

Esim. Pistemäärä-aineisto oli seuraava:

The decimal point is 1 digit(s) to the right of the |

```
0 | 9
1 | 0112233334444
1 | 566777799
2 | 001
```

Leikattu keskiarvo: Lasketaan esimerkiksi pistemäärän 10 % leikattu keskiarvo \bar{t}_{10} .

Lasketaan ensin se montako havaintoa 26:sta havainnosta on 10 %:

$$100 \cdot \frac{x}{26} = 10 \Rightarrow x = 2.6$$

\Rightarrow Poistetaan aineistosta kolme pienintä ja kolme suurinta havaintoa, joten $n = 26 - 3 - 3 = 20$.

$$\Rightarrow \bar{t}_{10} = \frac{1}{20} \cdot (11 + 12 + 12 + \dots + 17 + 19 + 19) \approx 14.8$$

Winsoroitu keskiarvo: Lasketaan esimerkiksi pistemäärän 10 % winsoroitu keskiarvo \bar{w}_{10} .

Korvataan leikattua keskiarvoa laskettaessa poisjätetyt havainnot seuraavasti:

- 3 pienintä havaintoa korvataan pienimmällä aineistoon jääneellä arvolla
- 3 suurinta havaintoa korvataan suurimmalla aineistoon jääneellä arvolla

Lasketaan keskiarvo "manipuloidulle" aineistolle, jonka $n = 26 - 3 - 3 + 3 + 3 = 26$.

$$\Rightarrow \bar{w}_{10} = \frac{1}{26} \cdot (11 + 11 + 11 + 11 + 12 + \dots + 17 + 19 + 19 + 19 + 19 + 19) \approx 14.8$$

Esimerkki harmonisen keskiarvon käytöstä:

Herra B tekee autollaan edestakaisen matkan 120 kilometrin päässä sijaitsevalle paikkakunnalle. Menomatalla keskinopeus oli 60 km/h ja paluumatkalla 120 km/h. Mikä oli hänen keskinopeutensa edestakaisella matkalla?

Jos keskinopeus lasketaan aritmeettisen keskiarvon avulla eli

$$\frac{60 + 120}{2} = 90 \text{ km / h}$$

saadaan virheellinen tulos!

Herra B:n ajoaika menomatalla oli $120/60 = 2$ h ja paluumatkalla $120/120 = 1$ h. Oikea keskinopeus on siis

$$\frac{120 + 120}{3} = 80 \text{ km / h}$$

Tulos on harmoninen keskiarvo nopeuksista 60 km/h ja 120 km/h, sillä

$$H = \frac{2}{\frac{1}{60} + \frac{1}{120}} = 80 \text{ km / h}$$

Esimerkki geometrisen keskiarvon käytöstä:

Yhtiön A tulos 3-kertaistui vuonna 2002, 4-kertaistui vuonna 2003 ja 1.5-kertaistui vuonna 2004. Koko tänä aikana tulos siis $3 \cdot 4 \cdot 1.5 = 18$ -kertaistui. Koska tutkittava ilmiö on kerrannainen, keskimääräinen vuotuinen kasvu saadaan geometrisena keskiarvona:

$$\sqrt[3]{3 \cdot 4 \cdot 1.5} = 2.62$$

Kutakin vuotta kohti tulos siis keskimäärin 2.62-kertaistui. (Aritmeettinen keskiarvo antaisi tässä virheellisen tuloksen 2.83).

YHDEN MUUTTUJAN TARKASTELU:

MITTA-ASTEIKKO			
LUOKITTELUASTEIKKO	JÄRJESTYSASTEIKKO	VÄLIMATKA-ASTEIKKO	SUHDEASTEIKKO
TAULUKOITA: - frekvenssijakauma (<i>frequency distribution</i>)	- frekvenssijakauma - summajakauma (<i>cumulative frequency distribution</i>)	- frekvenssijakauma - summajakauma	- frekvenssijakauma - summajakauma
GRAAFISIA ESITYKSIÄ: - pylväsdiagrammi (<i>bar chart</i>) - piirakkakuviio (<i>pie chart</i>)	- pylväsdiagrammi - piirakkakuviio	jatkuva / diskreetti - histogrammi (<i>histogram</i>) / janadiagrammi - frekvenssimonikulmio - summakäyrä / porraskuviio - runko-lehti -kuviio (<i>stem and leaf plot</i>) - laatikko-jana -kuviio (<i>boxplot</i>) - pistekuviio (<i>dot plot</i>)	
TUNNUSLUKUJA: - moodi (<i>mode</i>) (- entropiasuhde)	- moodi - mediaani (<i>median</i>) - vaihteluväli - kvartiiliväli	- moodi - mediaani (<i>median</i>) - aritmeettinen keskiarvo (<i>mean</i>) - leikattu keskiarvo (<i>trimmed mean</i>) - winsoroitu keskiarvo (<i>winsorized mean</i>) - geometrinen keskiarvo - harmoninen keskiarvo - vaihteluväli - vaihteluvälin pituus (<i>range</i>) - kvartiiliväli - kvartiilivälin pituus (<i>IQR=interquartile range</i>) - keskipoikkeama (<i>mean deviation</i>) - keskihajonta (<i>standard deviation</i>) - varianssi (<i>variance</i>) - variaatioeroin (<i>coefficient of variation, vain suhdeasteikolle</i>) - vinousmitta (<i>coefficient of skewness, skewness</i>) - huipukkuusmitta (<i>coefficient of kurtosis</i>)	

Tutkitaan muuttujien välistä riippuvuutta ehdollisten prosenttijakaumien avulla eri sukupuoli-

ESIMERKKI vakioinnista kun muuttujien välistä riippuvuutta on tutkittu ristintaulukoinnin avulla.

Tutkimuksessa selvitettiin erään ihmisjoukon kiinnostusta urheiluun (= x) ja politiikkaan (= y). Muuttujien välille saatiin seuraavaanlainen ristintaulukko:

	kiinnostus urheiluun		Yhteensä
	suuri	vähäinen	
kiinnostus suuri	55	36	91
politiikkaan vähäinen	51	54	105
Yhteensä	106	90	196

Tutkitaan muuttujien välistä riippuvuutta ehdollisten prosenttijakaumien avulla:

	kiinnostus urheiluun		Yhteensä
	suuri	vähäinen	
kiinnostus suuri	52 %	40 %	46 %
politiikkaan vähäinen	48 %	60 %	54 %
Yhteensä	100 %	100 %	100 %

Koska ehdolliset prosenttijakaumat poikkeavat toisistaan, näyttää tarkasteltavien muuttujien muuttujien x ja y välillä olevan riippuvuutta: urheilusta paljon kiinnostuneiden joukossa on suhteellisesti enemmän paljon politiikasta kiinnostuneita.

Sukupuolen (= z) epäillään kuitenkin vaikuttavan x:n ja y:n väliseen riippuvuuteen. Tämän takia kaksitoteiset jakaumat (eli ristintaulukot) lasketaan sukupuoli-muuttujan eri luokissa erikseen ts. *vakioidaan sukupuoli-muuttuja*.

miehet:

	kiinnostus urheiluun		Yhteensä
	suuri	vähäinen	
kiinnostus suuri	45	20	65
politiikkaan vähäinen	25	11	36
Yhteensä	70	31	101

naiset:

	kiinnostus urheiluun		Yhteensä
	suuri	vähäinen	
kiinnostus suuri	10	16	26
politiikkaan vähäinen	26	43	69
Yhteensä	36	59	95

miehet:

	kiinnostus urheiluun		Yhteensä
	suuri	vähäinen	
kiinnostus suuri	64 %	65 %	64 %
politiikkaan vähäinen	36 %	35 %	36 %
Yhteensä	100 %	100 %	100 %

- ehdolliset prosenttijakaumat ovat likimain samat, muuttujien välillä ei ole riippuvuutta.

naiset:

	kiinnostus urheiluun		Yhteensä
	suuri	vähäinen	
kiinnostus suuri	28 %	27 %	27 %
politiikkaan vähäinen	72 %	73 %	73 %
Yhteensä	100 %	100 %	100 %

- ehdolliset prosenttijakaumat ovat likimain samat, muuttujien välillä ei ole riippuvuutta.

Nyt nähdään, että alussa havaittu riippuvuusuhde muuttujien x ja y välillä on syntynyt niin, että miehet olivat usein kiinnostuneita sekä urheilusta että politiikasta, kun taas naisia ei kiinnostanut useimmiten kovin paljon kumpikaan.

Pearsonin tulomomenttikorrelaatiokertoimen r arvojen tulkinnassa voi käyttää (tässä vaiheessa kurssia) seuraavia adjektiveja:

- 1 ≤ r < -0.8: vahva (voimakas) negatiivinen lineaarinen riippuvuus
- 0.8 ≤ r < -0.6: melko vahva (melko voimakas) negatiivinen lineaarinen riippuvuus
- 0.6 ≤ r < -0.4: kohtalainen (melkoinen) negatiivinen lineaarinen riippuvuus
- 0.4 ≤ r < -0.2: melko heikko negatiivinen lineaarinen riippuvuus
- 0.2 ≤ r < -0.1: heikko negatiivinen lineaarinen riippuvuus
- 0.1 ≤ r ≤ 0.1: ei *lineaarista* riippuvuutta
- 0.1 < r ≤ 0.2: heikko positiivinen lineaarinen riippuvuus
- 0.2 < r ≤ 0.4: melko heikko positiivinen lineaarinen riippuvuus
- 0.4 < r ≤ 0.6: kohtalainen (melkoinen) positiivinen lineaarinen riippuvuus
- 0.6 < r ≤ 0.8: melko vahva (melko voimakas) positiivinen lineaarinen riippuvuus
- 0.8 < r ≤ 1: vahva (voimakas) positiivinen lineaarinen riippuvuus

Huom. Samat adjektiivit sopivat myös Spearmanin ρ :ta ja suhteellista kontingenssikerrointa C/C_{\max} tulkittaessa:

Esim.,

$\rho = 0.45$: muuttujien väiltä vallitsee kohtalainen positiivinen riippuvuus

$C/C_{\max} = 0.45$: muuttujien väiltä vallitsee kohtalainen riippuvuus

Esim. Erään tilastotieteen kurssin opiskelijoiden sukupuolen ja koulutusohjelman ristiintaulukko oli seuraava:

koulutusohjelma	sukupuoli		Yhteensä
	mies	nainen	
kansantaloustiede	9	9	18
laskentatoimi	8	15	23
markkinointi	2	16	18
jokin muu	4	3	7
yhteensä	23	43	66

Ovatko sukupuoli ja koulutusohjelma toisistaan riippumattomia?

Ehdolliset prosenttijakaumat ovat:

koulutusohjelma	sukupuoli		Yhteensä
	mies	nainen	
kansantaloustiede	39	21	27
laskentatoimi	35	35	35
markkinointi	9	37	27
jokin muu	17	7	11
yhteensä	100	100	100

Ehdolliset prosenttijakaumat näyttävät poikkeavan jonkin verran toisistaan, joten koulutusohjelman ja sukupuolen välillä vallitsee kohtalainen riippuvuus. Miehistä kansantaloustieteen koulutusohjelmaan kuuluu noin 39 % ja naisista noin 21 %. Vastaavasti esim. markkinoinnin koulutusohjelmaan kuuluu noin 9 % miehistä ja 37 % naisista.

Mikäli tarkasteltavien muuttujien välillä ei olisi riippuvuutta ollenkaan, olisivat ehdolliset prosenttijakaumat seuraavat:

koulutusohjelma	sukupuoli		Yhteensä
	mies	nainen	
kansantaloustiede	27	27	27
laskentatoimi	35	35	35
markkinointi	27	27	27
jokin muu	11	11	11
yhteensä	100	100	100

Edellä olevaan tilanteeseen päädyttäisiin, jos ristiintaulukon havaitut solufrekvenssit (f_{ij}) olisivat olleet alunperin seuraavat:

koulutusohjelma	sukupuoli		Yhteensä
	mies	nainen	
kansantaloustiede	6.27	11.73	18
laskentatoimi	8.02	14.98	23
markkinointi	6.27	11.73	18
jokin muu	2.44	4.56	7
yhteensä	23	43	66

Edellä olevassa taulukossa olevat teoreettiset frekvenssit ovat ns. odotettuja frekvenssejä (e_{ij}) tilanteessa, jossa muuttujat sukupuoli ja koulutusohjelma ovat toisistaan riippumattomia ja ristiintaulukon reunafrekvenssit ovat ne, mitkä on havaittu.

Lasketaan seuraavaksi muuttujien välistä riippuvuutta kuvaavan tunnusluvun (*suhteellisen kontingenssikertoimen*) arvo.

Edellisessä taulukossa esiintyvät riippumattomuus-oletuksen mukaiset odotetut frekvenssit e_{ij} saadaan laskettua kaavalla $e_{ij} = \frac{f_i \cdot f_j}{n}$, missä indeksi i viittaa ristiintaulukon i. riville ja indeksi j viittaa ristiintaulukon j. sarakkeelle. Esimerkkiaineistossa

$$\begin{aligned} e_{11} &= \frac{f_1 \cdot f_{.1}}{n} = \frac{18 \cdot 23}{66} \approx 6.27, & e_{12} &= \frac{f_1 \cdot f_{.2}}{n} = \frac{18 \cdot 43}{66} \approx 11.73, \\ e_{21} &= \frac{f_2 \cdot f_{.1}}{n} = \frac{23 \cdot 23}{66} \approx 8.02, & e_{22} &= \frac{f_2 \cdot f_{.2}}{n} = \frac{23 \cdot 43}{66} \approx 14.98, \\ e_{31} &= \frac{f_3 \cdot f_{.1}}{n} = \frac{18 \cdot 23}{66} \approx 6.27, & e_{32} &= \frac{f_3 \cdot f_{.2}}{n} = \frac{18 \cdot 43}{66} \approx 11.73, \\ e_{41} &= \frac{f_4 \cdot f_{.1}}{n} = \frac{7 \cdot 23}{66} \approx 2.44, & e_{42} &= \frac{f_4 \cdot f_{.2}}{n} = \frac{7 \cdot 43}{66} \approx 4.56 \end{aligned}$$

χ^2 -tunnusluvun havaittu arvo on

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^r \frac{(f_{ij} - e_{ij})^2}{e_{ij}} = \frac{(9 - 6.27)^2}{6.27} + \frac{(9 - 11.73)^2}{11.73} + \dots + \frac{(3 - 4.56)^2}{4.56} \approx 7.82$$

$$\Rightarrow C = \sqrt{\frac{\chi^2}{n + \chi^2}} = \sqrt{\frac{7.82}{66 + 7.82}} \approx 0.33 \quad \Rightarrow C_{MAX} = \sqrt{\frac{q-1}{q}} = \sqrt{\frac{2-1}{2}} = \sqrt{\frac{1}{2}}$$

Edellä $q = \min(m,r) = \min(4,2) = 2$.

$$\Rightarrow C/C_{MAX} \approx 0.33/\sqrt{0.5} \approx 0.47$$

Muuttujien välillä vallitsee kohtalainen riippuvuus.

Vastaavasti voitaisiin laskea ns. cramerin V:

$$V = \sqrt{\frac{\chi^2}{nu}} \approx \sqrt{7.82/66 \cdot 1} \approx 0.34$$

, missä $u = \min(m-1, r-1) = \min(4-1, 2-1) = 1$

Muuttujien välillä vallitsee melko heikko riippuvuus.

Esimerkki ristitulosuhteen käytöstä:

Tutkittaessa odottavan äidin runsaan tupakoinnin (yli 10 savuketta/päivä) vaikutusta lapsen syntymäpainoon prof. Paula Rantakallion havaintoaineiston avulla:

Lapsen syntymäpaino	Äidin tupakointi		Yhteensä
	yli 10 savuketta/päivä	ei ollenkaan	
alle 2500 g	23	345	368
vähintään 2500 g	253	9360	9613
Yhteensä	276	9705	9981

Tutkitaan muuttujien välistä riippuvuutta

- ehdollisten prosenttijakaumien avulla:

Ehdolliset prosenttijakaumat kannattaa laskea nyt sarakkeittain, sillä äidin tupakoinnin voidaan ajatella olevan riski (syy) lapsen alipainoisuudelle (seuraus). Ehdolliset prosenttijakaumat ovat:

Lapsen syntymäpaino	Äidin tupakointi		Yhteensä
	yli 10 savuketta/päivä	ei ollenkaan	
alle 2500 g	8.3 %	3.6 %	3.7 %
vähintään 2500 g	91.7 %	96.4 %	96.3 %
Yhteensä	100 %	100 %	100 %

Ehdolliset prosenttijakaumat poikkeavat toisistaan, joten äidin runsaan tupakoinnin ja lapsen syntymäpainon välillä näyttää olevan riippuvuutta. Paljon tupakoivien äitien lapsista noin 8.3 % on alipainoisia, kun taas tupakoimattomien äitien lapsista alipainoisia on noin 3.6 %. Paljon tupakoivat äidit saavat siis suhteellisesti enemmän alipainoisia lapsia kuin tupakoimattomat äidit.

- ristitulosuhteen (OR) ja riskisuhteen (RR) avulla:

$$OR = \frac{f_{11}f_{22}}{f_{12}f_{21}} = \frac{23 \cdot 9360}{345 \cdot 253} \approx 2.47$$

OR:n tulkinta: muuttujien välillä on (selvä) positiivinen riippuvuus.

$$RR = \frac{\pi_1}{\pi_0} = \frac{f_{11}/f_{1.}}{f_{12}/f_{2.}} = \frac{23/276}{345/9705} \approx 2.34$$

RR:n tulkinta: Paljon tupakoivilla äideillä on noin 2.34-kertainen riski saada alipainoinen lapsi tupakoimattomiin äiteihin verrattuna.

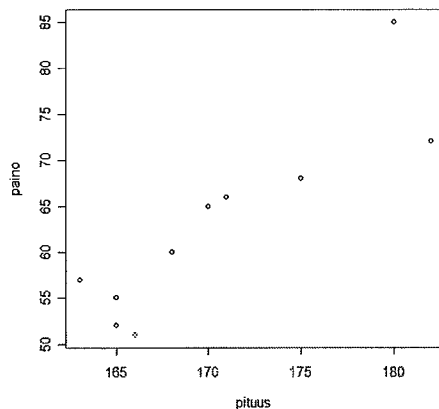
riskiero (RD):

$$RD = \pi_1 - \pi_0 = \frac{23}{276} - \frac{345}{9705} \approx 0.048.$$

Tarkastellaan seuraavaksi pituuden (=x, cm) ja painon (=y, kg) välistä riippuvuutta 10 henkilön aineistossa:

henkilö (i)	1	2	3	4	5	6	7	8	9	10
pituus (x _i)	175	180	165	170	168	166	165	171	163	182
paino (y _i)	68	85	52	65	60	51	55	66	57	72

Pituuden ja painon välinen korrelaatiodiagrammi on seuraava:



Muuttujien välillä näyttää vallitsevan lineaarista riippuvuutta, joten korrelaatiokerroin r on mielekästä laskea.

Nyt n=10,

$$\sum_{i=1}^{10} x_i = 1705, \bar{x} = 170.5, \sum_{i=1}^{10} x_i^2 = 291089$$

$$\sum_{i=1}^{10} y_i = 631, \bar{y} = 63.1, \sum_{i=1}^{10} y_i^2 = 40793,$$

$$\sum_{i=1}^{10} x_i y_i = (175 \cdot 68) + (180 \cdot 85) + \dots + (182 \cdot 72) = 108132$$

$$s_x = \sqrt{\frac{1}{n-1} \left(\sum_{i=1}^{10} x_i^2 - n\bar{x}^2 \right)} = \sqrt{\frac{1}{10-1} (291089 - 10 \cdot 170.5^2)} \approx 6.553,$$

$$s_y = \sqrt{\frac{1}{n-1} \left(\sum_{i=1}^{10} y_i^2 - n\bar{y}^2 \right)} = \sqrt{\frac{1}{10-1} (40793 - 10 \cdot 63.1^2)} \approx 10.418,$$

$$s_{xy} = \frac{1}{n-1} \left(\sum_{i=1}^{10} x_i y_i - n\bar{x}\bar{y} \right) = \frac{1}{10-1} (108132 - 10 \cdot 170.5 \cdot 63.1) \approx 60.722$$

$$\Rightarrow r_{xy} = \frac{s_{xy}}{s_x s_y} \approx \frac{60.722}{6.553 \cdot 10.418} \approx 0.889$$

r_{xy}:n tulkinta: Muuttujien välillä vallitsee voimakas positiivinen lineaarinen riippuvuus.

ESIMERKKI Sovitaan jo aiemmin esillä olleeseen (pituus,paino)-aineistoon regressiosuora $y = a + bx$, missä $y =$ paino ja $x =$ pituus. Aineisto oli seuraava:

henkilö (i)	1	2	3	4	5	6	7	8	9	10
pituus (x_i)	175	180	165	170	168	166	165	171	163	182
paino (y_i)	68	85	52	65	60	51	55	66	57	72

Edellisessä ko. aineiston käsitellessä esimerkiksi havaittiin korrelaatiogrammin avulla, että muuttujien välillä vallitsee voimakas lineaarinen riippuvuus. Lisäksi aineistosta on laskettu mm. seuraavat tunnusluvut:

$$n = 10, \quad \bar{x} = 170,5, \quad s_x \approx 6,553, \\ \bar{y} = 63,1, \quad s_y \approx 10,418, \quad s_{xy} \approx 60,722 \quad \text{ja} \quad r_{xy} \approx 0,889.$$

$$\Rightarrow b = \frac{s_{xy}}{s_x^2} \approx \frac{60,722}{6,553^2} \approx 1,414$$

$$a = \bar{y} - b\bar{x} \approx 63,1 - 1,414 \cdot 170,5 \approx -177,987$$

\Rightarrow Regressiosuora on

$$\hat{y} = -177,987 + 1,414 x.$$

Regressioyhtälön antama painon ennuste 178 cm pitkälle henkilölle on siis

$$\hat{y} = -177,987 + 1,414 \cdot 178 \approx 73,7 \text{ kg}$$

Aineiston ensimmäisellä tilastoyksiköllä muuttujien x ja y havaitut arvot olivat 175 cm ja 68 kg. Regressioyhtälön antama ennuste arvo \hat{y}_1 ko. tilastoyksikön arvoiksi muuttujalle y on

$$\hat{y}_1 = -177,987 + 1,414 \cdot 175 \approx 69,5 \text{ kg}$$

ja jäätännöstermi

$$e_1 = y_1 - \hat{y}_1 \approx 68 - 69,5 = -1,5 \text{ kg}$$

Mallissa on nyt vain yksi selittävä muuttuja, joten regressiomallin selityssaste

$$R^2 = r_{xy}^2 \approx 0,889^2 \approx 0,79.$$

Regressioerotointen tulkinna:

- b: Kun henkilön pituus kasvaa sentillä, kasvaa henkilön paino keskimäärin noin 1,414 kg.
a: Kun henkilön pituus on nolla senttiä, on henkilön paino keskimäärin noin -178 kg!

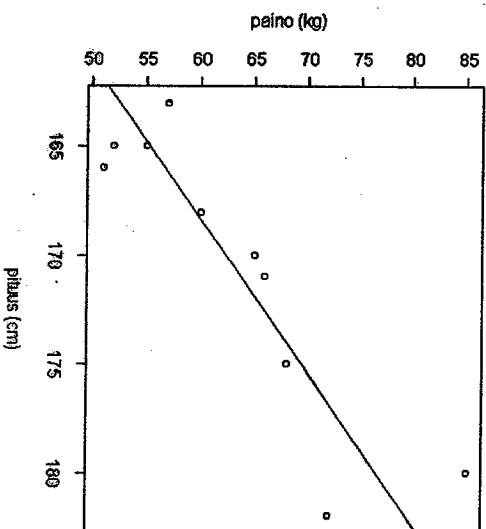
Determinaatioerotointen eli selityssasteen tulkinna:

R^2 : Regressiomallilla (henkilön pituudella) voidaan selittää noin 79 % painon kokonaisvaihtelusta.

Regressiosuoran piirtämiseksi korrelaatiogrammiin tarvitaan suoralla kaksi pistettä:

$$x = 165 : \hat{y} = -177,987 + 1,414 \cdot 165 \approx 55,3 \\ x = 180 : \hat{y} = -177,987 + 1,414 \cdot 180 \approx 76,5$$

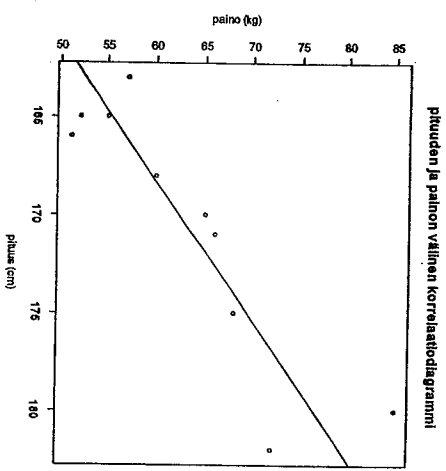
Pisteet (165,55.3) ja (180,76.5) merkitään korrelaatiogrammiin ja yhdistetään ne suoralla.



Neliösumma hajotelmassa SST = SS_{reg} + SS_{res} ;

Yhden selittävän muuttujan regressioanalyysi R-ohjelmalla:

SS_T = väitteen kokonaisvaihtelu
 SS_{reg} = regressiomallin selittävä osuus väitteen kokonaisvaihteluista
 SS_{res} = regressiomallille selittämättä jäänyt osuus väitteen kokonaisvaihteluista



```
> pituus <- c(175,180,155,170,168,166,165,171,163,182)
> paino <- c(68,85,52,65,60,51,55,66,57,72)
> lm(paino~pituus) ← komento muotoa: lm (värite ~ selittäjä)
```

```
Call:
lm(formula = paino ~ pituus)

Coefficients:
(Intercept) -177.982      pituus      1.414
=> y = -177.982 + 1.414 X

> abline(lm(paino~pituus)) ← saadaan lisää tietoja mallinnuksesta
```

```
Call:
lm(formula = paino ~ pituus)

Residuals:
    Min       1Q   Median       3Q      Max
-7.36067 -2.85809  0.05589  2.50349  8.46727

Coefficients:
(Intercept) -177.982      pituus      1.414
Estimate Std. Error t value Pr(>|t|)
(Intercept) -177.982    43.841   -4.060 0.003635 ***
pituus        1.414     0.257    5.503 0.000572 ***

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

Residual standard error: 5.052 on 8 degrees of freedom
Multiple R-Squared: 0.791
F-statistic: 30.28 on 1 and 8 DF, p-value: 0.000572
```

determinaattokerroin R²

R² = 0.6213 in tulkinna: Regressiomallin (eli muuttujien AGE15-64 ja LIT_MALE) voidaan selittää noin 62% väitteen eli BKT/asukas-muutuksen kokonaisvaihteluista aineiston Afrikan maat.

ESIMERKKI kahden selittäjän regressioanalyysistä, aineistona WORLD2005-aineiston Afrikan maat.

Muuttujat:
 Y = GDP_CAPI: BKT/asukas, US dollari
 X₁ = AGE15_64: 15-64-vuotiaiden prosentuaalinen osuus maan väestöstä
 X₂ = LIT_MALE: miesten lukutaitoprosentti

```
malli_1 <- lm(GDP_CAPI~AGE15_64+LIT_MALE, data=afrikka)
summary(malli_1) => y = b0 + b1X1 + b2X2
```

```
Call:
lm(formula = GDP_CAPI ~ AGE15_64 + LIT_MALE, data = afrikka)

Residuals:
    Min       1Q   Median       3Q      Max
-3741.6 -1056.8 -320.3   731.3  5478.8

Coefficients:
(Intercept) -22154.89      AGE15_64      2753.06      LIT_MALE      13.22
Estimate Std. Error t value Pr(>|t|)
(Intercept) -22154.89    2753.06   -8.047 1.41e-10 ***
AGE15_64     432.31     52.58    8.222 7.55e-11 ***
LIT_MALE     13.22     16.56    0.799 0.428

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

Residual standard error: 1799 on 50 degrees of freedom
Multiple R-Squared: 0.6213
F-statistic: 41.02 on 2 and 50 DF, p-value: 2.863e-11
```

Tulkinna: b₀ kun työlläisten suhteellinen osuus on nollla ja miesten lukutaitoprosentti on nollla, on maan BKT/asukas keskimäärin noin -22154 dollaria (ei käytännössä mielekäs tulos!)

b₁: Kun työlläisten suhteellinen osuus kasvaa yhdellä prosenttiyksiköllä ja samaan aikaan miesten lukutaitoprosentti ei muutu, kasvaa maan BKT/asukas keskimäärin noin 432 dollarilla.
 b₂: Kun miesten lukutaitoprosentti kasvaa yhdellä prosenttiyksiköllä ja samaan aikaan työlläisten suhteellinen osuus ei muutu, kasvaa maan BKT/asukas keskimäärin noin 13 dollarilla.