

Emne	Kompendium
1 Repetisjon	
2 Kontinuerlige stokastiske variabler	[I] 1.1
3 Forventning og varians	[I] 1.2 - 1.4
4 Median og andre prosentiler	[I] 1.5 - 1.6

### Oppgaver for Forelesning 7

Oppgaver fra arbeidsboken	[DA] 7.1 - 7.2
Oppgaver fra kompendiet	[I] 1.15, 1.19, 1.22 - 1.24

### ① Repetisjon

$$f(\underline{x}) = \underline{x}^T A \underline{x} \quad \text{kvadr. form} = \underline{u}^T D \underline{u} = \lambda_1 u_1^2 + \lambda_2 u_2^2 + \dots + \lambda_n u_n^2$$

Variableksifte:

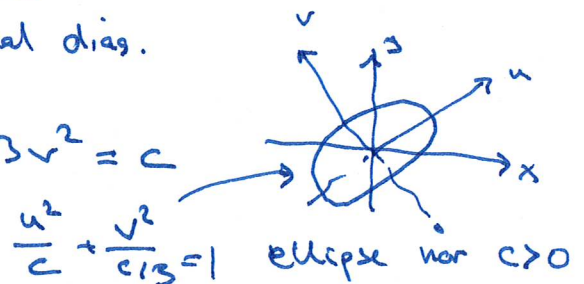
$$\underline{x} = P \cdot \underline{u}$$

der  $P^T A P = D$

er en ortogonal diag.

$$f(x,y) = 2x^2 - 2xy + 2y^2 = u^2 + 3v^2 = c$$

$$A = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \quad \lambda_1 = 1 \\ \lambda_2 = 3$$



Optimering:

$$f(\underline{x}) = \underline{x}^T A \underline{x} + B \underline{x} + C$$

$$f'(\underline{x}) = 2A \underline{x} + B^T$$

Stasjonære pkt:

$$f'(\underline{x}) = 2A \underline{x} + B^T = \underline{0}$$

A pos defn / pos semidefn: f konvex

A neg " neg " : f konkav

**6.8.** Gitt et datasett med  $N$  observasjoner av variabelen  $y$  som skal forklares og av de  $n$  forklaringsvariable  $x_1, x_2, \dots, x_n$ , så ønsker vi å finne den lineære likningen

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

som gir beste tilnærming til dataene i datasettet. Vi kan sjelden finne en tilnærming uten feil, og må derfor istedet å gjøre feilen  $\varepsilon$  gitt ved

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \varepsilon$$

minst mulig. Vi antar at datasettet er gitt ved følgende tabell (der hver linje svarer til en observasjon):

	$x_1$	$x_2$	$\dots$	$x_n$	$y$
1	$x_{11}$	$x_{21}$	$\dots$	$x_{n1}$	$y_1$
2	$x_{12}$	$x_{22}$	$\dots$	$x_{n2}$	$y_2$
3	$x_{13}$	$x_{23}$	$\dots$	$x_{n3}$	$y_3$
$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\vdots$
$N$	$x_{1N}$	$x_{2N}$	$\dots$	$x_{nN}$	$y_N$

Vi får det lineært likningssystem  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  der

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad X = \begin{pmatrix} 1 & x_{11} & x_{21} & \dots & x_{n1} \\ 1 & x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{1N} & x_{2N} & \dots & x_{nN} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{pmatrix}$$

Vi sier at feilen er minst mulig når  $\|\boldsymbol{\varepsilon}\|^2 = \varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_N^2$  er minst mulig, og metoden kalles derfor *minste kvadraters metode*. Vis at problemet har entydig løsning

$$\boldsymbol{\beta} = (X^T X)^{-1} \cdot X^T \mathbf{y}$$

når  $\det(X^T X) \neq 0$ .

**6.9.** Estimer  $\beta_0, \beta_1$  i den lineære regresjonsmodellen  $Y = \beta_0 + \beta_1 x + \varepsilon$  ut fra de tre observasjonene

$$(x_1, y_1) = (0, 1)$$

$$(x_2, y_2) = (1, 0)$$

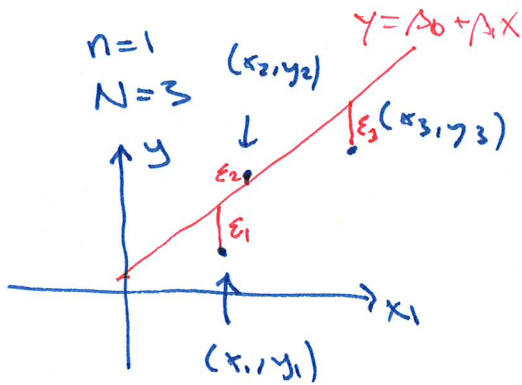
$$(x_3, y_3) = (1, 1)$$

og beregn feilleddene  $e_1, e_2$  og  $e_3$ .

$$\mathbf{y} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

LDA3 6.8

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$



Minste kvadraters metode

Velg  $\beta_0, \beta_1, \dots, \beta_n$  slik at

$$\epsilon_1^2 + \epsilon_2^2 + \dots + \epsilon_N^2 \text{ blir minst mulig.}$$

Dvs:

$$\min_{(\beta_0, \dots, \beta_n)} \epsilon_1^2 + \epsilon_2^2 + \dots + \epsilon_N^2$$

	$x_1$	$x_2$	...	$x_n$	$y$
1	$x_{11}$	$x_{21}$		$x_{n1}$	$y_1$
2	$x_{12}$	$x_{22}$		$x_{n2}$	$y_2$
⋮	⋮	⋮		⋮	⋮
N					

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{21} + \dots + \beta_n x_{n1} + \epsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_{12} + \beta_2 x_{22} + \dots + \beta_n x_{n2} + \epsilon_2$$

$$\vdots$$

$$\underline{y} = \underbrace{\begin{pmatrix} 1 & x_{11} & x_{21} & \dots \\ 1 & x_{12} & x_{22} & \dots \end{pmatrix}}_X \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{pmatrix}$$

Liten på matriseform:

$$\underline{y} = X \cdot \underline{\beta} + \underline{\epsilon} \quad \leftarrow \quad \underline{\epsilon} = \underline{y} - X \underline{\beta}$$

$$\begin{aligned} \|\underline{\epsilon}\|^2 &= \underline{\epsilon} \cdot \underline{\epsilon} = \underline{\epsilon}^T \underline{\epsilon} = (\underline{y} - X \underline{\beta})^T \cdot (\underline{y} - X \underline{\beta}) = (\underline{y}^T - \underline{\beta}^T X^T) (\underline{y} - X \underline{\beta}) \\ &= \underline{y}^T \cdot \underline{y} - \underline{\beta}^T X^T \underline{y} + \underline{y}^T (-X \underline{\beta}) + \underline{\beta}^T X^T X \underline{\beta} \\ &= \underbrace{\underline{\beta}^T (X^T X)}_A \underline{\beta} - \underbrace{\underline{y}^T X}_B \underline{\beta} - \underbrace{(\underline{\beta}^T X^T \underline{y})^T}_{\underline{y}^T X \underline{\beta}} - \underbrace{\underline{y}^T \underline{y}}_C \end{aligned}$$

$1 \times (n+1)$     $(n+1) \times N$     $N \times 1$   
 $1 \times 1$ -matrise

$$\|\underline{\varepsilon}\|^2 = \underline{\beta}^T (X^T X) \underline{\beta} - 2(Y^T X) \underline{\beta} + \mathbf{1}^T Y = \underbrace{\beta^T A \beta}_{X^T X} + \underbrace{b \beta}_{-2Y^T X} + \underbrace{c}_{Y^T Y}$$

$$\min_{(\beta_0, \dots, \beta_n)} \varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_N^2 \rightsquigarrow \min f(\beta) = \beta^T (X^T X) \beta - 2Y^T X \beta + \mathbf{1}^T Y$$

Stasjonære pkt:  $-f'(\underline{\beta}) = 2(X^T X) \underline{\beta} + (-2Y^T X)^T$

$$= 2X^T X \underline{\beta} - 2X^T Y = \underline{0} \quad | :2$$

$$X^T X \underline{\beta} - X^T Y = \underline{0}$$

$$(X^T X) \cdot \underline{\beta} = X^T \cdot Y$$

Brukes:

i)  $X^T X$  er pos. semidefn.

ii)  $|X^T X| \neq 0 \Rightarrow X^T X$  invertibel

$(n+1) \times (n+1)$

$(n+1) \times 1$

$$\underline{\beta} = (X^T X)^{-1} \cdot X^T Y$$

$f$  konveks  $\Rightarrow \underline{\beta} = (X^T X)^{-1} \cdot (X^T Y)$   
er global min  
for  $f(\underline{\beta}) = \varepsilon_1^2 + \dots + \varepsilon_N^2$



## ② Kontinuertlige stokastiske variabler

### Stokastiske forsøk:

- mengden av mulige utfall (utfallsrommet) er kjent
- det enkelte utfall er ukjent

Stokastisk variabel:  $X$  variabel med verdi som avhenger av utfallet i et stokastisk forsøk

Eks: Vi kaster to terninger

$X =$  summen  
(diskret)

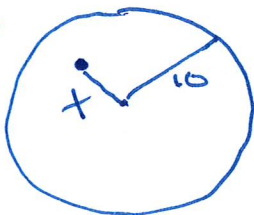
$Y =$  max

$U = \{(1,1), (1,2), \dots, (6,6)\}$

Mulige  $X$ -verdier:  $2, \dots, 12$

Ser på tilfellet med en kontinuerlig stokastisk variabel

Eks:



Vi kaster pil på en balle

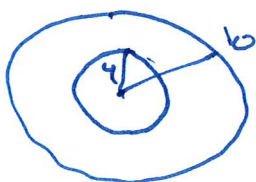
$X =$  avstand til sentrum

Mulige  $X$ -verdier:  $[0, 10]$

$$P(X \leq 4) = \frac{\pi \cdot 4^2}{\pi \cdot 10^2} = \frac{16}{100}$$

$$P(X \leq b) = \frac{\pi \cdot b^2}{\pi \cdot 10^2} = \frac{b^2}{100}$$

$(0 \leq b \leq 10)$



$$F(b) = b^2/100 \quad F(x) = x^2/100$$

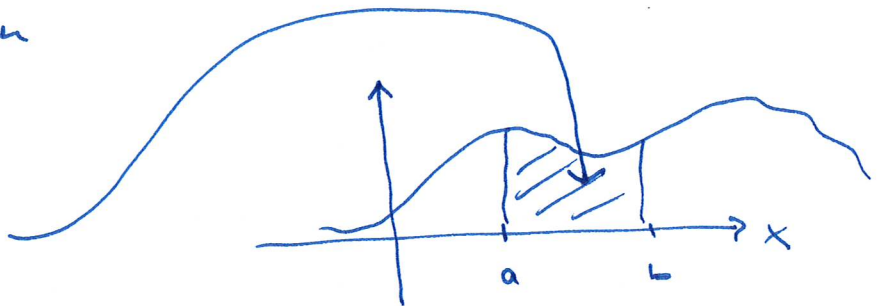
$$\Rightarrow f(x) = F'(x) = 2x/100 = x/50$$

For en kontinuerlig stokastisk variabel fins det en funksjon  $f(x)$  som kalles sannsynlighetstettheten til  $X$

Som har egenskapen

$$P(a \leq X \leq b)$$

$$= \int_a^b f(x) dx$$

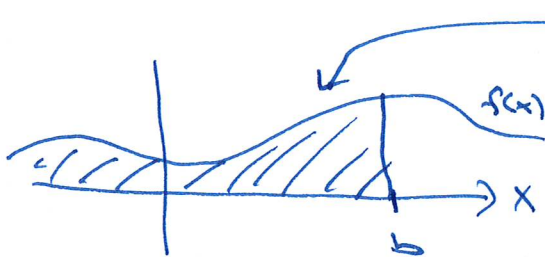


Krav til tetthetsfunksjon:

i)  $f(x) \geq 0$  for alle  $x$

ii)  $\int_{-\infty}^{\infty} f(x) dx = 1$

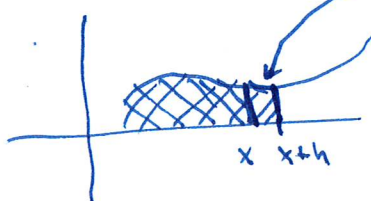
Kumulativ fordelingsfunksjon for en kontinuerlig stokastisk variabel  $X$ :  $F(x) = P(X \leq x)$       $F(b) = P(X \leq b)$



$$F(b) = \int_{-\infty}^b f(x) dx \quad F(x) = \int_{-\infty}^x f(x) dx$$

Integralregningens fundamentalsats:  $F'(x) = f(x)$

$$F'(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x)}{h}$$

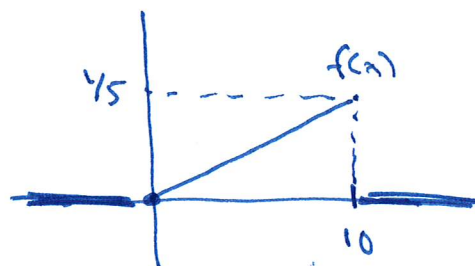
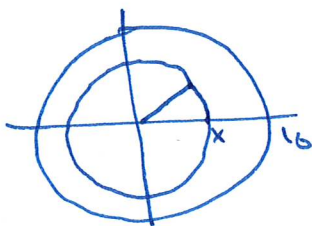


$$\frac{h \cdot f(x)}{h} = f(x)$$

Eks (pilkast)

$$F(x) = x^2/100$$

$$f(x) = \frac{2x}{100} = \frac{x}{50}$$

Krav til kumulativ fordelingsfunksjon:i)  $F$  vokendeii)  $\lim_{x \rightarrow \infty} F(x) = 1$  og  $\lim_{x \rightarrow -\infty} F(x) = 0$ Skrimemåte:

$$f(x) = f_X(x)$$

$$F(x) = F_X(x)$$

Sannsynlighetstettheten til  $X$   
kum. fordelingsfun. - " -3) Forventning og variansDefn:  $E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx$  forventningsverdien til  $X$ Eks: Pilkast

$$f(x) = x/50, \quad 0 \leq x \leq 10$$

$$E(X) = \int_0^{10} x \cdot f(x) dx = \int_0^{10} x \cdot x/50 dx = \left[ \frac{1}{50} \cdot \frac{1}{3} x^3 \right]_0^{10}$$

$$= \frac{1}{150} (10^3 - 0^3) = \frac{1000}{150} = \frac{100}{15} = \frac{20}{3} \approx \underline{\underline{6.67}}$$

Hvis  $h(x)$  er en avledet stokastisk variabel:

$$h(x) = x^2 \quad \text{eller} \quad h(x) = x^3 - x$$



Forventningsverdien til  $h(x)$ :

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

Ekso: Pikkost

$$\begin{aligned} E(x^2) &= \int_{-\infty}^{\infty} x^2 f(x) dx = \int_0^{10} x^2 \cdot \frac{1}{50} dx = \left[ \frac{1}{50} \cdot \frac{1}{4} x^4 \right]_0^{10} \\ &= \frac{1}{200} (10^4 - 0^4) = \frac{10000}{200} = \underline{\underline{50}} \end{aligned}$$

$$E(x)^2 = \left(\frac{20}{3}\right)^2 = \frac{400}{9} \neq E(x^2) = 50$$

Defn: Varians til  $X$ ,  $X$  stokastisk variabel med  $E(x) = \mu$

$$\begin{aligned} \text{Var}(x) &= E[(x - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx \\ &= \int_{-\infty}^{\infty} (x^2 - 2\mu x + \mu^2) f(x) dx = \int_{-\infty}^{\infty} x^2 f(x) dx - 2\mu \int_{-\infty}^{\infty} x f(x) dx \\ &\quad + \mu^2 \int_{-\infty}^{\infty} f(x) dx = E(x^2) - 2\mu E(x) + \mu^2 \cdot 1 \\ &= E(x^2) - 2\mu \cdot \mu + \mu^2 = E(x^2) - \mu^2 \end{aligned}$$

$$\boxed{\text{Var}(x) = E(x^2) - E(x)^2}$$

Ekso: Pikkost

$$\begin{aligned} E(x) = \frac{20}{3} \quad E(x^2) = 50 \quad \Rightarrow \quad \text{Var}(x) &= 50 - \left(\frac{20}{3}\right)^2 = \frac{50 \cdot 9 - 400}{9} \\ &= \underline{\underline{50/9}} \end{aligned}$$



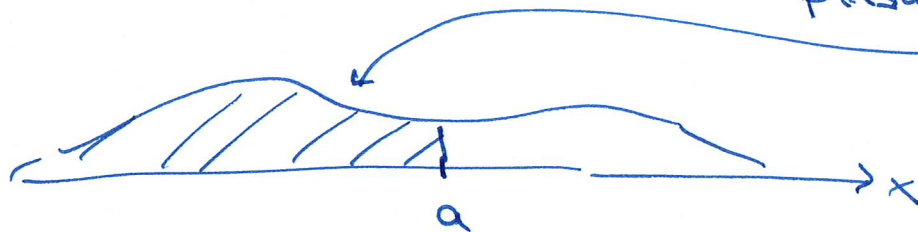
Generelt: 
$$\text{Var}(X) = \int_{-\infty}^{\infty} (x-\mu)^2 f(x) dx \geq 0$$

$$E(X^2) - E(X)^2 \geq 0$$

$\Rightarrow$  Standardavviket til X:  $\sigma_X = \sqrt{\text{Var}(X)}$

#### ④ Median:

Hvis X er en kontinuert stokastisk variabel med kumulativ fordelingsfunksjon  $F(x)$ , så er medianen til X det tallet  $a$  slik at  $F(a) = 1/2$ ,  $P(X \leq a) = 1/2$ .



$$P(X \leq a) = 1/2$$

$$F(a) = 1/2$$

$$\int_{-\infty}^a f(x) dx = 1/2$$

Eks: Pikkost

$$F(b) = b^2/100 = 1/2$$

$$b^2 = 100 \cdot 1/2 = 50$$

$$b = \sqrt{50}$$

$$\text{median: } \sqrt{50} \approx 7$$

$$E(X): 20/3 \approx 6.67$$

90% prosent:  $F(b) = b^2/100 = 0.9$

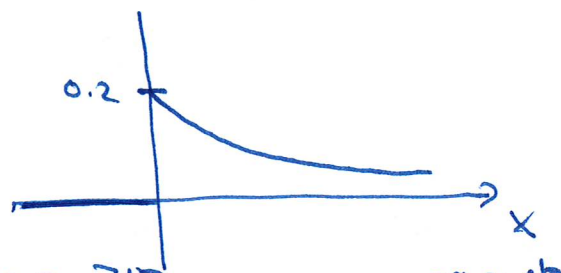
$$b^2 = 90$$

$$b = \sqrt{90}$$

Eks.  $X$  eksponensiell fordelt med  $f(x) = 0.2 \cdot e^{-0.2x}$ ,  $x \geq 0$

Skisse: i)  $f(x) \geq 0$  ok.

ii)  $\int_{-\infty}^{\infty} f(x) dx = 1$  ok.



$$\int_0^{\infty} 0.2 e^{-0.2x} dx = \left[ -e^{-0.2x} \right]_0^{\infty} = \lim_{b \rightarrow \infty} \left[ -e^{-0.2x} \right]_0^b$$

$$= \lim_{b \rightarrow \infty} \left( -e^{-0.2b} + 1 \right) = 1 \quad \underline{\text{ok}}$$

$$E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx = \int_0^{\infty} x \cdot 0.2 e^{-0.2x} dx$$

$$\int x e^{-0.2x} dx = \begin{array}{l} \boxed{\begin{array}{l} u = \frac{e^{-0.2x}}{-0.2} \\ u' = e^{-0.2x} \\ v = x \\ v' = 1 \end{array}} \\ = \frac{1}{-0.2} e^{-0.2x} \cdot x - \int \frac{1}{-0.2} e^{-0.2x} \cdot 1 dx \\ = -5x e^{-0.2x} + 5 \int e^{-0.2x} dx \\ = -5x e^{-0.2x} + 5 \left( \frac{1}{-0.2} \right) e^{-0.2x} + C \\ = -5x e^{-0.2x} - 25 e^{-0.2x} + C \end{array}$$

$$\begin{aligned} &= \lim_{b \rightarrow \infty} \left[ -5x e^{-0.2x} - 25 e^{-0.2x} \right]_0^b = \lim_{b \rightarrow \infty} \left( -5b e^{-0.2b} - 25 e^{-0.2b} + 25 \right) \\ &= \lim_{b \rightarrow \infty} \left( \frac{-5b}{e^{0.2b}} + \lim_{b \rightarrow \infty} \frac{-25}{e^{0.2b}} + 25 \right) = 25 \cdot 0.2 = \underline{\underline{5}} \end{aligned}$$