

Tentamenskrivning: TMS145 - Grundkurs i matematisk statistik och bioinformatik, 7,5 hp.

Tid: Onsdag den 19 augusti 2009, kl 08:30-12:30 Väg och vatten

Examinator: Olle Nerman, tel 7723565.

Jour: Frank Eriksson, tel 073-7263051.

Hjälpmaterial: Egen handskriven formelsamling (fyra A4-sidor) samt med skrivningen utdelade tabellsidor.

Maxpoäng: 32. För godkänt krävs minst 15 poäng totalt och minst 4 poäng på sannolikhetsteori- och statistik-delen vardera samt minst 3 poäng på bioinformatikdelen. För betyget 4 krävs 20 poäng och för betyget 5 25 poäng.

Sannolikhetsteori

1. a Vad menas med medianen för en kontinuerligt fördelad stokastisk variabel?

(2p)

- b Vad är medianen för en exponentialfördelad stokastisk variabel med väntevärde 3?

(2p)

2. Avrundningsfelen vid en viss typ av vägningar kan antas vara likformigt fördelade i intervallet -0,005 och 0,005. Om verkliga vikterna varierar mycket mellan 100 olika enheter som vägs, så kan de 100 mätfelen dessutom betraktas som oberoende.

- a Beräkna sannolikheten vid kast av 3 vanliga tärningar (med poängsidorna 1,2,3,4,5 eller 6) för händelsen $A = \text{"summan av kasten blir exakt 5"}$?

(2p)

- b Vad är den betingade sannolikheten vid kast av 3 vanliga tärningar för händelsen $B = \text{"minst en tärning visar poängen 3"}$ givet händelsen $A = \text{"summan av kasten är exakt 5"}$?

(2p)

Ledning: Räkna möjliga och gynnsamma fall

Vänd!

3. Ett visst flygplan kan på grund av säkerhetsskäl, när det är fullsatt, ha en obalans mellan stolsraderna på vänster och högersidan som är maximalt 250KG. På varje sida finns 50 platser. En typisk passagerare har en vikt som har väntevärdet 80 Kg och standardavvikelsen 10 Kg. Antag oberoende mellan resenärernas vikt på resp. plats och att planet är fullsatt.

- a Vilket väntevärde har totala viktskillnaden (med tecken) mellan passagerarna till höger och vänster i ett fullsatt plan?

(1p)

- b Vilken varians har totala viktskillnaden mellan passagerarna till höger och vänster i ett fullsatt plan?

(1p)

- c Vad är approximativt sannolikheten för att obalansen överskrider gränsen 250 kg.

(2p)

Statistik

4. a Tag fram Maximum Likelihood-skattaren för θ ur ett stickprov från en Rayleigh-fördelning, dvs. en fördelning med täthetsfunktion

$$f(x) = \theta x e^{-\theta x^2/2}, \quad x \geq 0.$$

(2p)

- b Ange två önskvärda egenskaper hos en punktskattare.

(2p)

5. a Vilken är fördelen med att använda intervallskattning istället för punktskattning?

(1p)

- b Ge definitionen av *p-värde*.

(1p)

- c Anslag att du har ett stickprov av storlek n från en normalfördelning med väntevärde μ och varians $\sigma^2 = 4$. Du vill testa $H_0 = \mu = 10$ mot $H_1 : \mu < 10$ på signifikansnivå 0.01. Du använder den vanliga teststatistikan, d.v.s. $(\bar{X} - \mu)/(\sigma/\sqrt{n})$. Antag att det sanna väntevärdet är 9. Hur stor bör stickprosstörleken n vara för att teststyrkan (power) blir minst 90%?

(2p)

Vänd!

6. Vikt X (pounds) och systoliskt blodtryck Y (blodtryck vid hjärtats sammardragning) hos 26 slumpmässigt utvalda (och oberoende) män i åldrarna 25-30 visas i tabellen nedan.

Patientnr	Vikt	Blodtryck	Patientnr	Vikt	Blodtryck
1	165	130	14	172	153
2	167	133	15	159	128
3	180	150	16	168	132
4	155	128	17	174	149
5	212	151	18	183	158
6	175	146	19	215	150
7	190	150	20	195	163
8	210	140	21	180	156
9	200	148	22	143	124
10	149	125	23	240	170
11	158	133	24	235	165
12	169	135	25	192	160
13	170	150	26	187	159

Antag nu att vikt och blodtryck är bivariat normalfördelade. Från värdena ovan får vi att

$$\sum_i x_i = 4743 \quad \sum_i y_i = 3786$$

$$\sum_i x_i^2 = 880545 \quad \sum_i y_i^2 = 555802 \quad \sum_i x_i y_i = 697076$$

$$\sum_i (x_i - \bar{x})^2 \approx 15312 \quad \sum_i (y_i - \bar{y})^2 \approx 4502 \quad \sum_i (x_i - \bar{x})(y_i - \bar{y}) \approx 6422$$

Vänd!

En (något beskuren) utskrift från programmet R med datamännen ovan ger

```
lm(formula = y ~ x)
```

Coefficients:

	Estimate	Std. Error	t value
(Intercept)	69.10437	12.91013	5.353
x	0.41942	0.07015	5.979

```
Residual standard error: 8.681 on 24 degrees of freedom  
Multiple R-squared: 0.5983, Adjusted R-squared: 0.5815  
F-statistic: 35.74 on 1 and 24 DF, p-value: 3.591e-06  
---
```

Analysis of Variance Table

Response: y

	Df	Sum Sq	Mean Sq
x	1	2693.58	2693.58
Residuals	24	1808.57	75.36

- a Ange uttrycket för en linjär regressionsmodell med blodtryck som svarsvariabel baserad på informationen ovan (där alla koefficienter samt regressionslinjens slutliga form anges).

(1p)

- b Tolka koefficienterna (uttryckt i termer relaterade till problemet).

(1p)

- c Baserat på informationen ovan, finns det ett statistiskt säkerställt samband mellan vikt och blodtryck? Motivera!

(1p)

- d Hur stor andel av variationen i männens blodtryck kan enligt modellen förklaras av personernas varierande vikt?

(1p)

Vänd!

1.a) Medianen för en stokastisk variabel \bar{X}
är det värde $\hat{\mu}$ så att

$$P(\bar{X} \leq \hat{\mu}) = \frac{1}{2}$$

TMS 195
090819

b) $\bar{X} \sim \text{exp}(\beta)$. Finn $\hat{\mu}$.

$$f_{\bar{X}(x)} = \frac{1}{\beta} e^{-\frac{x}{\beta}}, x > 0$$

$$0.5 = \int_0^{\hat{\mu}} \frac{1}{\beta} e^{-\frac{x}{\beta}} dx = \left[\frac{1}{\beta} (-\beta) e^{-\frac{x}{\beta}} \right]_0^{\hat{\mu}}$$

$$= 1 - e^{-\hat{\mu}/\beta}$$

$$\frac{1}{2} = 1 - e^{-\hat{\mu}/\beta} \text{ ger } \hat{\mu} = 3 \log 2$$

2. Kart av tre tärningar

a) A "summan av kartan 5"

Gynnsamma utfall: 115 122 131 212 221 311

Totalt antal möjliga fall: $6 \cdot 6 \cdot 6 = 6^3$

$$P(A) = \frac{\# \text{gynnsamma}}{\# \text{möjliga}} = \frac{6}{6^3} = \frac{1}{6^2}$$

b) B "minst en tärning visar 3"

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{3/6^3}{6/6^3} = \frac{3}{6} = \frac{1}{2}$$

$$A \cap B = \{113, 131, 311\}$$

3) Låt $X_1, \dots, X_{50}, Y_1, \dots, Y_{50}$ vara värten
av personerna på höger respektive vänster sida

$$\text{a) } E\left[\sum_{i=1}^{50} X_i - \sum_{j=1}^{50} Y_j\right] = \sum_{i=1}^{50} E[X_i] - \sum_{j=1}^{50} E[Y_j]$$

$$= 50 \cdot E[X_i] - 50 \cdot E[Y_j]$$

$$= 50 \cdot 80 - 50 \cdot 80 = 0$$

$$\text{b) } \text{Var}\left(\sum_{i=1}^{50} X_i - \sum_{j=1}^{50} Y_j\right) = \sum_{i=1}^{50} \text{Var}(X_i) + \sum_{j=1}^{50} \text{Var}(Y_j)$$

Observer

$$= 50 \cdot 10^2 + 50 \cdot 10^2$$

$$= 10000$$

c) 50+50 relativt stort, vi kan använda
normalapproximationen genom centrala
gränsvärdesatsen.

Låt $\sum_{i=1}^{50} X_i - \sum_{j=1}^{50} Y_j \sim \text{approximert } N(0, 10000)$

$$P\left(|\sum_{i=1}^{50} X_i - \sum_{j=1}^{50} Y_j| > 250\right)$$

$$= P\left(\frac{|\sum_{i=1}^{50} X_i - \sum_{j=1}^{50} Y_j| - 0}{\sqrt{10000}} > \frac{250 - 0}{\sqrt{10000}}\right)$$

$$= P(Z > 2) \sim N(0, 1)$$

$$= P(|Z| > 2,5) = 2 \cdot \Phi(-2,5) = 2 \cdot 0,062$$

$$= 0,124$$

ha) x_1, \dots, x_n obs. från $\mathbb{X}_1, \dots, \mathbb{X}_n$ som
har fönstret

$$f(x) = \theta x e^{-\theta x^2/2}, \quad x \geq 0$$

Likolikhoodfunktion

$$\mathcal{L}(\theta) = \prod_{i=1}^n f(x_i) = \prod_{i=1}^n \theta x_i e^{-\theta x_i^2/2}$$

$$\begin{aligned} \frac{\partial}{\partial \theta} \log \mathcal{L}(\theta) &= \frac{\partial}{\partial \theta} \left(n \cdot \log \theta - \frac{\theta}{2} \sum_{i=1}^n x_i^2 \right) \\ &= \frac{n}{\theta} - \frac{1}{2} \sum_{i=1}^n x_i^2 \end{aligned}$$

$$\frac{\partial}{\partial \theta} \log \mathcal{L}(\theta) \Big|_{\theta=\hat{\theta}} = \frac{n}{\hat{\theta}} - \frac{\sum_{i=1}^n x_i^2}{2} = 0$$

$$\text{gör } \hat{\theta} = \frac{2n}{\sum_{i=1}^n x_i^2}$$

$$\text{Kolla max! } \hat{\theta}^* = \frac{2n}{\sum_{i=1}^n x_i^2} \text{ MLE}$$

b) Önskvärda egenskaper: Väntevärdesnittig
Lag värde

5a Punktshatningen säger ingenting om
säkerheten i vår shatning

b Sannolikheten att observera ett lika eller
mer extremt värde än vi gjort givet
att nullhypotesen är sann. ~~Ett lägt~~
~~p-värde tyder på att innebär att antagande~~
~~om nullhypotesen är osannolikt.~~

$\Sigma \dots x_n$ observationer av
 $\bar{x}_1, \dots \bar{x}_n$ som är $N(\mu, \sigma^2 = n)$ -fördelade.
Vi testar $H_0: \mu = 10$ mot $H_1: \mu < 10$
på signifikansnivå $\alpha = 0.01$. Det sanna
värdet är 9. Hur stort bör n vara
för att styrkan $1 - \beta(\gamma) \geq 0.9$, dvs $\beta(\gamma) \leq 0.1$.

$$\beta(\gamma) = 1 - \Phi\left(-z_{\alpha} + \frac{\mu_0 - \mu}{\sigma/\sqrt{n}}\right) \text{ sid } 298$$

i vrt fall

$$0.1 > \beta(\gamma) = 1 - \Phi\left(-2.33 + \frac{10 - 9}{\sigma/\sqrt{n}}\right)$$

$$\text{gör, } \Phi\left(-2.33 + \frac{\sqrt{n}}{2}\right) > 0.9$$

$$-2.33 + \frac{\sqrt{n}}{2} > \Phi^{-1}(0.9) = 1.28$$

$$n > (2 \cdot (1.28 + 2.33))^2$$

Så n bör vara minst 53. $= 52,12$

a) $E[Y] = \beta_0 + \beta_1 X$

$\hat{\beta}_0 = 69,1$, $\hat{\beta}_1 = 0,02$ sättningar av β_0, β_1

b) En röntgening på en pound förväntas ha blodtrycket β_1 . Tolkningen är endast giltig för män som väger mellan ca 150-250 kg

c) Testar $H_0: \beta_1 = 0$
med
 $H_a: \beta_1 \neq 0$

her ett t-värde på $5,979 >> 1,7$ som ger ett p-värde på $0,05$ (2u frihetegrader). Vi kan förhålla ~~H₀~~ H₀ med stor säkerhet.

d) $R^2 = 0,59$, 59% av variationen förklaras av modellen.