

# SF1901: SANNOLIKHETSLÄRA OCH STATISTIK

## FÖRELÄSNING 15.

### ENKEL LINJÄR REGRESSION

Jan Grandell & Timo Koski

07.03.2016



- Enkel linjär regression, repetition
- Enkel linjär regression: linjär algebra & numerisk analys
  - Modellens antaganden, teoretisk regressionslinje  $y = \alpha + \beta x$
  - Minstakvadratskattningarna  $\alpha_{\text{obs}}^*$  och  $\beta_{\text{obs}}^*$
  - Skattade regressionslinjen, residualkvadratsumman
- Analys av enkel linjär regression
  - Extra formler för minstakvadratskattningarna  $\alpha_{\text{obs}}^*$  och  $\beta_{\text{obs}}^*$
  - Normalfördelning, väntevärde och varians för  $\alpha^*$  och  $\beta^*$
- Konfidensintervall
  - Konfidensintervall för  $\alpha$  och  $\beta$
  - Konfidensintervall för  $\alpha + \beta x$
  - Prediktionsintervall för en ny observation
- Multipel regression: Statistisk prediktion av vinets kvalitet
- Bilaga i summaformler för regressionsanalys

Teorin och praktiken handlar om följande: Det föreligger  $n$  par av värden

$$(x_1, y_1), \dots, (x_n, y_n)$$

där  $x_1, \dots, x_n$  är givna storheter och  $y_1, \dots, y_n$  är observationer av oberoende s.v.  $Y_1, \dots, Y_n$ , där  $Y_i \in N(\mu_i, \sigma)$ .

Observera att  $\sigma$  förutsätts att ej bero av  $x$ , vilket ofta är det kritiska antagandet. Varje väntevärde  $\mu_i$  är linjärt beroende av  $x_i$ , d.v.s.

$$\mu_i = \alpha + \beta x_i, \quad i = 1, \dots, n.$$

Linjen

$$y = \alpha + \beta x \tag{1}$$

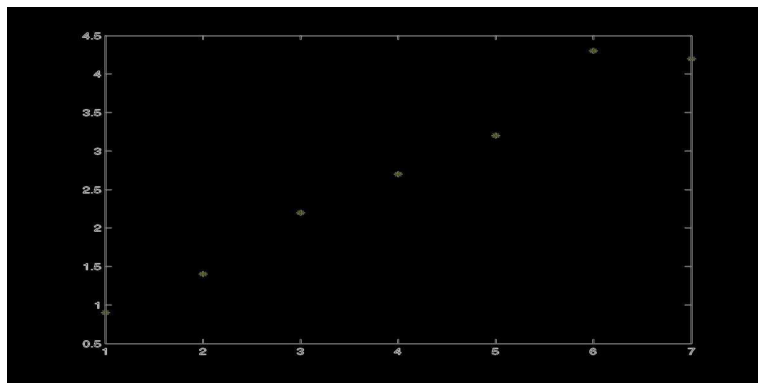
kallas den *teoretiska regressionslinjen*.



# ENKEL LINJÄR REGRESSION

Ett fiktivt exempel

$$\begin{array}{ccccccc} x = & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ y = & 0.9 & 1.4 & 2.2 & 2.7 & 3.2 & 4.3 & 4.2 \end{array}$$



Vi skattar parametrarna  $\alpha$  och  $\beta$  med Minsta-Kvadratmetoden, dvs. vi minimerar summan av de kvadrerade **lodräta** avstånden mellan den teoretiska linjen och  $y$ -värdena, d.v.s vi minimerar

$$Q(\alpha, \beta) = \sum_{i=1}^n (y_i - E(Y_i))^2 = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

m.a.p.  $\alpha$  och  $\beta$ . De värden  $\alpha_{\text{obs}}^*$  och  $\beta_{\text{obs}}^*$  som ger minimum kallas *MK-skattningarna* av  $\alpha$  och  $\beta$ .

Vi får

$$\alpha_{\text{obs}}^* = \bar{y} - \beta_{\text{obs}}^* \bar{x} \quad \text{och} \quad \beta_{\text{obs}}^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

# MINSTA KVADRAT VIA LINJÄR ALGEBRA & NUMERISK ANALYS

Låt oss skriva det ovanstående med en matrisformalism.

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \theta = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$n \times 1$  -vektor,  $n \times 2$  - matris,  $2 \times 1$  -vektor. Den teoretiska regressionmodellen är då

$$\mathbf{y} = \mathbf{X}\theta + \epsilon$$

där  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$  är en  $n \times 1$  vektor av normalfördelade variabler.



# MINSTA KVADRAT VIA LINJÄR ALGEBRA & NUMERISK ANALYS

För  $\mathbf{z}$  en  $n \times 1$  -vektor, definieras den euklidiska normen (i kvadrat) av  $\mathbf{z}$  som

$$\|\mathbf{z}\|^2 = \sum_{i=1}^n z_i^2.$$

Då gäller att

$$Q(\alpha, \beta) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = \|\mathbf{y} - \mathbf{X}\theta\|^2$$

och  $Q$  minimeras av ( $\mathbf{X}^T$  är den transponerade matrisen,  $^{-1}$  är matrisinversen)

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

I Matlab<sup>TM</sup> ges detta av

$$\hat{\theta} = \mathbf{X} \setminus \mathbf{y}$$





# MINSTA KVADRAT VIA LINJÄR ALGEBRA & NUMERISK ANALYS

Genom att multiplicera matriserna och utveckla inversen kan man verifiera att

$$\begin{pmatrix} \alpha_{\text{obs}}^* \\ \beta_{\text{obs}}^* \end{pmatrix} = \hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

d.v.s vi har samma lösning som ovan. Vi har att residualerna

$$\mathbf{y} - \mathbf{X}\hat{\theta}$$

är **ortogonala** mot kolonnrummet av matrisen  $\mathbf{X}$



Linjen

$$\hat{y} = \alpha_{\text{obs}}^* + \beta_{\text{obs}}^* x.$$

kallas den *skattade regressionslinjen*. De lodräta avstånden  $\varepsilon_i$  från  $y_i$  till den skattade regressionslinjen i  $x_i$ ,

$$\varepsilon_i = y_i - \alpha_{\text{obs}}^* - \beta_{\text{obs}}^* x_i$$

kallas *observerade residualer*.  $Q_0$  definieras som

$$Q_0 = Q(\alpha_{\text{obs}}^*, \beta_{\text{obs}}^*) = \sum_{i=1}^n \varepsilon_i^2.$$

och kallas residualkvadratsumman.

$\sigma^2$  skattas med

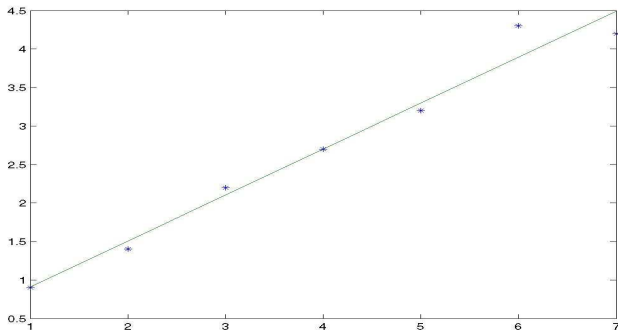
$$s^2 = \frac{Q_0}{n-2}.$$



# ENKEL LINJÄR REGRESSION

I exemplet ovan är skattade regressionslinjen, residualkvadratsumman och  $s^2$

$$y = 0.3143 + 0.5964x, \quad Q_0 = 0.2796, \quad s^2 = 0.0559$$



# FORMLERNÄ FÖR $\alpha_{\text{obs}}^*$ OCH $\beta_{\text{obs}}^*$ OMSKRIVNA

$$\alpha_{\text{obs}}^* = \bar{y} - \beta_{\text{obs}}^* \bar{x} \quad \text{och} \quad \beta_{\text{obs}}^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

$\Leftrightarrow$

$$\beta_{\text{obs}}^* = \frac{S_{xy}}{S_{xx}} \quad \text{och} \quad \alpha^* = \bar{y} - \beta_{\text{obs}}^* \bar{x}, \quad (2)$$

där vi sätter

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad \text{och} \quad S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2.$$

$$c_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

with  $S_x$  and  $S_y$ ,

$$S_x = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2}, S_y = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \bar{y})^2}$$

and get the *correlation coefficient* as

## DEFINITION

$$r_{xy} \stackrel{\text{def}}{=} \frac{c_{xy}}{S_x S_y}$$

$$|r_{xy}| \leq 1.$$

# RETURN TO THE MEAN (1)

$$\alpha_{\text{obs}}^* = \bar{y} - \beta_{\text{obs}}^* \bar{x}$$

and

$$\hat{y} = \alpha_{\text{obs}}^* + \beta_{\text{obs}}^* x = \bar{y} + \beta_{\text{obs}}^* (x - \bar{x})$$

i.e.,

$$\hat{y} - \bar{y} = \beta_{\text{obs}}^* (x - \bar{x}).$$



# RETURN TO THE MEAN (1)

$$\hat{y} - \bar{y} = \beta_{\text{obs}}^* (x - \bar{x}).$$

But by the above

$$\beta_{\text{obs}}^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Hence (factor in  $1/n$ )

$$\beta_{\text{obs}}^* = \frac{c_{xy}}{s_x s_y} \frac{s_y}{s_x} = r_{xy} \frac{s_y}{s_x}$$



# RETURN TO THE MEAN (2)

with

$$\beta_{\text{obs}}^* = r_{xy} \frac{S_y}{S_x}$$

we get

$$\hat{y} - \bar{y} = \beta_{\text{obs}}^* (x - \bar{x}) = r_{xy} \frac{S_y}{S_x} (x - \bar{x})$$

or

$$\frac{\hat{y} - \bar{y}}{S_y} = r_{xy} \frac{(x - \bar{x})}{S_x}$$





## RETURN TO THE MEAN (3)

If  $-1 < r_{xy} < 1$ , then we see by

$$\frac{\hat{y} - \bar{y}}{S_y} = r_{xy} \frac{(x - \bar{x})}{S_x}$$

that the predicted standardized value  $\hat{y}$  of  $y$  is closer to its mean  $\bar{y}$  than the standardized value of  $x$  is to its mean  $\bar{x}$ . Thus the data points  $(x_i, y_i)_{i=1}^n$  display regression toward the mean (as found and formulated by Francis Galton).



# RETURN TO THE MEAN (4)

The idea of regression toward the mean resolved an important difficulty of Darwinian selection:

*if offspring were always identical to parents, then evolution by natural selection was not possible. But, on the other hand, there was also intergenerational stability, as all experience under fairly constant environmental conditions showed that the range of variability on short time scales, as between two generations, was essentially constant. (Stephen Stigler)*



Vi kommer nu att ägna oss åt en rätt så minutiös exercis i omskrivning av formlerna för  $\alpha_{\text{obs}}^*$  och  $\beta_{\text{obs}}^*$ . De (i sig enkla) räkneregler som kommer till användning här ges i Bilagan till detta dokument (se nedan). Slutresultaten ingår avsnitt 13.3 i Formelsamlingen.

Vinsten med detta är att vi överför skattningarna på en form som gör det tydligt, hur vi får normalfördelningarna för de mot svarande  $\alpha_{\text{obs}}^*$  och  $\beta_{\text{obs}}^*$  stickprovsvariablerna  $\alpha^*$  och  $\beta^*$  (se avsnitt 13.1 i Formelsamlingen). Detta ger i sin tur s.g.s automatiskt ett antal konfidensintervall (se avsnitt 13.2 i Formelsamlingen) för enkel linjär regression.

Korssumman  $S_{xy}$  kan skrivas på flera sätt och vi har (se (7) & (8) i Bilagan)

$$\begin{aligned} S_{xy} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i \\ &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}. \end{aligned}$$

Vi har här utnyttjat att summorna  $\sum_{i=1}^n (x_i - \bar{x})$  (och  $\sum_{i=1}^n (y_i - \bar{y})$  båda) är 0. Se formel (6) i Bilagan .

Motsvarande likhet gäller för kvadratsumman  $S_{xx}$ , formel (9) i Bilagan och/eller avsnitt 13.3 i Formelsamlingen:

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2.$$

Det är viktigt att notera att båda skattningarna är linjära uttryck i observationerna  $y_i$ . Man har ju (använd formel (4) i Bilagan upprepade ggr)

$$\begin{aligned}\beta_{\text{obs}}^* &= \frac{S_{xy}}{S_{xx}} = \frac{1}{S_{xx}} \left( \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right) = \\ &= \frac{1}{S_{xx}} \left( \sum_{i=1}^n x_i y_i - \bar{x} \sum_{i=1}^n y_i \right) = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right) y_i = \sum_{i=1}^n c_i y_i\end{aligned}$$

$$\begin{aligned}\alpha^* &= \bar{y} - \beta_{\text{obs}}^* \bar{x} = \frac{1}{n} \sum_{i=1}^n y_i - \bar{x} \sum_{i=1}^n c_i y_i = \\ &= \sum_{i=1}^n \left( \frac{1}{n} - \bar{x} c_i \right) y_i = \sum_{i=1}^n d_i y_i.\end{aligned}$$

D.v.s.

$$\beta_{\text{obs}}^* = \sum_{i=1}^n c_i y_i \quad \text{och} \quad \alpha_{\text{obs}}^* = \sum_{i=1}^n d_i y_i$$

där

$$c_i = (x_i - \bar{x}) / S_{xx} \quad \text{och} \quad d_i = \frac{1}{n} - c_i \bar{x}. \quad (3)$$

De motsvarande stickprovsvariablerna är

$$\beta^* = \sum_{i=1}^n c_i Y_i \quad \text{och} \quad \alpha^* = \sum_{i=1}^n d_i Y_i$$



Därmed är stickprovsvariablerna  $\alpha^*$  och  $\beta^*$  linjära funktioner i  $Y$ -variablerna, kom ihåg att  $x$ en är givna tal.

$$\beta^* = \sum_{i=1}^n c_i Y_i \quad \text{och} \quad \alpha^* = \sum_{i=1}^n d_i Y_i$$

Därmed är  $\beta^*$  och  $\alpha^*$  normalfördelade ! Det gäller att beräkna de respektive väntevärdena och varianserna.

Vi inleder med

$$E(\alpha^*) = \sum_{i=1}^n d_i E(Y_i).$$

Men vi vet att  $Y_i \in N(\alpha + \beta x_i, \sigma)$ , så att

$$\sum_{i=1}^n d_i E(Y_i) = \sum_{i=1}^n d_i (\alpha + \beta x_i) = \alpha \sum_{i=1}^n d_i + \beta \sum_{i=1}^n d_i x_i$$

Vi kommer att visa att

- (1)  $\sum_{i=1}^n d_i = 1$
- (2)  $\sum_{i=1}^n d_i x_i = 0$

# VÄNTEVÄRDET PÅ $\alpha^*$ : (1) $\sum_{i=1}^n d_i = 1$

Håll i minnet att  $c_i = (x_i - \bar{x})/S_{xx}$ .

$$\begin{aligned}\sum_{i=1}^n d_i &= \sum_{i=1}^n \left( \frac{1}{n} - c_i \bar{x} \right) = \sum_{i=1}^n \frac{1}{n} - \sum_{i=1}^n c_i \bar{x} \\ &= 1 - \bar{x} \sum_{i=1}^n c_i = 1 - \bar{x} \frac{1}{S_{xx}} \sum_{i=1}^n (x_i - \bar{x}) \\ &= 1 - \bar{x} \frac{1}{S_{xx}} \cdot 0 = 1.\end{aligned}$$

ty  $\sum_{i=1}^n (x_i - \bar{x}) = 0$  enligt (6) i Bilagan. Observera att vi även fick

$$\sum_{i=1}^n c_i = 0.$$



# VÄNTEVÄRDET PÅ $\alpha^*$ : (2) $\sum_{i=1}^n d_i x_i = 0$

$$\begin{aligned}\sum_{i=1}^n d_i x_i &= \sum_{i=1}^n \left( \frac{1}{n} - c_i \bar{x} \right) x_i \\ &= \frac{1}{n} \sum_{i=1}^n x_i - \bar{x} \sum_{i=1}^n c_i x_i.\end{aligned}$$

Men

$$\begin{aligned}\sum_{i=1}^n c_i x_i &= (1/S_{xx}) \sum_{i=1}^n (x_i - \bar{x}) x_i = (1/S_{xx}) \sum_{i=1}^n (x_i^2 - \bar{x} x_i) \\ &= (1/S_{xx}) \left( \sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i \right) =\end{aligned}$$

och från ovan

$$= (1/S_{xx}) \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = (1/S_{xx}) \cdot S_{xx} = 1.$$

# VÄNTEVÄRDET PÅ $\alpha^*$ : (2) $\sum_{i=1}^n d_i x_i = 0$

Så vi fick att  $\sum_{i=1}^n c_i x_i = 1$  och detta ger

$$\begin{aligned}\sum_{i=1}^n d_i x_i &= \sum_{i=1}^n \left( \frac{1}{n} - c_i \bar{x} \right) x_i \\ &= \frac{1}{n} \sum_{i=1}^n x_i - \bar{x} \underbrace{\sum_{i=1}^n c_i x_i}_{=1} \\ &= \frac{1}{n} \sum_{i=1}^n x_i - \bar{x} = 0.\end{aligned}$$

# VÄNTEVÄRDET PÅ $\alpha^*$ :

Således har vi visat att

$$E(\alpha^*) = \alpha \sum_{i=1}^n d_i + \beta \sum_{i=1}^n d_i x_i = \alpha.$$

dvs.  $\alpha^*$  är väntevärdesriktig.



Vissa hjälpsresultat har vi dessutom visats:

(A)

$$\sum_{i=1}^n c_i x_i = 1$$

(B)

$$\sum_{i=1}^n c_i = 0,$$

(C)

$$\sum_{i=1}^n d_i x_i = 0$$

# VÄNTEVÄRDET PÅ $\beta^*$

Vi har  $\beta^* = \sum_{i=1}^n c_i Y_i$  och detta ger

$$E(\beta^*) = \sum_{i=1}^n c_i E(Y_i).$$

Men  $Y_i \in N(\alpha + \beta x_i, \sigma)$ , så att

$$\sum_{i=1}^n c_i E(Y_i) = \sum_{i=1}^n c_i (\alpha + \beta x_i) = \alpha \sum_{i=1}^n c_i + \beta \sum_{i=1}^n c_i x_i$$

Hjälpresultaten (A) och (B) ovan ger

$$E(\beta^*) = \beta$$

och  $\beta^*$  är väntevärdesriktig.



Vi ha p.g.a oberoendet mellan  $Y_i$ na att

$$V(\alpha^*) = \sum_{i=1}^n d_i^2 V(Y_i).$$

Men  $Y_i \in N(\alpha + \beta x_i, \sigma)$ , så att

$$\sum_{i=1}^n d_i^2 V(Y_i) = \sigma^2 \sum_{i=1}^n d_i^2$$

Näst visar vi att

- $\sum_{i=1}^n d_i^2 = \frac{1}{n} + \frac{\bar{x}}{S_{xx}}$

$$\begin{aligned}\sum_{i=1}^n d_i^2 &= \sum_{i=1}^n \left( \frac{1}{n} - c_i \bar{x} \right)^2 \\ &= \sum_{i=1}^n \frac{1}{n^2} - \frac{2\bar{x}}{n} \sum_{i=1}^n c_i + \bar{x}^2 \sum_{i=1}^n c_i^2\end{aligned}$$

Men hjälpresultatet (B) ovan ger  $\sum_{i=1}^n c_i = 0$ . Av definition på  $c_i$  och  $S_{xx}$  fås

$$\sum_{i=1}^n c_i^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / S_{xx}^2 = \frac{1}{S_{xx}}$$

och  $\sum_{i=1}^n \frac{1}{n^2} = 1/n$ . Sammanfattningsvis

$$\sum_{i=1}^n d_i^2 = \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}.$$

$$V(\alpha^*) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right).$$

I härledningen fick vi även hjälpresultatet

(D)

$$\sum_{i=1}^n c_i^2 = \frac{1}{S_{xx}}$$

Vi har  $\beta^* = \sum_{i=1}^n c_i Y_i$  och p.g.a oberoendet mellan  $Y_i$ na att

$$V(\beta^*) = \sum_{i=1}^n c_i^2 V(Y_i).$$

Men  $Y_i \in N(\alpha + \beta x_i, \sigma)$ , så att

$$\sum_{i=1}^n c_i^2 V(Y_i) = \sigma^2 \sum_{i=1}^n c_i^2 = \frac{\sigma^2}{S_{xx}}$$

enligt det senaste hjälpresultatet (D).

## SATS

$$\alpha^* \in N \left( \alpha, \sigma \sqrt{\left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right)$$
$$\beta^* \in N \left( \beta, \sigma \sqrt{\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right).$$

## SATS

Vidare gäller att, här betyder  $Q_0$  stickprovsvariabeln  $Q(\alpha^*, \beta^*)$ ,

$$\frac{Q_0}{\sigma^2} = \frac{(n-2)S^2}{\sigma^2}$$

är  $\chi^2(n-2)$ -fördelad och att  $S^2$  är oberoende av  $\alpha^*$  och  $\beta^*$ .

En algebraisk exercis i ovanstående stil visar även att

$$\begin{aligned} Q_0 &= Q(\alpha_{\text{obs}}^*, \beta_{\text{obs}}^*) = \sum_{i=1}^n \varepsilon_i^2 = \\ &= \sum_{i=1}^n (y_i - \alpha_{\text{obs}}^* - \beta_{\text{obs}}^* x_i)^2 = S_{yy} - \frac{S_{xy}}{S_{xx}}. \end{aligned}$$

Det ovanstående gör att vi kan konstruera konfidensintervall och test som förut, både då  $\sigma$  är känt och okänt. För att inte behöva skriva alla intervall två gånger så betraktar vi fallet då  $\sigma$  är okänt.





Eftersom  $\alpha$  nu förekommer som en parameter, så ger vi konfidensintervall med konfidensgrad 95%. Metoden är definiera ett  $\theta$  och sedan bilda

$$I_{\theta} = \theta_{\text{obs}}^* \pm t_{0.025}(n-2)d(\theta^*),$$

där  $d(\theta^*)$  är skattningen av  $D(\theta^*)$ .

Vi använder som skattning av  $\sigma$  i konfidensintervallen

$$s = \sqrt{\frac{Q_0}{n-2}}$$

(dvs.kvadratroten av residualkvadratsumman/ $n-2$ )

$$I_{\theta} = \theta_{\text{obs}}^* \pm t_{0.025}(n-2)d(\theta^*),$$

(1)  $\theta = \alpha$

Vi har  $V(\theta^*) = V(\alpha^*) = \sigma^2 \cdot \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)}$ , vilket ger

$$I_{\alpha} = \alpha^* \pm t_{0.025}(n-2)s\sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)}.$$

$$(2) \quad \underline{\theta = \beta}$$

Vi har  $V(\beta^*) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$ , vilket ger

$$I_\beta = \beta_{\text{obs}}^* \pm t_{0.025}(n-2) \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

# Hypotesprövning för $\beta$

Det kan hända att inget samband i form av en teoretisk regressionslinje finns, dvs.  $\beta = 0$ . Vi kan testa detta genom

$$H_0 : \beta = 0$$

mot

$$H_1 : \beta \neq 0$$

Förkasta  $H_0$  på nivån 0.05 om

$$0 \notin I_\beta = \beta_{\text{obs}}^* \pm t_{0.025}(n-2) \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

# Konfidensintervall för regressionslinjen

$$(3) \quad \theta = \alpha + \beta x_0$$

Detta är alltså ett konfidensintervall för den teoretiska regressionslinjen i punkten  $x = x_0$ .

Vi observerar först

$$E(\alpha^* + \beta^* x_0) = E(\alpha^*) + E(\beta^* x_0) = \alpha + \beta x_0.$$

$$\begin{aligned} V(\alpha^* + \beta^* x_0) &= V\left(\sum_{i=1}^n d_i Y_i + \sum_{i=1}^n c_i Y_i x_0\right) \\ &= V\left(\sum_{i=1}^n (d_i + c_i x_0) Y_i\right) \end{aligned}$$

och p.g.a. oberoendet

$$= \sum_{i=1}^n (d_i + c_i x_0)^2 V(Y_i) = \sigma^2 \sum_{i=1}^n (d_i + c_i x_0)^2.$$



$d_i = \frac{1}{n} - c_i \bar{x}$  så att

$$\begin{aligned}\sum_{i=1}^n (d_i + c_i x_0)^2 &= \sum_{i=1}^n \left( \frac{1}{n} + c_i (x_0 - \bar{x}) \right)^2 = \\ &= \sum_{i=1}^n \frac{1}{n^2} + 2 \frac{(x_0 - \bar{x})}{n} \sum_{i=1}^n c_i + (x_0 - \bar{x})^2 \sum_{i=1}^n c_i^2 = \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx}.\end{aligned}$$

ty hjälpresultatet (B) ovan ger  $\sum_{i=1}^n c_i = 0$  och hjälpresultatet (D) ger  $\sum_{i=1}^n c_i^2 = 1/S_{xx}$ .

Vi slår ihop från ovan

$$V(\alpha^* + \beta^* x_0) = \sigma^2 \cdot \left( \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx} \right).$$

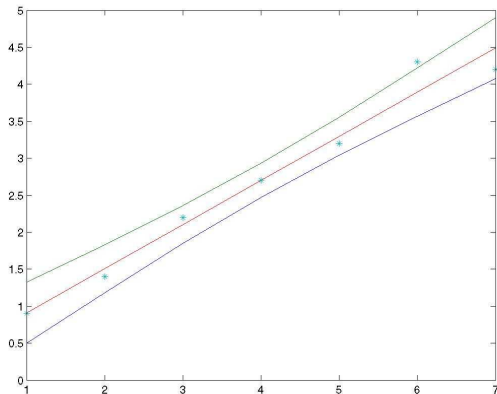
Detta ger det slutliga konfidensintervallet som

$$I_{\alpha+\beta x_0} = \alpha_{\text{obs}}^* + \beta_{\text{obs}}^* x_0 \pm t_{0.025}(n-2) s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$



# Konfidensintervall för linjen i exemplet ovan

Data, regressionslinje och  $\alpha_{\text{obs}}^* + \beta_{\text{obs}}^* x_0 \pm t_{0.025}(n-2)s\sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$



Låt  $Y_0$  vara en ny observation som svarar mot  $x_0$  dvs.  $Y_0 \in N(\alpha + \beta x_0, \sigma)$ .  
Vi predikterar  $Y_0$  genom  $\hat{Y}_0$  som ges av den skattade regressionslinjen

$$\hat{Y}_0 = \alpha_{\text{obs}}^* + \beta_{\text{obs}}^* x_0.$$

Prediktionsfelet

$$Y_0 - \hat{Y}_0$$

är normalfördelat med väntevärdet noll. Eftersom  $Y_0$  och  $\hat{Y}_0$  är oberoende

$$V(Y_0 - \hat{Y}_0) = V(Y_0) + V(\hat{Y}_0)$$

Ovan har vi i fallet (3) visat att

$$V(\hat{Y}_0) = V(\alpha^* + \beta^* x_0) = \sigma^2 \cdot \left( \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx} \right).$$

Alltså

$$V(Y_0 - \hat{Y}_0) = \sigma^2 \left( 1 + \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx} \right)$$

Vi skattar  $\sigma^2$  med  $s^2$ .

Man kan visa att kvoten

$$\frac{Y_0 - \hat{Y}_0}{s \cdot \sqrt{1 + \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx}}}$$

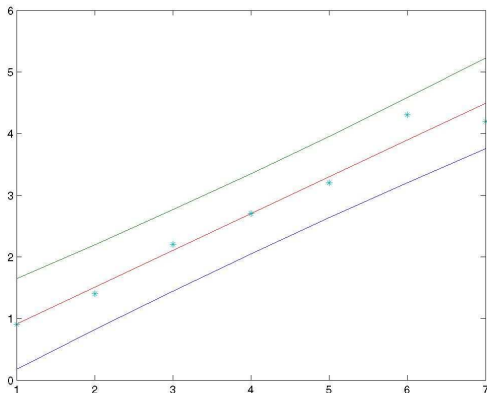
är t-fördelad med  $n - 2$  frihetsgrader. Då kan vi på bekant manér visa att intervallet

$$\hat{Y}_0 \pm t_{0.025}(n - 2)s \cdot \sqrt{1 + \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx}}$$

täcker  $Y_0$  med sannolikheten 0.95. Intervallet kallas ett 95%-igt prediktionsintervall för  $Y_0$

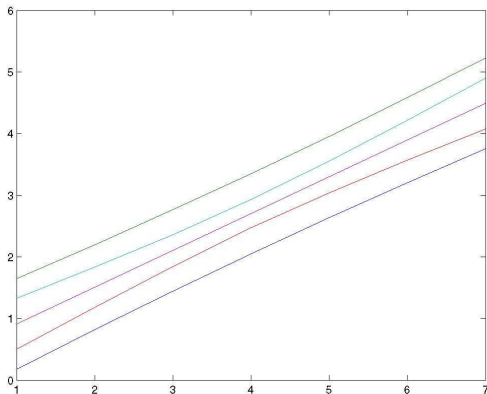
## Prediktionsintervall i exemplet ovan

Data, regressionslinje och  $\hat{Y}_0 \pm t_{0.025}(n-2)s \cdot \sqrt{1 + \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx}}$



# Prediktionsintervall och konfidensintervall i exemplet ovan

Regressionslinje,  $\alpha_{\text{obs}}^* + \beta_{\text{obs}}^* x_0 \pm t_{0.025}(n-2)s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$  och  
 $\hat{Y}_0 \pm t_{0.025}(n-2)s \cdot \sqrt{1 + \frac{1}{n} + (x_0 - \bar{x})^2 / S_{xx}}$



# MULTIPEL LINJÄR REGRESSIONSANALYS: VINETS KVALITET

Multipel regression med tre förklarande variabler  $x_1, x_2, x_3$  för statistisk prediktion (SPR) av (Bordeaux-)vinernas kvalitet:

$$\text{vinets kvalitet} = 12.145 + 0.00117x_1 + 0.0614x_2 - 0.00386x_3$$

där

$x_1 =$  nederbörd under vintern

$x_2 =$  genomsnittstemperaturen under växtsäsongen

$x_3 =$  nederbörd under skördeperioden

av Orley Ashenfelter, se <http://www.liquidasset.com/> .



# MULTIPEL LINJÄR REGRESSIONSANALYS: VINETS KVALITET

Statistisk prediktion (SPR) av (Bordeaux-)vinernas kvalitet:

$$\text{vinets kvalitet} = 12.145 + 0.00117x_1 + 0.0614x_2 - 0.00386x_3$$

Poängen här alltså att vi vill uttala oss om kvaliteten hos en årgång av (Bordeaux-)vin INNAN NÅGON smakat på vinet. Detta tycks fungera, men upplevs trots detta (eller kanske just p.g.a detta) som stötande av många goda vinexperter.



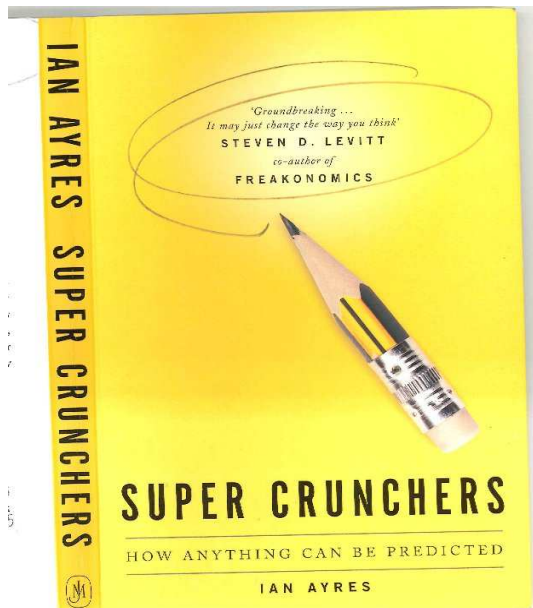


For a very wide range of prediction problems, statistical prediction rules (SPRs), often rules that are very easy to implement, make predictions that are as reliable as, and typically more reliable than, human experts. The success of SPRs forces us to reconsider our views about what is involved in understanding, explanation, and good reasoning.

Orley Ashenfelters SPR för vinets kvalitet är hämtad ur boken

I. Ayres: Super Crunchers. How anything can be predicted. John Murray (Publishers), Paperback edition 2008, London.





## DEFINITION

$$(1) \quad \underline{\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n.}$$

## SATS

$$(2) \quad \underline{\sum_{i=1}^n a \cdot x_i = a \sum_{i=1}^n x_i.}$$

## SATS

$$(3) \quad \underline{\sum_{i=1}^n (x_i + y_i) = \sum_{i=1}^n x_i + \sum_{i=1}^n y_i.}$$

## SATS

$$(4) \quad \underline{\sum_{i=1}^n (ax_i + by_i) = a \sum_{i=1}^n x_i + b \sum_{i=1}^n y_i}$$

## SATS

$$(5) \quad \underline{\sum_{i=1}^n (x_i + y_i)^2 = \sum_{i=1}^n x_i^2 + 2 \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2.}$$

Låt  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ . Då gäller

SATS

$$(6) \quad \underline{\sum_{i=1}^n (x_i - \bar{x}) = 0.}$$

SATS

$$(7) \quad \underline{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x}) y_i = \sum_{i=1}^n x_i (y_i - \bar{y}).}$$

## SATS

$$(8) \quad \underline{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}.}$$

## SATS

$$(9) \quad \underline{\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2.}$$