

## Estimering

### Viktige estimatoregenskaper:

Estimatoren  $\hat{\Theta}$  bør være **forventningsrett**, DVS  $E(\hat{\Theta}) = \Theta$ , hvor  $\Theta$  er et parameter du prøver å estiemere. **Variansen** til  $\hat{\Theta}$  bør være synkende med økende antall observasjoner.

Om du har to estimatorer  $\hat{\Theta}_1$ ,  $\hat{\Theta}_2$  er estimatoren med minst varians den mest effektive estimatoren for  $\Theta$ .

### Noen vanlige estimatorer, standardsituasjoner:

$\mu$ : For et tilfeldig utvalg av størrelse  $n$  fra en populasjon med forventning  $\mu$  og varians  $\sigma^2$  er en estimator for  $\mu$  gitt ved:

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$   $E(\bar{X}) = \mu$   $Var(\bar{X}) = \frac{\sigma^2}{n}$   
 $\sigma^2$ : For et tilfeldig utvalg av størrelse  $n$  fra en populasjon med forventning  $\mu$  og varians  $\sigma^2$  er en estimator for  $\sigma^2$  gitt ved:

$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$   $E(S^2) = \sigma^2$   $Var(S^2) = \frac{2\sigma^4}{n-1}$   
 $p$ : For et tilfeldig utvalg av størrelse  $n$  fra et binomisk forsøk (Bernulli-forsøksrekke) med sannsynlighet  $p$ . En estimator for  $p$  er gitt ved

$\hat{p} = \frac{\bar{X}}{n}$   $E(\hat{p}) = p$   $Var(\hat{p}) = \frac{p(1-p)}{n}$   
 $\mu_1 - \mu_2$ : For to uavhengige utvalg av størrelser  $n_1$ ,  $n_2$  fra populasjoner med forventning  $\mu_1$ ,  $\mu_2$  og varianser  $\sigma_1^2$ ,  $\sigma_2^2$  er en estimator for  $\mu_1 - \mu_2$  gitt ved

$\bar{X}_1 - \bar{X}_2$   $E(\bar{X}_1 - \bar{X}_2) = \mu_1 - \mu_2$   $Var(\bar{X}_1 - \bar{X}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$

$\sigma_1^2$ : For to tilfeldige utvalg av størrelser  $n_1$ ,  $n_2$  fra normalfordelte populasjoner med forventninger  $\mu_1$ ,  $\mu_2$  of varianser  $\sigma_1^2$ ,  $\sigma_2^2$  er en estimator for  $\frac{\sigma_1^2}{\sigma_2^2}$  gitt ved:

$p_1 - p_2$ : For to parvisse utvalg fra binomiske forsøk med sannsynligheter  $p_1$ ,  $p_2$  er en estimator for  $p_1 - p_2$  gitt ved:

$\hat{p}_1 - \hat{p}_2 = \frac{\bar{X}_1}{n_1} - \frac{\bar{X}_2}{n_2}$   $E(\hat{p}_1 - \hat{p}_2) = p_1 - p_2$   
 $Var(\hat{p}_1 - \hat{p}_2) = \frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}$

$\mu_D$ : For to parvisse tilfeldige utvalg av størrelse  $n$  der differansene fra populasjonene med forventning  $\mu_D$  og varians  $\sigma_D^2$  er en estimator for  $\mu_D$  gitt ved:

$\bar{D}$   $E(\bar{D}) = \mu_D$   $Var(\bar{D}) = \frac{\sigma_D^2}{n}$

## Utalgsfordelinger

En utvalgsfordeling er fordelinga for en observator (funksjon av de stokastiske variablene i utvalget) for et (tilfeldig) utvalg data. Vi er gjerne interresert i fordelinga til (de spesielle observatorene som er ) estimatorer for paramerere i populasjonen, ofte på standardisert form.

### Standardsituasjonene (Utalgsfordelinger)

$\bar{X}$ ,  $Z$ : For et tilfeldig utvalg av størrelse  $n$  fra en normalfordelt populasjon med forventning  $\mu$  og varians  $\sigma^2$  vil:

$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$   $Z = \frac{\bar{X}-E(\bar{X})}{\sqrt{Var(\bar{X})}} = \frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$   
**SGT**: Selv om populasjonen ikke er normalfordelt vil resultatet over gjelde når  $n \rightarrow \infty$ . Regner vanligvis tilnærminga for god når  $n \geq 30$

$T$ : For et tilfeldig utvalg av størrelser  $n$  fra en normalfordelt populasjon med forventning  $\mu$  og varians  $\sigma^2$ , der variansen estimeres ved  $S^2$  fra utvalget har vi at:

$T = \frac{\bar{X}-\mu}{\frac{S}{\sqrt{n}}} = \frac{\frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}}}{\sqrt{\frac{S^2}{\sigma^2}}} \sim t_{n-1}$

$S^2$ : For et tilfeldig utvalg av sørrelse  $n$  fra en normalfordelt populasjon med forventning  $\mu$  og varians  $\sigma^2$ :

$\frac{(n-1)S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi^2_{n-1}$   
 $\hat{p}$ : For et tilfeldig utvalg av størrelse  $n$  fra et binomisk forsøk med sannsynlighet  $p$  har vi tilnærmet at:  
 $Z = \frac{\hat{p}-E(\hat{p})}{\sqrt{Var(\hat{p})}} = \frac{\hat{p}-p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1)$

$\bar{X}_1 - \bar{X}_2$ : For to uavhengige tilfeldige utvalg  $n_1$ ,  $n_2$  (normalfordelt) med forventning  $\mu_1$ ,  $\mu_2$  og varianser  $\sigma_1^2$ ,  $\sigma_2^2$  vil:

$\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2, \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right)$

Eller  
 $Z = \frac{(\bar{X}_1-\bar{X}_2)-(\mu_1-\mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$

Gjelder også for tilnærma uten å forutsette normalfordelt (stor  $n_1$  og  $n_2$ )

**For ukjent men lik varians**  $\sigma_1^2 = \sigma_2^2 = \sigma_p^2$  estimeres med  $S_p^2$ :

$T = \frac{(\bar{X}_1-\bar{X}_2)-(\mu_1-\mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim T_{n_1+n_2-2}$

**For ukjent og ulik varians**  $\sigma_1^2 \neq \sigma_2^2$ , estimeres med  $S_1^2$ ,  $S_2^2$

$T = \frac{(\bar{X}_1-\bar{X}_2)-(\mu_1-\mu_2)}{\sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}} \sim T_\nu$   $\nu = \frac{\left(\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}\right)^2}{\frac{\hat{\sigma}_1^4}{n_1} + \frac{\hat{\sigma}_2^4}{n_2}}$

$F$ : For to uavhengige tilfeldige utvalg  $n_1$ ,  $n_2$  (normalfordelt) med forventning  $\mu_1$ ,  $\mu_2$  og varianser  $\sigma_1^2$ ,  $\sigma_2^2$  har vi at:

$F = \frac{\frac{\hat{\sigma}_1^2}{n_1}}{\frac{\hat{\sigma}_2^2}{n_2}} \sim F_{n_1-1, n_2-1}$

$\hat{p}_1 - \hat{p}_2$ : med to tilfeldige utvalg  $n_1$ ,  $n_2$  (binomisk) med sannsynligheter  $p_1$ ,  $p_2$  har vi tilnærma at (om  $n_1$ ,  $n_2$  stor nok):  
 $Z = \frac{(\hat{p}_1-\hat{p}_2)-(\mu_1-\mu_2)}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} \sim N(0, 1)$

$\bar{D}$ : for to parvisse tilfeldige utvalg av størrelse  $n$  der differansene er normalfordelte med forventning  $\mu_D$  og varians  $\sigma_D^2$  og der variansen estimeres ved  $S_D^2$  fra utvalget har vi at:

$T = \frac{\bar{D}-\mu_D}{\frac{S_D}{\sqrt{n}}} \sim t_{n-1}$

## Intervallestimering

**Konfidensintervall for  $\mu$  (KI)**  
KI for  $\mu$  med kjent  $\sigma^2$ :

$P\left(-z_{\frac{\alpha}{2}} < \frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}} < z_{\frac{\alpha}{2}}\right) = 1 - \alpha$   
 $\Rightarrow \mu = \bar{X} \pm z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$

KI for  $\mu$  med  $S^2$  som estimator for ukjent  $\sigma^2$  med  $n-1$  frihetsgrader:  $P\left(-t_{\frac{\alpha}{2}} < \frac{\bar{X}-\mu}{\frac{S}{\sqrt{n}}} < t_{\frac{\alpha}{2}}\right) = 1 - \alpha$   
 $\Rightarrow \mu = \bar{X} \pm t_{\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}$

**Konfidensintervall for  $\sigma^2$  med  $n-1$  (KI)**  
 $P\left(\chi^2_{1-\frac{\alpha}{2}} < \frac{(n-1)S^2}{\sigma^2} < \chi^2_{\frac{\alpha}{2}}\right)$   
 $\Rightarrow \sigma^2 \in \left[\frac{(n-1)S^2}{\chi^2_{\frac{\alpha}{2}}}, \frac{(n-1)S^2}{\chi^2_{1-\frac{\alpha}{2}}}\right]$

**Konfidensintervall for  $p$**   
KI for  $p$  med utvalg av størrelse  $n$  fra et binomisk forsøk med sannsynlighet  $p$  (OBS:  $n > 30$  /stor nok)

$P\left(-z_{\frac{\alpha}{2}} < \frac{\hat{p}-p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} < z_{\frac{\alpha}{2}}\right) = 1 - \alpha$   
 $\Rightarrow p = \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$  Erstatt  $p$  i uttrykket med  $\hat{p}$  for å kunne regne ut

**Konfidensintervall for  $\mu_1 - \mu_2$**   
For to uavhengige tilfelig utvalgte utvalg av størrelser  $n_1$  og  $n_2$  fra normalfordelte populasjoner med forventning  $\mu_1$ ,  $\mu_2$  og vareanser  $\sigma_1^2$ ,  $\sigma_2^2$ :  
For **kjente** vareanser  $\sigma_1^2$ ,  $\sigma_2^2$ :

$\mu_1 - \mu_2 = (\bar{X}_1 - \bar{X}_2) \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = 1 - \alpha$

For **ukjente men like** vareanser  $\sigma_1^2 = \sigma_2^2$  med antall frihetsgrader  $\nu = n_1 + n_2 - 2$ :

$\mu_1 - \mu_2 = (\bar{X}_1 - \bar{X}_2) \pm s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = 1 - \alpha$

For **ukjente men ulike** vareanser  $\sigma_1^2 \neq \sigma_2^2$  der antall frihetsgrader  $\nu$  er gitt fra utvalgsfordelंगा:

$\mu_1 - \mu_2 = (\bar{X}_1 - \bar{X}_2) \pm t_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = 1 - \alpha$

**Konfidensintervall for  $\frac{\sigma_1^2}{\sigma_2^2}$**   
Normalfordelt, med  $\mu_1$ ,  $\mu_2$  og  $\sigma^2$ ,  $\sigma_2^2$ :

$P\left(f_{1-\frac{\alpha}{2}, \nu_1, \nu_2} < \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} < f_{\frac{\alpha}{2}, \nu_1, \nu_2}\right)$   
 $\Rightarrow \frac{\sigma_1^2}{\sigma_2^2} \in \left[\frac{S_1^2}{S_2^2} \frac{1}{f_{\frac{\alpha}{2}, \nu_1, \nu_2}}, \frac{S_1^2}{S_2^2} f_{\frac{\alpha}{2}, \nu_1, \nu_2}\right]$

**Konfidensintervall for  $p_1 - p_2$**   
Uavhengig tilfeldig utvalg  $n_1$ ,  $n_2$  fra binomisk forsøk med sannsynlighet  $p_1$ ,  $p_2$  har vi (om  $n_1$ ,  $n_2$  store nok):

$(\hat{p}_1 - \hat{p}_2) \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} = 1 - \alpha$

**Konfidensintervall for  $\mu_D$**   
For to parvisse tilfeldige utvalg, størrelse  $n$  med normalfordelte differanse med forventning  $\mu_D$  og varians  $\sigma_D^2$  (Variansen estimeres med  $S_D^2$ )  
 $\bar{D} \pm t_{\frac{\alpha}{2}} \frac{S_D}{\sqrt{n}} = 1 - \alpha$

**Langden av KI**  
For en normalfordelt estimator med kjent varians  $\sigma^2$  kan vi regne ut hva lengden vil bli for et gitt valg av  $n$  og  $\alpha$   
 $L = 2 \cdot z_{\frac{\alpha}{2}} \cdot SE(\hat{\Theta})$ , to typiske tilfeller:

- Konstruere KI for  $\mu$  basert på estimatoren:  $\bar{X}$  blir lengda:  
 $L = 2 \cdot z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \Rightarrow n = \left(\frac{2z_{\frac{\alpha}{2}} \sigma}{L}\right)^2$   
Når  $\sigma^2$  må estimeres kan vi finne approx forventna lengde
- Når vi konstruerer KI for  $p$  basert på andelsestimatoren  $\hat{p}$  ( $l \approx L$ ):  
 $L = 2 \cdot z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{l}} \leq \frac{z_{\frac{\alpha}{2}}}{\sqrt{n}} \Rightarrow n = \left(\frac{2z_{\frac{\alpha}{2}} \sqrt{\hat{p}(1-\hat{p})}}{l}\right)^2$

## Estimeringsfeil

Når fordelinga til estimatoren er kjent (utvalgsfordelंगा) kan vi regne ut hvor stor feil vi gjør i estimeringa. Vil ofte ha en viss sannsynlighet for at feilen ikke skal overskride en verdi  $e$   
 $P(|\hat{\Theta} - \Theta| < e = 1 - \alpha)$  (vanlige tilfeller)  
 $\bar{X}$  (normalfordelt):

$P(|\bar{X} - \mu| < e = 1 - \alpha) \Rightarrow n = \left(\frac{z_{\frac{\alpha}{2}} \sigma}{e}\right)^2 \Rightarrow e = z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$   
 $\hat{p}$  (binomisk):

$P(|\hat{p} - p| < e = 1 - \alpha) \Rightarrow n \approx \left(\frac{z_{\frac{\alpha}{2}} \sqrt{\hat{p}(1-\hat{p})}}{e}\right)^2 \Rightarrow$   
 $e \approx z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

## Noen regneregler

$(u \cdot v)' = u'v + u \cdot v'$   
 $\left(\frac{u}{v}\right)' = \frac{u'v-u \cdot v'}{v^2}$   
 $g(f(x))' = g'(f(x)) \cdot f'(x)$   
 $\ln(x)' = \frac{1}{x}$   
 $\int \ln(x) dx = x \ln(x) - x + C$   
u er et utrykk med x  
 $\int a^x dx = \frac{a^x}{\ln(a)} + C$   
 $\int \frac{1}{u} dx = \frac{\ln(u)}{u} + C$   
 $\int u \cdot v' dx = u \cdot v - \int u'v dx$

## Kombinatorikk

**Multiplikasjonsregelen**  
Forsøk i k etapper, med  $m_1$ ,  $m_2$ ,  $\dots$ ,  $m_k$  mulige utfall i etappene. Totalt antall utfall:  $\prod_{i=1}^k m_i$

**Potensregelen**  
 $n$  merka enheter, velger  $k$  med tilbakelegging. Antall ordna utfall:  $n^k$

**Permutasjonsregelen (nPr)**  
 $n$  merka enheter, velger  $k$  uten tilbakelegging, antall ordna utfall:  $nPr = \frac{n!}{(n-k)!}$

**Kombinasjonsregelen**  
 $n$  merka enheter, velger  $k$  uten tilbakelegging, antall ikke-ordna utfall:  $nCr = \binom{n}{k} = \frac{n!}{(n-k)!k!}$

## Prediksjonsintervall

Et intervall som sier noe om neste verdi  $X_0$  i en normalfordelt populasjon:  $X \sim N(\mu, \sigma)$

**Kjent  $\mu$  og  $\sigma$**  (spredningsintervall)  
 $P\left(\mu - z_{\frac{\alpha}{2}} \cdot \sigma < X_0 < \mu + z_{\frac{\alpha}{2}} \cdot \sigma\right) = 1 - \alpha$   
 $\Rightarrow X_0 = \mu \pm z_{\frac{\alpha}{2}} \cdot \sigma$

**Ukjent  $\mu$ , kjent  $\sigma$**   
 $P\left(\bar{X} - z_{\frac{\alpha}{2}} \cdot \sigma \sqrt{1 + \frac{1}{n}} < X_0 < \bar{X} + z_{\frac{\alpha}{2}} \cdot \sigma \sqrt{1 + \frac{1}{n}}\right) = 1 - \alpha$   
 $\Rightarrow X_0 = \bar{x} \pm z_{\frac{\alpha}{2}} \cdot \sigma \sqrt{1 + \frac{1}{n}}$

**Ukjent  $\mu$ ,  $\sigma$**  frihetsgrad  $\nu = n - 1$

$X_0 = \bar{x} \pm z_{\frac{\alpha}{2}} \cdot S \sqrt{1 + \frac{1}{n}}$

**Formel for  $S_p^2$**  (Gjelder for alt på arket)

$S_p^2 = \frac{S_1^2(n_1-1) + S_2^2(n_2-1)}{n_1+n_2-2}$

## Foredelinger

**Binomisk fordeling**

- $n$  uavhengige delforsøk
- To utfall: suksess/ikke suksess
- Samme sannsynlighet for  $p = P(a)$  (suksess) i alle delforsøk

**Hypergeometrisk fordeling**  
1. Populasjon med  $N$  elementer  
2.  $k$  av disse regnes som "suksess",  $N - k$  som "fiasko"  
3. Trekker  $n$  elementer uten tilbakelegging  
Sannsynligheten  $p$  endrer seg mellom hvert delforsøk.

**Negativ binomisk fordeling**  
Antall forsøk du må gjøre for at hendelsen  $A$  (suksess) skal inntreffe  $k$  ganger

**Geometrisk fordeling**  
Antall forsøk du må gjøre for at hendelsen  $A$  (suksess) skal inntreffe første gang

**Poisson-fordeling**  
 $\mu = \lambda t$ ,  $\sigma^2 = Var(X) = \lambda t$   
 $f(x) = \frac{\mu^x}{x!} e^{-x}$

- Antallet av  $A$  disjunkte tidsintervall er uavhengige
- Forventa antall av  $A$  er konstant list  $\lambda$  (raten) per tidsenhet
- Kan ikke få to forekomster samtidig

**Gammalfordeling**  
En kontinuerlig variabel  $X$  er gammalfordelt med parameter  $\alpha > 0$  og  $\beta > 0$  dersom tetthetsfunksjonen er gitt ved (se blå tabell). Ventetida til hendelse nummer  $k$  i en Poisson-prosess vil være gammalfordelt med  $\alpha = k$  og  $\beta = \frac{1}{\lambda}$

**Eksponsiellfordeling**  
Ventetida til første hendelse (og mellom etterfølgende hendelser) i en Poisson-prosess følger en eksponensialfordeling. En kontinuerlig variabel  $X$  har eksponentialfordeling med parameter  $\beta > 0$  dersom tetthetsfunksjonen er gitt ved (se tabellbok:  $f(x)$ ). Eksponensialfordelंगा er en variant av gammalfordeling med  $\alpha = 1$

## Stokastisk variabel

**En variabel**

$\mu = E(X) = \begin{cases} \sum x \cdot x \cdot f(x), & \text{deskret} \\ \int_{-\infty}^{\infty} x \cdot f(x) dx, & \text{kontinuerlig} \end{cases}$

$\sigma^2 = Var(X) = E(X^2) - E(X)^2 = \begin{cases} \sum x(x - \mu)^2 f(x), & \text{deskret} \\ \int_{-\infty}^{\infty} (x - \mu)^2 \cdot f(x) dx, & \text{kontinuerlig} \end{cases}$

$F(x) = \int_0^x f(t) dt$   
**Funksjoner av stokastiske variabler**

$\mu_{g(X)} = E(g(X)) = \begin{cases} \sum g(x) \cdot f(x), & \text{deskret} \\ \int_{-\infty}^{\infty} g(x) \cdot f(x) dx, & \text{kontinuerlig} \end{cases}$

$\sigma_{g(X)}^2 = Var(g(X)) = \begin{cases} \sum x(g(x) - \mu_{g(X)})^2 \cdot f(x), & \text{deskret} \\ \int_{-\infty}^{\infty} (g(x) - \mu_{g(X)})^2 \cdot f(x) dx, & \text{kontinuerlig} \end{cases}$

$\mu_{g(X, Y)} = E(g(X, Y)) = \begin{cases} \sum x \sum y g(x, y) f(x, y), & \text{deskret} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dy dx & \text{kontinuerlig} \end{cases}$

**Simultanfordeling for to variabler**

$P(X, Y) = \begin{cases} \sum \sum f(x, y), & \text{deskret} \\ \int \int f(x, y) dy dx, & \text{kontinuerlig} \end{cases}$

**Marginale fordelinger** (to variabler)

$g(x) = \begin{cases} \sum_y f(x, y), & \text{deskret} \\ \int_{-\infty}^{\infty} f(x, y) dy & \text{kontinuerlig} \end{cases}$

$h(y) = \begin{cases} \sum_x f(x, y), & \text{deskret} \\ \int_{-\infty}^{\infty} f(x, y) dx & \text{kontinuerlig} \end{cases}$

**Betinga fordeling** (for  $Y$  gitt  $X$ )  
 $f(y|x) = \frac{f(x, y)}{g(x)}$ ,  $g(x) > 0$

**Kovarians**  
 $\sigma_{XY} = Cov(X, Y) = \begin{cases} \sum x \sum y (x - \mu_x)(y - \mu_y) f(x, y) \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) f(x, y) dy dx \end{cases}$   
Korrelasjon:

$\rho_{XY} = \frac{Cov(X, Y)}{\sqrt{Var(X) \cdot Var(Y)}}$ ,  $-1 < \rho_{XY} < 1$

## Sannsynlighetsregler

**Addisjonsregelen:**  
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

**Multiplikasjonsregelen:**  
 $P(A \cap B) = P(A) \cdot P(B|A) = P(B) \cdot P(A|B)$   
**Betinga sannsynlighet:**  $P(A|B) = \frac{P(A \cap B)}{P(B)}$

**Total sannsynlighet:**

$P(A) = \sum_{i=1}^n P(B_i \cap A) = \sum_{i=1}^k P(B_i) \cdot P(A|B_i)$

**Bayes setning:**  $P(B|A) = \frac{P(B)P(A|B)}{P(A)}$

**Tilfeller ved uavhengighet:**

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

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

Lager tester som tar stilling til påstander. Har da: nullhypotesen ( $H_0$ ) og en utfordrende hypotese ( $H_1$ ). Vi bruker testobservatorer for å kunne si noe om  $H_0$  er sann eller ikke.

**Normalfordelt estimator:**  
Dersom vi har en normalfordelt forventningsrett estimator  $\hat{\Theta}$  for parameter  $\Theta \sim N(\Theta, SE(\hat{\Theta}))$

Kan vi konstruere testobservator  $Z = \frac{\hat{\Theta}-\Theta_0}{SE(\hat{\Theta}_0)}$

(Måler hvor mange standardavvik  $\hat{\Theta}$  er fra  $\Theta_0$ .)  
-Om  $H_0$  er sann vil  $Z$  bli standardnormalfordelt  
-Om  $H_1$  er sann vil vi forvente at  $\hat{\Theta}$  blir større enn  $\Theta_0$ , og dermed at  $Z$  blir stor.

Forkastningsområdet velges slik at med sannsynlighet  $\alpha$  for å få en så stor verdi, dersom  $H_0$  er sann. Kritisk verdi blir da  $z_\alpha$ , om  $Z > z_\alpha$ .

**P-verdi:**  
Som et alternativ fil å finne et forkastningsområde for  $Z$  kan man finne en  $p$  – verdi som er det minste signifikansnivået som ville forkastet nullhypotesen.

**Feil av type I/II:**  
P(Feil av type I) = P(Forkaste en rett  $H_0$ )= $\alpha$   
P(Feil av type II) = P(ikke forkaste gal  $H_0$ )= $\beta$   
Det ideelle er å få  $\alpha$  og  $\beta$  minst mulig.

$1 - \beta$  kalles for styrken til en test og sier hvor sannsynlig det er å forkaste  $H_0$  som funksjon av den samme parameterverdien. Ved å finne denne for ulike verdier av

## Forenklinger for lineærkombinasjoner

*Var*(*aX* + *bY*) = *a*<sup>2</sup> *Var*(*X*) + *b*<sup>2</sup> *Var*(*Y*) + 2*ab* · *Cov*(*X*, *Y*)

Merk: hvis *X* og *Y* er uavhengige er *Cov*(*X*, *Y*) = 0

*E*(*aX* + *b*) = *aE*(*X*) + *b*

*E*(*X*<sub>*n*−1</sub><sup>2</sup>) = *n* − 1

*E*(*T*<sub>*v*</sub>) = 0

*Var*(*T*<sub>*v*</sub>) = 






ν



ν
−
2


,


 
ν
≥
3

**Tsjebysjeffs teorem:** Sannsynligheten for at verdien på en variabel *X* ligger innåfor *k* avstander fra forventningsverdien er minst 1 − 



1

k

2




{\displaystyle 1-{\frac {1}{k^{2}}}}

**Sannsynlighetsmaksimeringsfunksjonen:**

*L*(Θ) = ∏



i
=
1


n




{\displaystyle \prod \_{i=1}^{n}}

 *f*(*x*<sub>*i*</sub>, Θ)

Løses ved å:

1. Ta logaritmen av funksjonen ln *L*

ln



(


a
b



)


=
ln
(
a
)
−
ln
(
b
)
,


 
ln
(
a
b
)
=
ln
(
a
)
+
ln
(
b
)
,


 
ln
(

a

b


)
=
b
ln
(
a
)

2. Deriver uttrykket med hensyn på variabelen Θ oftest 






∂
ln
⁡
L


∂
Θ




{\displaystyle {\frac {\partial \ln L}{\partial \Theta }}}

 = 0

3. Sett uttrykket lik 0 og løs med hensyn på variabelen. 






∂
ln
⁡
L


∂
Θ




{\displaystyle {\frac {\partial \ln L}{\partial \Theta }}}

 = 0

**Tilfeldig utvalg:**

Når et utvalg er tilfeldig, kan vi se på det som en mengde identiske og uavhengige observasjoner.

**QQ-plot:**

- Plotter utvalgskvantiler (Observasjonene ordna etter størrelse) mot teoretiske kvantiler ("Ideelle observasjoner") fra aktuell fordeling

- Teoretiske kvantiler er gitt ver invers kumulativ fordeling "jevnt spredte" sannsynligheter mellom 0 og 1

- Om antatt fordeling stemmer skal plotter gj e en tilnerma rett linje *x* = *y*

### Enkel linær regresjon

Har observasjonspår (*x*<sub>1</sub>, *y*<sub>1</sub>), . . . , (*x*<sub>*n*</sub>, *y*<sub>*n*</sub>).

**Regresjonsmodellen**

*Y*<sub>*i*</sub> = *β*<sub>0</sub> + *β*<sub>1</sub>*x*<sub>*i*</sub> + *ε*<sub>*i*</sub>, 



 
i
=
1
,
.
.
.
,
n

Forutsetningene er at

*E*(*ε*<sub>*i*</sub>) = 0, 



 
Var
(

ε

i


)
=

σ

2


,


 
i
=
1
,
.
.
.
,
n

og *ε*<sub>1</sub>, . . . , *ε*<sub>*n*</sub> er inbyrdes uavhengige. Dette medfører at

*E*(*Y*<sub>*i*</sub>) = *μ*<sub>*y*|*x*<sub>*i*</sub></sub> = *β*<sub>0</sub> + *β*<sub>1</sub>*x*<sub>*i*</sub>, 



 
Var
(

Y

i


)
=

σ

2


y

i


|

x

i


=

σ

2


,


 
i
=
1
,
.
.
.
,
n

der *Y*<sub>1</sub>, . . . , *Y*<sub>*n*</sub> også er innbyrdes uavhengige.

**MKM**

Bruker miste kvadrat-metoden til å estimere *β*<sub>0</sub> og *β*<sub>1</sub> fra data for å få ei estimert (tilpassa) regresjonslinje for forventningslinja *β*<sub>0</sub> + *β*<sub>1</sub>*x*:

*ŷ* = *b*<sub>0</sub> + *b*<sub>1</sub>*x*

MKM baserer seg på å minimere

*SSE* = 




∑

n


i
=
1



(

y

i


−

ŷ

i


)

2


=

∑

n


i
=
1



(

y

i


−

b

0


−

b

1


x

i


)

2




{\displaystyle \sum \_{i=1}^{n}(y\_{i}-{\hat {y}}\_{i})^{2}=\sum \_{i=1}^{n}(y\_{i}-b\_{0}-b\_{1}x\_{i})^{2}}

Gir estimatorer

B

1


=



∑

n


i
=
1



(

x

i


−
x
¯
)
(

Y

i


−
Y
¯
)



∑

n


i
=
1



(

x

i


−
x
¯
)

2





=



∑

n


i
=
1



(

x

i


−
x
¯
)

Y

i




∑

n


i
=
1



(

x

i


−
x
¯
)

2





,


B

0


=
Y
¯
−

B

1


x
¯

der

E
(

B

1


)
=

β

1


,


 
Var
(

B

1


)
=



σ

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2

E
(

B

0


)
=

β

0


,


 
Var
(

B

0


)
=



σ

2




∑

n


i
=
1



x

i


2




∑

n


i
=
1



(

x

i


−
x
¯
)

2

I tillegg har vi (forventningsrett) estimator for *σ*<sup>2</sup>:

S

2


=



SSE


n
−
2




=



∑

n


i
=
1



(

Y

i


−

Ŷ

i


)

2




n
−
2




=



∑

n


i
=
1



(

Y

i


−

B

0


−

B

1


x

i


)

2




n
−
2

**Parameterinferens**

Om vi antar at feilledda er normalfordelte

*ε*<sub>*i*</sub> ~ *N*(0, *σ*)

vil

B

1


∼
N
⎛

β

1


,


σ

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





⎞
,


 
B

0


∼
N
⎛

β

0


,


σ

2


n
⋅

∑

n


i
=
1



(

x

i


−
x
¯
)

2





⎞

og KI og tester kan konstrueres ved å ta utgangspunkt i at

V
=



(
n
−
2

)

S

2




σ

2

er kikkvadratfordelt med *n* − 2 frihetsgrader, og dermed at

T
=



B

1


−

β

1




S
E
(

B

1


)


=



B

1


−

β

1




S


√

∑

n


i
=
1



(

x

i


−
x
¯
)

2





=



B

0


−

β

0




S
E
(

B

0


)


=



B

0


−

β

0




S


√



∑

n


i
=
1



x

i


2




n
⋅

∑

n


i
=
1



(

x

i


−
x
¯
)

2

vil være t-fordelte med *n* − 2 frihetsgrader. Dette gir oss da KI og tester ofr

*β*<sub>1</sub>, *β*<sub>0</sub>, *σ*

**Inferens for** 




μ

y

|

x

0




{\displaystyle \mu \_{y|x\_{0}}}

Dersom *x* er en gitt verdi, *x*<sub>0</sub>, vil en estimator for forventa verdi av responsen bli

*Ŷ*<sub>0</sub> = *B*<sub>0</sub> + *B*<sub>1</sub>*x*<sub>0</sub>

Fordelinga for *Ŷ*<sub>0</sub> blir:

Ŷ

0


∼
N
⎛

μ

y

|

x

0


,


σ

√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





⎞

Og da kan ta utgangspunkt i at

T
=



Ŷ

0


−

μ

y

|

x

0




S
√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2

er t-fordelt med *n* − 2 frihetsgrader til for eksempel å lage 100(1 − *α*)%-konfidensintervall for 




μ

y

|

x

0




{\displaystyle \mu \_{y|x\_{0}}}

 (forventa verdi av responsen *Y* når *x* = *x*<sub>0</sub>):

⎡

ŷ

0


−

t


z



σ


√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





,


ŷ

0


+

t


z



σ


√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





⎤

**PI for** *Y*<sub>0</sub>

Ønsker ofte intervall som med sannsynligheten 100(1 − *α*)% vil inneholde verdien av en ny observasjon *Y*<sub>0</sub> (når *x* = *x*<sub>0</sub>), altså et *prediksjonsintervall*. Tar utgangspunkt i forskjellen mellom estimatoren for forventa respons *Ŷ*<sub>0</sub> og verdien av en ny observasjon, *Y*<sub>0</sub>

Ŷ

0


−

Y

0


∼
N
⎛
0
,


σ

√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





⎞

Da kan vi bruke at

T
=



Ŷ

0


−

Y

0




S
√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2

er t-fordelt med *n* − 2 frihetsgrader.

Et 100(1 − *α*)%-PI for responsen *Y*<sub>0</sub> vil da finnes ved

⎡

ŷ

0


−

t


z



σ


√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





,


ŷ

0


+

t


z



σ


√



1


n


+



(

x

0


−
x
¯
)

2




∑

n


i
=
1



(

x

i


−
x
¯
)

2





⎤

**Betydning av lineær modell**

Total variasjon =	Forklart variasjon	+	Restvariasjon
<span><span>     ∑<!-- ∑ -->  n   i = 1    (  y  i   −<!-- − --> y ¯<!-- ¯ --> )  2     =  ∑<!-- ∑ -->  n   i = 1    (  ŷ<!-- ŷ -->  i   −<!-- − --> y ¯<!-- ¯ --> )  2      </span></span>	<span><span>     ∑<!-- ∑ -->  n   i = 1    (  y  i   −<!-- − --> y ¯<!-- ¯ --> )  2      </span></span>	+	<span><span>     ∑<!-- ∑ -->  n   i = 1    (  y  i   −<!-- − -->  ŷ<!-- ŷ -->  i   )  2      </span></span>
<i>SST</i>	<i>SSR</i>	+	<i>SSE</i>

**forklaringsgraden** er definert ved

R

2


=



Forklart variasjon
Total variasjon


=



SSR
SST


=



SST
− SSE
SST


=
1
−



SSE
SST


,


 
0
≤

R

2


≤
1

Dette er et vanlig mål på hvor stor andel av variasjonen i responsen som kan tilskrives den lineære sammenhengen 




μ

y

|

x


=

β

0


+

β

1


x


{\displaystyle \mu \_{y|x}=\beta \_{0}+\beta \_{1}x}

 (Resten tilskrives tilfeldig feil, *ε*)

## Løsningsforslag

### Oppgave

Vi antar at konsentrasjonen i en prøve *Y*<sub>1</sub>, av volum *l*<sub>*i*</sub> liter har følgende fordeling:

Y

i


∼
N
⎛
μ
,


σ

√

l

i




⎞

Nå er det tatt *n* = 5 uavhengige prøver av ulik størrelse.

Sett opp en likelihoodfunksjon (sannsynlighetsmaksimeringsfunksjon) for de ukjente parameterne *μ* og *σ*<sup>2</sup>.

Vis as sannsynlighetsmaksimeringsfunksjonen for *μ* og *σ*<sup>2</sup> blir

μ
=



∑

n


i
=
1



l

i


Y

i




∑

n


i
=
1



l

i




,


 
og


σ

2


=



∑

n


i
=
1



l

i


(

Y

i


−
μ
)

2




n

L
(
μ
,

σ

2


)
=
∏

i
=
1


n



f
(

x

i


;
μ
,

σ

i


)
=
∏

i
=
1


n



1


√
2
π
σ


√

l

i




exp
⁡
⎛
−



(

y

i


−
μ
)

2




2

σ

2


l

i





⎞
=
∏

i
=
1


n



√

l

i




√
2
π
σ


exp
⁡
⎛
−



l

i


(

y

i


−
μ
)

2




2
σ

2





⎞
=



1


(
√
2
π
)

n




1


(
√

σ

2


)

n




exp
⁡
⎛
−



1


2
σ

2




∑

i
=
1


n



l

i


(

y

i


−
μ
)

2





⎞
∏

i
=
1


n



√

l

i

Finner sannsynlighetsmaksimeringsfunksjonen på vanlig måte:

ln
⁡
L
=
−


n
2



ln
⁡
(
2
π
)
−


n
2



ln
⁡
(

σ

2


)
−


1


2
σ

2




∑

i
=
1


n



l

i


(

y

i


−
μ
)

2


+


1
2



ln
⁡
⎛

∑

i
=
1


n



l

i





⎞

∂
ln
⁡
L


∂
μ


=
−


1


2
σ

2




∑

i
=
1


n



2

l

i


(

y

i


−
μ
)
(
−
1
)
=
0


 
(kjerneregelen)


⇒
μ


∑

i
=
1


n



l

i


=
∑

i
=
1


n



l

i


y

i


⇒
μ
=



∑

i
=
1


n



l

i


Y

i




∑

i
=
1


n



l

i

∂
ln
⁡
L


∂

σ

2




=
−


n


2
σ

2


+


1


2
σ

4




∑

i
=
1


n



l

i


(

y

i


−
μ
)

2


=
0


⇒
1


σ

2




∑

i
=
1


n



l

i


(

y

i


−
μ
)

2


=
n


⇒
σ

2


=



∑

i
=
1


n



l

i


(

Y

i


−
μ
)

2




n

**Deloppgave**

Vis at *μ* er forventningsrett. Vis at variansen til *μ* = 






σ

2




∑

i
=
1


n



l

i

Vis at *σ*<sup>2</sup> ikke er forventningsrett. Du kan bruke at 






n
σ

2




{\displaystyle {\frac {n\sigma ^{2}}{\sigma ^{2}}}}

 ~ *χ*<sub>*n*−1</sub><sup>2</sup>.

Hva blir en forventningsrett estimator for *σ*<sup>2</sup> ?

E
(
μ
)
=
E
⎛



∑

n


i
=
1



l

i


Y

i




∑

n


i
=
1



l

i





⎞
=



∑

n


i
=
1



l

i


E
(

Y

i


)



∑

n


i
=
1



l

i





=



∑

n


i
=
1



l

i


μ



∑

n


i
=
1



l

i





=
μ



∑

n


i
=
1



l

i




∑

n


i
=
1



l

i





=
μ

Var
(
μ
)
=
Var
⎛



∑

n


i
=
1



l

i


Y

i




∑

n


i
=
1



l

i





⎞
=



∑

n


i
=
1



l

i


2


Var
(

Y

i


)



(
∑

n


i
=
1



l

i


)

2





=



∑

n


i
=
1



l

i


2


σ

2




l

i




(
∑

n


i
=
1



l

i


)

2





=



σ

2




∑

n


i
=
1



l

i




(
∑

n


i
=
1



l

i


)

2





=



σ

2




∑

n


i
=
1



l

i

E
(

σ

2


)
=
E
⎛



n

σ

2




σ

2





⎞
=



σ

2


n


E
⎛



n

σ

2




σ

2





⎞
=

σ

2


n
−
1

S

2


=



∑

n


i
=
1



l

i


(

Y

i


−
μ
)

2




n
−
1

## Oppgave

Vi skal se på rekkevidde for en elektrisk bilmodell (kjørelengde fra batteriet er fullt til tomt). Det blir påstått at bilprodusenter oppgir urealistisk lang rekkevidde, derfor lar vi *n* = 15 tilfeldig valgte sjåfører bruke bilen til normal kjøring til batteriet er tomt, for hver sjåfør er kjørelengda (km) registrert. Vi går ut ifra at disse kjørelengdene er et tilfeldig utvalg fra en populasjon med forventning *μ* (*μ* = 201.0667) og varians *σ*<sup>2</sup> (*σ*<sup>2</sup> = 111.0667). Produsenten av bilmodellen påstår er forventa rekkevidde er under tilsvarende forhold 210 (km). Vi vil undersøke om det er grunnlag for å påstå at en reell rekkevidde er lavere enn dette.

Sett opp hypoteser og testobservator, og utfør en test for problemstillinga over med signifikansnivå 5%. Finn både forkastningsområde og (tilnærma) *p*-verdi. Hva blir konklusjonen? Vi går ut ifra at observatorene i et tilfeldig utvalg *X*<sub>1</sub>, . . . , *X*<sub>15</sub> fra en normalfordelt populasjon med forventning *μ* og standardavvik *σ*. Vi vil teste hypotesene:

*H*<sub>0</sub> : *μ* = 210

*H*<sub>1</sub> : *μ* < 210

Ettersom *σ* er ukjent og maa estimeres fra data ved *S* blir dette en T-test, testobservator:

T
=



X
¯
−

μ

0




S


√
n

som er *t-fordelt* med *n* − 1 = 14 frihtsgrader når *H*<sub>0</sub> er sann. Insett data:

t
=



201.067
−
210


10.539


√
n





=
−
3.282

Forkaster *H*<sub>0</sub> om *t* < *t*<sub>0.02,14</sub> = −1.761. Så *H*<sub>0</sub> blir forkasta, og vi kan påstå at forventnia rekkevidde *μ* er under 210km.

Kan finne *p*-verdien fra:

p
−
verdi
=
P
(

T

14


<
−
3.282
)
=
⎧
<
0.005


>
0.001


⎫

Så fra tabellen kan vi se at *p*-verdien er mellom 0.1% og 0.5%. Da dette er lavere enn signifikansnivået på 0.5% blir *H*<sub>0</sub> forkasta.

**Deloppgave**

Vi er og interessert i hvor mye rekkevidda kan variere. Utled et 90%-konfidensintervall for variansen *σ*<sup>2</sup>.

Finn intervallestimatet fra observatorene. Hva kan du konkludere fra intervallet?

Utleder et 90%-konfidensintervall fra variansen, *σ*<sup>2</sup>, bruker at vi kjenner fordelinga for den standariserte observatoren

(
n
−
1

)

S

2




σ

2





∼

χ

2


n
−
1

P
⎛

χ

2


0.95


<



(
n
−
1

)

S

2




σ

2





<

χ

2


0.05


⎞
=
0.90

P
⎛



χ

2


0.95


(
n
−
1

)

S

2





<


1


σ

2





<



χ

2


0.05


(
n
−
1

)

S

2





⎞
=
0.90

P
⎛



(
n
−
1

)

S

2





χ

2


0.05





<

σ

2


<



(
n
−
1

)

S

2





χ

2


0.95





⎞
=
0.90

Oppgave

Ola samler spillekort. Ønsker et spesifikt spillekort (A). X er antall kort han kjøper før han får dette. Det er 10 forskjellige spillekort i serien. Han kjøper pakker av kort med  $n = 20$  kort i hver pakke. Hver pakke er uavhengig. Hva er fordelinga for antall kort i pakka som er med spiller A? Hva er sannsynligheten for at Ola får minst ett kort med spiller A i pakka? Hva er forventna antall ulike spillekort i pakka? Ola gjør  $n = 20$  kjøp. Resultatet av hvert kjøp er uavhengige, og sannsynligheten for suksess (kort A) er lik  $\frac{1}{10}$  for hvert kjøp. Da blir  $Y =$  tall på kort A, binomisk fordelt med  $n = 20$  delforsøk og sannsynlighet  $p = \frac{1}{10}$ :

$$Y \sim \text{binom} \left( n = 20, p = \frac{1}{10} \right)$$
$$P(Y \geq 1) = 1 - P(X = 0) = 1 - \binom{20}{0} \left( \frac{1}{10} \right)^0 \left( 1 - \frac{1}{10} \right)^{20} = 1 - 0.122 = 0.878$$

Definerer korttypene  $i \in \{A, B, \dots, J\}$ :

$$I_i = \begin{cases} 1, & \text{minst ett kort av typen } i. \\ 0, & \text{ingen kort av typen } i \end{cases}$$

Da blir

$$E(I_i) = 1 \cdot P(I_i = 1) + 0 \cdot P(I_i = 0) = P(I_i = 1) = P(Y \geq 1) = 0.878$$

La tallet på ulike kort være  $N = \sum_i I_i$ . Da blir

$$E(N) = E(\sum_i I_i) = \sum_i E(I_i) = \sum_i 0.878 = 10 \cdot 0.878 = 8.78$$

Oppgave

Fordelinga til en kontinuerlig stokastisk variabel X er gitt ved følgende sannsynlighetstetthetsfunksjon:

$$f(x) = \begin{cases} \frac{2}{9}x, & 0 < x < 3 \\ 0, & \text{ellers} \end{cases}$$

Vi er interresert i følgende hendelser:

$$X = \{X > 1\} \qquad \qquad \qquad \text{og} \qquad \qquad \qquad B = \{X > 2\}$$

er disse hendelsene disjunkte?

Hva er sannsynligheten for B?

Hva er sannsynligheten for B gitt A?

Er disse hendelsene uavhengige?

Nei, A og B er ikke disjunkte, det er mulig at  $X > 1$  samtidig som  $X > 2$  ( $X = 2.5$ )

$$P(B) = P(X > 2) = \int_2^3 f(x) \, dx = \int_2^3 \frac{2}{9}x \, dx \left[ \frac{1}{9}x^2 \right]_2^3 = 1 - \frac{4}{9} = \frac{5}{9}$$
$$P(B|A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A|B)P(B)}{P(A)} = \frac{1 \cdot P(B)}{P(A)} = \frac{\frac{5}{9}}{\int_1^3 \frac{1}{9}x \, dx} = \frac{5}{8}$$

Vi ser at  $P(B|A) \neq P(B)$ , så A og B er ikke uavhengige

Oppgave

Vi vil tilpasse en lineær regresjonsmodell for  $Y =$  dager, og  $x =$  år, det vil si

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \qquad \qquad \qquad i = 1, \dots, 48$$

der vi går ut ifra at feilledda  $\epsilon_1, \dots, \epsilon_{48}$  er uavhengige og normalfordelte med forventning 0 og varians  $\sigma^2$ .

Regresjonsanalysen er gjennomført i R:

```
summary(lm(y ~ x)) -> (Intercept) = b_0, x = b_1
mean(x) = 24.5, sum((x - mean(x))^2) = 9212
res = residuals(lm(y ~ x)), sum(res^2) / (48 - 2) = 150.1756
```

Skriv ned minste kvadrat-estimatoren for  $\beta_0, \beta_1$ , og finn estimata fra R (ish) utskrifta. Gi ei presis tolkning av hva de estimerte verdiene forteller deg om sammenhengen.

Bruk den estimerte regresjonslinja til å finne et predikert verdi for dager i år 49

$$B_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, \qquad \qquad \qquad B_0 = \bar{Y} - B_1 \cdot \bar{X}$$

Fra R-utskrifta kan vi finne MKM-estimata for  $\beta_1, \beta_0$ , og den estimerte regresjonslinja:

$$b_1 = -0.411, \qquad \qquad b_0 = 67.175 \Rightarrow \qquad \qquad \hat{y} = b_0 + b_1 x = 67.175 - 0.411x$$

$b_0$  = skjæring med y-aksen, estimert forventna dager ved  $x=0$  år.

$b_1 = -0.411$  er stigningstallet, estimert endring i firvnta dager ekstra for hvert år.

Predikert verdi for  $x_0 = 49$  er:

$$\hat{y}_0 = b_0 + b_1 x_0 = 67.175 - 0.411 \cdot x = 47.036$$

.

Deloppgave

Vi er spesielt interresert i den totale endringa i antall dager som har skjedd fra  $x = 1$  til  $x = 48$ . Som estimator for endring i forventna dager i løpet av de 47 åra skal vi nytte:

$$\Delta \hat{Y} = \hat{Y}_{48} - \hat{Y}_1$$

(Her er  $\hat{Y}_i = B_0 + B_1 x$ , altså estimert regresjonslinje år  $x_i$ )

Vis at  $\Delta \hat{Y} = (48 - 1) \cdot B_1$ , der  $B_1$  er estimatoren for stigningstallet.

Vis at  $\Delta \hat{Y}$  er en forventningsrett estimator for endringa over 47 år, og finn variansen til estimatoren.

Bruk dette som utgangspunkt til å utlede et 95%-konfidensintervall for forventna endring, og estimer intervallet ved hjelp av utskriftene over.

Kan du utifra dette intervallet påstå signifikant endring i antall dager?

Differansen i estimert regresjonslinje mellom  $x_2 = 48$  og  $x = 1$ :

$$\Delta \hat{Y} = \hat{Y}_{48} - \hat{Y}_1 = B_0 + B_1 \cdot 48 - (B_0 + B_1 \cdot 1) = 47B_1$$
$$E(\Delta \hat{Y}) = E(47 \cdot B_1) = 47 \cdot E(B_1) = 47 \cdot \beta_1$$

$$\text{Vet}(\Delta \hat{Y}) = \text{Var}(47 \cdot B_1) = 47^2 \cdot \text{Var}(B_1) = 47^2 \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Her er det vrukt at vi vet at for stigningstallestimatoren er

$$E(B_1) = \beta_1 \qquad \qquad \qquad \text{og} \qquad \qquad \qquad \text{Var}(B_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Merk at forventningsverdien til estimatoren er lik forventna endring over 47 år da

$$\Delta \mu_Y = \mu_{Y|48} - \mu_{Y|1} = \beta_0 + \beta_1 \cdot 48 - (\beta_0 + \beta_1 \cdot 1) = 47 \cdot \beta_1$$

Da vi vet at estimatoren  $B_1$  er normalfordelt vil og  $\Delta \hat{Y}$  bli det, og vi har at

$$\Delta \hat{Y} \sim N \left( \Delta \mu_Y, \frac{47 \sigma}{\sqrt{\sum_{i=1}^n 48(x_i - \bar{x})^2}} \right)$$

Dermed kan vi utlede et 95%-KI ved å ta utgangspunkt i en t-fordelt observator

$$P \left( -t_{0.025} < \frac{\Delta \hat{Y} - \Delta \mu_Y}{\frac{47 \cdot S}{\sqrt{\sum_{i=1}^n 48(x_i - \bar{x})^2}}} < t_{0.025} \right) = 0.95$$
$$P \left( \Delta \hat{Y} - t_{0.025} \frac{47 \cdot S}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} < \Delta \mu_Y < \Delta \hat{Y} + t_{0.025} \frac{47 \cdot S}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \right) = 0.95$$

. Med  $n - 2 = 46$  frihetsgrader (bruker  $t_{0.025, 40} = 2.021$ ) blir et 95%-intervallestimat for  $\Delta \mu_Y$ :

$$47b_1 \pm t_{0.025, 46} \frac{47s}{\sqrt{\sum_{i=1}^n 48(x_i - \bar{x})^2}} \approx 47 \cdot (-0.411) \pm 2.021 \cdot \frac{47 \cdot 12.25}{\sqrt{9212}} = (-31.440, -7.193)$$

Fra R:  $x = 12.25$  og  $\sum_{i=1}^{48} (x_i - \bar{x})^2 = 9212$ . Alternativt:  $\hat{SE}(B_1) = \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$ . Det er 95% sannsynlighet for at intervallet skal inneholde forventna endring over 47 år. Og da intervallet ikke dekker 0 kan vi påstå med 5% signifikansnivå at det har skjedd ei endring.

You can do it!

I belive in you!  
JUST DO IT!  
(;

*This page was intentionally left blank*