

STOCKHOLMS UNIVERSITET
Statistiska institutionen
Ellinor Fackle-Fornius

TENTAMEN I STATISTISK TEORI MED TILLÄMPNINGAR II
2017-04-18

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.

Genomgång av tentamen sker 2017-05-09 kl. 15 i sal D207.

OBS! Glöm inte att ange nödvändiga antaganden överallt.

Uppgift 1. (20 poäng)

Avgör om vart och ett av följande påståenden är sanna eller falska samt motivera ditt svar.

- i)* Om p-värdet vid ett hypotestest är 0.028 kan nollhypotesen förkastas på signifikansnivån $\alpha = 0.01$.
- ii)* Styrkan hos ett hypotestest beräknas alltid givet att nollhypotesen är sann.
- iii)* Signifikansnivån α är sannolikheten att nollhypotesen är falsk.
- iv)* Om p-värdet är väldigt litet vid ett hypotestest för att jämföra två populationsmedelvärden måste skillnaden mellan medelvärdena vara stor.
- v)* Vid ett likelihoodkvotest är det möjligt att $L(\hat{\Omega}_0) > L(\hat{\Omega})$.

Uppgift 2. (20 poäng)

En stor livsmedelskedja är intresserad av andelen kunder som är positivt inställda till att använda självscanning när de handlar. I ett slumpmässigt urval av 120 kunder svarade 27 st att de var positivt inställda till att använda självscanning. Av de 120 kunderna i urvalet var 84 st medlemmar i kedjans kundklubb, varav 20 st var positivt inställda. Av de resterande 36 kunderna i urvalet (som alltså var icke-medlemmar) var 7 st positivt inställda.

- Beräkna ett 95 %-igt konfidensintervall för andelen kunder som är positivt inställda till att använda självscanning när de handlar.
- Hur stort urval skulle krävas om man vill att längden av ett 95 %-igt konfidensintervall för andelen kunder som är positivt inställda till att använda självscanning när de handlar inte ska överstiga 0.1.
- Beräkna ett 95 %-igt konfidensintervall för skillnaden mellan andelen kundklubsmedlemmar och andelen icke-medlemmar som är positivt inställda till att använda självscanning när de handlar.

Uppgift 3. (20 poäng)

I en fabrik utvärderas två typer av maskiner, A och B, genom ett experiment. Tiden det tar att utföra en viss uppgift kontrolleras för varje typ av maskin. Åtta maskintekniker är involverade i experimentet, varje tekniker använder maskin A och maskin B i slumpmässig ordning. Tiden i sekunder som det tog att slutföra uppgiften i respektive maskin visas i följande tabell.

Tekniker	Maskin A	Maskin B
1	32	30
2	40	39
3	42	42
4	26	23
5	35	36
6	29	27
7	45	41
8	22	21

- Testa med ett lämpligt icke-parametriskt test om det går att påvisa någon skillnad mellan tiden det tar för maskin A och maskin B att utföra uppgiften.
- Testa med ett lämpligt parametriskt test om det går att påvisa någon skillnad mellan tiden det tar för maskin A och maskin B att utföra uppgiften.
- I det här experimentet användes samma maskintekniker för båda maskiner, är det en fördel eller nackdel?

Uppgift 4. (20 poäng)

Antal gånger som en student på en grundkurs i statistik behöver tentera för att få ett godkänt betyg på kursen antas följa den geometriska fördelningen

$$p(y) = p(1-p)^{y-1}, \quad y = 1, 2, \dots$$

där y är antal försök och p är sannolikheten att få ett godkänt betyg. I ett slumpmässigt urval av 8 studenter blev antalet gånger som var och en behövde tentera:

1, 1, 1, 2, 3, 1, 1, 1

- Härled momentskattningen \hat{p}_{MOM} och beräkna skattningen för urvalet.
- Härled maximumlikelihood-skattningen \hat{p}_{ML} och beräkna skattningen för urvalet.
- Vad blir maximumlikelihood-skattningen för väntevärdet för antal gånger som man behöver tentera för att få ett godkänt betyg.

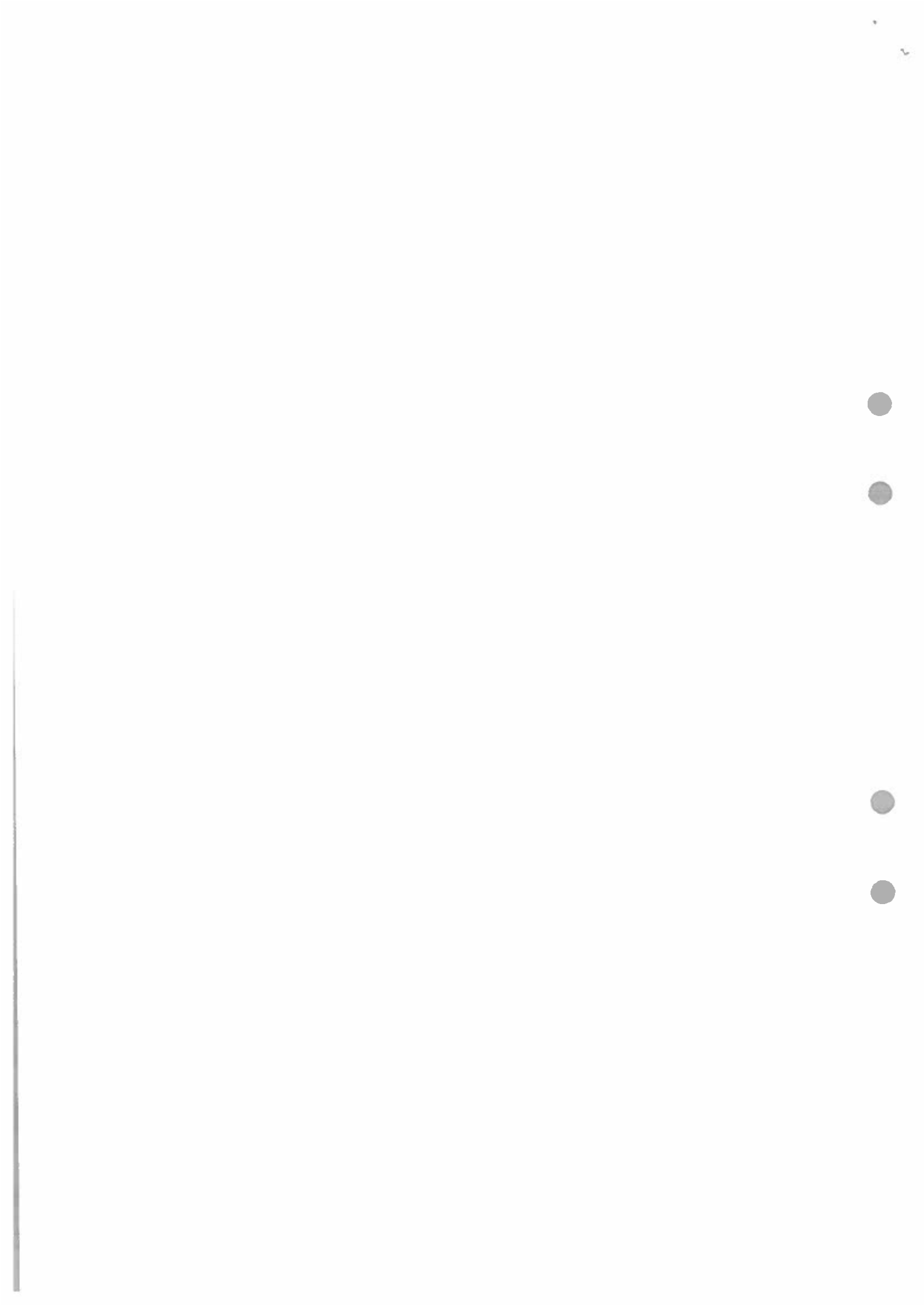
Uppgift 5. (20 poäng)

Låt Y_1, Y_2, Y_3 vara ett slumpmässigt urval från

$$f(y) = \begin{cases} \frac{2y}{\theta^2}, & 0 \leq y \leq \theta \\ 0, & \text{annars.} \end{cases}$$

$\hat{\theta} = \frac{Y_1 + Y_2 + Y_3}{2}$ är en estimator för θ .

- Bestäm väntevärdet för $\hat{\theta}$. Är $\hat{\theta}$ väntevärdesriktig?
- Bestäm variansen för $\hat{\theta}$.





Stockholms
universitet

Statistiska institutionen

Rättningsblad

Datum: 18/4-2017

Sal: Brunnsvikssalen

Tenta: Statistisk teori med tillämpningar II

Kurs: Statistisk teori med tillämpningar

ANONYMKOD:

STM-0032

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 82
Lär. ant. 20	20	19	19	20					

98 + 8 bonus

POÄNG	BETYG	Lärarens sign.
106	A	



① i.)

Falskt. P-värdet är sannolikheten att erhålla det observerade värdet eller ett mer extremt värde.

P-värdet är den lägsta signifikansnivån som nollhypotesen kan förkastas på. Eftersom $0,01 < 0,029$ kan nollhypotesen inte förkastas på signifikansnivån $\alpha = 0,01$.

Deremot kan nollhypotesen förkastas på signivån $0,05$ tex.

4

ii.)

Falskt. Styrkan = Power(Θ) är $1 - \beta$

β är sannolikheten för att ej förkasta H_0 givet att H_1 är sann $\beta = P(\text{ej förkasta } H_0 | H_1 \text{ sann})$

Styrkan är sannolikheten att förkasta H_0 om

H_1 är sann $\text{Power}(\Theta) = P(\text{förkasta } H_0 | H_1 \text{ sann})$

För att beräkna styrkan måste man veta H_1 och förkastelseområdet.

4

iii.)

Falskt. Signifikansnivån α är sannolikheten att

förkasta en sann nollhypotes: $P(\text{förkasta } H_0 | H_0 \text{ sann})$

och är samma sak som ett Typ I-fel

4

iv.

Falskt. Ett litet p-värde kan även fås om

man har många observationer eller liten

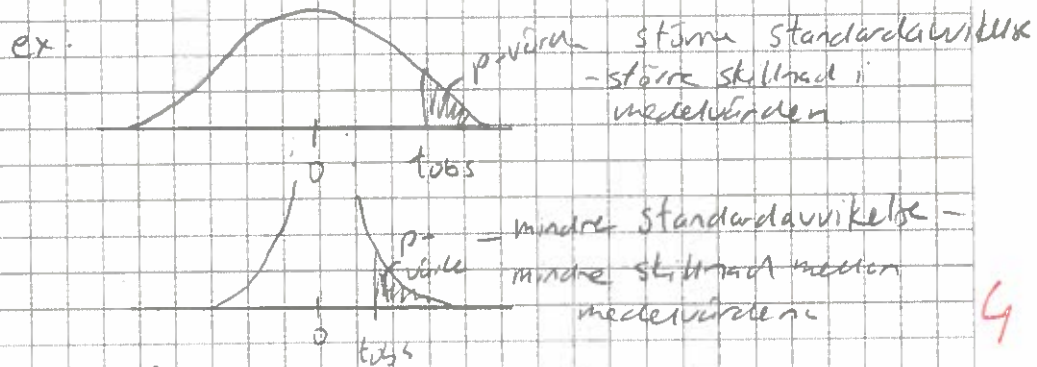
standardavvikelse. Med exempelvis ett t-test:

$$T = \frac{\bar{Y}_1 - \bar{Y}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Stor skillnad i medelvärden

ger ett stort T-värde,

med även en liten standardavvikelse och/eller stora
 urval ser ett stort T-värde och därmed mindre
 p-värde.



v) Falskt:

$$\lambda_{LR} = \frac{L(\hat{\theta}_0)}{L(\hat{\theta})}$$

$L(\hat{\theta}_0)$ är likelihooden för H_0 $\theta \in \Omega_0$

$L(\hat{\theta})$ är likelihooden för den bästa av alla

modeller. $0 < \lambda_{LR} < 1$ Om λ_{LR} är nära 1 är

$L(\hat{\theta}_0)$ lika stor som $L(\hat{\theta})$ och vi kan

ex förkasta H_0 , då H_0 är den bästa modellen.

Om λ_{LR} är liten betyder det att likelihooden

för $L(\hat{\theta})$ är större och det finns en

bättre modell än H_0 . Då förkastas H_0 .

(20) Bra!

(2)

p- "andelen kunder som är positivt inställda till spålvskanning"

$n = 120$ kunder

Totalt var 27 kunder positiva

84 st var medlemmar i kundklubben varav 20 positiva.

Av de 36 icke-medlemmarna var 7 positiva.

a) Beräkna 95% -igt konfidensintervall för andelen kunder som är positivt inställda till spålvskanning.

$$p = \frac{\text{antalet positiva}}{\text{totalt antal}} = \frac{27}{120} = 0,225$$

$$K.I. : \hat{p} \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

Vi kan anta att andelen är normalfördelad pga det stora urvalet enligt CLS. Därmed kan vi använda $Z_{\alpha/2}$ till konfidensintervallet.

(Tumregel $np > 5$, $n(1-p) > 5$) \hat{p}

Vi antar att urvalet består av oberoende observationer.

$$Z_{\alpha/2} = 1,96 \Rightarrow$$

$$0,225 \pm 1,96 \cdot \sqrt{\frac{0,225(1-0,225)}{120}} \Leftrightarrow 0,225 \pm 0,0747 \quad R$$

Svar: Ett 95% -igt konfidensintervall för p är $0,225 \pm 0,0747$ eller $[0,1503; 0,2997]$

för andelen kunder som är positivt inställda.

b) Hur stort urval skulle krävas om längden på ett 95%:igt k.i. $\leq 0,1$

$$\hat{p} \pm z_{\alpha/2} \underbrace{\sqrt{\frac{p(1-p)}{n}}}_B \quad R$$

Längden på konfidensintervallet är $2B = 0,1$

$$\Rightarrow B = 0,05$$

$$\Rightarrow z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} = 0,05 \quad R$$

Om vi har anledning att tro att p ska vara ett särskilt värde kan vi använda det, men för säkerhets skull använder vi $p = 0,5$ för att uppskatta variansen, då detta ger det största värdet på variansen och därmed kan vi vara säkra på att erhålla ett värde på B som stämmer med det givna villkoret för det n vi räknar ut. Ber

$$z_{\alpha/2} = z_{0,025} = 1,96 \quad p = 0,5$$

$$1,96 \cdot \sqrt{\frac{0,5 \cdot 0,5}{n}} = 0,05$$

$$\sqrt{\frac{0,25}{n}} = \frac{0,05}{1,96}$$

$$\sqrt{n} = \frac{1,96 \cdot \sqrt{0,25}}{0,05} = 19,6$$

$$\Rightarrow n = 384,16$$

Svar: Det skulle behövas 385 personer i urvalet. R

7

(2) c) Andelen kundklubsmedlemmar som är positivt inställda = $p_1 = \frac{20}{84} = 0,2381$ $n_1 = 84$

Andelen icke-medlemmar som är positivt inställda = $p_2 = \frac{7}{36} = 0,1944$ $n_2 = 36$

Samma antaganden som i a)

Skillnaden i andelen positiva: $p_1 - p_2 = 0,2381 - 0,1944 = 0,04365$

95% -igt konfidensintervall för skillnaden

$$\hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2} \sigma(\hat{p}_1 - \hat{p}_2) \Rightarrow \hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$0,0436 \pm 1,96 \sqrt{\frac{0,238(1-0,238)}{84} + \frac{0,1944(1-0,1944)}{36}}$$

$$0,0436 \pm 1,96 \cdot 0,0806$$

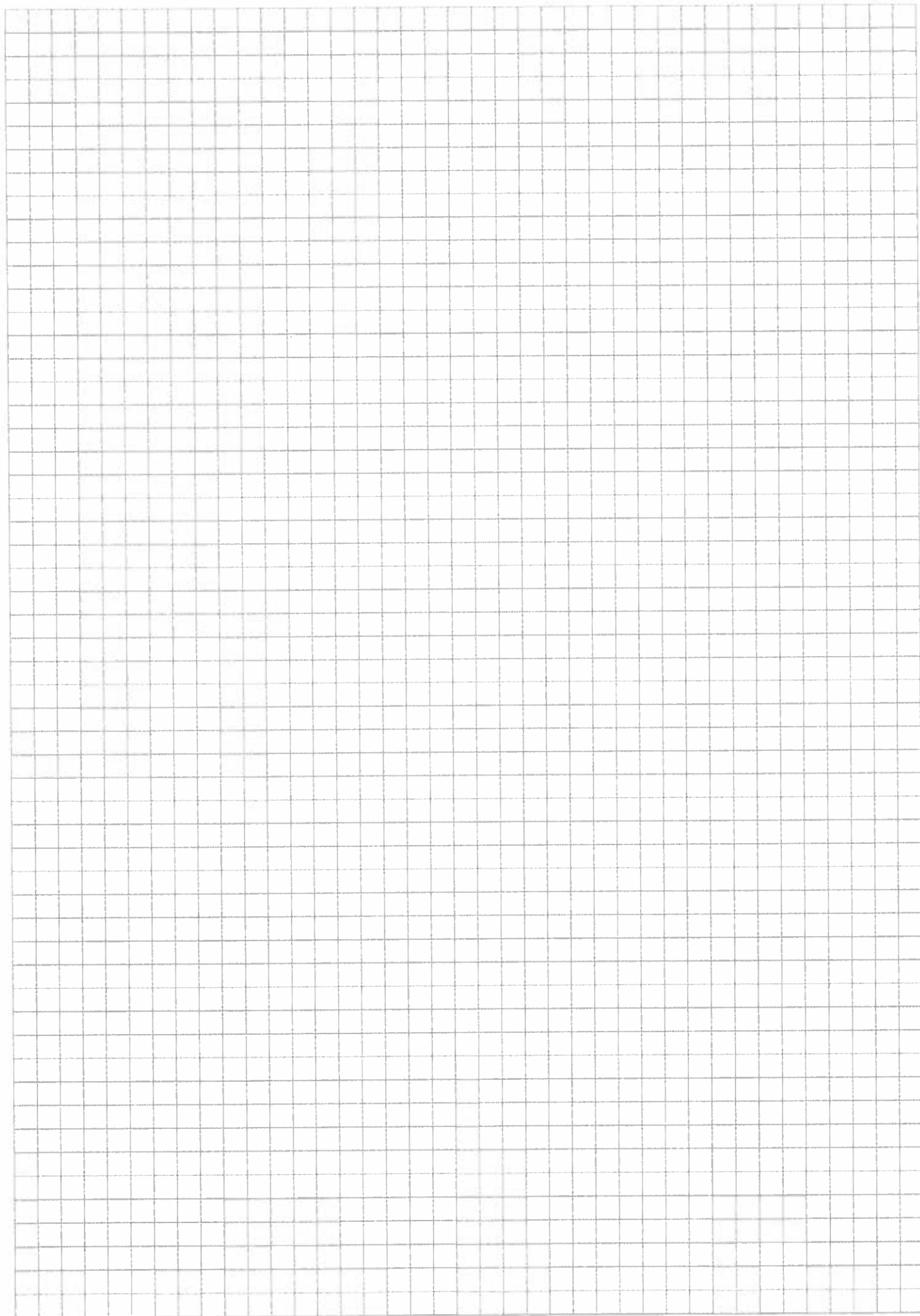
$$0,0436 \pm 0,1581$$

Svar: Ett 95% -igt konfidensintervall för skillnaden

i andelen som är positivt inställda, de båda sätten ges av $0,0436 \pm 0,158$ eller

$$[-0,114; 0,202]$$

7
20



- ③ Maskinerna A och B utvärderas med ett experiment. Man mäter tiden det tar att göra en viss uppsett. Alla maskintekniker utför ett test på maskin A och B.

Tekniker	Maskin A	Maskin B	D_i	Rang
1	32	30	2	4,5
2	40	39	1	2
3	42	42	0	tie
4	26	23	3	6
5	35	36	-1	②
6	29	27	2	4,5
7	45	41	4	7
8	22	21	1	2

$n=7$ pga
en tie

- a) Testa med icke-parametriskt test, om det går att påvisa någon skillnad mellan tiden med A eller B.

Jag väljer Wilcoxon Signed rank test.

Antaganden: Parvisa observationer, beroende observationer mellan varje försök som utförs av samma tekniker.

Oberoende observationer mellan de olika teknikernas försök.

Antas symmetriska fördelningar.

Hypoteser:

H_0 : Fördelningarna är lika

H_1 : Fördelningarna skiljer sig åt

Teststatistika T_{min} (T^-, T^+)

Jag väljer signifikansnivå $\alpha = 0,05$

Förkasta H_0 om $\bar{T} = \min(T^-, T^+) \leq T_0 = 2$ (tabell 9)
[Dubbelsidigt test, $p = 0,05$ $n = 7$ $= T_0 = 2$] R

$$T^- = 2$$

$$T^+ = 26$$

R

$$T = \min(2, 26) = 2$$

$$T^- = 2 \leq 2 = T_0$$

Vi kan precis förkasta H_0 på signifikansnivå 0,05.

Svar: Vi kan förkasta H_0 på signifikansnivå 0,05.

Fördelningarna skiljer sig åt i lägen, dvs det är skillnad mellan de olika maskinerna i tid att utföra uppgiften. 8

b) Utöver antagandena i a) antar vi att tiderna är normalfördelade och utför ett t-test då σ är okänt och det är litet urval (CGS gäller inte)

Parvisa observationer.

Hypoteser: $H_0: \bar{D} = 0$ $\alpha = 0,05$

$$H_A: \bar{D} \neq 0 \quad \checkmark \quad \underline{\underline{H_0}}$$

teststatistika: $T = \frac{\bar{D} - 0}{S/\sqrt{n}} \sim t(n-1) = t(7)$ R

$$\bar{D} = 1,5$$
 R

$$S^2 = \frac{\sum_{i=1}^8 D_i^2 - n \cdot \bar{D}^2}{n-1} = \frac{36 - 8 \cdot 1,5^2}{7} = 2,5714$$
 R

Förkasta H_0 om $|T_{\text{oss}}| > T_{\text{crit}} = 2,36$ (tabell 3, $f_5 = 7, \alpha = 0,025$)
R

forts. på
nästa blad

③ 6 forts.)

$$T_{0.05} = \frac{1,5}{\sqrt{257/8}} = 2,64597$$

$$T_{0.05} = 2,65 > 2,36 = T_{krit}$$

\Rightarrow Förkasta H_0

Svar H_0 förkastas på signifikansnivå 0,05

Det finns bevis för att det är skillnad mellan tiden det tar att utföra utgången med de olika maskinerna.

7

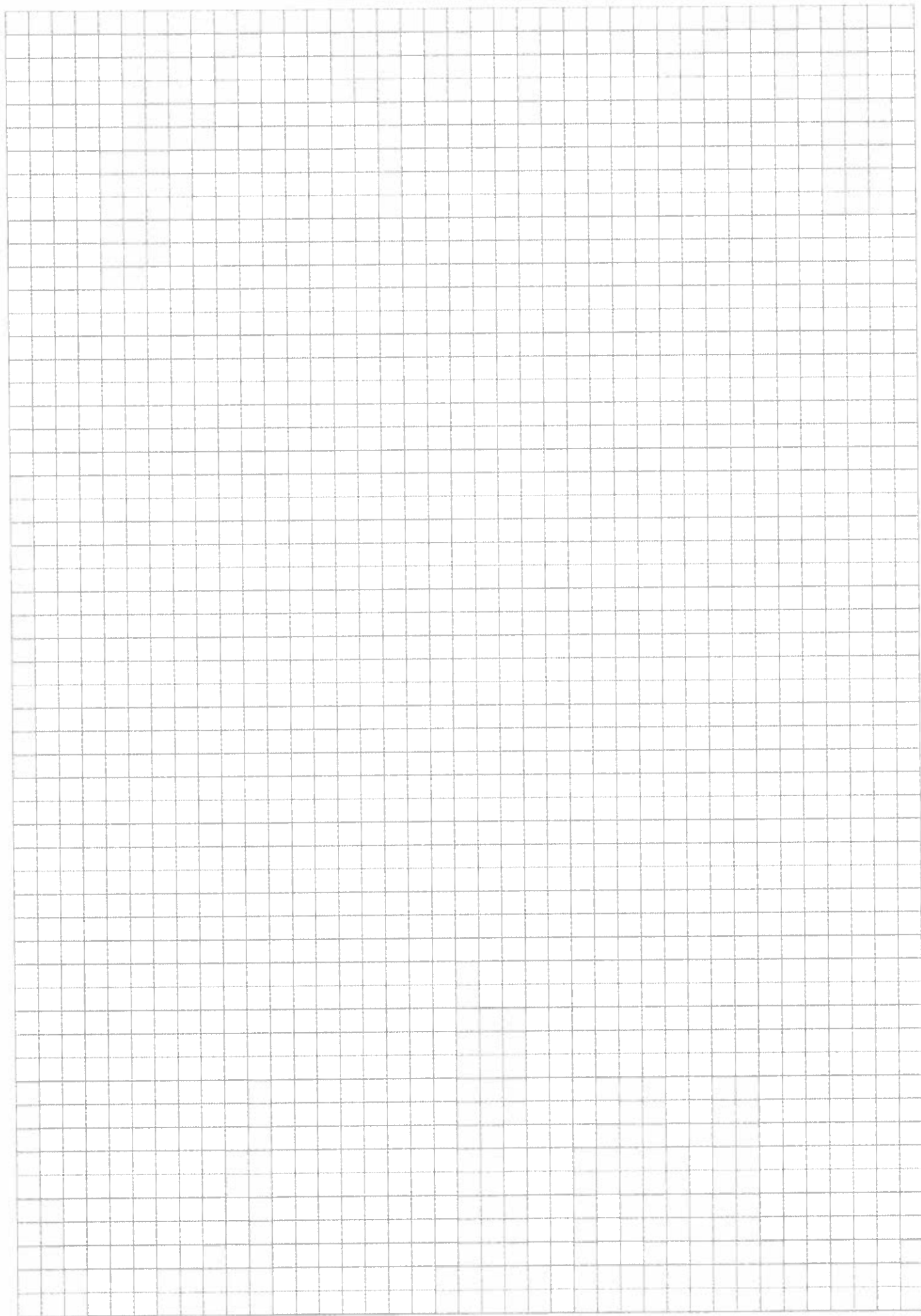
c)

För att utföra ett parvis test behöver det vara samma person som utför testen. Olika personer arbetar ju olika snabbt. Eftersom det skiljer sig ganska mycket åt mellan de olika teknikerna, hur lång tid processen tar, från 22 till 45, så påverkar teknikerna väldigt mycket.

Genom att samma tekniker gör processen på både maskin A och B är det endast skillnaden, maskinerna som utsör skillnaden, tiden för varje parvis test. Man skulle också kunna göra ett test där samma tekniker gör många test på maskin A och B och göra ett Mann-Whitney U-test / olika tekniker

4

19



(4)

Y - "antal gånger en student behöver tentera för att få godkänt"

$$p(y) = p(1-p)^{y-1} \quad y = 1, 2, \dots$$

Delta är en geometrisk fördelning

p - sannolikheten att få godkänt

Väntevärde för Y i en geometrisk fördelning: $\frac{1}{p}$

$n = 8$ studenter Resultat $Y = 1, 1, 1, 2, 3, 1, 1, 1$

a) Härled momentskattningen \hat{p}_{mom} och beräkna skattningen för urvalet

Momentskattningen sörs genom att populationsmomentet μ'_1 antas kunna uppskattas med urvalsmomentet m'_1

$$\mu'_1 = \frac{1}{p} \quad m'_1 = \frac{\sum Y_i}{n} = \bar{Y} \quad R \quad R$$

$$\mu'_1 = m'_1 = \frac{\sum Y_i}{n} = \frac{1}{p} \Rightarrow p = \frac{n}{\sum Y_i} = \frac{1}{\bar{Y}}$$

$$\hat{p}_{mom} = \frac{1}{\bar{Y}} \quad \bar{Y} = \frac{1+1+1+2+3+1+1+1}{8} = 1,375 \quad R$$

$$\hat{p}_{mom} = \frac{1}{1,375} = 0,727 \quad R$$

Så $\hat{p}_{mom} = \frac{1}{\bar{Y}}$ och blir 0,727 med detta urval.

b) Härled maximum likelihood-estimation \hat{p}_{ML}

$$p(y) = p(1-p)^{y-1}$$

$$L(p) = \prod_{i=1}^n p(1-p)^{y_i-1} = \{ \text{pga oberoende} \} =$$

$$p^n \cdot (1-p)^{\sum_{i=1}^n y_i - n} \quad \left(\sum_{i=1}^n y_i \right) \quad R$$

logaritmera:

$$l(p) = \ln \left(p^n (1-p)^{\sum_{i=1}^n y_i} \right) = n \cdot \ln p + \left(\sum_{i=1}^n y_i - n \right) \ln(1-p)$$

$$= n \ln p + \sum_{i=1}^n y_i \cdot \ln(1-p) - n \ln(1-p) \quad R$$

Derivera med avseende på p :

$$\frac{dl(p)}{dp} = \frac{n}{p} + \frac{\sum_{i=1}^n y_i}{1-p} (-1) - \frac{n}{1-p} (-1) = \frac{n}{p} - \frac{\sum_{i=1}^n y_i}{1-p} + \frac{n}{1-p} \quad R$$

Sätt derivatan = 0 för att hitta maxvärdet

$$\frac{n}{p} - \frac{\sum_{i=1}^n y_i}{1-p} + \frac{n}{1-p} = 0 \quad \frac{n(1-p) - \sum_{i=1}^n y_i + n \cdot p}{p(1-p)} = 0$$

$$n(1-p) - \sum_{i=1}^n y_i + n \cdot p = 0 \quad n - n \cdot p - \sum_{i=1}^n y_i + n \cdot p = 0$$

$$n = \sum_{i=1}^n y_i \Rightarrow p = \frac{n}{\sum_{i=1}^n y_i} = \frac{1}{\bar{y}} \quad R$$

$$\hat{p}_{ML} = \frac{1}{\bar{y}} \quad \text{Med det samma urvalet blir skattningen}$$

$$\hat{p}_{ML} = 0,727 \quad R$$

12

Svar $\hat{p}_{ML} = \frac{1}{\bar{y}} = 0,727$

c) $E(\hat{Y}) = \frac{1}{\hat{p}_{ML}} = \frac{1}{0,727} = 1,3755$

Svar 1,376 sgr behövs
man testern i genomsnitt 4

$$(5) \quad a) \quad f(y) = \begin{cases} \frac{2y}{\theta^2} & 0 \leq y < \theta \\ 0 & \text{annars} \end{cases}$$

Y_1, Y_2, Y_3 s.u. från $f(y)$ y , kont. variabel
Antag oberoende observationer.

$$\hat{\theta} = \frac{Y_1 + Y_2 + Y_3}{2}$$

$$E(\hat{\theta}) = E\left(\frac{Y_1 + Y_2 + Y_3}{2}\right) = \frac{E(Y_1) + E(Y_2) + E(Y_3)}{2} \quad R$$

$$E(Y) = \int_{-\infty}^{\infty} y \cdot f(y) dy = \int_0^{\theta} y \cdot \frac{2y}{\theta^2} dy = \int_0^{\theta} \frac{2y^2}{\theta^2} dy =$$

$$\left[\frac{2y^3}{3\theta^2} \right]_0^{\theta} = \frac{2\theta^3}{3\theta^2} = \frac{2\theta}{3} \quad R$$

$$E(\hat{\theta}) = \frac{\frac{2\theta}{3} + \frac{2\theta}{3} + \frac{2\theta}{3}}{2} = \frac{6\theta}{6} = \frac{6\theta}{6} = \theta \quad R$$

Svar $E(\hat{\theta}) = \theta$ Ja, θ är väntevärdesskattis.

R

8

$$b) \quad V(Y) = \int_{-\infty}^{\infty} y^2 f(y) dy - [E(Y)]^2$$

$$\int_{-\infty}^{\infty} \frac{y^2 \cdot 2y}{\theta^2} dy = \int_0^{\theta} \frac{2y^3}{\theta^2} dy = \left[\frac{2y^4}{4\theta^2} \right]_0^{\theta} = \frac{2\theta^4}{4\theta^2} = \frac{\theta^2}{2} \quad R$$

$$V(Y) = \frac{\theta^2}{2} - \left(\frac{2\theta}{3}\right)^2 = \frac{\theta^2}{2} - \frac{4\theta^2}{9} = \frac{9\theta^2 - 8\theta^2}{18} = \frac{\theta^2}{18} \quad R$$

$$V(\hat{\theta}) = V\left(\frac{Y_1 + Y_2 + Y_3}{3}\right) = \frac{1}{9} \cdot V(Y_1 + Y_2 + Y_3) = \text{pga oserveni} \quad R$$

$$\frac{1}{9} (V(Y_1) + V(Y_2) + V(Y_3)) = \frac{1}{9} \left(\frac{\theta^2}{18} + \frac{\theta^2}{18} + \frac{\theta^2}{18}\right) = \frac{3\theta^2}{18} \cdot \frac{1}{9} =$$

$$\frac{3\theta^2}{72} = \frac{\theta^2}{24} \quad R$$

Svar: Variansen $V(\hat{\theta})$ är $\frac{\theta^2}{24}$

12

20