



Norges teknisk-naturvitenskapelige universitet

Institutt for matematiske fag

Løsningsforslag - Eksamen desember 2011

ST1201 Statistiske metoder

Oppgave 1

a)

Dette er en ANOVA-tabell for k -utvalg med $k = 4$ og $n_j = 6$ for $j = 1, 2, 3, 4$. Den fullstendige ANOVA-tabellen blir

Kilde	df	SS	MS	F
Betong	$k - 1 = 3$	47203.13	15734.38	2.90
Error	$6 \cdot 4 - 4 = 20$	108671.50	5433.58	
Total	$6 \cdot 4 - 1 = 23$	155874.63		

der

$$SSTR = MSTR \cdot 3 = 47203.14,$$

$$MSE = \frac{SSE}{20} = \frac{10861.50}{20} = 5433.58,$$

$$SSTOT = SSTR + SSE = 47203.13 + 108671.50 = 155874.63$$

og

$$F = \frac{MSTR}{MSE} = \frac{15734.38}{5433.58} = 2.90.$$

Testobservatoren F relaterer seg til hypotesene

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 \quad \text{mot} \quad H_1 : \text{ikke slik,}$$

der μ_i , for $i = 1, 2, 3, 4$, er forventet opptatt fuktighet for betong av type nummer i .

Når H_0 er riktig er F Fisher fordelt med 3 og 20 frihetsgrader. Finner kritisk verdi for $\alpha = 0.05$ fra tabell til å være $f_{0.05,3,20} = 3.10$. Beslutningsregelen blir dermed at vi skal forkaste H_0 når $F > 3.10$. Betongdataene gav $F = 2.90 < 3.10$ slik at konklusjonen blir at vi skal ikke forkaste H_0 .

b)

En to-utvalg t -test baserer seg på at man har observasjoner av stokastiske variabler X_1, \dots, X_n og Y_1, \dots, Y_m der alle X_i -er og Y_i er uavhengige av hverandre,

$$X_i \sim N(\mu_X, \sigma^2) \quad \text{og} \quad Y_i \sim N(\mu_Y, \sigma^2).$$

Man benytter da testobservatoren

$$T = \frac{\bar{X} - \bar{Y}}{\sqrt{S_p^2 \left(\frac{1}{n} + \frac{1}{m}\right)}}$$

som er Student t -fordelt med $n + m - 2$ frihetsgrader når $H_0 : \mu_X = \mu_Y$ er riktig. Varians-estimatoren S_p^2 er gitt ved formelen

$$S_p^2 = \frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}.$$

Vi lar X_i -ene og Y_i -ene være henholdsvis data for betong av type 3 og 4. Ved å benytte oppgitte verdier for S_X^2 og S_Y^2 i tabellen på første side av oppgavesettet får man

$$s_p^2 = \frac{5 \cdot 3593.50 + 5 \cdot 3704.27}{10} = 3648.89, \quad t = \frac{\frac{3663}{6} - \frac{2924}{6}}{\sqrt{3648.89 \left(\frac{1}{6} + \frac{1}{6}\right)}} = 3.53.$$

Man må her benytte en tosidig test slik at kritisk verdi blir $t_{\frac{\alpha}{2}, n+m-2} = t_{0.025, 10} = 2.228$. Beslutningsregelen blir dermed at man skal forkaste H_0 dersom $T < -2.228$ eller $T > 2.228$. Vi observerte $t = 3.53 > 2.228$, slik at konklusjonen blir at vi forkaster H_0 .

Det er ikke urimelig at vi i punkt **a)** konkluderer med at det ikke er signifikant forskjell mellom forventningsverdiene til de fire utvalgene, mens vi her i punkt **b)** konkluderer med at det er signifikant forskjell mellom forventningsverdiene til utvalg 3 og 4. Vi kan spesielt legge merke til at vi her i punkt **b)** sammenligner de to av de fire utvalgene som har størst avvik i gjennomsnittsverdi. Vi kan også legge merke til at empirisk varians for utvalg nummer 1 er betydelig større enn for de andre tre utvalgene. ANOVA-analysen baserer seg som kjent på antagelsen om lik varians for alle utvalg. Den store empiriske variansen for utvalg nummer 1 vil dermed føre til at estimert ("pooled") varians i ANOVA-analysen blir betydelig større enn tilsvarende størrelse i t -testen.

Oppgave 2

a)

Rimelighetsfunksjonen blir her

$$L(p) = f_Y(y; p) = \binom{m}{y} p^y (1-p)^{m-y}.$$

Log-likelihoodfunksjonen blir dermed

$$l(p) = \ln[L(p)] = \ln \binom{m}{y} + y \cdot \ln p + (m - y) \ln(1 - p).$$

Deriverer og setter lik null:

$$l'(p) = 0 + \frac{y}{p} + \frac{m-y}{1-p} \cdot (-1) = \frac{y}{p} - \frac{m-y}{1-p} = 0 \Rightarrow y(1-p) = p(m-y) \Rightarrow y - yp = pm - py \Rightarrow p = \frac{y}{m}.$$

Dermed får vi at SME blir

$$\hat{p} = \frac{Y}{m}.$$

Siden vi kun har en parameter vi skal estimere, og kun har en observasjon, finner man momentestimatoren ved å sette forventet verdi for Y lik observert verdi for Y . Siden $E[Y] = mp$ får vi

$$m\hat{p} = Y \Rightarrow \hat{p} = \frac{Y}{m}.$$

b)

En estimator $\hat{\theta}$ for en parameter θ er en *beste* estimator hvis den er forventningsrett og har minst like liten varians som enhver annen forventningsrett estimator.

For å vise at en estimator er en *beste* estimator kan man sjekke at den er forventningsrett og sjekke at variansen til estimatoren er lik Cramér-Raos nedre grense for forventningsrette estimatorene.

For \hat{p} har vi

$$E[\hat{p}] = E\left[\frac{Y}{m}\right] = \frac{E[Y]}{m} = \frac{mp}{m} = p,$$

og

$$\text{Var}[\hat{p}] = \text{Var}\left[\frac{Y}{m}\right] = \frac{\text{Var}[Y]}{m^2} = \frac{mp(1-p)}{m^2} = \frac{p(1-p)}{m}.$$

Regner ut Cramér-Raos nedre grense, Starter med

$$\ln f_Y(y; p) = \ln \binom{m}{y} + y \cdot \ln p + (m - y) \ln(1 - p).$$

Deriverer to ganger med hensyn på p :

$$\begin{aligned} \frac{\partial \ln f_Y(y; p)}{\partial p} &= \frac{y}{p} + \frac{m-y}{1-p} \cdot (-1) = \frac{y}{p} - \frac{m-y}{1-p}, \\ \frac{\partial^2 \ln f_Y(y; p)}{\partial p^2} &= -\frac{y}{p^2} - \frac{m-y}{(1-p)^2} \cdot (-1) = \frac{m-y}{(1-p)^2} - \frac{y}{p^2}. \end{aligned}$$

Tar forventningsverdien,

$$\begin{aligned} E \left[\frac{\partial \ln f_Y(y; p)}{\partial p} \right] &= E \left[\frac{m - y}{(1 - p)^2} - \frac{y}{p^2} \right] = \frac{m - E[Y]}{(1 - p)^2} - \frac{E[Y]}{p^2} = \frac{m - mp}{(1 - p)^2} - \frac{mp}{p^2} \\ &= \frac{m}{1 - p} - \frac{m}{p} = \frac{m(1 - p + p)}{p(1 - p)} = \frac{m}{p(1 - p)}. \end{aligned}$$

Cramé-Raos nedre grense for varians av forventningsrette estimatorene blir dermed (hvor vi benytter at \hat{p} er basert på kun $n = 1$ stokastiske variabler),

$$\left\{ -n E \left[\frac{\partial^2 \ln f_Y(Y; \theta)}{\partial \theta^2} \right] \right\}^{-1} = \frac{p(1 - p)}{m}.$$

Vi ser dermed at \hat{p} var forventningsrett og at variansen for \hat{p} er lik Cramér-Raos nedre grense. Dermed er \hat{p} en beste estimator for p .

c)

Vi ser at $\hat{\theta}$ er forventningsskjev fordi

$$\begin{aligned} E[\hat{\theta}] &= E \left[Y - \frac{Y^2}{m} \right] = E[Y] - \frac{E[Y^2]}{m} = mp - \frac{\text{Var}[Y] + E[Y]^2}{m} \\ &= mp - \frac{mp(1 - p) + (mp)^2}{m} = mp - p(1 - p) + mp^2 = (m - 1)p(1 - p) = \frac{m - 1}{m} \theta \neq \theta. \end{aligned}$$

Vi ser at $\hat{\theta}$ er assymptotisk forventningsrett fordi

$$\lim_{m \rightarrow \infty} E[\hat{\theta}] = \lim_{m \rightarrow \infty} \left[\frac{m - 1}{m} \theta \right] = \theta.$$

Vi ser at vi får en forventningsrett estimator med å dele $\hat{\theta}$ på $(m - 1)/m$. Den forventningsrette estimatoren blir dermed

$$\hat{\theta} = \frac{\hat{\theta}}{\frac{m-1}{m}} = \frac{mY}{m-1} \left(1 - \frac{Y}{m} \right).$$

Oppgave 3

a)

Rimelighetsfunksjonen blir

$$L(\alpha, \beta) = \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi} \sigma_0} \exp \left\{ -\frac{1}{2\sigma_0^2} (y_i - \alpha - \beta(x_i - \bar{x}))^2 \right\} \right].$$

Logrimelighetsfunksjonen blir

$$\begin{aligned} l(\alpha, \beta) &= \sum_{i=1}^n \left[-\frac{1}{2} \ln(2\pi) - \ln \sigma_0 - \frac{1}{2\sigma_0^2} (y_i - \alpha - \beta(x_i - \bar{x}))^2 \right] \\ &= -\frac{n}{2} \ln(2\pi) - n \ln \sigma_0 - \frac{1}{2\sigma_0^2} \sum_{i=1}^n (y_i - \alpha - \beta(x_i - \bar{x}))^2. \end{aligned}$$

Partiellderiverer med hensyn på α og setter lik null,

$$\begin{aligned} \frac{\partial l}{\partial \alpha} &= -\frac{1}{2\sigma_0^2} \sum_{i=1}^n 2(y_i - \alpha - \beta(x_i - \bar{x})) \cdot (-1) = 0 \\ \Rightarrow \sum_{i=1}^n y_i - n\alpha - \beta \sum_{i=1}^n (x_i - \bar{x}) &= 0 \Rightarrow \alpha = \frac{1}{n} \sum_{i=1}^n y_i, \end{aligned}$$

hvor vi har benyttet at $\sum_{i=1}^n (x_i - \bar{x}) = 0$. Partiellderiverer så med hensyn på β og setter lik null,

$$\begin{aligned} \frac{\partial l}{\partial \beta} &= -\frac{1}{2\sigma_0^2} \sum_{i=1}^n 2(y_i - \alpha - \beta(x_i - \bar{x})) \cdot (-(x_i - \bar{x})) = 0 \\ \Rightarrow \sum_{i=1}^n y_i(x_i - \bar{x}) - \alpha \sum_{i=1}^n (x_i - \bar{x}) - \beta \sum_{i=1}^n (x_i - \bar{x})^2 &= 0 \\ \Rightarrow \beta &= \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, \end{aligned}$$

hvor vi igjen har benyttet at $\sum_{i=1}^n (x_i - \bar{x}) = 0$.

Sannsynlighetsmaksimeringsestimatorene for α og β er dermed gitt ved

$$\hat{\alpha} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{og} \quad \hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Ved å benytte at Y_1, \dots, Y_n er uavhengige får vi

$$\begin{aligned} \text{Var}[\hat{\beta}] &= \text{Var} \left[\frac{\sum_{i=1}^n (x_i - \bar{x})Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] = \frac{\text{Var}[\sum_{i=1}^n (x_i - \bar{x})Y_i]}{(\sum_{i=1}^n (x_i - \bar{x})^2)^2} = \frac{\sum_{i=1}^n \text{Var}[(x_i - \bar{x})Y_i]}{(\sum_{i=1}^n (x_i - \bar{x})^2)^2} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}[Y_i]}{(\sum_{i=1}^n (x_i - \bar{x})^2)^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_0^2}{(\sum_{i=1}^n (x_i - \bar{x})^2)^2} = \frac{\sigma_0^2 \sum_{i=1}^n (x_i - \bar{x})^2}{(\sum_{i=1}^n (x_i - \bar{x})^2)^2} \\ &= \frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}. \end{aligned}$$

b)

$\hat{\beta}$ er en lineærkombinasjon av Y_i -ene som er uavhengige og normalfordelte variabler. Dermed blir også $\hat{\beta}$ normalfordelt, dvs.

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right).$$

Dermed har vi også at

$$\frac{\hat{\beta} - \beta}{\sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim N(0, 1)$$

slik at

$$P\left(-z_{\frac{a}{2}} \leq \frac{\hat{\beta} - \beta}{\sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \leq z_{\frac{a}{2}}\right) = 1 - a$$

Løser hver ulikhet med hensyn på β . Starter med den første,

$$-z_{\frac{a}{2}} \leq \frac{\hat{\beta} - \beta}{\sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \Leftrightarrow -\hat{\beta} - z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \leq -\beta \Leftrightarrow \hat{\beta} + z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \geq \beta$$

$$\frac{\hat{\beta} - \beta}{\sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \leq z_{\frac{a}{2}} \Leftrightarrow -\beta \leq -\hat{\beta} + z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \Leftrightarrow \beta \geq \hat{\beta} - z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Må dermed også ha at

$$P\left(\hat{\beta} - z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \leq \beta \leq \hat{\beta} + z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}\right) = 1 - a.$$

Et $(1 - a) \cdot 100\%$ konfidensintervall for β blir dermed

$$\left[\hat{\beta} - z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}, \hat{\beta} + z_{\frac{a}{2}} \sqrt{\frac{\sigma_0^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}\right].$$

c)

For å finne et prediksjonsintervall tar vi utgangspunkt i

$$\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0.$$

Denne vil være normalfordelt fordi det er en lineærkombinasjon av uavhengige normalfordelte variabler, nemlig Y_1, \dots, Y_n og Y_0 . Dette kan vi se ved å sette inn i uttrykket over hva vi i punkt **a)** fant for $\hat{\alpha}$ og $\hat{\beta}$.

Forventningsverdien til dette uttrykket blir

$$E \left[\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0 \right] = E[\hat{\alpha}] + E[\hat{\beta}] (x_0 - \bar{x}) - E[Y_0] = \alpha + \beta(x_0 - \bar{x}) - (\alpha + \beta(x_0 - \bar{x})) = 0.$$

Siden Y_0 åpenbart er uavhengig av $\hat{\alpha}$ og $\hat{\beta}$ får vi for tilhørende varians

$$\begin{aligned} \text{Var} \left[\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0 \right] &= \text{Var} \left[\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) \right] + \text{Var}[Y_0] \\ &= \text{Var}[\hat{\alpha}] + \text{Var} \left[\hat{\beta}(x_0 - \bar{x}) \right] + 2\text{Cov} \left[\hat{\alpha}, \hat{\beta}(x_0 - \bar{x}) \right] + \text{Var}[Y_0] \\ &= \text{Var}[\hat{\alpha}] + (x_0 - \bar{x})^2 \text{Var}[\hat{\beta}] + 2(x_0 - \bar{x}) \text{Cov}[\hat{\alpha}, \hat{\beta}] + \text{Var}[Y_0] \end{aligned}$$

Fra oppgaveteksten har vi uttrykk for $\text{Var}[\hat{\alpha}]$ og $\text{Var}[\hat{\beta}]$, og $\text{Var}[Y_0] = \sigma_0^2$. Trenger å regne ut $\text{Cov}[\hat{\alpha}, \hat{\beta}]$. Ved å benytte at $\text{Cov}[Y_i, Y_j] = \text{Var}[Y_i] = \sigma_0^2$ dersom $i = j$ og lik 0 hvis $i \neq j$ får vi at

$$\begin{aligned} \text{Cov}[\hat{\alpha}, \hat{\beta}] &= \text{Cov} \left[\frac{1}{n} \sum_{i=1}^n Y_i, \frac{\sum_{j=1}^n (x_j - \bar{x}) Y_j}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov}[Y_i, (x_j - \bar{x}) Y_j] \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n \sum_{j=1}^n (x_j - \bar{x}) \text{Cov}[Y_i, Y_j] \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \sigma_0^2 = \frac{\sigma_0^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) = 0. \end{aligned}$$

Dermed får vi at

$$\text{Var} \left[\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0 \right] = \frac{\sigma_0^2}{n} + \frac{\sigma_0^2 (x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} + \sigma_0^2 = \sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right),$$

og dermed også

$$\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0 \sim N \left(0, \sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \right)$$

og

$$\frac{\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0}{\sqrt{\sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}} \sim N(0, 1),$$

slik at

$$P \left(-z_{\frac{a}{2}} \leq \frac{\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - Y_0}{\sqrt{\sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}} \leq z_{\frac{a}{2}} \right) = 1 - a.$$

Løser så hver ulikhet hver for seg med hensyn på Y_0 og setter deretter ulikhetene sammen igjen med Y_0 alene i midten, og får

$$\begin{aligned} P \left(\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - z_{\frac{a}{2}} \sqrt{\sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \leq Y_0 \right. \\ \left. \leq \hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) + z_{\frac{a}{2}} \sqrt{\sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right) = 1 - a. \end{aligned}$$

Et $(1 - a) \cdot 100\%$ prediksjonsintervall for Y_0 blir dermed

$$\left[\hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) - z_{\frac{a}{2}} \sqrt{\sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}, \hat{\alpha} + \hat{\beta}(x_0 - \bar{x}) + z_{\frac{a}{2}} \sqrt{\sigma_0^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right].$$

Vi ser at prediksjonsintervallet blir kortest når $x_0 - \bar{x} = 0$, dvs. når $x_0 = \bar{x}$.