

---

TAMPEREEN YLIOPISTO  
Tilastollinen mallintaminen

---

Kari Tokola ja Jouni Kirjola

Aivojen koon  
riippuvuus ruumiinpainosta

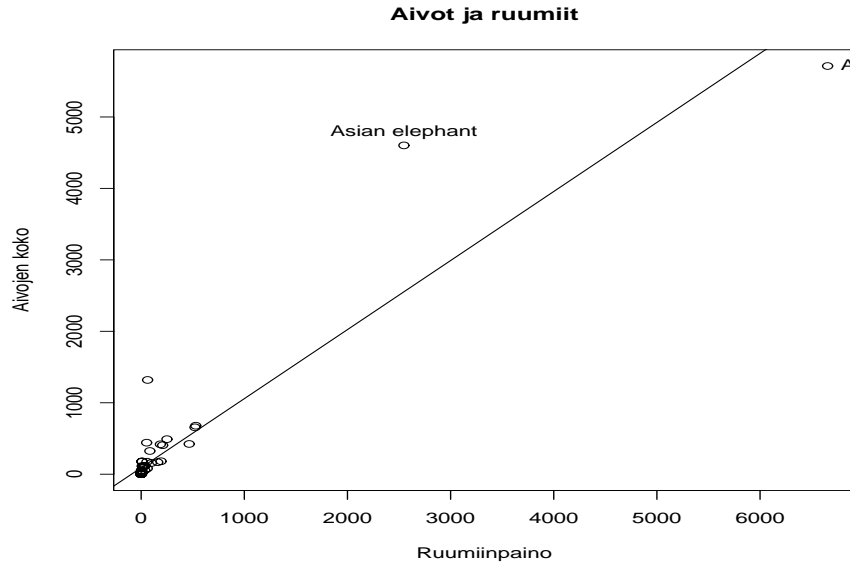
---

Matematiikan, tilastotieteen ja filosofian laitos  
Tilastotiede  
15. lokakuuta 2000

---

# 1 Johdanto

Aineistona meillä oli nisäkkäiden aivojen koko ja ruumiinpaino muuttujat. Tilastotyksiköitä aineistossa oli reilut kuusikymmentä. Eläinten koko vaihteli hiirestä norsuun. Tehtävänäimme oli tutkia eläinten painon riippuvuutta aivojen kokoon.



Kuva 1: Aivojen koon riippuvuus ruumiinpainosta

## 2 Perusregressiomalli

Aloitimme tarkastelun yksinkertaisesta regressiomallista. Kuten kuvasta 1 on helppo havaita, että  $y$  (aivojen koko) kasvaa  $x$ :n (ruumiin paino) kasvaessa. Valitulla asteikolla muutamat havainnot erottuvat selvästi muusta aineistosta ja loput muodostavat yhden "möykyn". Nämä poikkeavat havainnot ovat afrikanelefantti ja intiannorsu. Niiden paino on aivan eri luokkaa vertattuna muiden eläinten painoon. Tämän vuoksi niillä on todella suuri vaikutus mallin parametreihin. Ja mallihan oli jokaiselle havainnolle  $y_i$ ,

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (1)$$

ja virhetermit  $\epsilon_i, i = 1, \dots, n$  ovat toisistaan riippumattomia ja noudattavat normaalijakaumaa  $N(0, \sigma^2)$ . Ajamalla seuraava R-komento saadaan muodostettua malli (1). Jossa siis objektiin species.lm "sijoitetaan" malli.

```
species.lm <- lm(y ~ x)
```

Tulokseksi saadaan:

Residuals:

Min	1Q	Median	3Q	Max
-810.07	-88.52	-79.64	-13.02	2050.33

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	91.00440	43.55258	2.09	0.0409 *
x	0.96650	0.04766	20.28	<2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 334.7 on 60 degrees of freedom

Multiple R-Squared: 0.8727, Adjusted R-squared: 0.8705

F-statistic: 411.2 on 1 and 60 degrees of freedom, p-value: 0

Aineestoon sovitettu suora on siis,

$$y_i = 91,0 + 0,97 * x_i, \quad (2)$$

Regressiokertoimen ( $\beta_1$ ) estimaatin  $\hat{\beta}_1$   $t$ -arvo on erittäin merkitsevä (p-arvo lähellä nollaa). Suora näyttäisi sopivan aineestoon. Tämä ei kuitenkaan vakuuttanut vielä meitä. Suoritimme graafisia diagnostisia tarkistuksia kome-nolla

`plot(species.lm, ask=T)`

Komento tulostaa seuraavat kuviot: *residuals vs. fitted*<sup>1</sup>, *Q-Q-diagramma*<sup>2</sup>, *Scale-location plot*<sup>3</sup> ja *Cook's distance*<sup>4</sup>. Näistä enemmän seuraavassa alakappaleessa.

## 2.1 Mallin (1) diagnostiset tarkastelut

Mallin oletusten mukaan virhetermit  $\epsilon_i, i = 1, \dots, n$  noudattavat normaali-jakaumaa. Koska virhetermejä ei voida havaita, niiden käyttäytymistä tutkimme residuaalien  $y - \hat{y}$  (havaittu arvo - sovite) avulla. *Q - Q*-diagramma antoiinkin viitteitä, etteivät residuaalit noudattaisi normaalijakaumaa. Lisäksi ulkoisesti Studentoidut residuaalit, jotka saadaan kaavalla,

$$t_i = \frac{y_i - \hat{y}_i}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}}, \quad h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_x} \quad (3)$$

ovat testisuureita havaintojen vierauden (Outlier) tilastolliseksi testaamiseksi. Sillä ne noudattavat  $t$ -jakaumaa vapausastein 59.

<sup>1</sup>Kuva residuaaleista ja sovitetuista arvoista

<sup>2</sup>Residuaalien normaalisuus

<sup>3</sup>Studentoitujen residuaalien itseisarvojen neliöjuuri ja sovite

<sup>4</sup>Cookin etäisyys

Laskemalla Studentoidut residuaalit käskyllä

```
rstudent(species.lm)
```

saadaan jokaiselle havainnolle p-arvo (t-jakaumasta) hypoteesille

$$H_0 : \text{"Havainto ei ole vieras"}$$

Nämä arvot laskemalla saadaan, että havainnot 19,32 ja 33 (intiannorsu, ihmisen ja afrikanelefantti) ovat vieraita havaintoja. Lisäksi samaan tulokseen pääsimme käyttämällä komentoa,

```
influence.measures(species.lm)
```

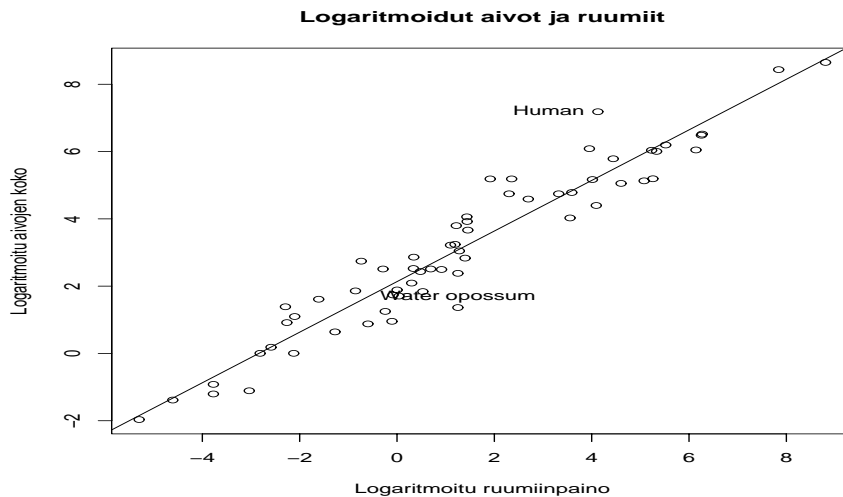
joka paljasti saman tuloksen. Kyseinen komento laskee muunmuassa Cookin etäisyydet, joihin jo aikaisemmin viittasimme.

### 3 Logaritmoitu malli

Halusimme lähteä tarkastelemaan logaritmoitua mallia, koska alkuperäisessä mallissa kolme havainnoista vetivät regressiosuoraa voimakkaasti puoleensa. Logaritmoimalla havaintojen arvojen suhteet muuttuvat "linearisemmiksi".

$$\frac{y_i}{x_i} \implies \frac{\ln(y_i)}{\ln(x_i)} = \ln(y_i) - \ln(x_i) \quad (4)$$

Havaintojen arvojen muuttumisen voimme huomata seuraavasta kuvasta.



Kuva 2: Logaritmoidut muuttujat

Logaritmoitu teoreettinen malli on siis seuraavanlainen,

$$\ln(y_i) = \beta_0 + \beta_1 \ln(x_i) + \epsilon_i \quad (5)$$

jossa on siis samat oletukset kuin mallissa (1). Tässä kohtaa on siis hyvä huomata, ettei logaritmin kannalla ole väliä. Tämän varmistaminen jätetään lukijan tehtäväksi. Nyt jääkin siis muodostettavaksi malli R:llä. Annamme mallille (objektille) nimeksi `species.llm`, kuten seuraavasta komennosta huomaamme.

```
species.llm <- lm(ly ~ lx)
```

*Huom. muuttujille on tehty logaritointi ennen mallin luomista. Ajettaessa yhteenveto mallista saadaan seuraavanlaiset tulokset.*

Residuals:

Min	1Q	Median	3Q	Max
-1.71550	-0.49228	-0.06162	0.43597	1.94829

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	2.13479	0.09604	22.23	<2e-16 ***
lx	0.75169	0.02846	26.41	<2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6943 on 60 degrees of freedom

Multiple R-Squared: 0.9208, Adjusted R-squared: 0.9195

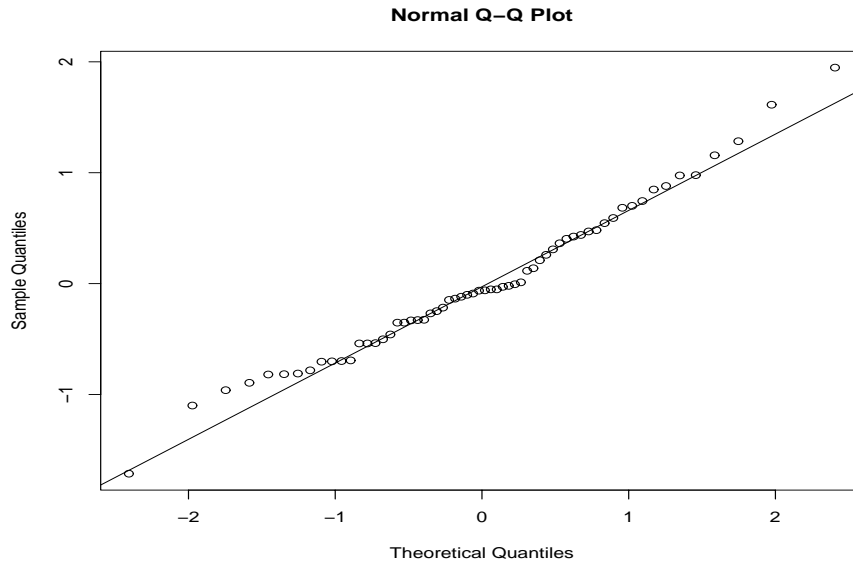
F-statistic: 697.4 on 1 and 60 degrees of freedom, p-value: 0

Vertaamalla näitä tuloksia mallin (1) antamiin tuloksiin huomasimme seuraavaa. Ensinnäkin huomattava muutos on tapahtunut estimaattorien hajonnoissa  $\text{var}(\hat{\beta})$ . Mallissa (1)  $\hat{\beta}_0$ :n hajonta on huomattavasti suurempi kuin mallissa (4). Toinen huomattava asia on kohonnut selitysaste. Syy tähän löytynee tarkemmista tarkasteluista.

### 3.1 Mallin (4) tarkemmat tarkastelut

Lähdimme tarkastelemaan "uuden"mallin hyvyttä kuten aikaisemman mallin. Piirsimmekin ensin  $Q-Q$ -diagramman, joka antoi viitteitä että residuaalit olisivat jakautuneet normaalimmin kuin mallissa (1). Kuten huomaamme kuvasta 4. Testattaessa Studentin residuaaleja kolme havaintoa tuli merkit-

seväksi. Tämä siis tarkoittaa, että havaintojen p-arvot olivat  $\leq 0.05$ .



Kuva 3: Q-Q -diagramma mallissa (4)

## 4 Mallien vertailu

On itsestään selvää, että logaritmoitaessa  $\hat{\beta}$ -vektorin arvot muuttuvat. Myös estimaattoreiden hajonnat pienenevät huomattavasti logaritmoidessa. Tämä johtuu siitä, että havaintojen arvojen pienessä huomattavasti pienenee myös hajonta. Logaritmoidussa mallissa jäännösneliösumma (SSE) on paljon pienempi ja se johtaa seuraaviin tuloksiin. Logaritmoidun mallin estimaattorien p-arvot pienenevät. Myös selitysaste tulee korkeammaksi. Tehtäessä tämä muunnos, alkuperäinen  $y$  on seuraavaa muotoa:

$$y_i = \beta_0 + e^{x_i} \quad (6)$$

## 5 Yhteenveto

Lukijalle saattaa tulla käsitys, että logaritmointiin päättymisen on itsestään selvyyys. Näin ei kuitenkaan ole. Joskus saattaa muuttujien tulkinnan kannalta olla hankalaa logaritmoida ne. Toisin sanoen niiden merkitys saattaa muuttua. Meille autuaana tapahtanut logaritmointi, johti meidät parempaan malliin. Siis meidän mielestä logaritmoitu malli on todellakin parempi kuin alkuperäinen, verrattaessa esim. F-testien p-arvoja, estimaattorien hajontoja ja selitysasteita.

## 5.1 Vieraista havainnoista (outlier)

Vieraiden havaintojen havaitsemiseen voidaan käyttää edellä mainitsemiamme testejä, esim. Studentoidut residuaalit ja Cookin etäisyydet. Alkuperäisessä aineistossa havaittiin kolme vierasta havaintoja, kuten olemme kappaleessa 2.1. todenneet. Syy miksi emme niitä poistaneet on se, että ne vetävät regressiosuoraa puoleensa. Logaritmoidussa mallissa vieraita havaintoja olivat ihminen ja *water opussum*<sup>5</sup> ja *rhesus monkey*. Näillä havainnoilla oli itseisarvoiltaan suurimmat Studentoidut residuaalit. Haluttaessa voitaisiin näitä arvoja käyttää ns. *viisausmittarina*. Negatiivinen arvo vastaisi normaalia tyhmempää ja päinvastoin. Ihmisen ja *rhesus monkey*'n arvot olivat positiivisia. Sensijaan *water opussumin* arvo oli negatiivinen. Edellisen perusteella jokainen vetäköön johtopäätöksensä. Riippuen tutkimuksen tarkoituksesta jatkettaisiin vieraiden havaintojen poistamisella tai niiden tarkemmalla tutkimisella.

## 6 R-koodi

```
library(mass)
data(mammals)
attach(mammals)
species <- row.names(mammals)
x <- body; y <- brain
plot(x,y,main="Aivot ja ruumiit",xlab="Ruumiinpaino",ylab="Aivojen koko")
species.lm <- lm(y ~ x)
abline(species.lm)
identify(x,y,species)
summary(species.lm)
plot(species.lm,ask=T)
rstudent(species.lm)
influence.measures(species.lm)
lx <- log(x)
ly <- log(y)
species.llm <- lm(ly ~ lx)
plot(log(x),log(y),main="Logaritmoidut aivot ja ruumiit",xlab="Logaritmoitu
ruumiinpaino",ylab="Logaritmoitu aivojen koko")
abline(species.llm)
identify(log(x),log(y),species)
summary(species.llm)
plot(species.llm,ask=T)
rstudent(species.llm)
influence.measures(species.llm)
```

---

<sup>5</sup>Emme tienneet suomenkielisiä nimiä :)