

STOCKHOLMS UNIVERSITET
Statistiska institutionen
Jessica Franzén

TENTAMEN I STATISTISK TEORI MED TILLÄMPNINGAR II
2018-11-01

Skrivtid: 10.00-15.00

Godkända hjälpmedel: Miniräknare, språklexikon.

Tentamen består av fem uppgifter. För full poäng på en uppgift krävs tydliga, utförliga och väl motiverade lösningar.

OBS! Glöm inte att ange nödvändiga antaganden där det behövs.

Uppgift 1. (20 poäng)

Förklara följande begrepp:

- a) Konsistens
- b) Apriori- och aposteriorifördelning
- c) Medelkvadratfel (MSE)
- d) Neyman-Pearsons lemma
- e) Runs test

Uppgift 2. (20 poäng)

Innan valet gjordes en telefonintervju-undersökning där ett slumpmässigt urval av 1000 väljare fick svara på frågan "Hur stort förtroende har du för Kristdemokraternas partiledare Ebba Busch Thor?" 189 stycken av dessa svarade att de har stort eller mycket stort förtroende för partiledaren.

- a) Bestäm ett 99 %-igt konfidensintervall för andelen väljare som har stort eller mycket stort förtroende för partiledaren.
- b) Genomför ett hypotestest för att avgöra om undersökningen ger stöd för att påstå att partiledaren har stort eller mycket stort förtroende hos mindre än 20 % av väljarna. Använd signifikansnivån 5 %.
- c) Hur stort urval skulle krävas för testet i b) för att sannolikheten för typ II-fel ska vara högst 0.03 om den sanna andelen är 0.18.

Uppgift 3. (20 poäng)

En skotillverkare har gjort ett experiment för att jämföra slitstyrkan hos två olika material till klacken på en herrsko. Ett slumpmässigt urval av 7 män har fått ett par nya skor där ena klacken är tillverkad av material A och den andra av material B. Efter tre månaders daglig användning mättes klackhöjden på båda skorna med följande resultat (i mm):

Material A	6.6	5.5	8.3	8.2	5.2	9.1	6.3
Material B	7.4	7.0	8.8	8.0	6.8	9.3	7.9

- Använd ett lämpligt icke-parametriskt test för att pröva om det finns någon skillnad mellan slitstyrkan hos de två materialen.
- Använd ett lämpligt parametriskt test för att pröva om det finns någon skillnad mellan slitstyrkan hos de två materialen.
- Beskriv eventuella skillnader och likheter mellan de båda testmetoderna när det gäller testens förutsättningar och resultat.

Uppgift 4. (20 poäng)

Y är en stokastisk variabel för en persons inkomst. Enligt den s.k. Paretolagen gäller att $P(Y \geq y) = \left(\frac{k}{y}\right)^\theta$ där k är den lägsta inkomsten i hela populationen. Då är $F(y) = 1 - \left(\frac{k}{y}\right)^\theta$ och vidare följer att

$$f(y) = \theta k^\theta \left(\frac{1}{y}\right)^{\theta+1} \quad y \geq k \quad \theta \geq 1$$

- Visa med utgångspunkt från $F(y)$ att $f(y) = \theta k^\theta \left(\frac{1}{y}\right)^{\theta+1}$.
- Antag att k är känd. Härled maximum-likelihoodskattningen av θ .
- Antag att den lägsta inkomsten k i en populationen är 15 000 kronor/månad. Ett slumpmässigt urval av 5 personers inkomster från denna population ges nedan. Beräkna med hjälp av maximum-likelihoodskattningen sannolikheten att en slumpmässigt vald person från populationen har en inkomst som överstiger 50 000 kronor.

Inkomster (kronor/månad) från 5 slumpmässigt valda personer:
26 500, 45 100, 58 000, 22 300, 19 500

Uppgift 5. (20 poäng)

Antag att ett slumpmässigt urval av en enda observation tas från en population där täthetsfunktionen är

$$f(y) = \left(\frac{1}{\lambda}\right) e^{-y/\lambda} \quad y > 0$$

med syftet att testa

$$H_0 : \lambda = 1$$

$$H_A : \lambda > 1$$

Nollhypotesen förkastas om observationen $y \geq 3.20$.

- a) Beräkna sannolikheten för ett typ I-fel d.v.s. α .
- b) Beräkna sannolikheten för ett typ II-fel d.v.s. β när $\lambda = \frac{4}{3}$
- c) Rita en bild över fördelningarna under H_0 och H_A ($\lambda = \frac{4}{3}$) som visar hur α och β uträknade i a) och b) förhåller sig till varandra både vad gäller läge och storlek.
- d) Härled momentskattningen av λ för ett slumpmässigt urval av n observationer.

Statistiska institutionen



Rättningsblad

Datum: 1/11-2018

Sal: Brunnsvikssalen

Tenta: Statistisk teori med tillämpningar 2

Kurs: Statistisk teori med tillämpningar

ANONYMKOD:

0039-LGB

Jag godkänner att min tenta får läggas ut anonymt på hemsidan som studentsvar.

OBS! SKRIV ÄVEN PÅ BAKSIDAN AV SKRIVBLADEN

Markera besvarade uppgifter med kryss

1	2	3	4	5	6	7	8	9	Antal inl. blad
X	X	X	X	X					7
Lär.ant. 20	12	19	20	18					

POÄNG	BETYG	Lärarens sign.
89+6=95	A	JF

Uppgift 1.

a) En estimator $\hat{\theta}$ av en populationsparameter θ är konsistent om det gäller att $\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| < \varepsilon) = 1$ för alla reella $\varepsilon > 0$

Om $\hat{\theta}$ är väntevärdesriktig gäller att $\hat{\theta}$ är konsistent om $\lim_{n \rightarrow \infty} V(\hat{\theta}) = 0$. Om n betecknar n antalet observationer som estimatorn $\hat{\theta}$ beräknas utifrån.

b) Om en okänd populationsparameter θ enligt ett Bayesianskt synsätt betraktas som en stokastisk variabel så är apriorifördelningen $g(\theta)$ täthetsfunktionen (eller i diskreta fall frekvensfunktionen) för denna. Apriorifördelningens utseende representerar vår (fram tills nu) erhållna information om θ samt eventuella antaganden vi gör om θ .

Efter erhållna observationer ur populationen har ny information erhållits som leder till en revidering av täthetsfunktionen för θ och vi får en s.k. aposteriorifördelning g^* som tar hänsyn både till tidigare information/antaganden samt till de nya data som observationerna utgör.

c) Medelkvadratkylet för en estimator $\hat{\theta}$ av en populationsparameter θ definieras som $MSE(\hat{\theta}) = E((\hat{\theta} - \theta)^2)$, d.v.s som väntevärdet av kvadraten på skillnaden mellan parameter och estimator. Om $\hat{\theta}$ är väntevärdesriktig är medelkvadratkylet samma som variansen hos $\hat{\theta}$ ($\hat{\theta}$ har då ingen bias).

4

d) Enligt Neyman-Pearsons lemma är det starkaste test vi kan använda för att testa en nollhypotes $H_0: \theta = \theta_0$ gentemot en alternativhypotes $H_a: \theta = \theta_1$ på formen
$$\frac{L(y_1, y_2, \dots, y_n, \theta = \theta_0)}{L(y_1, y_2, \dots, y_n, \theta = \theta_1)} < k$$
 där L är likelihood-

funktionen och k är någon positiv konstant.

Om villkoret ovan uppfylls förkastas H_0 .

4

e) För att undersöka om en sekvens av identiskt Bernoulli fördelade observationer är slumpmässig eller resultatet av någon påverkan som gruppering eller motverka gruppering kan ett Runstest användas.

Om exempelvis 3 av 10 observationer är ettor och resten nollor, är det sannoliktare att det bildas 5 sekvenser

än att det bildas 2: $\underline{0010001100}$ (5 sekvenser)

$\underline{0000000111}$ (2 sekvenser)

Runstestet räknar antalet "runs" (sekvenser med identiska värden) och undersöker sannolikheten för observerat

antal runs jämfört uppskattat \hat{p} (andel ettor t.ex.)

4

samt $\hat{p} \cdot n > 5$
och $(1 - \hat{p}) \cdot n > 5$

Uppgift 2

a) Eftersom det här rör sig om så pass många observationer ($n = 1000$), så kan vi anta att $\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$ (där $y_i \in \{0, 1\}$ för $1 \leq i \leq 1000$)

är approximativt normalfördelad enligt centrala gränsvärdesatsen.

För variansen av \hat{p} gäller:

$$V(\hat{p}) = V\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(y_i) = \frac{1}{n^2} \cdot n V(y) = \frac{1}{n} V(y) = \frac{p(1-p)}{n}$$

Vi uppskattar $V(\hat{p})$ genom att använda vår skattning \hat{p} av p :

$$\widehat{V(\hat{p})} = \frac{\hat{p}(1-\hat{p})}{n}$$

Vårt 99% konfidensintervall kommer att vara $\hat{p} \pm Z_{0,005} \cdot \sqrt{\widehat{V(\hat{p})}}$
Med insatta värden ($\hat{p} = 0,189$, $n = 1000$, $Z_{0,005} = 2,5758$) får vi:

$$0,189 \pm 2,5758 \cdot \sqrt{\frac{0,189 \cdot 0,811}{1000}}$$

$$0,189 \pm 2,5758 \cdot 0,01238059$$

$$0,189 \pm 0,031889923$$

Detta kan uttryckas ($0,1571; 0,2209$)

Svar: Det sökta 99%-iga konfidensintervallet är $(0,1571; 0,2209)$.

ⓐ

Uppgift 2; forts.

b) För att undersöka om partiledaren har stort eller mycket stort förtroende hos mindre än 20% av röstarna ställer vi upp följande hypoteser:

$$H_0: p = 0,20 \quad \checkmark \quad (\text{här } p \text{ är andelen röstare som har stort/m mycket stort förtroende})$$

$$H_a: p < 0,20$$

Vi antar att samtliga observationer är oberoende av varandra.

Under förutsättning att H_0 är sann gäller:

$$\sigma_{\hat{p}}^2 = V(\hat{p}) = V\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(y_i) = \{y_i \text{ Bernoulli-fördelade med } p=0,20\} =$$

$$= \frac{1}{n^2} \cdot n(0,2 + (1-0,2)) = \{n=1000\} = 0,00016$$

$$\text{Inför teststatistikan } Z = \frac{\hat{p} - 0,20}{\sqrt{V(\hat{p})}} = \frac{\hat{p} - 0,20}{\sigma_{\hat{p}}}$$

eftersom antalet observationer är så pass stort, och $\hat{p} \cdot n = 189 > 5$ samt $(1-\hat{p}) \cdot n = 811 > 5$ kan vi anta att statistikan Z är standardiserat normalfördelad enligt centrala gränsvärdesatsen om H_0 är sann.

För att testa H_0 gentemot H_a på signifikansnivån 5% använder vi den följande beslutsregel:

• H_0 förkastas om $Z < -\lambda_{0,05} = \{\text{ur tabell}\} = -1,6449$, annars

vänt förkastelseområde RR är således $\{Z < -1,6449\}$

Låt oss nu beräkna värdet av vår teststatistika:

$$Z_{\text{obs}} = \frac{0,189 - 0,20}{\sqrt{0,00016}} = \frac{-0,011}{0,012649111} \approx -0,8696 \neq -\lambda_{0,05} \quad \checkmark$$

Enligt vår beslutsregel kan vi inte förkasta H_0 .

Svar: Vi har, på 5% signifikansnivå, inte stöd för att påstå att partiledaren har stort eller mycket stort förtroende hos mindre än 20% av röstarna.

Uppgift 2; forts.

⊖) Typ II-fel innebär att vi inte förkastar H_0 trots att H_0 är falsk.

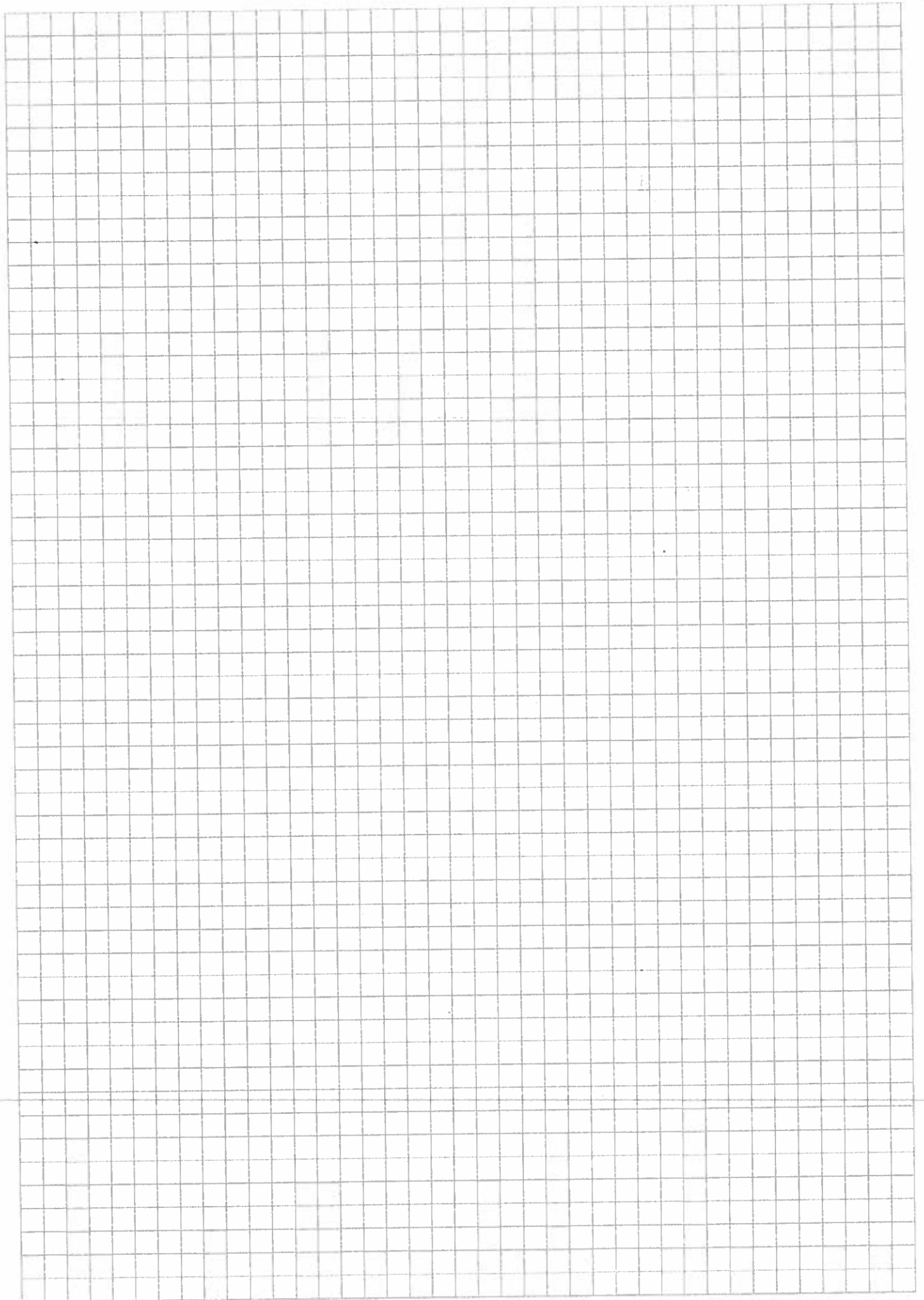
I detta fall, om $p = 0,18$, gäller att

$$\begin{aligned} V(\hat{p}) &= V\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(y_i) = \frac{1}{n^2} \cdot n \cdot p(1-p) = \frac{p(1-p)}{n} \\ &= \{\text{med insatta värden}\} = \frac{0,1476}{n} \end{aligned}$$

Ett motsvarande test som i b) skulle ha teststatistikan

$$Z = \frac{\hat{p} - 0,20}{\sqrt{\quad}}$$





Uppgift 3

a) Låt oss utföra Wilcoxon's teckenrangtest. Vi rankar data:

	Material A	Material B	Differens:	Rang:	Tecken:
Man 1	6,6	7,4	-0,8	4	-
Man 2	5,5	7,0	-1,5	5	-
Man 3	8,3	8,8	-0,5	3	-
Man 4	8,2	8,0	0,2	1,5	+
Man 5	5,2	6,8	-1,6	6,5	-
Man 6	9,1	9,3	-0,2	1,5	-
Man 7	6,3	7,9	-1,6	6,5	-

Ur rankingarna får att $T^- = 4 + 5 + 3 + 6,5 + 1,5 + 6,5 = 26,5$
och att $T^+ = 1,5$

Nollhypotes H_0 : Det finns ingen skillnad i slitstyrkan mellan de två materialen.

Alternativhypotes H_a : Det finns en skillnad i slitstyrkan mellan de två materialen.

Antaganden: Samtliga observerade differenser är oberoende av varandra. Beröranden förekommer endast mellan data i respektive observationspar.

P.g.a. hypotesens natur utförs ett t-äsidigt test. Vår teststatistiska blir då $T = \min(T^+, T^-)$.

Beslutregel: H_0 förkastas om $T \leq 2$

Vårt förkastelseområde är alltså $RR: \{T \leq 2\}$

Enligt tabell i formelsamlingen uppnås en signifikansnivå på $\alpha = 0,05$ på detta sätt.

v.g.v \Rightarrow

Uppgift 3; forts.

a) forts.

Vi observerar att $T = \min(T; T^*) = \min(26,5; 1,5) = 1,5$

Dämed gäller $T \leq 2$ och H_0 förkastas, enligt vår beslutsregel.

Slutsats: Wilcoxon's teckenrangtest med 5% signifikansnivå visar att slitstyrkan hos de två materialen skiljer sig åt.

b) För vårt parametriska test antar vi att "klackhöjdsdifferensen" är normalfördelad och att de n st. observationerna är oberoende av varandra.

Differens i klackhöjd för ett observerat par skor är en stokastisk variabel $Y \sim N(\mu, \sigma)$.

$$H_0: \mu = 0$$

$$H_a: \mu \neq 0$$

Teststatistika: $T = \frac{\bar{y} - 0}{\left(\frac{s}{\sqrt{n}}\right)}$ där $S^2 = \frac{\sum_{i=1}^n y_i^2 - n(\bar{y})^2}{n-1}$ är urvals- variansten.

Beslutsregel: H_0 förkastas om $|T| > t_{0,025}(n-1)$.

Detta ger ett test med 5% signifikansnivå.

$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$: beräknas ur tabellen på följande sätt till

$$\bar{y} = \frac{1}{7} (-0,8 - 1,5 - 0,5 + 0,2 - 1,6 - 0,2 - 1,6) = -\frac{5}{7} \approx -0,85714286$$

$$\sum_{i=1}^7 y_i^2 = 0,8^2 + 1,5^2 + 0,5^2 + 0,2^2 + 1,6^2 + 0,2^2 + 1,6^2 = 8,34$$

$$S^2 = \frac{\sum_{i=1}^7 y_i^2 - 7(\bar{y})^2}{7-1} = \frac{8,34 - 7 \cdot (-0,85714286)^2}{6} \approx 0,532857$$

$$S = \sqrt{S^2} \approx 0,7299706 \quad (= \hat{\sigma})$$

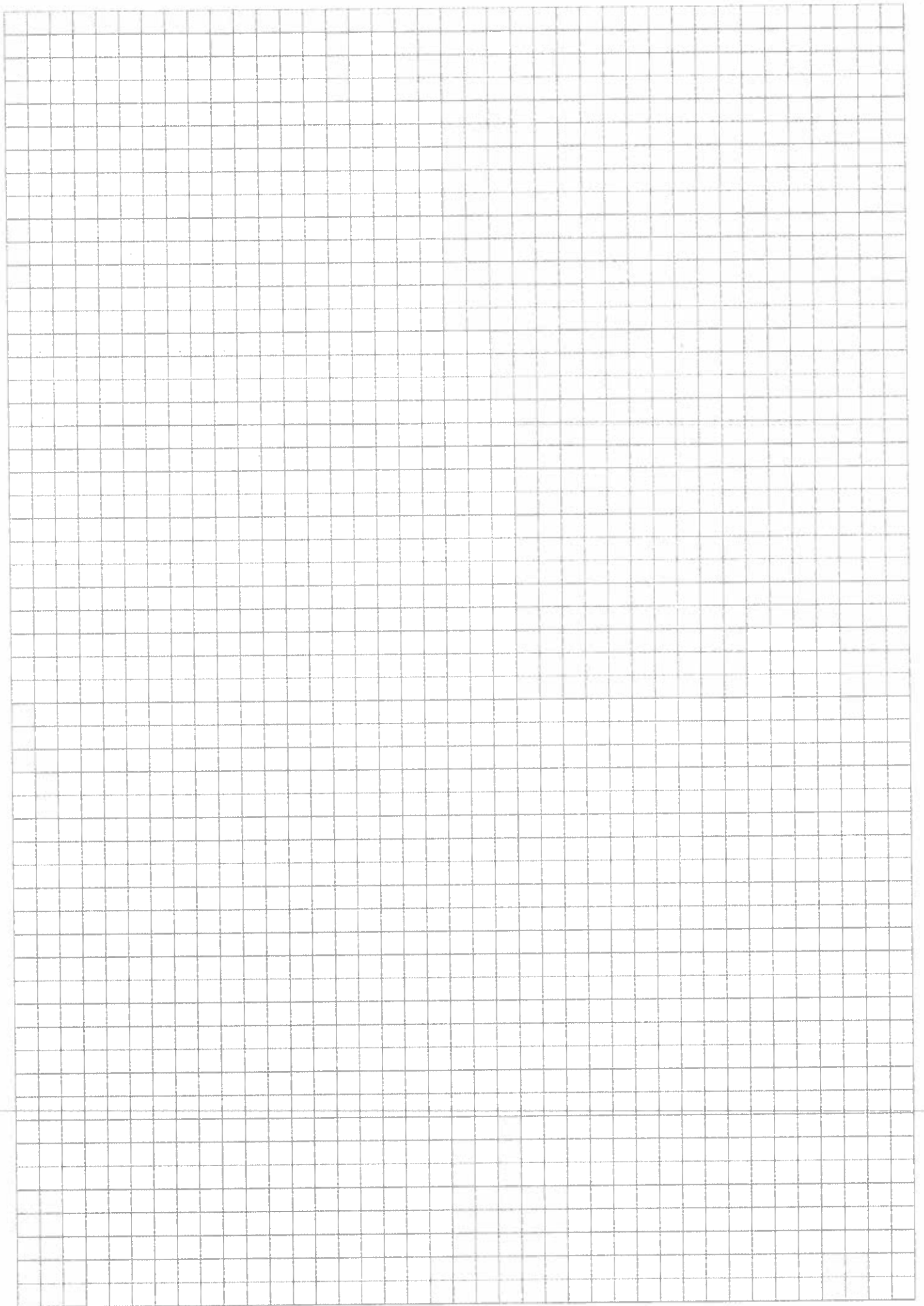
Vår teststatistika värde observeras således till $T_{obs} = \frac{-0,85714286 - 0}{\left(\frac{0,7299706}{\sqrt{7}}\right)} \approx -3,10668 \Rightarrow |T_{obs}| = 3,10668 > t_{0,025}(7-1) \approx 2,45$

Enligt beslutsregeln förkastas H_0 , och vi drar slutsatsen att det finns en skillnad i slitstyrkan mellan de två materialen.

Uppgift 3: forts.

- c) • Teckenrangtestet förutsätter inte att differensen i
klockhöjd är normalfördelade, vilket det använda
parametriska testet gör.
- Upprättat p-värde i vårt parametriska test tycks vara
nästan 0,01 ($\chi_{0,01} = 3,14$ enligt tabell för t-fördelningen)
medan upprättat p-värde i teckenrangtestet enligt tabell
ligger mellan 0,02 och 0,05. Det extra antagandet
om normalfördelning leder som väntat till lägre
p-värde (den nästa signifikansnivån på vilken vi
kunnat förkasta H_0).

(3)



Uppgift 4.

- a) För täthetsfunktionen $f(y)$ och fördelningsfunktionen $F(y)$ gäller sambandet $f(y) = F'(y)$ för alla y där F är deriverbar. Således:

$$f(y) = F'(y) = \frac{d}{dy} \left(1 - \left(\frac{k}{y} \right)^\theta \right) = \frac{d}{dy} (1 - k^\theta y^{-\theta}) = -k^\theta \cdot (-\theta) y^{-\theta-1}$$

$$= k^\theta \theta \cdot \left(\frac{1}{y} \right)^{\theta+1} = \theta k^\theta \left(\frac{1}{y} \right)^{\theta+1} \quad \text{v.s.v.} \quad (4)$$

- b) Vi antar att den lägsta inkomsten i populationen, k , är känd. Likelihoodfunktionen givet n st observationer y_1, y_2, \dots, y_n är:

$$L(y_1, y_2, \dots, y_n, \theta) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \left(\theta k^\theta \left(\frac{1}{y_i} \right)^{\theta+1} \right) = \theta^n k^{n\theta} \left(\prod_{i=1}^n y_i \right)^{-\theta-1}$$

Maximum-likelihoodskattningen av θ får vi genom att lokalisera det värde $\hat{\theta}_{ML}$ på θ som maximerar $L(y_1, y_2, \dots, y_n, \theta)$.

Da log-likelihoodfunktionen $l(y_1, y_2, \dots, y_n, \theta)$ har maximum på samma ställe som $L(y_1, y_2, \dots, y_n, \theta)$ beräknar vi denna och sätter

partiella derivatan n.a.p. θ till 0 (vilket ju gäller i maximum (i denna kontext!)):

$$L(y_1, y_2, \dots, y_n, \theta) = \ln(L(y_1, y_2, \dots, y_n, \theta)) = n \cdot \ln \theta + n\theta \cdot \ln k + (-\theta-1) \sum_{i=1}^n \ln y_i$$

$$\frac{dL(y_1, y_2, \dots, y_n, \theta)}{d\theta} = \frac{n}{\theta} + n \ln k - \sum_{i=1}^n \ln y_i$$

$$\frac{dL(y_1, y_2, \dots, y_n, \theta = \hat{\theta}_{ML})}{d\theta} = 0 \Rightarrow \frac{n}{\hat{\theta}_{ML}} + n \ln k - \sum_{i=1}^n \ln y_i = 0 \Rightarrow$$

$$\frac{n}{\hat{\theta}_{ML}} = \sum_{i=1}^n \ln y_i - n \ln k \Rightarrow \hat{\theta}_{ML} = \frac{n}{\sum_{i=1}^n \ln y_i - n \ln k}$$

Svar: $\hat{\theta}_{ML} = \frac{n}{\sum_{i=1}^n \ln y_i - n \ln k}$

(12)

Uppgift 4, forts.

c) Inkomster (kr/månad) för $n=5$ slumpmässigt valda personer:

$$y_1 = 26500$$

$$y_2 = 45100$$

$$y_3 = 58000$$

$$y_4 = 22300$$

$$y_5 = 19500$$

Lägst inkomsten i populationen: $k = 15000$ (kr/månad)

Enligt resultat i b) ges $\hat{\theta}_{ML}$ av:

$$\hat{\theta}_{ML} = \frac{n}{\sum_{i=1}^n \ln y_i - n \ln k} =$$

$$= 5 \left(\ln(26500) + \ln(45100) + \ln(58000) + \ln(22300) + \ln(19500) - 5 \cdot \ln(15000) \right)^{-1} \approx$$

$$\approx 1,358245$$

Sannolikheten för att en slumpmässigt vald person från populationen har en inkomst som överstiger 50000 kr/månad kan nu uppskattas med hjälp av föregående ML-skattning av parametern θ enligt:

$$P(Y > 50000) = \left(\frac{k}{50000} \right)^{\hat{\theta}_{ML}} \approx \left\{ \text{efter insättning av aktuella värden} \right\} \approx$$
$$\approx \left(\frac{15000}{50000} \right)^{1,358245} \approx 0,195$$

Svar: Sannolikheten att en slumpmässigt vald person från populationen har en inkomst som överstiger 50000 kr är ungefär 0,195.

4

Uppgift 5

Populationens täthetsfunktion är

$$f(y) = \left(\frac{1}{\lambda}\right) e^{-y/\lambda} \quad y > 0$$

Om vi gör en enda observation y_{obs} , så förkastar $H_0: \lambda = 1$ om $y_{obs} \geq 3,20$.

a) Att göra ett typ I-fel innebär att förkasta H_0 givet att den är sann.

Anta att H_0 är sann. Då gäller att täthetsfunktion för populationen

$$\text{är: } f_{\lambda=1}(y) = \left(\frac{1}{1}\right) e^{-y/1} = e^{-y} \quad y > 0 \quad \text{eftersom } \lambda = 1 \text{ enligt } H_0.$$

Sannolikheten att vår enda observation y_{obs} uppfyller $y_{obs} \geq 3,20$

$$\text{blir då } \alpha = \int_{3,20}^{\infty} f_{\lambda=1}(y) dy = \int_{3,20}^{\infty} e^{-y} dy = \left[-e^{-y}\right]_{3,20}^{\infty} =$$

$$= 0 - (-e^{-3,20}) = e^{-3,20} \approx 0,0407622$$

Eftersom $y_{obs} \geq 3,20$ just är vår beslutsregel för att förkasta H_0 , så är sannolikheten för ett typ I-fel $\alpha \approx 0,0408$

Svar: $\alpha \approx 0,0408$.

b) När $\lambda = \frac{4}{3}$ gäller att populationens täthetsfunktion är

$$f_{\lambda=\frac{4}{3}}(y) = \frac{1}{\left(\frac{4}{3}\right)} e^{-y/\left(\frac{4}{3}\right)} = \frac{3}{4} e^{-3y/4} \quad y > 0$$

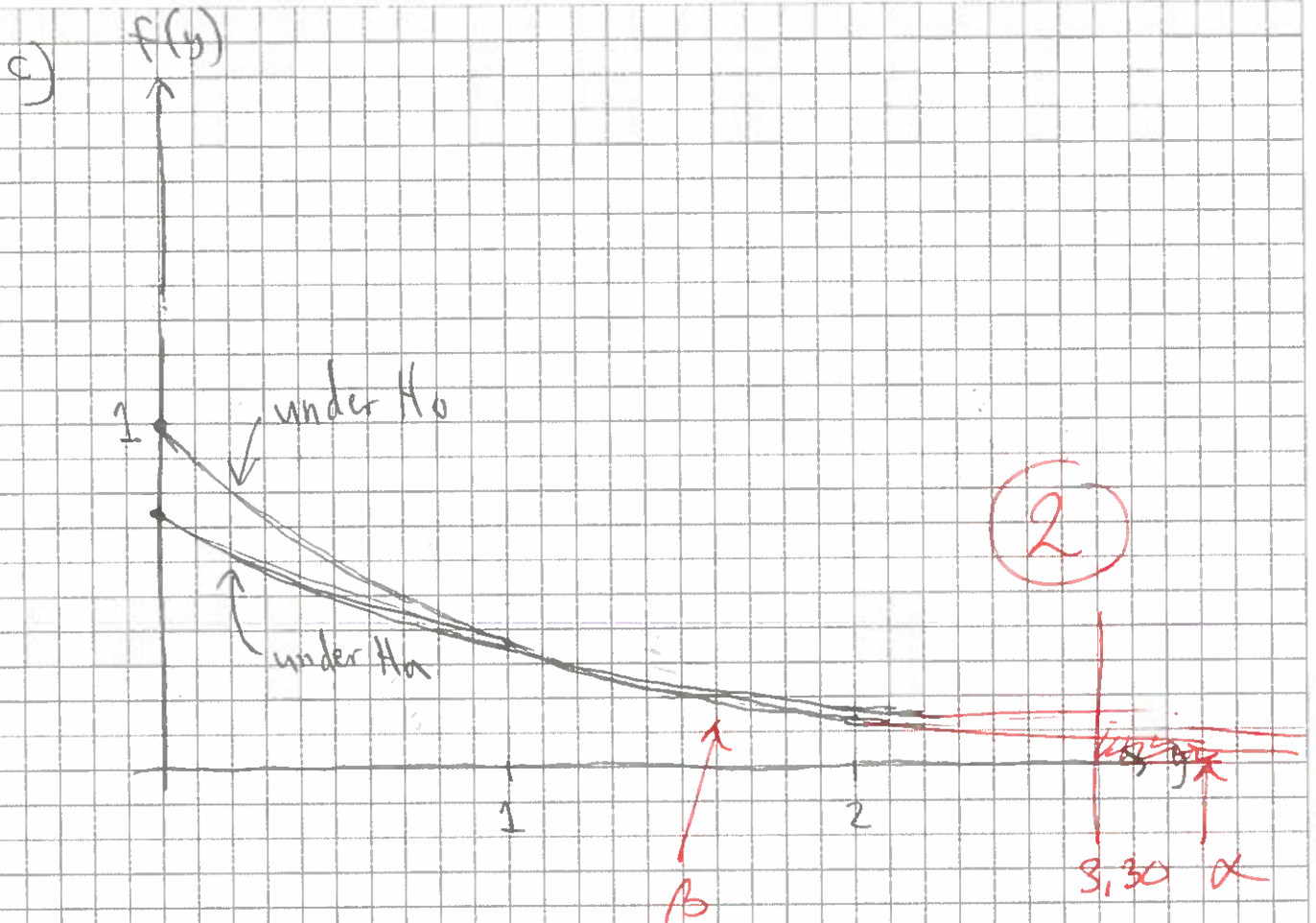
Vi kommer nu, då H_0 är falsk, att göra ett typ II-fel om vi inte förkastar H_0 . Denna situation inträffar när $y_{obs} < 3,20$.

Sannolikheten för detta, givet $\lambda = \frac{4}{3}$, är

$$\beta = \int_0^{3,20} f_{\lambda=\frac{4}{3}}(y) dy = \int_0^{3,20} \left(\frac{3}{4} e^{-3y/4}\right) dy = \left[-e^{-3y/4}\right]_0^{3,20} =$$

$$= (-e^{-3 \cdot 3,20/4}) - (-e^0) = 1 - e^{-2,4} \approx 0,909282$$

Svar: När $\lambda = 4/3$ gäller att $\beta \approx 0,909$



d) Vi ser att täthetsfunktionen $f(y)$ svarar mot en exponentialfördelning. Vi vet att för en sådan gäller $E(y) = \lambda$ (kan visas genom att t.ex. beräkna $E(y) = \int_0^{\infty} y \cdot \frac{1}{x} e^{-y/x} dy$ med partiell integrering; beräkningar utelämnas).

Anta att vi har n observationer y_1, y_2, \dots, y_n .

Momentestimatorn av λ fås genom att sätta $m'_1 = \frac{1}{n} \sum_{i=1}^n y_i$ till lika $\mu'_1 = E(y | \lambda = \hat{\lambda}_{\text{mom}})$, d.v.s.:

$$\frac{1}{n} \sum_{i=1}^n y_i = \hat{\lambda}_{\text{mom}}$$

Svar:
$$\hat{\lambda}_{\text{mom}} = \frac{1}{n} \sum_{i=1}^n y_i$$