

# Sisällysluettelo

ALKUSANAT .....	4
ALKUSANAT E-KIRJA –VERSIOON .....	5
SISÄLLYSLUETTELO .....	6
LYHYT SANASTO VASTA-ALKAJILLE .....	7
<b>1. MONIMUUTTUJAMENETELMÄT IHMISTIEETEISSÄ</b> .....	<b>9</b>
1.1 MONIMUUTTUJA-AINEISTON ERITYISPIIRTEITÄ .....	10
1.2 AINEISTON ALUSTAVA TARKASTELU .....	13
1.2.1 Korrelaatio ja käyräviivainen yhteys .....	13
1.2.2 Outlierit .....	14
1.2.3 Normaalisuus .....	16
1.2.4 Multikollinearisuus ja singulaarisuus .....	17
1.3 KIRJAN RAKENTEESTA .....	18
<b>2. REGRESSIOANALYYSI</b> .....	<b>20</b>
2.1 PERINTEINEN REGRESSIOANALYYSI (RA) .....	21
2.1.1 Missä tilanteessa toimii parhaiten .....	21
2.1.2 Rajoitukset ja oletukset .....	22
2.1.3 Lyhyesti teoriasta ja käsitteistä .....	22
2.1.4 Lisätestit ja jatkoanalyysit .....	29
2.1.5 Tekninen suoritus SPSS-ohjelmistolla ja tulkinta .....	32
2.2 LOGISTINEN REGRESSIOANALYYSI (LRA) .....	40
2.2.1 Missä tilanteessa toimii parhaiten .....	40
2.2.2 Rajoitukset ja oletukset .....	41
2.2.3 Lyhyesti teoriasta ja käsitteistä .....	42
2.2.4 Lisätestit ja jatkoanalyysit .....	47
2.2.5 Tekninen suoritus SPSS-ohjelmistolla ja tulkinta .....	48
2.3 KANONINEN KORRELAATIO (CC) .....	56
2.3.1 Missä tilanteessa toimii parhaiten .....	56
2.3.2 Rajoitukset ja oletukset .....	56
2.3.3 Lyhyesti teoriasta ja käsitteistä .....	57
2.3.4 Lisätestit ja jatkoanalyysit .....	60
2.3.5 Tekninen suoritus SPSS-ympäristössä ja tulkinta .....	60
<b>3. LOPUKSI</b> .....	<b>67</b>
<b>LIITE A. AINEISTOSSA KÄYTETYT ALKUPERÄISET MUUTTUJAT</b> .....	<b>68</b>
<b>LIITE B. MATRIISILASKENNASTA KEVYESTI</b> .....	<b>60</b>
<b>LÄHTEET</b> .....	<b>63</b>
<b>ASIA- JA HENKILÖHAKEMISTO</b> .....	<b>65</b>

Joissain tapauksissa havaitsemme, että residuaalit eivät olekaan homoskedastisia selitettävän muuttujan suhteen. Huomaamme, että myös **selitettävää muuttujaa saatetaan joutua korjaamaan ja muuntamaan**. Mainitut nyrkkisäännönmaiset muunnokset soveltuvat myös tähän tapaukseen. Mainittujen muunnosten lisäksi on olemassa myös sofistikoitumpi tapa etsiä oikeaa muunnosta. Box-Cox -muunnoksella voidaan **etsiä sitä Y-muuttujan potenssia, joka toisi parhaan mahdollisen korjauksen heteroskedastiseen tilanteeseen**. Mikäli nimittäin Y:n potenssi on 1, ei muunnosta tarvita, mutta jos Y:n potenssi olisi esimerkiksi  $\frac{1}{2}$ , tarvittaisiin neliöjuurimuunnos. Box ja Cox esittivät (1964), että parasta potentiaalista potenssia voidaan etsiä mallittamalla Y:n potenssin  $Y^{(1-b)}$  parametri  $b$  (*Slope*). Jo aiemmin mainittuja muunnoksia vastaavat eri  $b$ :n arvot seuraavasti:

<b>b:n arvo</b>	<b>Kaavan potenssi</b>	<b>Suosittelava muunnos</b>
-1	2	neliöön korottaminen ( $Y^2$ )
0	1	ei muunnosta (Y)
$\frac{1}{2}$	$\frac{1}{2}$	neliöjuuri ( $\sqrt{Y}$ )
1	0	logaritmi (LOG(Y) tai LN(Y))
$1\frac{1}{2}$	$-\frac{1}{2}$	käänteisluku ja neliöjuuri ( $1/Y^2$ )
2	-1	käänteisluku ( $1/Y$ )

SPSS-ohjelmistossa ei pystytä suoraan laskemaan tätä teknistä muunnosehdotusta, mutta samaa ideaa käyttäen pystytään muunnosta arvioimaan ns. *Spread-versus-level* -kuvalla. Tämä tapahtuu *Analyze*-valikon *Descriptives*-alavalikon *Explore*-valinnan avulla. Riippuviksi muuttujiksi (*Dependent List*) valitaan malliin mukaan tulleet selittävät muuttujat. Nämä tekijät asetetaan selitettävää muuttujaa (*Factor List*) vastaan. *Plots*-valinnasta valitaan *Spread vs. Level with Levene Test* -kuva ja siellä erityisesti *Power estimation*. Kuvan alareunaan tulee aineiston pohjalta laskettu Y:n potenssin  $b$ -arvo (*Slope*) sekä ehdotettu muunnos (*Power for transformation*). Harvoin ehdotettu muunnos on kuitenkin juuri täsmälleen edeltäneen taulukon mukainen. Muunnosta kokeillaan sillä potenssilla, joka lähinnä vastaa taulukon arvoa. Jos siis ehdotettu muunnoksen potenssi on 0.189, voidaan kokeilla logaritmimuunnosta, sillä 0.189 on lähempänä nollaa kuin puolta. Uudella muunnetulla muuttujalla tehdään joko edellä kuvattu *Spread-versus-level* -kuva tai tavallinen histogrammi ja tarkistetaan, olisiko nyt ehdotettu muunnoksen potenssi lähellä arvoa 1 eli ettei tarvittaisi muunnosta.

## Ennustearvo, residuaalit ja ristiinvaldointi

**Lisäanalyysina voidaan laskea ennustearvo** (*Predicted value*). Kyseistä arvoa voidaan hyödyntää esimerkiksi ennustettaessa muiden kuin otoksessa tai laskelmissa mukana olleiden havaintojen arvoja. Lisäanalyysina voidaan tehdä **analyyseja myös residuaaleille**. Voimme pyrkiä selittämään sitä, miksi osa havainnoista ei esimerkiksi saavuta ennustearvoa (residuaalit ovat negatiivisia) ja osa saa ennustetta paremman arvon (residuaalit positiivisia). **Ristiinvaldoinnissa** (*Cross-validation*) enkä kolmasosa tai jopa puolet havainnoista otetaan mallin rakentamiseen ja malli testataan lopuilla. On nimittäin niin, että tekninen mallittaminen pystyy kyllä löytämään hyvän mallin – jopa liian hyvän mallin – siihen aineistoon, joka oli mukana analyysissa. Emme vain tiedä, miten luotu malli ennustaa niitä havaintoja, jotka eivät ole mukana itse analyysissa. **Yleinen tapa tutkia mallin hyvyttä tai toimivuutta onkin ottaa mukaan analyysiin vain osa havainnoista ja tutkia mallin osuvuutta lopuilla havainnoilla**. Tämä vaatii kuitenkin huomattavan määrän havaintoja.

## 2.1.5 Tekninen suoritus SPSS-ohjelmistolla ja tulkinta

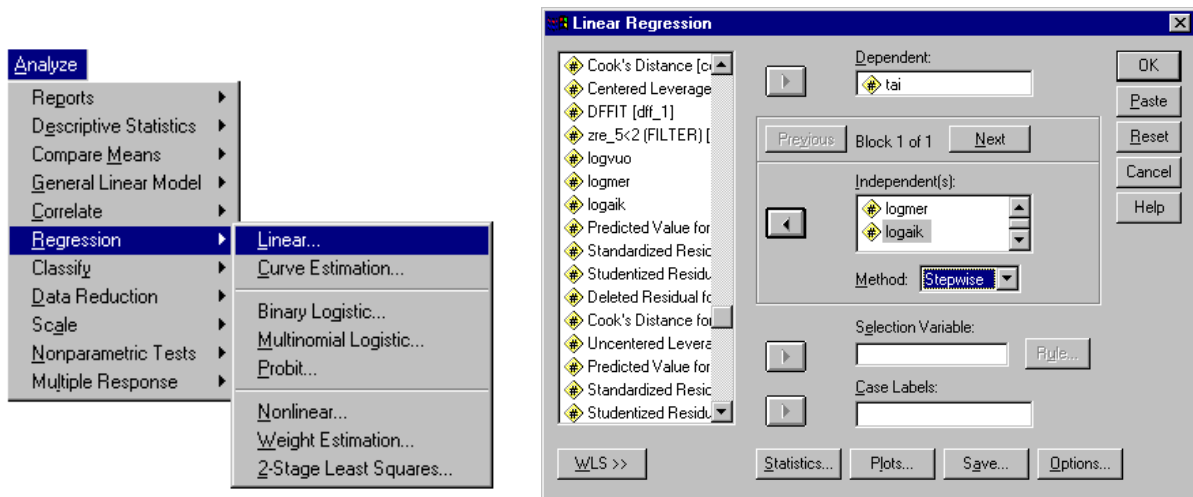
Tarkoituksemme on regressioanalyysin avulla selvittää, millä sitoutumiseen liittyvillä tekijöillä voidaan selittää sitä, että kyselyyn vastaaja kokee olevansa taitava oman alansa harrastajana. Korrelaatiomatriisin perusteella mukaan on valittu vain sellaisia tekijöitä, jotka korreloivat taitotason kanssa (Taulukko 21). Yleisesti ottaen korrelaatiot eivät ole korkeita (vaihtelevat välillä 0.29 - 0.50), mikä indikoi sitä, että mallin selitysaste ei tule olemaan korkea. Koska otoskoko on kohtuullinen (n. 700), korrelaatiot eroavat nolasta tilastollisesti merkitsevästi. Ennen regressioanalyysia on tutkittu muuttujien jakaumia ja havaittu, että erityisesti kaikki selittävät muuttujat ovat voimakkaasti vinoja. Kaikille selittäville muuttujille on tehty logaritmi-muunnokset. Muuttujat AIK ja VUO ovat voimakkaasti painottuneita pieniin arvoihin, joten niille on tehty suora luonnollinen logaritointi. Logaritointia varten muuttuja pitää olla positiivinen; koska jotkut vastaajista olivat harrastaneet vain vähän aikaa (vuosina 0) ja toisaalta jotkut ilmoittivat käyttävänsä 0 tuntia viikossa harrastukseensa, on muuttujiin ensin lisätty 1. Muuttuja MER puolestaan on painottuu suuriin arvoihin, joten se on ensin käännetty ja sitten otettu luonnollinen logaritmi. Näillä muunnoksilla pyritään siihen, että jakaumat olisivat normaalimmat, eikä jälkikäteen tarvitsisi tehdä muunnoksia residuaalien heteroskedastisuuden tai poikkeavien havaintojen vuoksi.

**Taulukko 2.1** Analyysiin mukaan tulevat muuttujat

Mja	selite	Merkitys mallissa	Korrelaatio Y:n kanssa
TAI	Itse koettu taitotaso	Y	
LOGMER	Harrastuksen koettu merkityksellisyys (logaritmi-muunnos [LN(6-MER)])	X	-.29
LOGAIK	Harrastukseen käytetty aika kuukaudessa (logaritmi-muunnos [LN(AIK+1)])	X	.29
LOGVUO	Harrastuksen kesto vuosina (logaritmi-muunnos [LN(AIK+1)])	X	.50

### Perusnäkökulma

SPSS-ympäristössä regressioanalyysi alkaa valinnoilla *Analyze – Regression – Linear...* Regressioanalyysin päävalikko näyttää seuraavalta:



## Valinnat

Muistamme, että regressioanalyysi voidaan tehdä usealla eri menettelyllä. Oletuksena SPSS-ohjelmistossa on, että halutaan tehdä pakotettu mallitus (*Method*-valinnan vaihtoehto *Enter*). Nut valitaan kuitenkin askeltava eli *Stepwise*-menettely, sillä käytössä ei ole teoriaa, joka väittäisi, juuri kyseiset valitut selittävät muuttujat muodostaisivat taitotason. Aiemmin jo todettiin, että on varmempaa kokeilla erilaisia regressiomenettelyjä luotettavimman tuloksen löytämiseksi. Päävalikossa on muuttujien ja regressiomenettelyn lisäksi valittavana neljä erilaista seikkaa, joilla voidaan vaikuttaa analyysin sisältöön ja laajuuteen: *Statistics*, *Plots*, *Save* ja *Options*. Tehdään oletusten lisäksi seuraavat valinnat:

- **Statistics**

R squared change,  Collinearity diagnostics

Residuals:  Casewise Diagnostics

[Haluamme tutkia multikollineaarisuutta ja pyrimme löytämään outlierieita.]

- **Plots**

Produce all Partial plot

Valitaan kuvat: Y:SDRESID [Standardoitut residuaalit]

X:ZPRED [Ennustearvo]

ja

Y:ZPRED [Ennustearvo]

X:DEPENDNT [Selitettävä muuttuja Y]

Standardized Residual plots:  Normal probability plot

[Voimme graafisesti tarkastella, ovatko residuaalit normaalisia.]

- **Save**

Predicted:  Standardized

Residuals:  Studentized,  Studentized deleted

Distance:  Mahalanobis,  Cook's,  Leverage value

Influence statistics:  DfFit

[Saamme hyvän käsityksen mallin oletusten toteutumisesta.]

- **Options**

[Annamme oletusten olla voimassa. Täällä määrätään millä sisäänotto- ja ulosheittokriteereillä muuttujia käsitellään askeltavassa valinnassa.]

## Tulokset ja niiden tulkinta

Seuraavassa käymme läpi regressioanalyysin tulokset tilanteessa, että meillä on kolme selitettävää muuttujaa ja yksi selitettävä muuttuja. Regressiomenetelmänä käytetään tässä askeltavaa menettelyä (*Stepwise*). Logistisen regression yhteydessä luvussa 3.2 teemme pakotetun mallituksen. Taulujen selitykset eivät oleellisesti poikkea erilaisia regressiomenettelyjä käytettäessä. Halukas lukija voi konsultoida esimerkiksi SPSS-manuaalia (1999a, 205 - 230) saadakseen laajemman kuvan erilaisista tulkintavaihtoehdoista.

Variables Entered/Removed<sup>d</sup>

Model	Variables Entered	Variables Removed	Method
1	LOGVUO	,	Stepwise (Criteria: Probability-of-F-to-enter <= ,050, Probability-of-F-to-remove >= ,100).
2	LOGMER	,	Stepwise (Criteria: Probability-of-F-to-enter <= ,050, Probability-of-F-to-remove >= ,100).
3	LOGAIK	,	Stepwise (Criteria: Probability-of-F-to-enter <= ,050, Probability-of-F-to-remove >= ,100).

a. Dependent Variable: TAI

Ensimmäinen taulu kertoo, missä **järjestyksessä muuttujat on otettu mukaan malliin** (*Variables Entered/Removed*). Osoittautuu, että kaikki mukaan valitut muuttujat ovat tilastollisessa mielessä hyviä selittäjiä – tähän olimme tietenkin varmistaneet valitsemalla mukaan muuttujia, jotka korreloivat muuttujan TAI kanssa. SPSS-ohjelmisto käsittelee kaikkia löytämiään erilaisia regressiomalleja erillisinä malleina (*Model*); viimeinen malli numero 3 on siis  $TAI = C(\text{vakio}) + \beta_1 \text{LOGVUO} + \beta_2 \text{LOGMER} + \beta_3 \text{LOGAIK} + \varepsilon$ . *Method* -sarakeessa kerrotaan, millä sisäänottokriteerillä (F-testin p-arvo  $\leq 0.05$ ) ja poisheittokriteerillä ( $p \geq 0.1$ ) menettely on toiminut.

Model Summary<sup>d</sup>

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Change Statistics				
					R Square Change	F Change	df1	df2	Sig. F Change
1	,616 <sup>a</sup>	,379	,378	,701	,379	440,547	1	722	,000
2	,637 <sup>b</sup>	,406	,405	,685	,027	33,317	1	721	,000
3	,643 <sup>c</sup>	,414	,411	,682	,007	8,913	1	720	,003

a. Predictors: (Constant), LOGVUO

b. Predictors: (Constant), LOGVUO, LOGMER

c. Predictors: (Constant), LOGVUO, LOGMER, LOGAIK

d. Dependent Variable: TAI

Mallien hyvyttä kuvataan **mallien yhteenveto** -taulukossa (*Model Summary*). Jokaisessa mallissa on mukana vakiotermi (*Constant*) sekä erinäinen joukko muita muuttujia. R-sarake (multippelikorrelaatiokerroin) kertoo selittävien muuttujien muodostaman muuttujajoukon ja selitettävän muuttujan välisen yhteiskorrelaation. Tätä tärkeämpi on kuitenkin multippelikorrelaatiokertoimen neliö  $R^2$  (*R Square*), joka kertoo mallin selitysasteen eli sen, kuinka monta prosenttia muuttujat yhdessä selittävät TAI-muuttujasta. Nyt kaikki kolme tekijää selittävät taitotasosta yhteensä vain n. 40 % (tarkalleen 41.4 %), mitä ei voi pitää kovin suurena osuutena. Koetusta taitotasosta jää siis selittymättä 59.6 %. Otoskoolla ja selittäjien määrällä korjattu  $R^2$  (*Adjusted R Square*) on hieman pienempi (koko mallille 0.411) kuin korjaamaton. Olimme lisäksi pyytäneet R:n muutosta koskevia tietoja. Taulun loppuosa kertoo, että muuttujien lisääminen on tilastollisessa mielessä vakuuttavasti  $R^2$ :n arvoa kasvattavaa. Vaikka viimeisen AIK-muuttujan lisääminen ei kasvata selitysastetta kuin 0.7 %, on muutos kuitenkin tilastollisesti merkitsevä ( $p = 0.003$ ).