

Föreläsning 13, Matematisk statistik 7.5 hp för E

Multipel linjär regression

Stas Volkov

Linjär regression

Vi har n st par av mätvärden (x_ℓ, y_ℓ) , $\ell = 1, \dots, n$ där y_ℓ är observationer av

$$Y_\ell = \alpha + \beta x_\ell + \varepsilon_\ell$$

där ε_ℓ är oberoende av varandra, och $\varepsilon_\ell \in \mathbf{N}(0, \sigma)$.

Parameterskattningarna

Skattningarna av α^* , β^* och $(\sigma^2)^*$ är

$$\alpha^* = \bar{y} - \beta^* \cdot \bar{x}, \quad \beta^* = \frac{\sum_{\ell=1}^n (x_\ell - \bar{x})(y_\ell - \bar{y})}{\sum_{\ell=1}^n (x_\ell - \bar{x})^2} = \frac{S_{xy}}{S_{xx}},$$

$$(\sigma^2)^* = s^2 = \frac{Q_0}{n-2}, \quad Q_0 = \sum_{\ell=1}^n (y_\ell - \alpha^* - \beta^* x_\ell)^2 = S_{yy} - \frac{S_{xy}^2}{S_{xx}}$$

Skattningarnas fördelning:

$$\alpha^* \in \mathbf{N} \left(\alpha, \sigma \cdot \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} \right), \quad \beta^* \in \mathbf{N} \left(\beta, \frac{\sigma}{\sqrt{S_{xx}}} \right)$$

Men de är **inte oberoende** av varandra.

Konfidens-, prediktions- och kalibreringsintervall ($f = n - 2$):

$$I_\alpha = \alpha^* \pm t_{a/2}(f) \cdot s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}},$$

$$I_\beta = \beta^* \pm t_{a/2}(f) \cdot \frac{s}{\sqrt{S_{xx}}},$$

$$I_{\mu(x_0)} = \alpha^* + \beta^* x_0 \pm t_{a/2}(f) \cdot s \cdot \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}},$$

$$I_{Y(x_0)} = \alpha^* + \beta^* x_0 \pm t_{a/2}(f) \cdot s \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}},$$

$$I_{x_0} = x_0^* \pm t_{a/2}(f) \cdot \frac{s}{|\beta^*|} \cdot \sqrt{1 + \frac{1}{n} + \frac{(y_0 - \bar{y})^2}{(\beta^*)^2 S_{xx}}}.$$

Linjärisering av exponentiella samband (vid behov)

För att få ett linjärt samband

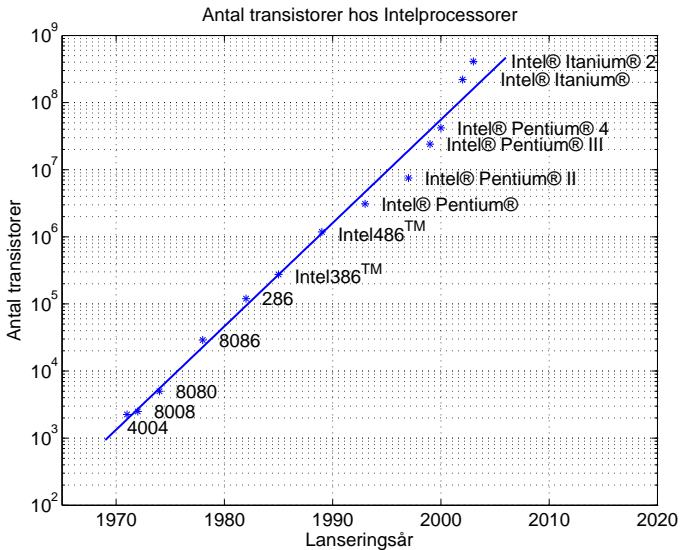
$$y_l = \alpha + \beta x_l + \varepsilon_l$$

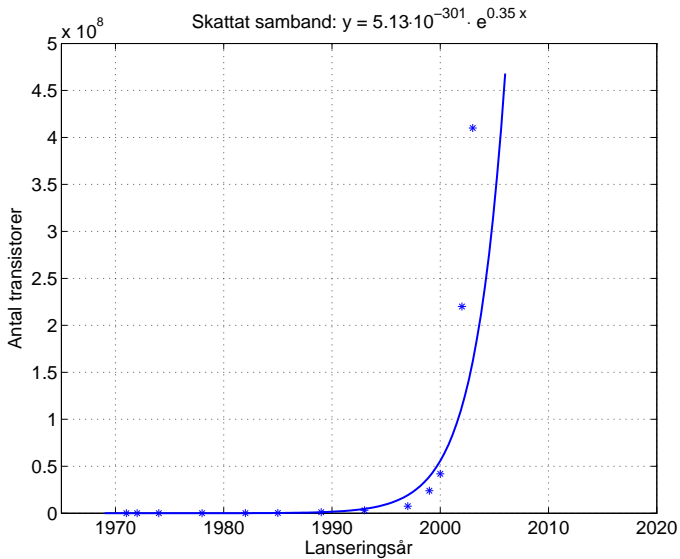
kan vissa exponent- och potenssamband logaritmeras.

$$z_l = a \cdot e^{\beta x_l} \cdot \varepsilon'_l \quad \xrightarrow{\ln} \quad \underbrace{\ln z_l}_{y_l} = \underbrace{\ln a}_{\alpha} + \beta \cdot x_l + \underbrace{\ln \varepsilon'_l}_{\varepsilon_l}$$

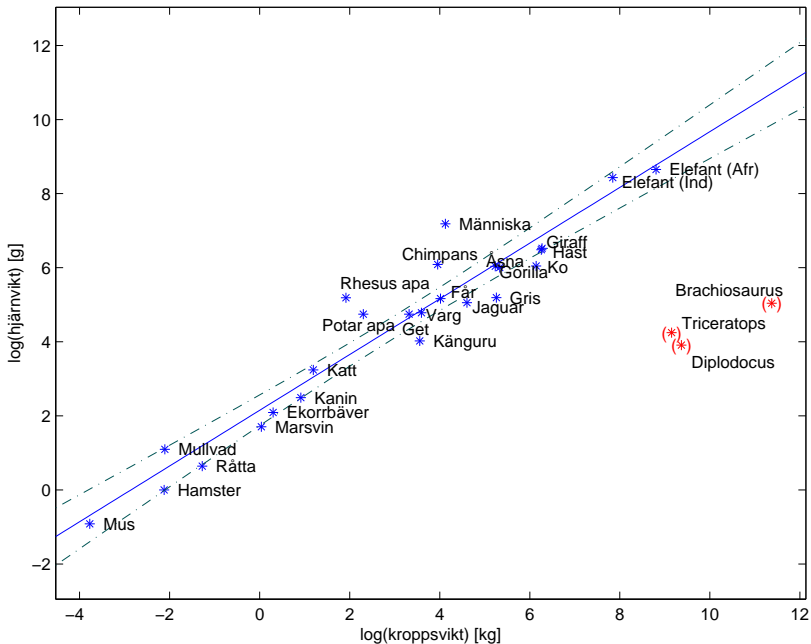
$$z_l = a \cdot t_l^\beta \cdot \varepsilon'_l \quad \xrightarrow{\ln} \quad \underbrace{\ln z_l}_{y_l} = \underbrace{\ln a}_{\alpha} + \beta \underbrace{\ln t_l}_{x_l} + \underbrace{\ln \varepsilon'_l}_{\varepsilon_l}$$

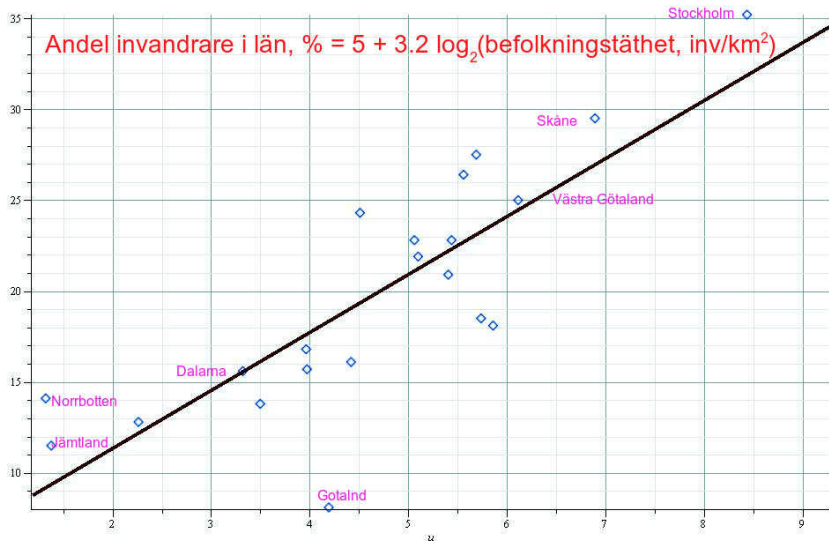
Om de multiplikativa felen, ε'_l , är lognormalfördelade blir då $\ln \varepsilon'_l \in \mathbf{N}$ och vi kan använda linjär regression för att skatta $\ln \alpha$ och β .





Samband vikt och hjärnstorlek





Multipel linjär regression

Modellen kan utökas med flera \mathbf{x} -variabler:

$$y_\ell = \beta_0 + \beta_1 x_{\ell 1} + \dots + \beta_k x_{\ell k} + \varepsilon_\ell, \quad \ell = 1, \dots, n, \quad \varepsilon_\ell \in \mathbf{N}(0, \sigma)$$

som kan skrivas på matrisform som

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

där \mathbf{y} och $\boldsymbol{\varepsilon}$ är $n \times 1$ -vektorer, $\boldsymbol{\beta}$ en $(k+1) \times 1$ -vektor och \mathbf{X} en $n \times (k+1)$ -matris

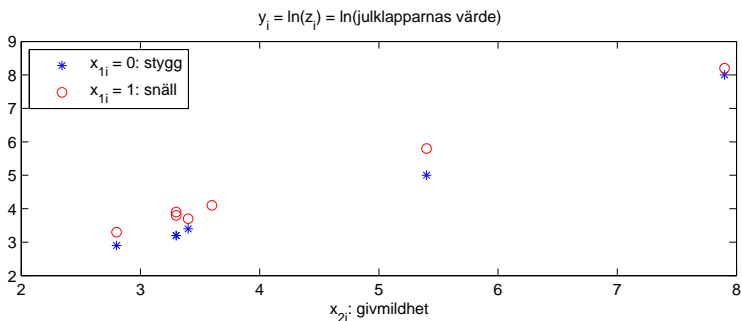
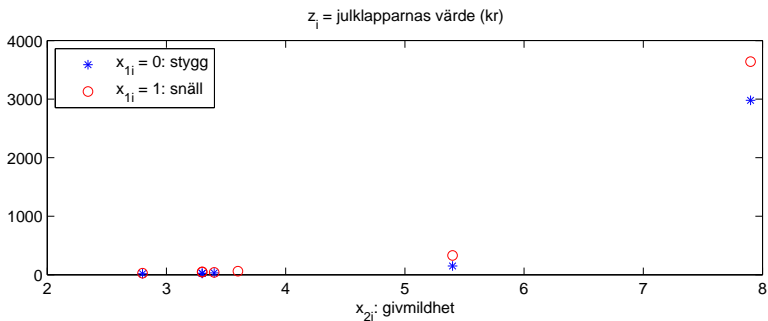
dvs.

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Exempel – Julklappar:

En liten flicka vill undersöka om det lönar sig att vara **snäll**. Hon har därför noterat värdet på de julklappar hon fick från olika släktingar i år, när hon varit **snäll**, och i fjol då hon var **styggt**. Hon har insett att värdet på julklapparna också till stor del beror på givarens ekonomi och allmänna generositet. Hon räknar därför också ut ett lämpligt mått på **givmildhet**.

Släkting	värde (kr.)		ln(värde)		givmildhet
	i fjol	i år	i fjol	i år	
Storebror	24.5	49.5	3.2	3.9	3.3
Lillebror	18.	27.	2.9	3.3	2.8
Mormor och morfar	2981.	3641.	8.0	8.2	7.9
Farmor och farfar	30.	40.	3.4	3.7	3.4
Mamma och pappa	148.	329.50	5.0	5.8	5.4
Moster	24.5	44.5	3.2	3.8	3.3
Kusin	?	62.	?	4.1	3.6



Lämplig regressionsmodell:

$$\ln z_\ell = y_\ell = \alpha + \beta_1 \cdot x_{1\ell} + \beta_2 \cdot x_{2\ell} + \varepsilon_\ell, \quad \ell = 1, \dots, 13$$

där e^{β_1} = relativa ökningen i värde när flickan är snäll.

Responsvariabeln:

z_ℓ = värdet (i kronor) av julklapp ℓ ,

$\implies y_\ell = \ln z_\ell$ = logaritmerat värde på julklapp ℓ

Förklarande variablerna:

$$x_{1\ell} = \begin{cases} 0 & \text{för alla fjolårets julklappar (då hon varit stygg)} \\ 1 & \text{för alla årets julklappar (då hon varit snäll)} \end{cases},$$

$x_{2\ell}$ = givmildheten hos givaren av julklapp ℓ

$\varepsilon_\ell \in N(0, \sigma)$ oberoende

- ▶ Testa på nivån **5%**, om det lönar sig att vara snäll, dvs om β_1 är signifikant större än **0**.
- ▶ Gör ett tväsidigt **95%** prediktionsintervall för värdet på Kusinens julklapp i fjol, dvs. då den lilla flickan varit stygg.

Modell med matriser: $Y = X\beta + \varepsilon$ där $n = 13$ samt $k = 2$ och

$$Y = \begin{bmatrix} 3.2 \\ 3.9 \\ 2.9 \\ 3.3 \\ 8.0 \\ 8.2 \\ 3.4 \\ 3.7 \\ 5.0 \\ 5.8 \\ 3.2 \\ 3.8 \\ 4.1 \end{bmatrix} \quad X = \begin{bmatrix} 1 & 0 & 3.3 \\ 1 & 1 & 3.3 \\ 1 & 0 & 2.8 \\ 1 & 1 & 2.8 \\ 1 & 0 & 7.9 \\ 1 & 1 & 7.9 \\ 1 & 0 & 3.4 \\ 1 & 1 & 3.4 \\ 1 & 0 & 5.4 \\ 1 & 1 & 5.4 \\ 1 & 0 & 3.3 \\ 1 & 1 & 3.3 \\ 1 & 1 & 3.6 \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_{12} \\ \varepsilon_{13} \end{bmatrix}$$

Skattning av parametrarna

Skattning av β

ML- och MK-skattningar av β_0, \dots, β_k (elementen i β) blir

$$\beta^* = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{(k+1) \times (k+1)} \underbrace{\mathbf{X}^T}_{k \times n} \underbrace{\mathbf{y}}_{n \times 1}$$

En väntevärdesriktig skattning av σ^2 ges av (korrigerad ML)

$$s^2 = \frac{Q_0}{n - (k + 1)} \quad \text{där } Q_0 = (\mathbf{y} - \mathbf{X}\beta^*)^T (\mathbf{y} - \mathbf{X}\beta^*)$$

Q_0 är alltså residualkvadratsumman och $k + 1$ är antalet skattade parametrar i Q_0 (obs. tumregeln!).

Skattningar:

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 13 & 7 & 55.8 \\ 7 & 7 & 29.7 \\ 55.8 & 29.7 & 278.5 \end{bmatrix} \quad \mathbf{X}^T \mathbf{y} = \begin{bmatrix} 58.5 \\ 32.8 \\ 289.09 \end{bmatrix}$$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 0.6530 & -0.1786 & -0.1118 \\ -0.1786 & 0.3098 & 0.0028 \\ -0.1118 & 0.0028 & 0.0257 \end{bmatrix}$$

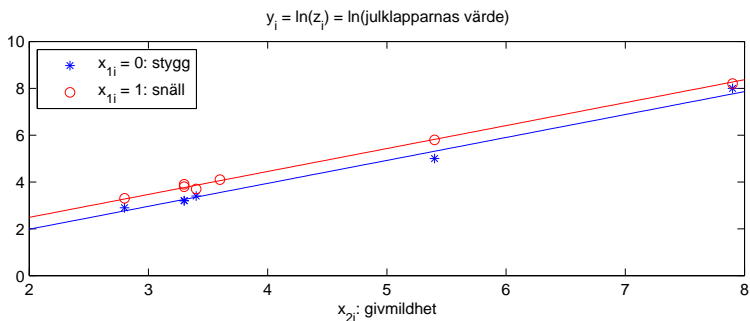
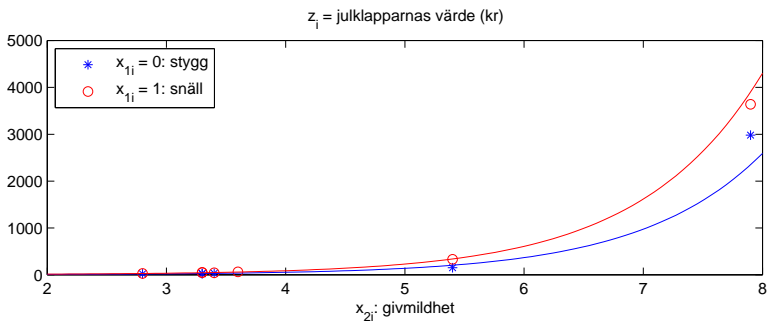
så

$$\beta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \begin{bmatrix} \beta_0^* \\ \beta_1^* \\ \beta_2^* \end{bmatrix} = \begin{bmatrix} 0.0208 \\ 0.5074 \\ 0.9799 \end{bmatrix}$$

$$Q_0 = (\mathbf{y} - \mathbf{X}\beta^*)^T (\mathbf{y} - \mathbf{X}\beta^*) = 0.2347$$

$$\sigma^* = s = \sqrt{\frac{Q_0}{f}} = 0.1532$$

$$\text{samt } f = n - (k + 1) = 13 - 3 = 10$$



Skattningarnas fördelning

Skattningarna av β är **linjära funktioner** av \mathbf{Y} och är därmed normalfördelade

$$\beta_l^* \in \mathbf{N}(\beta_l, \mathbf{D}(\beta_l^*)),$$

där $\mathbf{D}(\beta_l^*)$ ges av roten ur diagonalelementen i **kovariansmatrisen**

$$\mathbf{V}(\beta^*) = \underbrace{\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}}_{\substack{(k+1) \times (k+1) \\ \text{matris}}} = \begin{bmatrix} \mathbf{V}(\beta_0^*) & \mathbf{C}(\beta_0^*, \beta_1^*) & \cdots & \mathbf{C}(\beta_0^*, \beta_k^*) \\ \mathbf{C}(\beta_1^*, \beta_0^*) & \mathbf{V}(\beta_1^*) & \cdots & \mathbf{C}(\beta_1^*, \beta_k^*) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{C}(\beta_k^*, \beta_0^*) & \mathbf{C}(\beta_k^*, \beta_1^*) & \cdots & \mathbf{V}(\beta_k^*) \end{bmatrix}.$$

För residualkvadratsumman gäller

$$\frac{Q_0}{\sigma^2} \in \chi^2(n - (k + 1))$$

Konfidsensintervall och hypotestest för β_ℓ

Konfidsensintervall för β_ℓ blir alltså

$$I_{\beta_\ell} = \beta_\ell^* \pm t_{\alpha/2}(f) \cdot d(\beta_\ell^*) = \beta_\ell^* \pm t_{\alpha/2}(f) \cdot s \cdot \sqrt{[(\mathbf{X}^T \mathbf{X})^{-1}]_{\ell,\ell}}$$

där $[(\mathbf{X}^T \mathbf{X})^{-1}]_{\ell,\ell}$ är diagonalelement nr ℓ , $f = n - k - 1$.

Obs! det första elementet har nummer $\ell = 0$.

Intervallet kan användas för att testa hypotesen

$$H_0 : \beta_\ell = 0 \quad \text{kontra} \quad H_1 : \beta_\ell \neq 0$$

Alternativt kan man naturligtvis använda

$$T = \frac{\beta_\ell^* - 0}{d(\beta_\ell^*)} \quad \text{och förkasta } H_0 \text{ om } |T| > t_{\alpha/2}(n - (k + 1)).$$

(1) Vi vill testa $H_0: \beta_1 = 0$ mot $H_1: \beta_1 > 0$ på signifikansnivån $\alpha = 0.05$. Medelfelet blir

$$d(\beta_1^*) = s\sqrt{[(\mathbf{X}^T\mathbf{X})^{-1}]_{1,1}} = 0.1532 \cdot \sqrt{0.3098} = 0.0853$$

således (obs! ensidig...)

$$T = \frac{\beta_1^* - 0}{d(\beta_1^*)} = \frac{0.5074}{0.0853} = 5.9496 > t_{\alpha}(f) = t_{0.05}(10) = 1.81$$

(dvs. ”för mycket”) kan H_0 förkastas. **Ja**, det lönar sig att vara snäll.
Hur mycket lönar det sig? Ett tvåsidigt konfidensintervall för β_1 blir

$$\begin{aligned} I_{\beta_1} &= \beta_1^* \pm t_{\alpha/2}(f) \cdot d(\beta_1^*) = 0.5074 \pm \underbrace{t_{0.025}(10)}_{2.23} \cdot 0.0853 \\ &= (0.3174, 0.6974) \Rightarrow I_{e^{\beta_1}} = (e^{0.3174}, e^{0.6974}) = (1.37, 2.01) \end{aligned}$$

Att vara snäll ökar värdet på julklapparna med i genomsnitt **37-101%**!

Skattning av punkt på ”planet”

\mathbf{Y} -s väntevärde i en punkt $\mathbf{x}_0 = [\mathbf{1} \quad x_{01} \quad x_{02} \quad \cdots \quad x_{0k}]$ ges nu av

$$\mu^*(\mathbf{x}_0) = \beta_0^* + \sum_{\ell=1}^k \beta_\ell^* x_{0\ell} = \mathbf{x}_0 \beta^*$$

$$\text{med } V(\mu^*(\mathbf{x}_0)) = \mathbf{x}_0 V(\beta^*) \mathbf{x}_0^T = \sigma^2 \underbrace{\mathbf{x}_0}_{1 \times k} \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{k \times k} \underbrace{\mathbf{x}_0^T}_{k \times 1}$$

Ett **konfidensintervall** för $\mu^*(\mathbf{x}_0)$ blir således (med $f = n - (k + 1)$)

$$I_{\mu^*(\mathbf{x}_0)} = \underbrace{\mathbf{x}_0}_{1 \times k} \underbrace{\beta^*}_{k \times 1} \pm t_{\alpha/2}(f) \cdot s \cdot \sqrt{\mathbf{x}_0 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0^T}$$

För **prediktionsintervallet** får man **lägga till en etta** under kvadratroten, som tidigare

$$I_{Y(\mathbf{x}_0)} = \mathbf{x}_0 \beta^* \pm t_{\alpha/2}(f) \cdot s \cdot \sqrt{\mathbf{1} + \mathbf{x}_0 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0^T}$$

(2) Prediktionsintervall för Kusinens julklapp:

Vi har $\mathbf{x}_0 = [1 \quad 0 \quad 3.6]$ och skattningen

$$\mu^*(\mathbf{x}_0) = \mathbf{x}_0\boldsymbol{\beta}^* = 1 \cdot \beta_0^* + 0 \cdot \beta_1^* + 3.6 \cdot \beta_2^* = 3.5484,$$

$$e^{\mu^*(\mathbf{x}_0)} = e^{3.5484} = 35.76 \text{ kr},$$

$$\mathbf{x}_0(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{x}_0^T = 0.18,$$

$$\text{därför } I_{Y(\mathbf{x}_0)} = 3.5484 \pm 2.23 \cdot 0.1532\sqrt{1 + 0.18} = (3.21, 3.89)$$

Omräknat till kronor blir det

$$I_{e^{Y(\mathbf{x}_0)}} = (e^{3.21}, e^{3.89}) = (25.69, 48.94) \text{ kr}$$

Modellvalidering

Precis som för enkel regression bör man undersöka **residualerna**

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*,$$

och förvissa sig om att de verkar vara oberoende och $\mathbf{N}(\mathbf{0}, \sigma)$ -fördelade.
Plotta residualerna

- ▶ ”Som de kommer”, dvs mot $1, 2, \dots, n$. Ev. ett histogram
- ▶ Mot var och en av \mathbf{x}_ℓ -dataserierna
- ▶ I en normalfördelningsplot

För var och en av β_1, \dots, β_k (obs! i regel $\text{ej } \beta_0$) bör man kunna förkasta H_0 i testet

$$H_0: \beta_\ell = 0$$

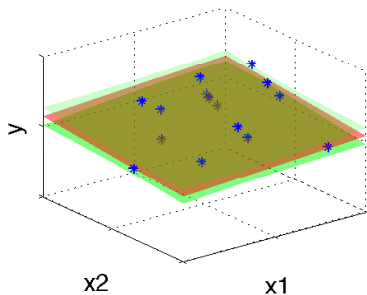
$$H_1: \beta_\ell \neq 0$$

eftersom β_ℓ anger ”hur mycket \mathbf{y} ändrar sig när vi ändrar \mathbf{x}_ℓ ”.

Kolinjäritet (ex. två variabler, motsv. för fler)

Man bör om möjligt välja sina (x_{1l}, x_{2l}) -värden så att de blir utspridda i (x_1, x_2) -planet och inte "klumpar ihop" sig längs en linje. Detta ger "en mer stabil grund" åt regressionsplanet.

Låg kolinjäritet



Hög kolinjäritet

