



Norges teknisk-naturvitenskapelige universitet
Institutt for matematiske fag

SIF 5062 Statistikk

Løsningsforslag - Eksamen mai 2003

Oppgave 1

a)

$$P(X \leq 120) = P\left(\frac{X - 115}{5} \leq \frac{120 - 115}{5}\right) = \Phi(1) = 0.841.$$

$$\begin{aligned} P(120 < X \leq 125) &= P(X \leq 125) - P(X \leq 120) \\ &= P\left(\frac{X - 115}{5} \leq \frac{125 - 115}{5}\right) - P(X \leq 120) \\ &= \Phi(2) - 0.841 = 0.977 - 0.841 = 0.136. \end{aligned}$$

A og B er disjunkte siden de to hendelsene ikke kan inntreffe samtidig.

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(\emptyset)}{P(B)} = \frac{0}{P(B)} = 0.$$

Siden $P(A) = 0.841 \neq P(A | B)$ er A og B ikke uavhengige.

b)

Sjekker kriteriene for Bernoulli-forsøk:

- Vi undersøker navnet til n jenter, dvs. gjør n forsøk
- I hvert forsøk er navnet enten Maud eller ikke Maud
- Sannsynligheten for at navnet er Maud er $p = 0.02$ og lik i alle forsøk
- Det er rimelig å anta at forsøkene er uavhengige

Dermed er de fire kriteriene for et Bernoulli-forsøk oppfylt, og det er rimelig å anta at Z er binomisk fordelt.

$$P(C) = P(Z \geq 1) = 1 - P(Z = 0) = 1 - \binom{15}{0} 0.02^0 (1 - 0.02)^{15} = 1 - 0.739 = 0.261$$

$$\begin{aligned} P(D | C) &= P(Z = 2 | Z \geq 1) \\ &= \frac{P(Z = 2 \cap Z \geq 1)}{P(Z \geq 1)} = \frac{P(Z = 2)}{P(Z \geq 1)} = \frac{\binom{15}{2} 0.02^2 (1 - 0.02)^{13}}{P(Z \geq 1)} = \frac{0.032}{0.261} = 0.123. \end{aligned}$$

Definer hendelsene

M : navnet på tilfeldig valgt elev er Maud

J : tilfeldig valgt elev er jente

G : tilfeldig valgt elev er gutt

Bruker lov om total sannsynlighet:

$$P(M) = P(M \cap J) + P(M \cap G) = P(M | J)P(J) + P(M | G)P(G) = 0.02 \cdot \frac{15}{25} + 0 \cdot \frac{15}{25} = 0.012.$$

Oppgave 2

a)

Det er mest rimelig med en venstresidig hypotesetest:

$$H_0 : \quad \mu = 16,$$

$$H_1 : \quad \mu < 16.$$

Begrunnelse: forhandleren sier at bilen kan forventes å kjøre **minst** 16 km pr liter. Vi vil avsløre ev. feil i markedsføringen.

NB: Hypotesetesten skal være uavhengig av målingene. En bør altså ikke velge alternativ hypotese på grunnlag av \bar{x} .

\bar{X} er normalfordelt med forventning μ og varians σ^2/n . Variansen er ukjent, derfor kreves T-fordeling med $\nu = n - 1 = 19$ frihetsgrader. Gjennomfører testen med $\alpha = 0.05$. Testobservator: $T_{\text{obs}} = \frac{\bar{X} - \mu}{S/\sqrt{n}}$. Observert verdi:

$$t_{\text{obs}} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{15.56 - 16}{0.94/\sqrt{20}} = -2.093.$$

Fra tabell over kvantiler i T-fordelingen; $-t_{0.05,19} = -1.729$. Altså: $t < -t_{0.05,19}$, dermed skal H_0 forkastes.

Hvis vi hadde valgt å bruke en normalfordelingshypotese, ville kvantilen $-z_{0.05} = -1.645$ gitt samme konklusjon. Imidlertid bør vi da argumentere for at avstanden til denne kvantilen er så stor at høyere varians i T-fordelingen ikke ville påvirket resultatet. Å sammenlikne med denne kvantilen kan ikke regnes som fullgodt svar.

b)

P-verdien finnes ved å lete opp verdien på testobservatoren fra a) i tabell. For T-fordeling med $\nu = 19$, finner vi $t_{0.025,19} = 2.093$. Ettersom T-fordelingen er symmetrisk, har vi at $P(T > t_{\alpha,\nu}) = P(T < -t_{\alpha,\nu})$. Dermed; $p = \alpha = 0.025 = 2.5\%$.

Testobservatoren er normalfordelt hvis $\sigma = s$. Dette bør være tilnærmet oppfylt for å bruke normalfordeling. Hvis en ikke har ekstra informasjon om σ , er det ikke anbefalt å tilnærme student-fordelingen med en normalfordeling når $n < 30$, da s ikke er et godt nok estimat.

Under normalfordelingen får vi p-verdi

$$P(z \leq -2.09) = P(z \geq 2.09) = 1 - P(z < 2.09) = 1 - \Phi(2.09) = 1 - 0.9817 = 0.0183 = 1.8\%.$$

c)

Antar $H'_1 : \mu = \mu_1 = 15.5$ og $\sigma = s$. Teststyrken er sannsynligheten for å forkaste H_0 under H'_1 , dvs

$$P\left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} < -1.645 \mid \mu = \mu_1 = 15.5\right).$$

For å få en normalfordelt variabel, flytter vi alt utenom \bar{X} , som er stokastisk, over på høyre side. Deretter trekker vi fra sann forventningsverdi μ_1 og dividerer med standardavviket på begge sider.

$$\begin{aligned} & P(\bar{X} < -1.645 \cdot \sigma/\sqrt{n} + \mu_0) \\ &= P\left(\frac{\bar{X} - \mu_1}{\sigma/\sqrt{n}} < \frac{-1.645 \cdot \sigma/\sqrt{n} + \mu_0 - \mu_1}{\sigma/\sqrt{n}}\right) \\ &= P(Z < 0.7338) \approx 0.767. \end{aligned}$$

Hvis vi ikke kunne bruke normalfordelingsantakelsen, ville teststyrken blitt svakere. Her er det forutsatt at vi er ganske sikre på variansen, f.eks. på grunnlag av data fra produsent.

Generelt må antall observasjoner økes for å oppnå økt teststyrke. (Dette er fullgodt svar.) Mulig tillegg: Hvis en har mulighet til å gjennomføre forsøket på en måte slik at variansen blir mindre, f.eks. kjøre bilene under mer kontrollerte former i et laboratorium, ville også teststyrken økes. Eventuelt kan en øke signifikansnivået α f.eks. til 0.1, og dermed øke teststyrken, men dette er sjelden aktuelt i praksis.

Oppgave 3

a)

Kumulativ fordelingsfunksjon finnes ved å integrere sannsynlighetstettheten.

$$F(t) = \int_{-\infty}^t f(\tau) d\tau = \int_0^t \frac{1}{\beta} e^{-\tau/\beta} d\tau = 1 - e^{-t/\beta},$$

for $t \geq 0$, og $F(t) = 0$ for $t < 0$.

For $\beta = 5$:

$$P(T_1 < 3) = F(3) = 1 - e^{-3/5} = 0.4512.$$

$$P(2 < T_1 < 4) = P(T_1 < 4) - P(T_1 \leq 2) = F(4) - F(2) = e^{-2/5} - e^{-4/5} = 0.2210.$$

b)

Apparatet fungerer bare hvis k uavhengige komponenter fungerer. Altså er levetiden X lik levetiden til komponenten som bryter sammen først. Dvs.

$$X = \text{MIN}(T_1, T_2, T_3, \dots, T_k).$$

Dermed er

$$\begin{aligned} P(X > x) &= P(T_1 > x \cap T_2 > x \cap T_3 > x \dots \cap T_k > x) \\ &= P(T_1 > x)P(T_2 > x)P(T_3 > x) \dots P(T_k > x). \end{aligned}$$

Den siste likheten kommer av at T_1, T_2, \dots, T_k er uavhengige.

Videre, $P(T_j > x) = 1 - F(x) = e^{-x/\beta}$ for alle $j = 1, 2, 3, \dots, k$. Dermed har vi

$$P(X > x) = \left(e^{-x/\beta}\right)^k = e^{-kx/\beta} = e^{-\frac{x}{\beta/k}}.$$

Dette tilsvarer en kumulativ eksponensialfordeling med parameter β/k , dvs

$$F_X(x) = 1 - e^{-\frac{x}{\beta/k}}, \quad x \geq 0.$$

Forventningsverdien til en eksponensialfordeling med parametrisering som gitt i formelsamlingen, er parameteren selv. Altså er forventningsverdien til X lik β/k_i . Med $k = 4$ og $\beta = 5$, er apparatets forventede levealder $5/4 = 1.25$. (Tidsenhet ikke oppgitt i oppgaven.)

c)

Forventningsverdi $\hat{\beta}$:

$$\begin{aligned} E(\hat{\beta}) &= E\left(\frac{1}{n} \sum_{i=1}^n X_i k_i\right) \\ &= \frac{1}{n} \sum_{i=1}^n E(X_i) k_i \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\beta}{k_i} k_i = \beta. \end{aligned}$$

Forventningsverdi $\tilde{\beta}$:

$$\begin{aligned} E(\tilde{\beta}) &= E\left(\frac{\sum_{i=1}^n X_i}{\sum_{i=1}^n k_i^{-1}}\right) \\ &= \frac{\sum_{i=1}^n E(X_i)}{\sum_{i=1}^n k_i^{-1}} \\ &= \frac{\sum_{i=1}^n \frac{\beta}{k_i}}{\sum_{i=1}^n k_i^{-1}} = \beta. \end{aligned}$$

Dermed er begge estimatorene forventningsrette. Beregner variansene:

Varians til $\hat{\beta}$:

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \frac{1}{n^2} \sum_{i=1}^n k_i^2 \text{Var}(X_i) \quad (\text{pga. uafhængighed}) \\ &= \frac{1}{n^2} \sum_{i=1}^n k_i^2 \left(\frac{\beta}{k_i}\right)^2 \\ &= \frac{1}{n} \beta^2. \end{aligned}$$

Varians til $\tilde{\beta}$:

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= \frac{\sum_{i=1}^n \text{Var}(X_i)}{\left(\sum_{i=1}^n k_i^{-1}\right)^2} \\ &= \frac{\sum_{i=1}^n \left(\frac{\beta}{k_i}\right)^2}{\left(\sum_{i=1}^n k_i^{-1}\right)^2} \\ &= \beta^2 \frac{\sum_{i=1}^n (k_i^{-1})^2}{\left(\sum_{i=1}^n k_i^{-1}\right)^2}. \end{aligned}$$

d)

Sannsynlighetsmaksimeringsestimator (SME):

Begynner med likelihood-funksjonen

$$L(X_1, X_2, \dots, X_n) = \prod_{i=1}^n f_X(X_i) = \prod_{i=1}^n \left(\frac{k_i}{\beta} \right) e^{-k_i X_i / \beta} = \frac{1}{\beta^n} e^{-\sum_{i=1}^n X_i k_i / \beta} \prod_{i=1}^n k_i.$$

Tar logaritmen til denne og deriverer mhp β .

$$\Lambda = \log(L) = -n \log(\beta) - \beta^{-1} \sum_{i=1}^n X_i k_i + \log\left(\prod_{i=1}^n k_i\right).$$

$$\frac{d\Lambda}{d\beta} = -\frac{n}{\beta} + \beta^{-2} \sum_{i=1}^n X_i k_i.$$

Setter til slutt den deriverte lik null for $\beta = \bar{\beta}$ og løser ut for $\bar{\beta}$.

$$\begin{aligned} -n\bar{\beta}^{-1} + \bar{\beta}^{-2} \sum_{i=1}^n X_i k_i &= 0 \\ \bar{\beta} &= \frac{1}{n} \sum_{i=1}^n X_i k_i = \hat{\beta}. \end{aligned}$$

Så sjekker vi påstanden at $\hat{\beta}$ har minst varians, ved å sette $r_i = 1/k_i$ som foreslått i oppgaveteksten.

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \frac{1}{n} \beta^2. \\ \text{Var}(\tilde{\beta}) &= \frac{\sum_{i=1}^n r_i^2}{\left(\sum_{i=1}^n r_i\right)^2} \beta^2. \end{aligned}$$

SME har minst varians hvis $\frac{1}{n} \leq \frac{\sum_{i=1}^n r_i^2}{\left(\sum_{i=1}^n r_i\right)^2}$, dvs.

$$\begin{aligned} \frac{1}{n} \left(\sum_{i=1}^n r_i \right)^2 &\leq \sum_{i=1}^n r_i^2 \\ \left(\frac{1}{n} \sum_{i=1}^n r_i \right)^2 &\leq \frac{1}{n} \sum_{i=1}^n r_i^2 \\ \frac{1}{n} \sum_{i=1}^n r_i^2 - \left(\frac{1}{n} \sum_{i=1}^n r_i \right)^2 &\geq 0. \end{aligned}$$

Det er oppgitt at nederste linje stemmer, dermed har SME alltid varians mindre enn eller lik variansen til den andre estimatoren.

e)

Setter $Y = 2k_i X_i / \beta$. Vi kan benytte variabelskifte til å sette inn i sannsynlighetstettheten til X_i , som er eksponensialfordelt med parameter β/k_i . La $X_i = w(y) = \frac{\beta}{2k_i} y$. Da er

$$\begin{aligned} f_Y(y) &= w'(y) f_X(w(y)) = \frac{\beta}{2k_i} f_X\left(\frac{\beta y}{2k_i}\right) \\ &= \frac{\beta}{2k_i} \frac{k_i}{\beta} e^{-\frac{\beta y}{2k_i} \frac{k_i}{\beta}} = \frac{1}{2} e^{-y/2}. \end{aligned}$$

Dette er en χ^2 -fordeling med 2 frihetsgrader.

Alternativt, sett inn i momentgenererende funksjon for eksponensialfordelingen:

$$\begin{aligned} M_X(t) &= \frac{1}{1 - \beta t / k_i}, \quad t < \frac{1}{\beta / k_i}. \\ M_Y(t) &= M_X\left(\frac{\beta}{2k_i} t\right) = \frac{1}{1 - 2t}, \quad t < \frac{1}{2}. \end{aligned}$$

Dette er momentgenererende funksjon for χ^2 -fordelingen med $\nu = 2$ frihetsgrader.

Ser deretter at $2n\hat{\beta}/\beta = \sum_{i=1}^n 2k_i X_i / \beta$ er en lineærkombinasjon av χ^2 -fordelte variable, som også er χ^2 -fordelt. Med to frihetsgrader fra hver X_i , er det totalt $\nu = 2n$ frihetsgrader.

Et $(1 - \alpha) \cdot 100\%$ konfidensintervall for $\chi^2 = 2n\hat{\beta}/\beta$ er $\chi_{1-\alpha/2, \nu}^2 < \chi^2 < \chi_{\alpha/2, \nu}^2$, altså

$$\begin{aligned} \chi_{1-\alpha/2, \nu}^2 &< 2n\hat{\beta}/\beta < \chi_{\alpha/2, \nu}^2 \\ \frac{\chi_{1-\alpha/2, \nu}^2}{2n\hat{\beta}} &< \frac{1}{\beta} < \frac{\chi_{\alpha/2, \nu}^2}{2n\hat{\beta}} \\ \frac{2n\hat{\beta}}{\chi_{\alpha/2, \nu}^2} &< \beta < \frac{2n\hat{\beta}}{\chi_{1-\alpha/2, \nu}^2}. \end{aligned}$$

For $n = 8$ og $\alpha = 0.05$, finner vi kvantilene i tabell: $\chi_{0.975, 16}^2 = 6.908$ og $\chi_{0.025, 16}^2 = 28.845$.

Setter inn dette og $\hat{\beta} = 8.3$ og får konfidensintervall

$$\begin{aligned} \frac{2 \cdot 8 \cdot 8.3}{28.845} &< \beta < \frac{2 \cdot 8 \cdot 8.3}{6.908} \\ 4.604 &< \beta < 19.224. \end{aligned}$$