

Plan:

- ① Lineær regresjon
- ② Repetisjon av hele kursen

Ekstra forelesning: 31/05

Repetisjon av kurset + gjennomgang av prøve-eksamen. (Kap 6-7)

Kontortid: Bortsett hvis alle  
Ellers kontortid 10-18 ca.

Kursevaluering : Husk å svare på kursevalueringen  
når den kommer!

Jeg kommer til å legge ut en prøve-eksamen til i  
god tid før forelesningen 31/05!

Prøve-eksamen I = Oppsattsett 15

Prøve-eksamen II (kommer senere)

① Linear regresjon

X: forklaringsvariabel

Y: responsvariabel

} Under å studere.  
Sammenheng mellom  
X og Y

(a) Korrelasjonskoeffisient

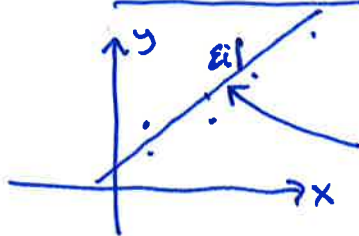
$$R = \frac{S_{XY}}{S_X \cdot S_Y} = \frac{\frac{1}{n-1} \cdot \sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \cdot \sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2}}$$

Husk:  $-1 \leq R \leq 1$  $R > 0$  positiv sammenheng $R < 0$  negativ " " " $|R|$  : hvor sterk sammenheng er

$x_i$	$y_i$
$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	$\vdots$
$x_n$	$y_n$

(b) Linear regresjon:

$$Y = \alpha + \beta X + \varepsilon \quad \hat{y} = \alpha + \beta X$$



Spredingsdiagram

regresjons-  
linjen:  
passer best  
mulig med  
dataene

Forutsetninger:

- $\varepsilon$  er  $N(0, \sigma)$  for en konstant  $\sigma$
- Når vi trekker  $n$  punkter (datasettet), er  $\varepsilon_1, \dots, \varepsilon_n$  uavhengige.

Beste estimat for regresjonslinjen:

$$y = \hat{\alpha} + \hat{\beta} X$$

$$Y_1 = \alpha + \beta x_1 + \varepsilon_1$$

 $\vdots$ 

$$Y_n = \alpha + \beta x_n + \varepsilon_n$$

$$\hat{\beta} = \frac{S_{XY}}{S_X^2} = R \cdot \frac{S_Y}{S_X}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \cdot \bar{X}$$

Beste estimat ved minste kvadraters metode, dvs  $\varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2$  er minst mulig.

Statistiske egenskaper: $\hat{\alpha}, \hat{\beta}$  estimerer: Stokastiske variable

$$SE(\hat{\alpha})^2 = \frac{\sigma^2 \cdot \sum x_i^2}{n (\sum x_i^2 - n\bar{x}^2)}$$

$$SE(\hat{\beta})^2 = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2}$$

Estimat for  $\sigma^2$ :

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{der } \hat{y}_i = \hat{\alpha} + x_i \hat{\beta}$$

Tolkning av  $R^2$ :

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\left\{ \begin{aligned} y_i - \bar{y} &= (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}) \\ &= \hat{\epsilon}_i + (\hat{y}_i - \bar{y}) \end{aligned} \right.$$

↑  
feil-ledd,  
residual

↑  
forklart  
av  
regresjonslinjen

Man kan vise:

$$1) \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SST = SSE + SSR$$

$$\begin{aligned} 2) R^2 &= \frac{SSR}{SST} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \\ &= \frac{\frac{1}{n-1} \sum (\hat{y}_i - \bar{y})^2}{\frac{1}{n-1} \sum (y_i - \bar{y})^2} = \frac{SSR}{s_y^2} \end{aligned}$$

Hvis  $\hat{\theta}$  er en  
estimator for  $\theta$ ,  
så er:

$$E(\hat{\theta}) = \theta \quad (\text{foru. rett})$$

$$SE(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$$

meist mulig

$$SE(\hat{\theta}) = \text{std. avvik til estimatoren } \hat{\theta}$$

Ex:  $\bar{X}$  estimator for  $\mu$

$$E(\bar{X}) = \mu$$

$$SE(\bar{X}) = \sigma/\sqrt{n}$$

når  $x_1, \dots, x_n \sim N(\mu, \sigma)$

Tolkning:

$R^2$  er andelen av  
variansen i  $y$   
som er forklart  
av regresjonslinjen

Utledning av  $S^2$ :

$$S^2 = \frac{1}{n-2} \cdot \sum_{i=1}^n \varepsilon_i^2 = \frac{1}{n-2} SSE$$

$$R^2 = \frac{SSR}{SST}$$

~~U~~

$$S^2 = \frac{1}{n-2} SSE = \frac{1}{n-2} (SST - SSR) \quad \downarrow$$

$$= \frac{1}{n-2} \cdot (SST - R^2 SST) = \frac{1}{n-2} SST \cdot (1 - R^2)$$

$$= \frac{1}{n-2} \cdot (n-1) S_y^2 \cdot (1 - R^2) = \frac{n-1}{n-2} S_y^2 (1 - R^2)$$

Ekse: Hypotese test for  $\beta$   $\left\{ \begin{array}{l} H_0: \beta = 0 \quad (\beta_0 = 0) \\ H_1: \beta \neq 0 \end{array} \right.$

Testobservator:  $T = \frac{\hat{\beta} - \beta_0}{SE(\hat{\beta})} = \frac{\hat{\beta}}{SE(\hat{\beta})}$

t-fordelt med  $n-2$  dk hvis  $\beta = \beta_0$

Forkastningsområde:  $|T| > t_{\alpha/2}^{n-2}$

Beregning av  $SE(\hat{\beta})$ :

$$SE(\hat{\beta})^2 = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2} \rightarrow \frac{S^2}{\sum x_i^2 - n\bar{x}^2} =$$

estimer  $\sigma^2$  vha  $S^2$

$$\rightarrow \frac{\frac{n-1}{n-2} S_y^2 (1 - R^2)}{(n-1) S_x^2} = \frac{1}{n-2} \frac{S_y^2}{S_x^2} (1 - R^2)$$

$$SE(\hat{\beta}) \text{ estimeres ved } \sqrt{\frac{1}{n-2} \frac{S_y^2}{S_x^2} (1 - R^2)}$$

$$SE(\hat{\alpha})^2 = \frac{\sigma^2 \sum x_i^2}{n \cdot (\sum x_i^2 - n\bar{x}^2)}$$

estimeres med

$$\frac{S^2 \sum x_i^2}{n (\sum x_i^2 - n\bar{x}^2)} = \frac{\frac{n-1}{n-2} S_y^2 (1-R^2) \sum x_i^2}{n \cdot (n-1) S_x^2} = \frac{1}{n(n-2)} \frac{S_y^2}{S_x^2} (1-R^2) \sum x_i^2$$

Konfiansintervall:

$$\text{For } \beta: \hat{\beta} \pm t_{\alpha/2}^{n-2} \cdot SE(\hat{\beta})$$

$$\text{For } \alpha: \hat{\alpha} \pm t_{\alpha/2}^{n-2} \cdot SE(\hat{\alpha})$$

② Repetisjon:① Sannsynlighetsregning

Kap 3 - 5

Stokastisk forsøkutfallsrom  $U = \{\omega_1, \omega_2, \dots\}$ Sannsynlighetsmål: funksjon  $p$  som til enhver hendelse  $A$  (dvs delmengde av  $U$ )tilordner et tall  $p(A)$  slik ati)  $p(A) \geq 0$  for enhver hendelse  $A$ ii)  $p(U) = 1$ iii)  $p(A_1 \cup A_2 \cup \dots) = p(A_1) + p(A_2) + \dots$   
når  $A_i, A_j$  er disjunkte  
(uten felles utfall)  
for alle  $i, j$ .Konsekvenser:

i)  $p(A^c) = 1 - p(A)$

ii)  $p(A \cup B) = p(A) + p(B) - p(A \cap B)$

Spesialtilfelle: Uniform sannsynlighet = alle utfall har samme sanns. het.

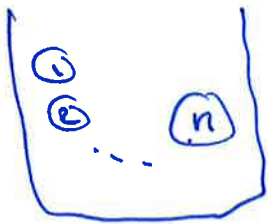
Da har vi

$$p(A) = \frac{\text{antall gunstige utfall}}{\text{antall mulige utfall}}$$



Kombinatorikk: Å telle opp ting

Urne modell: trekker  $r$  blant  $n$  kuler fra en urne



i) ordnet, med tilbakelegging:  $n^r$   
 ii) ordnet, uten tilbakelegging:  $n \cdot (n-1) \cdot \dots \cdot (n-r+1)$

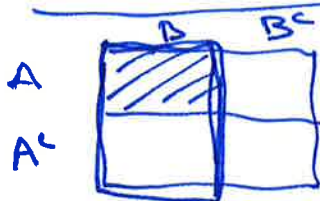
$$\boxed{nPr}$$

iii) uordnet, uten tilbakelegging:  $\frac{n(n-1) \cdot \dots \cdot (n-r+1)}{r!}$

på kalkulator:  $\boxed{nCr}$

skrivebrev:  $\rightarrow \binom{n}{r}$

Betinget sannsynlighet:



$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

sannsynligheten for A, gitt at B inntreffer

Chavensighet:

A og B er uavhengige hvis

$$P(A \cap B) = P(A) \cdot P(B)$$

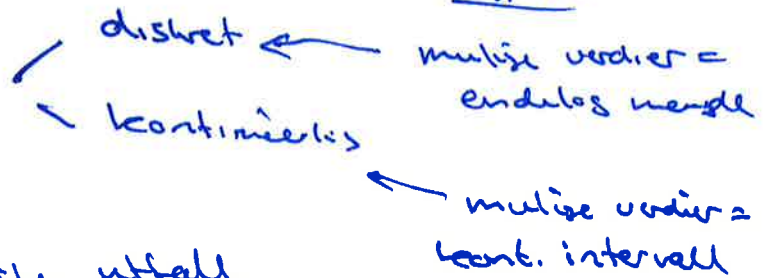
$\Leftrightarrow$

$$\frac{P(A \cap B)}{P(B)} = P(A)$$

$$P(A|B) = P(A)$$

Stokastiske variabler:  $X$  som tar verdier avhengig av utfallet til et stokastisk forsøk

Typisk:



i) Diskret:

fordelingsfunksjon =  $P$   
Sannsynligheter for ulike utfall

$x$	$p(x=x)$
$x_1$	$p(x_1)$
$x_2$	$p(x_2)$
$\vdots$	$\vdots$
$x_n$	$p(x_n)$

Forventningsverdi:

$$E(x) = \sum_x x \cdot p(x) \quad \text{forventningsverdi}$$

$$E(g(x)) = \sum_x g(x) \cdot p(x)$$

$$\text{Var}(x) = E[(x - E(x))^2] \quad \text{varians}$$

Simultan fordeling:  $X$  og  $Y$  simultant fordelt  $p(x,y) = p(X=x, Y=y)$

	$x_1$	$x_2$	$\dots$	$x_n$
$y_1$	$p(x_1, y_1)$	$\dots$		$p(Y=y_1)$
$y_2$	$\vdots$			
$\vdots$				
$y_n$	$\vdots$			
	$\underline{\hspace{2cm}}$			
	$p(X=x_1)$			

kolonne- og radsummer = marginale fordelingsfun.

Forventningsverdi:

$$E(x) = \sum_x x \cdot p_X(x) \quad E(y) = \sum_y y \cdot p_Y(y)$$

$$\begin{aligned} \text{Cov}(x,y) &= \sum_{x,y} (x - E(x))(y - E(y)) \\ &= E(xy) - E(x) \cdot E(y) \end{aligned}$$



Regne regler:

i)  $E(ax + b) = aE(x) + b$

ii)  $Var(x) = E(x^2) - E(x)^2$

iii)  $Var(ax + b) = a^2 \cdot Var(x)$

iv)  $Cov(x, y) = E(xy) - E(x) \cdot E(y)$

v)  $Cov(x_1 + x_2, y_1 + y_2) = Cov(x_1, y_1) + Cov(x_1, y_2) + Cov(x_2, y_1) + Cov(x_2, y_2)$

$Cov(x, y) = Cov(y, x)$

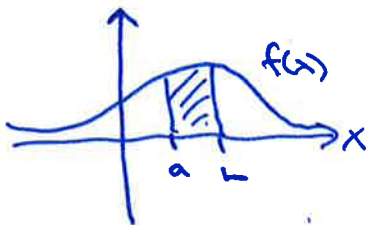
$Cov(cX, Y) = Cov(X, cY) = c \cdot Cov(X, Y)$

vi)  $Var(x) = Cov(x, x)$

ii) kontinuerlig:Tetthetsfunksjonen  $f(x)$ 

$$P(a \leq x \leq b) = \int_a^b f(x) dx$$

punkt sannsynlig:  $P(x=a) = \underline{0}$

Førutsetning:  $\infty$ 

$$E(x) = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

$$E(g(x)) = \int_{-\infty}^{\infty} g(x) \cdot f(x) dx$$

$$Var(x) = E[(x - E(x))^2]$$

Alle regneregler overfor gjelder også for kontinuerlige stokastiske variable.

Sannsynlighetsfordelinger:

Binomisk:  $X =$  antall forekster av en hendelse  $A$   
 når vi gir et forsøk der  $P(A) = p$   
 $n$  ganger, og disse er uavhengige

diskret  
variabelparametre:  $n, p$   $q = 1 - p$ Mulige  
verdier: $0, 1, 2, \dots, n$ 

$$P(X=i) = \binom{n}{i} p^i q^{n-i}, \quad i=0, 1, \dots, n$$

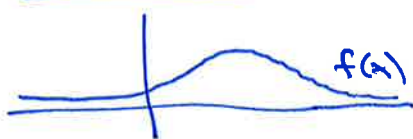
$$E(X) = np$$

$$\text{Var}(X) = npq$$

Normalfordeling: $N(\mu, \sigma)$ normalfordeling med  
forventning  $\mu$ , std. avvik  $\sigma$ 

$$E(X) = \mu, \quad \text{Var}(X) = \sigma^2$$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$



i)  $X$  normalfordelt  $\Rightarrow ax+b$  normalfordelt  
 $X, Y$  normalfordelt  $X+Y$  — | —

ii)  $X_1, \dots, X_n, \dots$   
uavhengige, identisk  
fordelte  $\Rightarrow \bar{X}_n = \frac{1}{n}(X_1 + X_2 + \dots + X_n)$   
 $\rightarrow$  normalfordelt  
 når  $n \rightarrow \infty$

Normaltilnærning til binomisk fordeling:Hvis  $X$  er binomisk fordelt med parametre  $n, p$   $0 <$  $n$  er stor, så er  $X$  tilnærmet  $N(\mu, \sigma)$  $(npq \geq 5)$ 

$$\mu = E(X) = np \quad \sigma^2 = \text{Var}(X) = npq$$

For å få en mest mulig nøyaktig tilnærming, bruker vi helkallskorrelasjon:

$$p(x \leq a) \approx \Phi\left(\frac{a + 0,5 - \mu}{\sigma}\right) \quad \text{hvor } \begin{cases} \mu = np \\ \sigma^2 = npq \end{cases}$$

↑  
korrelasjon = 0,5