
7	0,99	1,24	1,69	2,17	2,83	4,25	6,35	9,04	12,02	14,07	16,01	18,48	20,28
8	1,34	1,65	2,18	2,73	3,49	5,07	7,34	10,22	13,36	15,51	17,53	20,09	21,95
9	1,73	2,09	2,70	3,33	4,17	5,90	8,34	11,39	14,68	16,92	19,02	21,67	23,59
10	2,16	2,56	3,25	3,94	4,87	6,74	9,34	12,55	15,99	18,31	20,48	23,21	25,19
11	2,60	3,05	3,82	4,57	5,58	7,58	10,34	13,70	17,28	19,68	21,92	24,73	26,76
12	3,07	3,57	4,40	5,23	6,30	8,44	11,34	14,85	18,55	21,03	23,34	26,22	28,30
13	3,57	4,11	5,01	5,89	7,04	9,30	12,34	15,98	19,81	22,36	24,74	27,69	29,82
14	4,07	4,66	5,63	6,57	7,79	10,17	13,34	17,12	21,06	23,68	26,12	29,14	31,32
15	4,60	5,23	6,26	7,26	8,55	11,04	14,34	18,25	22,31	25,00	27,49	30,58	32,80
16	5,14	5,81	6,91	7,96	9,31	11,91	15,34	19,37	23,54	26,30	28,85	32,00	34,27
17	5,70	6,41	7,56	8,67	10,09	12,79	16,34	20,49	24,77	27,59	30,19	33,41	35,72
18	6,26	7,01	8,23	9,39	10,86	13,68	17,34	21,60	25,99	28,87	31,53	34,81	37,16
19	6,84	7,63	8,91	10,12	11,65	14,56	18,34	22,72	27,20	30,14	32,85	36,19	38,58
20	7,43	8,26	9,59	10,85	12,44	15,45	19,34	23,83	28,41	31,41	34,17	37,57	40,00
21	8,03	8,90	10,28	11,59	13,24	16,34	20,34	24,93	29,62	32,67	35,48	38,93	41,40
22	8,64	9,54	10,98	12,34	14,04	17,24	21,34	26,04	30,81	33,92	36,78	40,29	42,80
23	9,26	10,20	11,69	13,09	14,85	18,14	22,34	27,14	32,01	35,17	38,08	41,64	44,18
24	9,89	10,86	12,40	13,85	15,66	19,04	23,34	28,24	33,20	36,42	39,36	42,98	45,56
25	10,52	11,52	13,12	14,61	16,47	19,94	24,34	29,34	34,38	37,65	40,65	44,31	46,93
26	11,16	12,20	13,84	15,38	17,29	20,84	25,34	30,43	35,56	38,89	41,92	45,64	48,29
27	11,81	12,88	14,57	16,15	18,11	21,75	26,34	31,53	36,74	40,11	43,19	46,96	49,65
28	12,46	13,56	15,31	16,93	18,94	22,66	27,34	32,62	37,92	41,34	44,46	48,28	50,99
29	13,12	14,26	16,05	17,71	19,77	23,57	28,34	33,71	39,09	42,56	45,72	49,59	52,34
30	13,79	14,95	16,79	18,49	20,60	24,48	29,34	34,80	40,26	43,77	46,98	50,89	53,67
40	20,71	22,16	24,43	26,51	29,05	33,66	39,34	45,62	51,81	55,76	59,34	63,69	66,77
50	27,99	29,71	32,36	34,76	37,69	42,94	49,33	56,33	63,17	67,50	71,42	76,15	79,49
60	35,53	37,48	40,48	43,19	46,46	52,29	59,33	66,98	74,40	79,08	83,30	88,38	91,95
70	43,28	45,44	48,76	51,74	55,33	61,70	69,33	77,58	85,53	90,53	95,02	100,43	104,21
80	51,17	53,54	57,15	60,39	64,28	71,14	79,33	88,13	96,58	101,88	106,63	112,33	116,32
90	59,20	61,75	65,65	69,13	73,29	80,62	89,33	98,65	107,57	113,15	118,14	124,12	128,30
100	67,33	70,06	74,22	77,93	82,36	90,13	99,33	109,14	118,50	124,34	129,56	135,81	140,17

LIITE 6: F-jakauman kriittisiä arvoja



Taulukko 1: $\alpha = 0,05$

		Osoittajan vapausasteet																
		1	2	3	4	5	6	7	8	9	10	12	15	20	25	30	40	50
Nimitäjän vapausasteet	1	161,4	199,5	215,7	224,6	230,2	234,0	236,8	238,9	240,5	241,9	243,9	245,9	248,0	249,3	250,1	251,1	251,8
	2	18,51	19,00	19,16	19,25	19,30	19,33	19,35	19,37	19,38	19,40	19,41	19,43	19,45	19,46	19,46	19,47	19,48
	3	10,13	9,55	9,28	9,12	9,01	8,94	8,89	8,85	8,81	8,79	8,74	8,70	8,66	8,63	8,62	8,59	8,58
	4	7,71	6,94	6,59	6,39	6,26	6,16	6,09	6,04	6,00	5,96	5,91	5,86	5,80	5,77	5,75	5,72	5,70
	5	6,61	5,79	5,41	5,19	5,05	4,95	4,88	4,82	4,77	4,74	4,68	4,62	4,56	4,52	4,50	4,46	4,44
	6	5,99	5,14	4,76	4,53	4,39	4,28	4,21	4,15	4,10	4,06	4,00	3,94	3,87	3,83	3,81	3,77	3,75
	7	5,59	4,74	4,35	4,12	3,97	3,87	3,79	3,73	3,68	3,64	3,57	3,51	3,44	3,40	3,38	3,34	3,32
	8	5,32	4,46	4,07	3,84	3,69	3,58	3,50	3,44	3,39	3,35	3,28	3,22	3,15	3,11	3,08	3,04	3,02
	9	5,12	4,26	3,86	3,63	3,48	3,37	3,29	3,23	3,18	3,14	3,07	3,01	2,94	2,89	2,86	2,83	2,80
	10	4,96	4,10	3,71	3,48	3,33	3,22	3,14	3,07	3,02	2,98	2,91	2,85	2,77	2,73	2,70	2,66	2,64
	11	4,84	3,98	3,59	3,36	3,20	3,09	3,01	2,95	2,90	2,85	2,79	2,72	2,65	2,60	2,57	2,53	2,51
	12	4,75	3,89	3,49	3,26	3,11	3,00	2,91	2,85	2,80	2,75	2,69	2,62	2,54	2,50	2,47	2,43	2,40
	13	4,67	3,81	3,41	3,18	3,03	2,92	2,83	2,77	2,71	2,67	2,60	2,53	2,46	2,41	2,38	2,34	2,31
	14	4,60	3,74	3,34	3,11	2,96	2,85	2,76	2,70	2,65	2,60	2,53	2,46	2,39	2,34	2,31	2,27	2,24
	15	4,54	3,68	3,29	3,06	2,90	2,79	2,71	2,64	2,59	2,54	2,48	2,40	2,33	2,28	2,25	2,20	2,18
	16	4,49	3,63	3,24	3,01	2,85	2,74	2,66	2,59	2,54	2,49	2,42	2,35	2,28	2,23	2,19	2,15	2,12
	17	4,45	3,59	3,20	2,96	2,81	2,70	2,61	2,55	2,49	2,45	2,38	2,31	2,23	2,18	2,15	2,10	2,08
	18	4,41	3,55	3,16	2,93	2,77	2,66	2,58	2,51	2,46	2,41	2,34	2,27	2,19	2,14	2,11	2,06	2,04
	19	4,38	3,52	3,13	2,90	2,74	2,63	2,54	2,48	2,42	2,38	2,31	2,23	2,16	2,11	2,07	2,03	2,00
	20	4,35	3,49	3,10	2,87	2,71	2,60	2,51	2,45	2,39	2,35	2,28	2,20	2,12	2,07	2,04	1,99	1,97
21	4,32	3,47	3,07	2,84	2,68	2,57	2,49	2,42	2,37	2,32	2,25	2,18	2,10	2,05	2,01	1,96	1,94	
22	4,30	3,44	3,05	2,82	2,66	2,55	2,46	2,40	2,34	2,30	2,23	2,15	2,07	2,02	1,98	1,94	1,91	
23	4,28	3,42	3,03	2,80	2,64	2,53	2,44	2,37	2,32	2,27	2,20	2,13	2,05	2,00	1,96	1,91	1,88	
24	4,26	3,40	3,01	2,78	2,62	2,51	2,42	2,36	2,30	2,25	2,18	2,11	2,03	1,97	1,94	1,89	1,86	
25	4,24	3,39	2,99	2,76	2,60	2,49	2,40	2,34	2,28	2,24	2,16	2,09	2,01	1,96	1,92	1,87	1,84	
26	4,23	3,37	2,98	2,74	2,59	2,47	2,39	2,32	2,27	2,22	2,15	2,07	1,99	1,94	1,90	1,85	1,82	
28	4,20	3,34	2,95	2,71	2,56	2,45	2,36	2,29	2,24	2,19	2,12	2,04	1,96	1,91	1,87	1,82	1,79	
30	4,17	3,32	2,92	2,69	2,53	2,42	2,33	2,27	2,21	2,16	2,09	2,01	1,93	1,88	1,84	1,79	1,76	
40	4,08	3,23	2,84	2,61	2,45	2,34	2,25	2,18	2,12	2,08	2,00	1,92	1,84	1,78	1,74	1,69	1,66	
50	4,03	3,18	2,79	2,56	2,40	2,29	2,20	2,13	2,07	2,03	1,95	1,87	1,78	1,73	1,69	1,63	1,60	
60	4,00	3,15	2,76	2,53	2,37	2,25	2,17	2,10	2,04	1,99	1,92	1,84	1,75	1,69	1,65	1,59	1,56	
120	3,92	3,07	2,68	2,45	2,29	2,18	2,09	2,02	1,96	1,91	1,83	1,75	1,66	1,60	1,55	1,50	1,46	
200	3,89	3,04	2,65	2,42	2,26	2,14	2,06	1,98	1,93	1,88	1,80	1,72	1,62	1,56	1,52	1,46	1,41	
500	3,86	3,01	2,62	2,39	2,23	2,12	2,03	1,96	1,90	1,85	1,77	1,69	1,59	1,53	1,48	1,42	1,38	
750	3,85	3,01	2,62	2,38	2,23	2,11	2,02	1,95	1,89	1,84	1,77	1,68	1,58	1,52	1,47	1,41	1,37	
1000	3,85	3,00	2,61	2,38	2,22	2,11	2,02	1,95	1,89	1,84	1,76	1,68	1,58	1,52	1,47	1,41	1,36	

Taulukko 2: $\alpha = 0,01$

		Osoittajan vapausasteet																
		1	2	3	4	5	6	7	8	9	10	12	15	20	25	30	40	50
Nimitäjän vapausasteet	1	4052	4999	5404	5624	5764	5859	5928	5981	6022	6056	6107	6157	6209	6240	6260	6286	6302
	2	98,50	99,00	99,16	99,25	99,30	99,33	99,36	99,38	99,39	99,40	99,42	99,43	99,45	99,46	99,47	99,48	99,48
	3	34,12	30,82	29,46	28,71	28,24	27,91	27,67	27,49	27,34	27,23	27,05	26,87	26,69	26,58	26,50	26,41	26,35
	4	21,20	18,00	16,69	15,98	15,52	15,21	14,98	14,80	14,66	14,55	14,37	14,20	14,02	13,91	13,84	13,75	13,69
	5	16,26	13,27	12,06	11,39	10,97	10,67	10,46	10,29	10,16	10,05	9,89	9,72	9,55	9,45	9,38	9,29	9,24
	6	13,75	10,92	9,78	9,15	8,75	8,47	8,26	8,10	7,98	7,87	7,72	7,56	7,40	7,30	7,23	7,14	7,09
	7	12,25	9,55	8,45	7,85	7,46	7,19	6,99	6,84	6,72	6,62	6,47	6,31	6,16	6,06	5,99	5,91	5,86
	8	11,26	8,65	7,59	7,01	6,63	6,37	6,18	6,03	5,91	5,81	5,67	5,52	5,36	5,26	5,20	5,12	5,07
	9	10,56	8,02	6,99	6,42	6,06	5,80	5,61	5,47	5,35	5,26	5,11	4,96	4,81	4,71	4,65	4,57	4,52
	10	10,04	7,56	6,55	5,99	5,64	5,39	5,20	5,06	4,94	4,85	4,71	4,56	4,41	4,31	4,25	4,17	4,12
	11	9,65	7,21	6,22	5,67	5,32	5,07	4,89	4,74	4,63	4,54	4,40	4,25	4,10	4,01	3,94	3,86	3,81
	12	9,33	6,93	5,95	5,41	5,06	4,82	4,64	4,50	4,39	4,30	4,16	4,01	3,86	3,76	3,70	3,62	3,57
	13	9,07	6,70	5,74	5,21	4,86	4,62	4,44	4,30	4,19	4,10	3,96	3,82	3,66	3,57	3,51	3,43	3,38
	14	8,86	6,51	5,56	5,04	4,69	4,46	4,28	4,14	4,03	3,94	3,80	3,66	3,51	3,41	3,35	3,27	3,22
	15	8,68	6,36	5,42	4,89	4,56	4,32	4,14	4,00	3,89	3,80	3,67	3,52	3,37	3,28	3,21	3,13	3,08
	16	8,53	6,23	5,29	4,77	4,44	4,20	4,03	3,89	3,78	3,69	3,55	3,41	3,26	3,16	3,10	3,02	2,97
	17	8,40	6,11	5,19	4,67	4,34	4,10	3,93	3,79	3,68	3,59	3,46	3,31	3,16	3,07	3,00	2,92	2,87
	18	8,29	6,01	5,09	4,58	4,25	4,01	3,84	3,71	3,60	3,51	3,37	3,23	3,08	2,98	2,92	2,84	2,78
	19	8,18	5,93	5,01	4,50	4,17	3,94	3,77	3,63	3,52	3,43	3,30	3,15	3,00	2,91	2,84	2,76	2,71
	20	8,10	5,85	4,94	4,43	4,10	3,87	3,70	3,56	3,46	3,37	3,23	3,09	2,94	2,84	2,78	2,69	2,64
21	8,02	5,78	4,87	4,37	4,04	3,81	3,64	3,51	3,40	3,31	3,17	3,03	2,88	2,79	2,72	2,64	2,58	
22	7,95	5,72	4,82	4,31	3,99	3,76	3,59	3,45	3,35	3,26	3,12	2,98	2,83	2,73	2,67	2,58	2,53	
23	7,88	5,66	4,76	4,26	3,94	3,71	3,54	3,41	3,30	3,21	3,07	2,93	2,78	2,69	2,62	2,54	2,48	
24	7,82	5,61	4,72	4,22	3,90	3,67	3,50	3,36	3,26	3,17	3,03	2,89	2,74	2,64	2,58	2,49	2,44	
25	7,77	5,57	4,68	4,18	3,85	3,63	3,46	3,32	3,22	3,13	2,99	2,85	2,70	2,60	2,54	2,45	2,40	
26	7,72	5,53	4,64	4,14	3,82	3,59	3,42	3,29	3,18	3,09	2,96	2,81	2,66	2,57	2,50	2,42	2,36	
28	7,64	5,45	4,57	4,07	3,75	3,53	3,36	3,23	3,12	3,03	2,90	2,75	2,60	2,51	2,44	2,35	2,30	
30	7,56	5,39	4,51	4,02	3,70	3,47	3,30	3,17	3,07	2,98	2,84	2,70	2,55	2,45	2,39	2,30	2,25	
40	7,31	5,18	4,31	3,83	3,51	3,29	3,12	2,99	2,89	2,80	2,66	2,52	2,37	2,27	2,20	2,11	2,06	
50	7,17	5,06	4,20	3,72	3,41	3,19	3,02	2,89	2,78	2,70	2,56	2,42	2,27	2,17	2,10	2,01	1,95	
60	7,08	4,98	4,13	3,65	3,34	3,12	2,95	2,82	2,72	2,63	2,50	2,35	2,20	2,10	2,03	1,94	1,88	
120	6,85	4,79	3,95	3,48	3,17	2,96	2,79	2,66	2,56	2,47	2,34	2,19	2,03	1,93	1,86	1,76	1,70	
200	6,76	4,71	3,88	3,41	3,11	2,89	2,73	2,60	2,50	2,41	2,27	2,13	1,97	1,87	1,79	1,69	1,63	
500	6,69	4,65	3,82	3,36	3,05	2,84	2,68	2,55	2,44	2,36	2,22	2,07	1,92	1,81	1,74	1,63	1,57	
750	6,67	4,63	3,81	3,34	3,04	2,83	2,66	2,53	2,43	2,34	2,21	2,06	1,90	1,80	1,72	1,62	1,55	
1000	6,66	4,63	3,80	3,34	3,04	2,82	2,66	2,53	2,43	2,34	2,20	2,06	1,90	1,79	1,72	1,61	1,54	

LIITE 7: Pearsonin korrelaatiokertoimen kriittisten arvojen itseisarvoja

$$|r_{\alpha/2}| = \frac{t_{\alpha/2}}{\sqrt{n-2 + (t_{\alpha/2})^2}}$$

n	$\alpha=0,05$ (5%)		$\alpha=0,01$ (1%)		$\alpha=0,001$ (0,1%)	
	1-suunt.	2-suunt.	1-suunt.	2-suunt.	1-suunt.	2-suunt.
4	0.900	0.950	0.980	0.990	0.998	0.999
5	0.805	0.878	0.934	0.959	0.986	0.991
6	0.729	0.811	0.882	0.917	0.963	0.974
7	0.669	0.754	0.833	0.875	0.935	0.951
8	0.621	0.707	0.789	0.834	0.905	0.925
9	0.582	0.666	0.750	0.798	0.875	0.898
10	0.549	0.632	0.715	0.765	0.847	0.872
11	0.521	0.602	0.685	0.735	0.820	0.847
12	0.497	0.576	0.658	0.708	0.795	0.823
13	0.476	0.553	0.634	0.684	0.772	0.801
14	0.458	0.532	0.612	0.661	0.750	0.780
15	0.441	0.514	0.592	0.641	0.730	0.760
16	0.426	0.497	0.574	0.623	0.711	0.742
17	0.412	0.482	0.558	0.606	0.694	0.725
18	0.400	0.468	0.543	0.590	0.678	0.708
19	0.389	0.456	0.529	0.575	0.662	0.693
20	0.378	0.444	0.516	0.561	0.648	0.679
21	0.369	0.433	0.503	0.549	0.635	0.665
22	0.360	0.423	0.492	0.537	0.622	0.652
23	0.352	0.413	0.482	0.526	0.610	0.640
24	0.344	0.404	0.472	0.515	0.599	0.629
25	0.337	0.396	0.462	0.505	0.588	0.618
26	0.330	0.388	0.453	0.496	0.578	0.607
27	0.323	0.381	0.445	0.487	0.568	0.597
28	0.317	0.374	0.437	0.479	0.559	0.588
29	0.311	0.367	0.430	0.471	0.550	0.579
30	0.306	0.361	0.423	0.463	0.541	0.570
31	0.301	0.355	0.416	0.456	0.533	0.562
32	0.296	0.349	0.409	0.449	0.526	0.554
33	0.291	0.344	0.403	0.442	0.518	0.547
34	0.287	0.339	0.397	0.436	0.511	0.539
35	0.283	0.334	0.392	0.430	0.504	0.532
36	0.279	0.329	0.386	0.424	0.498	0.525
37	0.275	0.325	0.381	0.418	0.492	0.519
38	0.271	0.320	0.376	0.413	0.486	0.513
39	0.267	0.316	0.371	0.408	0.480	0.507
40	0.264	0.312	0.367	0.403	0.474	0.501
50	0.235	0.279	0.328	0.361	0.427	0.451
60	0.214	0.254	0.300	0.330	0.391	0.414
70	0.198	0.235	0.278	0.306	0.363	0.385
80	0.185	0.220	0.260	0.286	0.340	0.361
90	0.174	0.207	0.245	0.270	0.322	0.341
100	0.165	0.197	0.232	0.256	0.305	0.324
150	0.135	0.160	0.190	0.210	0.250	0.266
200	0.117	0.139	0.164	0.182	0.217	0.231
250	0.104	0.124	0.147	0.163	0.195	0.207
300	0.095	0.113	0.134	0.149	0.178	0.189
400	0.082	0.098	0.116	0.129	0.154	0.164
500	0.074	0.088	0.104	0.115	0.138	0.147
600	0.067	0.080	0.095	0.105	0.126	0.134
700	0.062	0.074	0.088	0.097	0.117	0.124
800	0.058	0.069	0.082	0.091	0.109	0.116
900	0.055	0.065	0.078	0.086	0.103	0.110
1000	0,052	0,062	0,074	0,081	0,098	0,104

LIITE 8: Asteikkovälin lineaarinen muunnos

Empiirisessä tutkimuksessa tulee silloin tällöin vastaan tilanne, jossa tarvitsisi kyetä muuntamaan tietyllä asteikkovälillä (esim. 1-5) mitatun muuttujan arvot vastaamaan jollakin toisella asteikkovälillä mitatun muuttujan arvoja (kuten esim. 1-7). Toisin sanoen haluttaisiin saada aikaan asteikkovälin lineaarinen muunnos. Samalla voitaisiin haluta esim. tietää, mikä muuttujan keskiarvo olisi, jos se olisi mitattu tällä toisella asteikkovälillä.

Yleisesti, kun tunnetaan lähtöasteikon X minimi ja maksimi sekä tavoiteasteikon Y minimi ja maksimi, lineaariset asteikkovälien muunnokset suoritetaan ratkaisemalla yhtälöparista vakiot a ja b , jolloin löydetään suoran yhtälölle sellainen kulmakerroin b ja sellainen vakio a , jolloin ko. suora kulkee sekä minimien että maksimien muodostamien pisteiden kautta xy -koordinaatistossa.

Eli esimerkiksi, jos lähtöasteikon minimi on 1 ja maksimi 6 sekä tavoiteasteikon minimi 1 ja maksimi 7, lähtee matemaattinen ratkaisu liikkeelle seuraavasti

$$\begin{array}{l} (Y) \quad (X) \\ \left\{ \begin{array}{l} Y = a + bX \\ Y = a + bX \end{array} \right\} \Rightarrow \begin{array}{l} (\min) 1 = a + b1 \\ (\max) 7 = a + b6 \end{array} \Rightarrow \end{array}$$

$a = 1 - b1$ ja sijoittamalla a :n lauseke toiseen yhtälöön saadaan

$$b = \frac{6}{5}, \text{ jolloin } a = 1 - \frac{6}{5}.$$

Esimerkiksi, kun X :n arvot ovat 6, 3 ja 1, saadaan suoran yhtälöstä $Y = a + bX$ ko. vakion arvoilla lineaarisesti muunnetut Y :n arvot

$$-\frac{1}{5} + \frac{6}{5} \times 6 = 7; \quad -\frac{1}{5} + \frac{6}{5} \times 3 = 3,4; \quad -\frac{1}{5} + \frac{6}{5} \times 1 = 1.$$

Kun asteikoiden ääripisteet tunnetaan, voidaan yhtälöparin ratkaisusta johtaa yleinen kaava lineaarisille asteikkovälin muunnoksille $X \rightarrow Y$, joka on seuraavanlainen

$$Y = Y_{\min} - \left[\frac{Y_{\max} - Y_{\min}}{X_{\max} - X_{\min}} \right] \cdot X_{\min} + \left[\frac{Y_{\max} - Y_{\min}}{X_{\max} - X_{\min}} \right] \cdot X.$$

Sisältö

1	JOHDANTO	1
1.1	Tilastomenetelmät luotettavan tutkimuksen perustana	1
1.1.1	Otos vs. näyte	1
1.1.2	Tilastollinen päättely ja tieteellisyden kriteerit	2
2	TEOREETTINEN JAKAUMA	3
2.1	Satunnaismuuttuja	3
2.1.1	Merkintätavat	3
2.2	Tilastollinen päättely ja teoreettinen jakauma	3
2.3	Empiirinen jakauma vs. teoreettinen jakauma	4
2.4	Teoreettisen jakauman peruskäsitteet	6
2.4.1	Todennäköisyys (p)	6
2.4.1.1	Satunnaismalli vs. deterministinen malli	6
2.4.1.2	Klassinen todennäköisyys	7
2.4.2	Tiheysfunktio $f(X)$	7
2.4.2.1	Diskreetti vs. jatkuva jakauma	8
2.4.2.2	Tasajakauma	8
2.4.2.3	Jatkuvan tasajakauman tiheysfunktio	9
2.4.2.4	Todennäköisyyksien määrääminen tiheysfunktion avulla	10
2.4.2.5	Tasajakauman odotusarvo ja varianssi	11
2.4.3	Kertymäfunktio $F(x)$	11
2.4.3.1	Tasajakauman kertymäfunktio	13
2.4.3.2	Kertymäfunktion ominaisuuksia	13
2.4.4	Diskreetti tasajakauma	15
2.5	Teoreettiset tunnusluvut	16
2.5.1	Odotusarvo ja varianssi	16
2.5.2	Muita tunnuslukuja	17
2.6	Normaalijakauma	20
2.6.1	Keskeinen raja-arvolause (central limit theorem)	20
2.6.2	Muuttujatyypit jotka noudattavat likimäärin normaalijakaumaa	20
2.6.2.1	Virhemuuttujat	20
2.6.2.2	Ominaisuudet	21
2.6.2.3	Otoskeskiarvon otantajakauma	21
2.6.3	Normaalijakauman tiheysfunktio	21
2.6.4	Normaalijakauman ominaisuuksista	23

2.6.5	Parametrien μ ja σ^2 vaikutus tiheysfunktioon	24
2.6.6	Taulukon korvaaminen Excelin funktioilla	24
2.6.7	Todennäköisyyksien määrittäminen (Z-jakaumassa)	25
3	LINEAARINEN MUUNNOS	27
3.1	Yleistä	27
3.1.1	Lineaarisen muunnoksen vaikutus muuttujan keskiarvoon ja varianssiin	27
3.2	Muuttujan standardoiminen	28
3.3	Lineaarisen muunnoksen vaikutus muuttujan fraktiilin määrittämiseen	30
3.3.1	Lineaarisen muunnoksen vaikutus muuttujan todennäköisyyksien määrittämiseen	30
4	OTANTAJAKAUMA	32
4.1	Yleistä	32
4.2	Otoskeskiarvon otantajakauma	32
4.3	Prosenttiluvun otantajakauma	34
4.4	Binomijakauma	36
4.4.1	Binomijakauma ja otosprosenttiluvun otantajakauma	37
4.4.2	Binomijakauman pistetodennäköisyysfunktio	37
4.4.3	Binomijakauman muodon riippuvuus parametreista n ja p	39
4.4.4	Binomijakauman kertymäfunktio	39
4.4.5	Binomijakauman odotusarvo ja varianssi	40
4.4.6	Binomi- ja normaalijakauman välinen yhteys	41
4.4.7	Excelin funktiot, jotka liittyvät binomijakaumaan	41
4.5	t-jakauma	42
4.5.1	t-jakauman yhteydet muihin jakaumiin	42
4.5.2	t-jakauman fraktiilit	43
4.5.3	t-jakauman tiheysfunktio	44
4.5.4	Excelin t-jakaumaan liittyvät funktiot	44
4.6	χ^2-jakauma — eli otosvarianssin s^2 otantajakauma	45
4.6.1	χ^2 -jakauman tiheysfunktio	45
4.6.2	χ^2 -jakauman odotusarvo ja varianssi	46
4.6.3	χ^2 -jakauman kriittiset arvot	47
4.6.4	Otosvarianssin s^2 ja χ^2 -jakauman välinen yhteys	48
4.6.5	Excelin funktiot, jotka liittyvät χ^2 -jakaumaan	49

4.7	F-jakauma — eli varianssitestien otantajakauma	50
4.7.1	F-jakauman kriittiset arvot	51
4.7.2	Excelin funktiot, jotka liittyvät F-jakaumaan	53
5	ESTIMOINTI	54
5.1	Johdanto	54
5.2	Piste-estimointi	54
5.2.1	Estimaattorin ”hyvyys”	54
5.3	Väliestimointi	56
6	LUOTTAMUSVÄLIEN LASKEMINEN	57
6.1	Luottamusvälin määrittäminen populaation keskiarvolle μ , kun σ^2 tunnetaan	57
6.2	Populaation keskiarvon μ luottamusvälin määrittäminen, kun σ^2 on tuntematon	59
6.3	Luottamusvälin määrittäminen populaation prosenttiluvulle θ	61
6.4	Luottamusvälin määrittäminen populaation varianssille σ^2	62
6.5	Luottamusvälin määrittäminen populaation korrelaatiokertoimelle ρ	63
7	TILASTOLLINEN HYPOTEESIN TESTAUS	66
7.1	Tilastollisen testauksen lähtökohta	66
7.2	Hypoteesit	66
7.2.1	Esimerkkejä tilastollisista hypoteeseista	67
7.3	Tilastollisen testin mahdolliset virheet	68
7.4	Testisuure	69
7.5	Kriittinen alue	70
7.6	Esimerkki tilastollisesta testauksesta	71
7.7	Tilastollinen testaus käytännössä	72
7.8	Testin oletukset eli vaatimukset	73

8	TILASTOLLISET TESTIT	73
8.1	Testien ryhmittely	73
8.2	Tixelin testipohja	76
8.3	Prosenttilukutestit	77
8.3.1	Yhden otoksen prosenttilukutesti	77
8.3.1.1	Pieni otoskoko	79
8.3.2	Kahden riippumattoman otoksen prosenttilukutesti	80
8.3.3	Kahden riippuvan otoksen prosenttilukutesti	81
8.4	Keskiarvotestit	84
8.4.1	Yhden otoksen keskiarvotesti	84
8.4.2	Kahden riippumattoman otoksen keskiarvotesti	86
8.4.3	Perusjoukon varianssit tunnetaan	87
8.5	Kahden riippuvan otoksen keskiarvotesti	88
8.5.1	Erotusmuuttujan keskihajonta tunnuslukujen avulla	89
8.5.2	Excelin keskiarvotestit	90
8.5.3	Yksisuuntainen varianssianalyysi	91
8.6	Riippuvuuslukutestit	91
8.6.1	Korrelaatiokerroin	91
8.6.2	Korrelaatiokertoimen kriittiset arvot	93
8.6.3	Korrelaatiokerroin nollihypoteesin mukaan eri suuri kuin nolla	94
8.6.4	Kahden riippumattoman korrelaatiokertoimen yhtäsuuruuden testaus	95
8.6.5	χ^2 -riippumattomuustesti	96
8.6.5.1	Ristiintaulukon riippuvuusluvut	99
8.7	χ^2 -yhteensopivuustesti	100
8.8	Kahden varianssin yhtäsuuruuden testi	101
9	TESTIN VOIMAKKUUSFUNKTIO (POWER)	105
9.1	Hyväksymisvirheen todennäköisyys β	105
9.2	Testin voimakkuuteen vaikuttavia tekijöitä	108
9.2.1	α :n vaikutus testin voimakkuuteen	108
9.2.2	μ :n arvo	109
9.2.3	Otoskoko (n)	110
9.2.4	σ^2 :n arvo	111
9.2.5	Käytetty testityyppi	112
9.3	Testin voimakkuusfunktion arvon laskeminen tietylle vaihtoehdoiselle hypoteesille	112

10	OTOSKOON MÄÄRÄÄMINEN	114
10.1	Otoskoon määrääminen testin voimakkuusfunktioita ($1-\beta$) käyttäen	114
10.1.1	Keskiarvo	115
10.1.2	Prosenttiluku	117
10.2	Otoskoon määrääminen luottamusvälin avulla	118
10.2.1	Keskiarvo	118
10.2.2	Prosenttiluku	119
10.2.3	Korrelaatiokerroin	121
11	OTANTAMENETELMÄT	123
11.1	Otantamenetelmien käytön perusta	123
11.2	Otantamenetelmät	124
11.3	Yksinkertainen satunnaisotanta	124
11.4	Systemaattinen otanta	126
11.5	Estimointi yksinkertaisessa satunnaisotannassa	127
11.5.1	Luottamusväli	128
11.6	Ositettu otanta	128
11.6.1	Estimointi ositetussa otannassa	129
11.6.2	Otoksen kiintiöinti	130
11.7	Ryväsotanta	132
11.8	Otoskoon määrääminen äärellisen perusjoukon tapauksessa	132
12	LIITTEET	134

LIITE 1: Standardoidun normaalijakauman kertymäfunktion arvoja

LIITE 2: t-jakauman kriittisiä arvoja

LIITE 3: Binomijakauman kertymäfunktion arvoja

LIITE 4: Binomijakauman kriittisiä arvoja

LIITE 5: χ^2 -jakauman kriittisiä arvoja

LIITE 6: F-jakauman kriittisiä arvoja

LIITE 7: Pearsonin korrelaatiokertoimen kriittisten arvojen itseisarvoja

LIITE 8: Asteikkovälin lineaarinen muunnos

Aihehakemisto

B

Bernoullin koe · 40
 binomijakauma · 40
 diskreetti · 40
 Excelin funktiot · 45
 kertymäfunktio · 44
 kriittinen arvo · 44
 odotusarvo · 45
 otantajakauma · 41
 parametrit · 40, 43
 pistetodennäköisyysfunktio · 41, 42
 todennäköisyys · 41
 todennäköisyysjakauma · 40
 varianssi · 45
 yhteys normaalijakaumaan · 41
 yhteys otosprosenttilukuun · 41

C

Cramerin phi · 109

D

deterministinen malli · 7
 diskreetti · 9, 16

E

empiirinen jakauma · 4, 7
 kertymäfunktio · 14
 empiirinen muuttuja · 4
 empiiriset (otos-)tunnusluvut · 18
 epäjatkuva jakauma · 9
 estimaatti · 35, 73
 estimaattori · 3, 35
 satunnaiskäyttäytyminen · 4
 estimointi · 1, 35, 60
 estimaatti · 60

parametri · 60
 piste-estimointi · 60
 estimaattori · 60, 61
 harhattomuus · 60, 61
 keskivirhe · 61
 tarkkuus · 61
 satunnaiskäyttäytyminen · 60
 väliestimointi · 62
 luottamustaso · 62
 luottamusvälin ylä- ja alaraja · 62
 satunnaisväli · 62

F

Fisherin transformaatio · 71, 104, 105, 131
 F-jakauma · 55
 Excelin-funktiot · 58
 kertymäfunktio · 58
 kriittinen arvo · 56, 57
 otantajakauma · 55
 satunnaismuuttuja · 55
 tiheysfunktio · 56
 todennäköisyys · 58
 vapausasteet · 55
 yhteys χ^2 -jakaumaan · 55
 fraktiili · 20, 33
 frekvenssijakauma · 7, 8

H

havaittu frekvenssi · 107
 huipukkuusluvut · 21
 hypoteesi · 85, 93, 94, 112
 hypoteesin testaus · 1, 73
 estimaatti · 73
 johtopäätös · 74
 nollahypoteesi · 73
 otos · 73
 parametri · 73, 74
 riskitodennäköisyys · 74
 testisuure · 74
 tilastollinen hypoteesi · 75
 tutkimushypoteesi · 73
 vaihtoehtoinen hypoteesi · 73

yksi- tai kaksisuuntainen · 74, 79

K

kertymäfunktio · 7, 13, 16
 ominaisuudet · 15
 porrasfunktio · 17
 todennäköisyys · 16
 keskeinen raja-arvolause · 22, 38, 68
 keskiarvotestit · 92
 erotusmuuttuja · 97
 Excelissä · 99
 hypoteesi · 93, 94
 t-testi · 93
 keskivirhe · 2
 khi^2 -riippumattomuustesti · 106
 khi^2 -riippumattomuustesti
 kriittinen alue · 109
 testisuure · 108, 109
 vapausasteet · 109
 khi^2 -yhteensopivuustesti · 110
 havaittu frekvenssi · 110
 odotettu frekvenssi · 110
 testisuure · 110
 vapausasteet · 111
 khi^2 -jakauma
 Excelin funktiot · 53
 kriittinen arvo · 51
 odotusarvo · 50
 otantajakauma · 49
 otosvarienssi · 49
 tiheysfunktio · 49
 vapausasteet · 49
 varienssi · 50
 khi^2 -jakauma · 49, 53
 yhteys normaalijakaumaan · 50
 yhteys otosvarienssiin · 52
 Klassinen todennäköisyys · 8
 suhteellinen frekvenssi · 8
 kontingenssikerroin · 109
 korrelaatiokerroin · 100
 H_0 :n mukaan eri suuri kuin nolla · 104
 kaksi riippumatonta otosta · 105
 kriittinen arvo · 102

merkitsevyytestaus · 101
 nollahypoteesi · 101
 otoskorrelaatiokerroin · 100
 Pearsonin tulomomentti-
 korrelaatiokerroin · 101
 perusjoukon
 korrelaatiokerroin · 101
 kovarianssi · 100
 kriittinen alue · 78
 kumulatiivinen jakauma · 7

L

lineaarinen muunnos · 30
 vaikutus fraktiilin
 määräämiseen · 33
 vaikutus keskiarvoon ja
 varianssiin · 30
 vaikutus todennäköisyyksien
 määräämiseen · 33
 luokkaväli · 8
 luottamusvälien laskeminen · 63
 populaation keskiarvo · 63
 populaation
 korrelaatiokerroin · 70
 populaation prosenttiluku ·
 68
 populaation varianssi · 69

M

mallintaminen · 5, 6
 matemaattinen malli · 4, 5, 41
 matemaattinen muuttuja · 3
 McNemarin testi · 90
 mediaani · 19
 merkitsevyytaso · 77
 mitta-asteikko · 3
 mittauksen virhelähteet · 22
 mittausvirhe · 4
 moodi · 19
 muuttujan standardoiminen · 31

N

non-parametriset testit · 3, 82,
 122
 normaalijakauma · 5, 22
 Excelin funktiot · 25, 27
 keskeinen raja-arvolause · 22

keskiarvo · 24
 kriittinen arvo · 64
 luottamusrajat · 25
 parametrit · 26, 27
 standardoitu · 24, 33, 64, 79
 symmetrisyys · 23, 25, 28
 tiheysfunktio · 23, 24, 25
 varianssi · 22, 24
 näyte · 1

O

odotettu frekvenssi · 107
 otantajakauma · 35, 81
 otoskeskiarvo · 23, 35, 36
 otoskoko · 36
 otoskorrelaatiokerroin · 71
 otosprosenttiluku · 38, 39, 68
 otosvarienssi · 53
 satunnaiskäyttäytyminen · 35
 varianssi · 36
 otantajakauma · 77
 otantamenetelmät · 133
 estimaatti · 133
 estimointi · 133
 kato · 134
 mittausvirhe · 134
 todennäköisyys · 134
 otantavirheen
 satunnaiskäyttäytyminen ·
 133
 kokonaistutkimus · 133
 luottamusväli · 138
 ositettu otanta · 138
 kiintiöinti · 140
 ositepaino · 139
 suhteellinen kiintiöinti ·
 139
 otantasuhde · 137
 otantatutkimus · 133
 otantavirhe · 133
 otos · 133
 otoskoon määrääminen
 äärellisen perusjoukon
 tapauksessa · 143
 ryväsotanta · 143
 satunnaisotanta · 134
 systemaattinen otanta · 136
 äärellinen perusjoukko · 133
 äärellisen perusjoukon
 korjaustermi · 133, 137
 otantavirhe · 2, 35, 73
 otos · 1

Otoskoon määrääminen · 124
 luottamusväli
 keskiarvo · 128
 korrelaatiokerroin · 131
 prosenttiluku · 129
 testin voimakkuusfunktio
 keskiarvo · 125
 prosenttiluku · 127
 otosprosenttiluku · 38
 keskivirhe · 40
 odotusarvo · 39
 varianssi · 39
 otostunnusluku · 2, 3, 35

P

parametri · 73
 parametriset testit · 3, 4, 82, 122
 perusjoukon tunnusluku · 35
 pistetodennäköisyys · 16
 pistetodennäköisyysfunktio · 9,
 17
 populaation prosenttiluku · 38
 populaation tunnusluku · 2
 power · 115
 prosenttifrekvenssi · 10
 prosenttilukutestit · 85
 binomijakauma · 87
 hypoteesi · 85
 McNemarin testi · 90
 pieni otoskoko · 87

R

riippuvuuslukutestit
 χ^2 -riippumattomuustesti
 hypoteesit · 108
 oletukset · 108
 Cramerin phi · 100
 havaittu frekvenssi · 107
 kontingenssikerroin · 100
 korrelaatiokerroin · 100
 mitta-asteikko · 106
 odotettu frekvenssi · 107
 ristiintaulukko · 106, 107
 riippuvuuslukutestit · 100
 riskitodennäköisyys · 78
 riskitaso · 77, 84
 ristiintaulukko · 107
 riippuvuusluvut · 109

S

satunnaisilmiö · 3
 satunnaiskäyttäytyminen · 6
 satunnaismalli · 7
 stokastinen · 7
 satunnaismuuttuja
 varianssi · 18
 satunnaismuuttuja · 3, 9, 10, 11, 60
 diskreetti · 19
 odotusarvo · 18
 satunnaisotos · 35
 satunnaisvaihtelu · 4
 standardoiminen · 31
 suhteellinen frekvenssi · 7

T

tasajakauma
 tiheysfunktio · 10
 tasajakauma · 10
 diskreetti · 17
 jatkuva · 15
 kertymäfunktio · 15
 odotusarvo · 12, 13, 19
 pistetodennäköisyysfunktio · 17
 tiheysfunktio · 11
 todennäköisyys · 11
 todennäköisyysjakauma · 11
 varianssi · 12, 13, 19, 118
 tasajakuma
 kertymäfunktio · 15
 teoreettinen jakauma
 odotusarvo · 19
 teoreettinen jakauma · 3, 4, 5, 7
 kertymäfunktio · 14
 normaalijakauma · 22
 tasajakauma · 9, 10, 11
 todennäköisyysjakauma · 9
 todennäköisyysmassa · 14
 varianssi · 19
 teoreettinen tunnusluku · 3
 teoreettiset tunnusluvut · 18
 testin voimakkuusfunktio · 115

1. ja 2. lajin virheen sekä powerin suhde toisiinsa · 116
 2. lajin virheen todennäköisyys · 115
 2.lajin virhe · 115
 estimaatti · 124
 hyväksymisvirhe · 115
 hyväksymisvirheen todennäköisyys · 115
 keskivirhe · 120
 kriittinen arvo · 115, 122
 laskeminen
 kriittinen arvo · 122
 laskeminen 2. lajin virheen todennäköisyys · 123
 laskeminen tietyllä H_1 lle · 122
 nollahypoteesi · 115
 vaihtoehtoinen hypoteesi · 115
 vaikuttavat tekijät
 käytetty testityyppi · 122
 otoskoko · 120
 α :n ja β n suhde · 119
 α :n vaikutus · 118
 μ :n arvo · 119
 σ^2 n arvo · 121
 vaikuttavat tekijät · 118
 yksi- tai kaksisuuntaisen H_1 n tapauksessa · 121
 testisuure · 5, 6, 74, 77, 78, 80, 84, 86
 tiheysfunktio · 7, 8, 9, 16
 tilastollinen päättely · 2
 tilastollinen testaus
 keskiarvotestit · 92
 oletukset · 93, 94
 keskivirhe · 78
 non-parametriset testit · 82
 parametriset testit · 82
 tilastollinen testaus · 73
 kriittinen alue · 78
 merkitsevyystaso · 77
 nollahypoteesi · 76, 83
 oletukset · 81, 83
 prosenttilukutestit · 85
 oletukset · 85
 päätöksenteko · 76
 riskitaso · 77
 riskitodennäköisyys · 77, 81
 testisuure · 83
 testityypit · 81
 ennen-jälkeen asetelma · 82
 eri muuttujien vertailu · 82
 riippuvat vs. riippumattomat otokset · 82
 vastinparimenettely · 82
 vaihtoehtoinen hypoteesi · 76, 83
 virhetyypit · 76
 1.lajin virhe · 76
 2.lajin virhe · 76
 tilastoyksikkö · 4, 8
 t-jakauma
 satunnaismuuttuja · 46
 t-jakauma · 46
 Excelin funktiot · 48
 fraktiili · 47
 otantajakauma · 48
 tiheysfunktio · 48
 vapausasteet · 48
 varianssi · 47
 yhteys normaalijakaumaan · 46
 todellinen arvo · 23
 todennäköisyys · 5, 7, 13, 14, 17, 118
 todennäköisyysjakauma · 3, 5
 todennäköisyysjakaumaa · 6
 t-testi · 93
 tunnusluvut · 18

V

varianssi · 5
 varianssien yhtäsuuruustesti · 112
 hypoteesi · 112
 oletukset · 112
 testisuure · 113
 vapausasteet · 113
 vinousluvut · 21
 virhemuuttuja · 22