

## Tentamen i Undersökningsmetodik (4,5 hp)

### Kurs: Regressionsanalys och undersökningsmetodik

2020-02-13

---

<b>Skrivtid:</b>	kl. 16.00 - 21.00 (5 timmar)
<b>Godkända hjälpmedel:</b>	Miniräknare utan lagrade formler och text
<b>Vidhäftade hjälpmedel:</b>	Formelsamling och Statistiska tabeller (endast de tabeller som krävs)

- Tentamen består av 5 uppgifter, i förekommande fall uppdelade i deluppgifter. Maximalt antal poäng anges per deluppgift.
- Svar med fullständiga redovisningar ska lämnas.
  - Använd endast skrivpapper som tillhandahålls i skrivsalen.
  - För full poäng på en uppgift krävs tydliga, utförliga och väl motiverade lösningar.
  - Kontrollera alltid dina beräkningar och lösningar! Slarvfel kan också ge poängavdrag!
  - Använd minst fem värdesiffror i dina beräkningar (1,2345 och 1234,5 är exempel på tal med fem värdesiffror). I förekommande fall är det inte möjligt pga. avrundning i t.ex. SAS-utskrifter men utgå då ifrån det som är givet. Du kan dock avrunda ditt slutliga svar.
- Tentamen kan maximalt ge 100 poäng och för godkänt resultat krävs minst 50.
- Betygsgränser:
  - A: 90 – 100 p
  - B: 80 – 89 p
  - C: 70 – 79 p
  - D: 60 – 69 p
  - E: 50 – 59 p
  - Fx: 40 – 49 p
  - F: 0 – 40 p

OBS! Fx och F är underkända betyg som kräver omexamination. Studenter som får betyget Fx kan alltså inte komplettera för högre betyg.

- Lösningförslag läggs ut på Athena kort efter tentamen.

**LYCKA TILL!**

### Uppgift 1. (20p)

En nyanställd statistiker fick i uppdrag att skatta den genomsnittliga kommunalskatten bland Sveriges  $N = 290$  kommuner. Ett obundet slumpmässigt urval utan återläggning av storlek  $n = 30$  drogs och kommunalskatten för de dragna kommunerna registrerades och vid en första sammanställning fick man följande resultat:

$$\sum y_k = 661.73 \quad \sum y_k^2 = 14751.37$$

- Skatta den genomsnittliga kommunalskatten och ange ett 90 % konfidensintervall. (8p)
- Med undersökningen och de framräknade resultaten som underlag vill man nu gå ut och meddela att den genomsnittliga kommunalskatt som svenskar betalar är den skattning som du fick i a) ovan (med felmarginal förstås). Vad är det för fel med den slutsatsen? (4p)
- Anta att man vill ha exakt samma felmarginal som du fick i a) ovan men öka konfidensgraden till 95% istället. Hur stort stickprov skulle du behöva för att få exakt samma felmarginal? Använd att

$$Z_{0.05} \cdot D^* = Z_{0.025} \cdot D$$

där  $D^*$  = standardfelet som du fick i a). (8p)

### Uppgift 2. (25p)

En läkare vill skatta andelen med för högt blodtryck hos förvärvsarbetande kvinnor i en mindre region. Läkaren gör ett stratifierat urval efter ålder och får följande resultat:

Stratum	Antal i populationen	Urvalsstorlek	Antal med för högt blodtryck
- 25	2510	25	1
26 - 35	5270	75	15
36 - 45	5310	75	39
46 - 55	4960	100	18
56 -	4400	125	36

I dina beräkningar ska ändlighetskorrektion användas.

- Skatta andelen i hela populationen som har för högt blodtryck och beräkna ett 95% konfidensintervall för andelen. (10p)
- Anta att läkaren istället hade gjort ett lika stort obundet slumpmässigt urval utan återläggning. Hur stor hade felmarginalen då varit? Blir det bättre eller sämre med stratifieringen? (7p)
- Läkaren planerar att upprepa undersökningen nästa år. Utgå ifrån de data som du har fått och föreslå läkaren en alternativ stickprovssallokering som kan ge en mindre felmarginal med samma totala stickprovstorlek. (8p)

### Uppgift 3. (25p)

En parkförvaltning vill uppskatta medelåldern på träd i en större stadspark. Att bestämma ålder är besvärligt eftersom man behöver räkna antalet åldersringar som man får genom att ta ett borrhov genom trädet. Dels är det tidskrävande men det kan också vara skadligt för trädet. Men, allmänt gäller att ju äldre trädet, desto större är dess diameter mätt över stammen och diameter är lättare att mäta.

Man mäter diametern på samtliga  $N = 1132$  träd i parken och finner att populationsmedelvärdet är  $\mu_x = 10.3$ . De väljer sedan slumpmässigt  $n = 20$  träd för åldersmätning. Låt alltså  $x_i =$  diametern mätt i cm och  $y_i =$  åldern mätt i år för träd nummer  $i$ . Följande summor beräknade från stickprovet finns nu tillgängliga:

$$\sum_{k \in S} x_k = 478$$

$$\sum_{k \in S} y_k = 2148$$

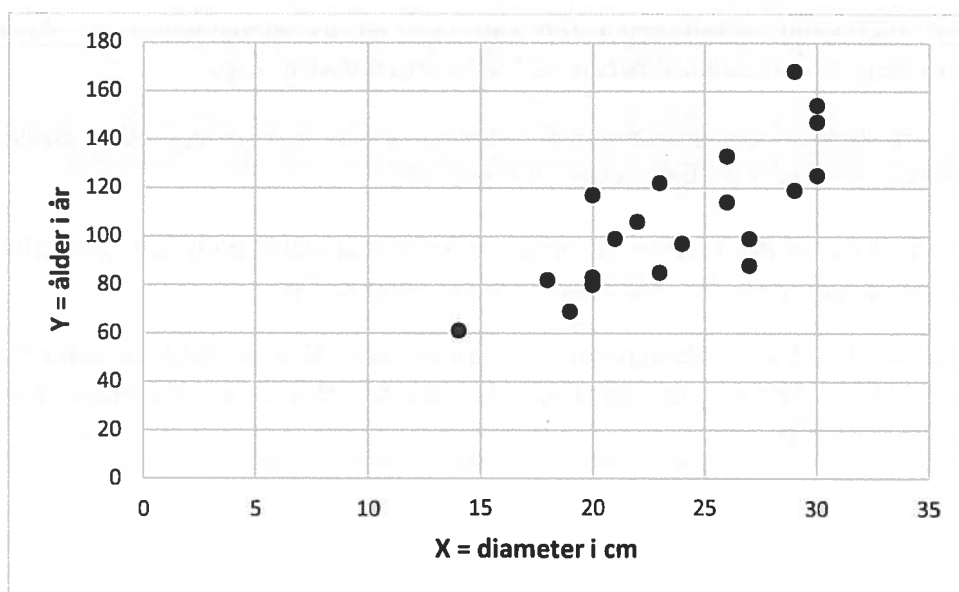
$$\sum_{k \in S} x_k^2 = 11832$$

$$\sum_{k \in S} y_k^2 = 246304$$

$$\sum_{k \in S} x_k y_k = 53315$$

Du ombeds skatta den genomsnittliga åldern för samtliga träd i parken. I beräkningarna ska ändlighetskorrektur användas.

- Skatta  $\mu_y =$  den genomsnittliga åldern med en regressionsestimator med diametern som hjälpvariabel. Beräkna sedan standardfelet för skattningen. (10p)
- Skatta  $\mu_y =$  den genomsnittliga åldern med en kvotestimator med diametern som hjälpvariabel. Beräkna sedan standardfelet för skattningen. (10p)
- Kommentera dina resultat i a) och b) ovan. Är det stora skillnader? Diskutera kortfattat vilken av de två skattningsmetoderna som det är mest lämplig i denna situation. Utgå ifrån definitionerna av variablerna, de bakomliggande modeller som skattningsmetoderna baseras på och diagrammet nedan som visar stickprovets spridning i bägge variabler. (5p)



#### Uppgift 4. (20p)

För var och en av följande deluppgifter ska du svara kortfattat. Hela uppgiften bör kunna redovisas på maximalt ca två A4-sidor. Du får gärna komplettera med bilder och skisser.

- Förklara vad teleskopeffekten är för något. Vad är bakåt respektive framåt effekter i detta sammanhang? Vad medför det för problem? (5p)
- Anta att man i en undersökning har en fråga av typen "Vilka av följande alternativ gäller för dig? Fler än ett alternativ kan väljas." Vad ska man tänka på när man definierar de olika svarsalternativen och hur kan svaren koda? (5p)
- Ange två olika klassifikationsstandarder som används i svensk statistik och vad de beskriver. Nämn en väsentlig fördel med klassifikationsstandarder. (5p)
- Beskriv översiktligt vad ett snöbollsurval är och när är en snöbollsdesign kan vara ett användbart alternativ. Är snöbollsurval ett sannolikhetsurval? (5p)

#### Uppgift 5. (10p)

Man har en population  $U$  bestående av  $N = 6$  element med följande värden på en variabel  $Y$ :

$$U = \{2, 6, 8, 14, 18, 24\}$$

Man föreslår en förenklad urvalsdesign där man endast tillåter tre möjliga stickprov  $S_k$ :

$$S_1 = \{2, 14\} \quad S_2 = \{6, 18\} \quad S_3 = \{8, 24\}$$

Vilket stickprov som dras är slumpmässigt och sannolikheten för var och en är lika stor dvs.  $1/3$ .

- Ange inklusionssannolikheten för vart och ett av elementen i  $U$ . Är den föreslagna urvalsdesignen ett sannolikhetsurval? Motivera ditt svar. (3p)
- Beräkna stickprovsmedelvärdena för vart och ett av stickproven ovan. Beräkna sedan  $\sigma_{\bar{y}}^2 =$  variansen för dessa stickprovsmedelvärden. (2p)
- Beräkna  $V(\bar{y}) =$  den teoretiska variansen för  $\bar{y}$  som skulle gälla om du istället hade valt OSU utan återläggning med  $n = 2$  som din urvalsdesign. (2p)
- Vilken av de två urvalsdesignerna ger lägst varians? Kan du förklara varför? TIPS: Vad är den förenklade urvalsdesignen egentligen för design? Hur är listan (ramen) för elementen i  $U$  organiserad? (3p)

Formulário de Avaliação

Nome do Aluno: \_\_\_\_\_

Data: \_\_\_\_\_

Assunto: \_\_\_\_\_

Descrição da atividade realizada:

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

# Formel- och tabellsamling

## DESKRIPTIV STATISTIK

Notation:  $U$  = populationen  
 $S$  = stickprov (stort  $S$ );  $\subseteq U$

Medelvärde:	$\mu = \frac{1}{N} \sum_{k \in U} y_k$	Varians:	$\sigma^2 = \frac{\sum_{k \in U} (y_k - \mu_y)^2}{N} = \frac{\sum_{k \in U} y_k^2 - N\mu_y^2}{N}$
	$\bar{y} = \frac{1}{n} \sum_{k \in S} y_k$		$s^2 = \frac{\sum_{k \in S} (y_k - \bar{y})^2}{n-1} = \frac{\sum_{k \in S} y_k^2 - n\bar{y}^2}{n-1}$
Andel:	$P = \frac{1}{N} \sum_{k \in U} y_k$		$\sigma^2 = P(1-P)$
( $y_k = 0$ eller 1)	$\hat{p} = \frac{1}{n} \sum_{k \in S} y_k$		$s^2 = \frac{n}{n-1} \hat{p}(1-\hat{p})$
Kovarians:	$\sigma_{xy} = Cov(x, y) = \frac{\sum_{k \in U} (x_k - \mu_x)(y_k - \mu_y)}{n-1} = \frac{\sum_{k \in U} x_k y_k - n\bar{x}\bar{y}}{n-1}$		
	$s_{xy} = Cov(x, y) = \frac{\sum_{k \in U} (x_k - \bar{x})(y_k - \bar{y})}{n-1} = \frac{\sum_{k \in U} x_k y_k - n\bar{x}\bar{y}}{n-1}$		
Korrelation:	$r_{xy} = Corr(x, y) = \frac{s_{xy}}{s_x \cdot s_y} = \frac{s_{xy}}{\sqrt{s_x^2 \cdot s_y^2}}$		

## Beräkningsformler för VARIANSER och REGRESSIONSKOEFFICIENT

$s^2 = \frac{n \sum y_k^2 - (\sum y_k)^2}{n(n-1)} = \frac{\sum y_k^2 - \frac{(\sum y_k)^2}{n}}{n-1} = \frac{\sum y_k^2 - n\bar{y}^2}{n-1} = \frac{\sum (y_k - \bar{y})^2}{n-1}$
$b = \frac{n \sum x_k y_k - (\sum x_k)(\sum y_k)}{n \sum x_k^2 - (\sum x_k)^2} = \frac{\sum x_k y_k - \frac{(\sum x_k)(\sum y_k)}{n}}{\sum x_k^2 - \frac{(\sum x_k)^2}{n}} = \frac{\sum x_k y_k - n\bar{x}\bar{y}}{\sum x_k^2 - n\bar{x}^2}$
$= \frac{\sum (x_k - \bar{x})(y_k - \bar{y})}{\sum (x_k - \bar{x})^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y}) / (n-1)}{\sum (x_i - \bar{x})^2 / (n-1)}$
$= \frac{s_{xy}}{s_x^2} = \frac{s_{xy}}{s_x^2} \cdot \frac{s_x s_y}{s_x s_y} = \frac{s_{xy}}{s_x s_y} \cdot \frac{s_y}{s_x} = r_{xy} \cdot \frac{s_y}{s_x}$

OBS! Notationen har förenklats ovan, summationsindex är alltid  $k$ : ex.  $\sum y_k = \sum_{k \in S} y_k$

## OBUNDET SLUMPMÄSSIGT URVAL u.å.

Parameter:	Estimator:	Varians $V(\cdot)$ :	Variansskattning $\hat{V}(\cdot)$ :
$\mu$	$\bar{y} = \frac{1}{n} \sum_{k \in S} y_k$	$V(\bar{y}) = \left( \frac{N-n}{N-1} \right) \frac{\sigma^2}{n}$	$\hat{V}(\bar{y}) = \left( 1 - \frac{n}{N} \right) \frac{s^2}{n}$
$\tau$	$\hat{\tau} = N\bar{y}$	$V(\hat{\tau}) = N^2 V(\bar{y})$	$\hat{V}(\hat{\tau}) = N^2 \cdot \hat{V}(\bar{y})$
$P$	$\hat{p} = \frac{1}{n} \sum_{k \in S} y_k$	$V(\hat{p}) = \left( \frac{N-n}{N-1} \right) \frac{P(1-P)}{n}$	$\hat{V}(\hat{p}) = \left( 1 - \frac{n}{N} \right) \frac{\hat{p}(1-\hat{p})}{n-1}$
$A$	$\hat{A} = N\hat{p}$	$V(\hat{A}) = N^2 V(\hat{p})$	$\hat{V}(\hat{A}) = N^2 \cdot \hat{V}(\hat{p})$

Stickprovsstorlek: 
$$n \geq \frac{N\sigma^2}{D^2(N-1) + \sigma^2}$$

## STRATIFIERAT URVAL u.å.

Notation:  $L =$  antal strata

$N_k =$  populationsstorleken för stratum  $k = 1, \dots, L$

$n_k =$  stickprovets storlek i stratum  $k = 1, \dots, L$

$W_k = N_k/N$

$\bar{y}_k =$  stickprovsmedelvärde i stratum  $k = 1, \dots, L$

$s_k^2 =$  stickprovsvarians i stratum  $k = 1, \dots, L$

Parameter	Estimator	Varians $V(\cdot)$	Variansskattning $\hat{V}(\cdot)$
$\mu$	$\bar{y}_{\text{str}} = \sum_{k=1}^L W_k \bar{y}_k$	$\sum_{k=1}^L W_k^2 \left( \frac{N_k - n_k}{N_k - 1} \right) \frac{\sigma_k^2}{n_k}$	$\sum_{k=1}^L W_k^2 \left( 1 - \frac{n_k}{N_k} \right) \frac{s_k^2}{n_k}$
$\tau$	$\hat{\tau}_{\text{str}} = N \bar{y}_{\text{str}}$	$\sum_{k=1}^L N_k^2 \left( \frac{N_k - n_k}{N_k - 1} \right) \frac{\sigma_k^2}{n_k}$	$\sum_{k=1}^L N_k^2 \left( 1 - \frac{n_k}{N_k} \right) \frac{s_k^2}{n_k}$
$P$	$\hat{p}_{\text{str}} = \sum_{k=1}^L W_k \hat{p}_k$	$\sum_{k=1}^L W_k^2 \left( \frac{N_k - n_k}{N_k - 1} \right) \frac{P_k(1-P_k)}{n_k}$	$\sum_{k=1}^L W_k^2 \left( 1 - \frac{n_k}{N_k} \right) \frac{\hat{p}_k(1-\hat{p}_k)}{n_k - 1}$
$A$	$\hat{A}_{\text{str}} = N \hat{p}_{\text{str}}$	$\sum_{k=1}^L N_k^2 \left( \frac{N_k - n_k}{N_k - 1} \right) \frac{P_k(1-P_k)}{n_k}$	$\sum_{k=1}^L N_k^2 \left( 1 - \frac{n_k}{N_k} \right) \frac{\hat{p}_k(1-\hat{p}_k)}{n_k - 1}$

Optimal allokering: 
$$n_k = n \cdot \frac{N_k \sigma_k}{\sum_{j=1}^L N_j \sigma_j}$$

## KLUSTERURVAL - OSU u.å.

Notation:  $U$  = population av kluster

$S$  = stickprov av kluster

$N$  = antal kluster totalt

$n$  = antal kluster i stickprovet

$M$  = totalt antal element

$m_i$  = antal element i kluster nr  $i = 1, 2, \dots, N$

$\bar{m}$  = stickprovsmedelvärde av klusterstorlekarna  $m_i$

$s_{m_i}^2$  = stickprovsvariansen av klusterstorlekarna  $m_i$

$\tau = \sum_{k \in U} y_k$  = totalvärdet för  $y$  i hela populationen

$\mu = \tau/M$  = populationsmedelvärde av  $y$

$\tau_i = \sum_{k \in C_i} y_k$  = totalvärdet för kluster nr  $i = 1, 2, \dots, N$

$\bar{\tau}$  = stickprovsmedelvärde av totalvärdena  $\tau_i$

$s_{\tau_i}^2$  = stickprovsvariansen av totalvärdena  $\tau_i$

$A = \sum_{k \in U} y_k$  = antalet ettor i hela populationen; ( $y_k = 0$  eller  $1$ )

$P = A/M$  = andelen ettor i hela populationen; ( $y_k = 0$  eller  $1$ )

Parameter	Estimator	Variansskattning
$M$	$\hat{M}_{vvr} = N \cdot \bar{m}$	$\hat{V}(\hat{M}_{vvr}) = N^2 \cdot \left(1 - \frac{n}{N}\right) \cdot \frac{s_{m_i}^2}{n}$
$\mu$	$\bar{y}_{vvr} = \frac{\hat{t}_{vvr}}{M} = \frac{N\bar{\tau}}{M}$	$\hat{V}(\bar{y}_{vvr}) = \frac{N^2}{M^2} \cdot \left(1 - \frac{n}{N}\right) \cdot \frac{s_{\tau_i}^2}{n}$
	$\bar{y}_{kvot} = \frac{\hat{t}_{vvr}}{\hat{M}} = \frac{\sum_{i \in S} \tau_i}{\sum_{i \in S} m_i}$	$\hat{V}(\bar{y}_{kvot}) = \left(\frac{1}{\bar{m}}\right)^2 \left(1 - \frac{n}{N}\right) \frac{\sum_{i \in S} (\tau_i - \bar{y}_{kvot} m_i)^2}{n(n-1)}$
$\tau$	$\hat{t}_{vvr} = N\bar{\tau}$	$\hat{V}(\hat{t}_{vvr}) = N^2 \cdot \left(1 - \frac{n}{N}\right) \cdot \frac{s_{\tau_i}^2}{n}$
	$\hat{t}_{kvot} = M\bar{y}_{kvot} = \frac{M}{\hat{M}} \hat{t}_{vvr}$	$\hat{V}(\hat{t}_{kvot}) = \left(\frac{M}{\bar{m}}\right)^2 \left(1 - \frac{n}{N}\right) \frac{\sum_{i \in S} (\tau_i - \bar{y}_{kvot} m_i)^2}{n(n-1)}$
$P$	<i>formler utgår</i>	
$A$	<i>formler utgår</i>	



## SKATTNINGSMETODER

Notation:  $\tau_y$  = totalvärdet för variabeln  $y$  för hela populationen  
 $\hat{\tau}_y$  = skattningen av  $\tau_y$  under OSU  
 $\mu_y$  = populationsmedelvärdet av för variabeln  $y$

### Kvotskattning under OSU u.å.:

Parameter      Punkt- resp. variansskattning

$\tau_y$	$\hat{\tau}_{\text{kvot}} = \hat{R} \cdot \tau_x = \frac{\sum_{k \in S} y_k}{\sum_{k \in S} x_k} \cdot \tau_x = \frac{\tau_x}{\hat{\tau}_x} \cdot \hat{\tau}_y \quad \text{där} \quad \hat{R} = \frac{\sum_{k \in S} y_k}{\sum_{k \in S} x_k} = \frac{\hat{\tau}_y}{\hat{\tau}_x}$ $\hat{V}(\hat{\tau}_{\text{kvot}}) = N^2 \cdot \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \left(\frac{\sum_{k \in S} (y_k - \hat{R}x_k)^2}{n-1}\right)$ <p>där <math>\sum_{k \in S} (y_k - \hat{R}x_k)^2 = \sum_{k \in S} y_k^2 - 2\hat{R} \sum_{k \in S} x_k y_k + \hat{R}^2 \sum_{k \in S} x_k^2</math></p>
$\mu_y$	$\hat{\mu}_{\text{kvot}} = \hat{R} \cdot \mu_x = \frac{\sum_{k \in S} y_k}{\sum_{k \in S} x_k} \cdot \mu_x = \frac{\mu_x}{\bar{x}} \cdot \bar{y}$ $\hat{V}(\hat{\mu}_{\text{kvot}}) = \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \left(\frac{\sum_{k \in S} (y_k - \hat{R}x_k)^2}{n-1}\right)$

### Regressionskattning under OSU u.å.:

Parameter      Punkt- och variansskattning

$\mu_y$	$\hat{\mu}_{\text{reg}} = \bar{y} + b(\mu_x - \bar{x}) \quad \text{där} \quad b = \frac{\sum_{k \in S} (y_k - \bar{y})(x_k - \bar{x})}{\sum_{k \in S} (x_k - \bar{x})^2}$ $\hat{V}(\hat{\mu}_{\text{reg}}) = \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \left(\frac{\sum_{k \in S} (y_k - \bar{y})^2 - b^2 \sum_{k \in S} (x_k - \bar{x})^2}{n-2}\right)$ <p>där <math>\sum_{k \in S} (y_k - \bar{y})^2 = \sum_{k \in S} y_k^2 - n\bar{y}^2</math></p>
$\tau_y$	$\hat{\tau}_{\text{reg}} = N \cdot \hat{\mu}_{\text{reg}}$ $\hat{V}(\hat{\tau}_{\text{reg}}) = N^2 \cdot \hat{V}(\hat{\mu}_{\text{reg}})$

### Poststratifiering under OSU u.å.:

Parametrar och estimatorer - se under **Stratifierat urval** ovan

OBS! Populationsvikterna  $W_k$  måste vara kända.

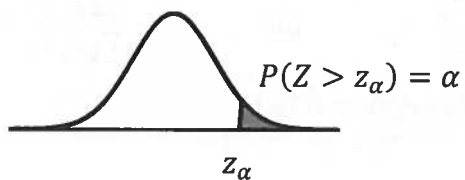
Variansskattning - *formler utgår*

## Från tabellsamlingen

**TABELL 2.** Normalfördelningens kvantiler, standardiserad

$Z \in N(0, 1)$ . Vilket värde har  $z_\alpha$  om  $P(Z > z_\alpha) = \alpha$  där  $\alpha$  är en given sannolikhet.

Utnyttja även  $\Phi(-z) = 1 - \Phi(z)$  för  $P(Z \leq -z_\alpha)$ .



$\alpha$	$z_\alpha$
0,25	0,6745
0,10	1,2816
0,05	1,6449
0,025	1,9600
0,010	2,3263
0,005	2,5758
0,0025	2,8070
0,0010	3,0902
0,0005	3,2905
0,00025	3,4808
0,00010	3,7190
0,00005	3,8906
0,000025	4,0556
0,000010	4,2649
0,000005	4,4172

Statistiska institutionen



Stockholms  
universitet

## Rättningsblad

**Datum:** 13/2-2020

**Sal:** Ugglevikssalen

**Tenta:** Undersökningsmetodik

**Kurs:** Regressionsanalys och undersökningsmetodik

**ANONYMKOD:**

0016-DU D

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.									
12	21	22	8	4					

POÄNG	BETYG	Lärarens sign.
67	D	ME



I  $N=290$  kommuner  
 $n=30$  kommuner

$$\sum Y_k = 661,73$$

$$\sum Y_k^2 = 14751,37$$

a) KI = 90%  $\Rightarrow z_{\frac{\alpha}{2}} = 1,6449$

$$\bar{Y} = \frac{\sum Y_k}{n} = \frac{661,73}{30} = 22,05766667 \approx 22,1 \text{ R}$$

$$s^2 = \frac{\sum (Y_k - \bar{Y})^2}{n-1} = \frac{\sum Y_k^2 - n\bar{Y}^2}{n-1} = \frac{14751,37 - 30 \cdot 22,05766667^2}{30-1} =$$

$$= 5,350008009 \text{ R}$$

$$D(Y) = \left(1 - \frac{n}{N}\right) \frac{s^2}{n} = \left(1 - \frac{30}{290}\right) \frac{5,350008009}{30} = 0,1598852968 \text{ R}$$

$$22,05766667 \pm 1,6449 \sqrt{0,1598852968}$$

$$= 22,05766667 \pm 0,6577241136 =$$

$$= [21,39994256; 22,71539078] \approx [21,4; 22,7] \text{ R}$$

Svaret Den genomsnittliga kommunalskatten  
 är 22,1 och ett 90% konfiansintervall  
 är [21,4; 22,7]

b) Svårt och med att det endast är en  
 skattning kan man inte vara  
 säker på att det svaret man  
 fick fram gäller för hela  
 Sverige. Man har ju endast  
 använt sig av 30 kommuner och  
 det är inte säkert att dessa  
 30 kommuner speglar hela landet.  
 För att ta fram en helt riktig  
 skattning bör man skatta med  
 hjälp av alla 290 kommuner.

och för sig rätt men inte  
 problemet här //

~~20,05 = 20,05  
 $1,6499 \cdot 10,1598852968 = 20,05$   
 D i uppgiften är  $10,1598852968$   
 då  $D = \sqrt{0,17}$   
 $\frac{20,05}{10,1598852968} = 1,9737241136$~~

Se sista bladet

8 + 1 + 3 = 12  
 ↗  
 sista bladet

SU, STATISTIK

Skrivsal: Aggleri KSSyan

Anonymkod: 0016-DVD

Blad nr: 2

$y$	$(N_k)$ Antal i populationen	$(n_k)$ Urvalsstorlek	$(x_k)$ Antal i urval
25	2510	25	7
26-35	5270	75	15
36-45	5310	75	39
46-55	4960	100	18
56-	4900	125	36
	22450	400	109

a)

$$\hat{P}_{str} = \sum W_k \cdot \hat{p}_k = \frac{2510}{22450} \cdot \frac{7}{25} + \frac{5270}{22450} \cdot \frac{15}{75} + \frac{5310}{22450} \cdot \frac{39}{75} + \frac{4960}{22450} \cdot \frac{18}{100} + \frac{4900}{22450} \cdot \frac{36}{125} = 0,0094721604 +$$

$$+ 0,0469487751 + 0,1229933185 + 0,0397687792 + 0,0564459343 = 0,270567929 \approx 0,27 \text{ ok}$$

$$D(P) = \sum W_k^2 \left(1 - \frac{n_k}{N_k}\right) \frac{\hat{p}_k(1 - \hat{p}_k)}{n_k - 1} = 0,1118040089^2 \left(1 - \frac{25}{2510}\right) \cdot$$

$$\cdot \frac{0,04(1-0,04)}{25-1} + 0,2347438753^2 \left(1 - \frac{75}{5270}\right) \cdot \frac{0,2(1-0,2)}{75-1}$$

$$+ 0,2365256125^2 \left(1 - \frac{75}{5310}\right) \cdot \frac{0,52(1-0,52)}{75-1} + 0,220935912^2 \cdot$$

$$\left(1 - \frac{100}{4960}\right) \cdot \frac{0,18(1-0,18)}{100-1} + 0,1959910913^2 \cdot \left(1 - \frac{125}{4900}\right) \cdot \frac{0,288(1-0,288)}{125-1}$$

$$= 4,585797367 \cdot 10^{-4} \text{ ok}$$

$X_k$ (M <sub>k</sub> /N)	$\hat{\phi}_k$ ( $\frac{1}{n} \cdot Y_k$ )
$\frac{2510}{22450} = 0,1118040089$	$\frac{1}{25} = 0,04$
$\frac{5270}{22450} = 0,234743973$	$\frac{15}{75} = 0,2$
$\frac{5310}{22450} = 0,236525625$	$\frac{39}{75} = 0,52$
$\frac{4960}{22450} = 0,220935412$	$\frac{18}{100} = 0,18$
$\frac{4400}{22450} = 0,1959910913$	$\frac{36}{125} = 0,288$

$KI = 95\% \Rightarrow z_{\alpha} = 1,96$

$$0,270567929 \pm 1,96 \sqrt{4,585297367 \cdot 10^{-4}} =$$

$$= 0,270567929 \pm 0,0411725 = 0,2294; 0,31196$$

Svar: Antalet med för högt blodtryck är ca 0,27 och ett 95% KI ger  $[0,2296; 0,3125]$

10



b) OSU

$$\hat{p} = \frac{1}{n} \cdot \sum y_k = \frac{109}{400} = 0,2725 \quad R$$

$$V(\hat{p}) = \left( \frac{N-n}{N-1} \right) \frac{p(1-p)}{n} = \left( \frac{22450-400}{22450-1} \right) \cdot \frac{0,2725(1-0,2725)}{400}$$

$$= 9,86800622 \cdot 10^{-4} \quad R$$

= variansen

$$D = \sqrt{9,86800622 \cdot 10^{-4}}$$

för OSU är felmarginalen  $\sqrt{9,86800622 \cdot 10^{-4}}$   
 och för stratifiering är den  
 $4,585797367 \cdot 10^{-4}$ ,  $D_{str}^2 < D_{osu}^2$  men rätt

svaret blir något bättre med  
 stratifiering då stratifierings felmarginal  
 är lite mindre.

$$\sqrt{\text{Varians}} = \text{standard fel} = D$$

$$\text{felmarginal} = Z_{\alpha/2} \cdot D = 1,96 \cdot \sqrt{\text{varians}}$$

7

$$C) n_k = n \cdot \frac{N_k \cdot \sigma_k}{\sum N_j \cdot \sigma_j}$$

$$\sigma_k = P(T-P)$$

$$0,09 \cdot 0,96 = 0,0384$$

$$0,2 \cdot 0,8 = 0,16$$

$$0,52 \cdot 0,48 = 0,2496$$

$$0,18 \cdot 0,82 = 0,1476$$

$$0,288 \cdot 0,712 = 0,205056$$

$$10 + 7 + 4$$

$$= 21$$

$$\sum N_j \cdot \sigma_j = 2510 \cdot 0,0384 + 5270 \cdot 0,16 + 5310 \cdot 0,2496$$

$$+ 4960 \cdot 0,1476 + 4500 \cdot 0,205056 = 3087,52 \text{ feli r'hand}$$

$$n_1 = 400 \cdot \frac{2510 \cdot 0,0384}{3087,52} = 12,48611506 \approx 12 \text{ pers}$$

$$n_2 = 400 \cdot \frac{5270 \cdot 0,16}{3087,52} = 109,2397782 \approx 109 \text{ pers}$$

$$n_3 = 400 \cdot \frac{5310 \cdot 0,2496}{3087,52} = 171,707593 \approx 172 \text{ pers}$$

$$n_4 = 400 \cdot \frac{4960 \cdot 0,1476}{3087,52} = 94,84583096 \approx 95 \text{ pers}$$

$$n_5 = 400 \cdot \frac{4500 \cdot 0,205056}{3087,52} = 116,8894647 \approx 117 \text{ pers}$$

Svari dan alternatifna st. k. prorsalokeringin / 4  
 dvi  $n_1 = 12 \text{ pers}$   $n_2 = 109 \text{ pers}$   $n_3 = 172 \text{ pers}$

$$n_4 = 95 \text{ pers}$$

7 w

$$\text{Summa} = 505, \text{ ej } 400$$

3)  $N=1132$  träd  $n=20$  träd

$$\mu_x = 10,3$$

$x_i$  = diameter i cm

$y_i$  = åltern mått i år

$$\sum x_k = 478$$

$$\sum y_k = 2148$$

$$\bar{x} = 23,9$$

$$\sum x_k^2 = 11832$$

$$\sum y_k^2 = 46304$$

$$\bar{y} = 107,4$$

$$\sum x_k y_k = 53315$$

a)  $\mu_y$  = genomsnittlig ålder  
diameter = självvariabel =  $x$

$$\hat{\mu}_{reg} = \bar{y} + b(\mu_x - \bar{x})$$

$$b = \frac{n \sum x_k y_k - (\sum x_k)(\sum y_k)}{n \sum x_k^2 - (\sum x_k)^2} = \frac{20 \cdot 53315 - 478 \cdot 2148}{20 \cdot 11832 - (478)^2} =$$

$$= \frac{39556}{8156} = 4,849926435 \text{ R}$$

$$\hat{\mu}_{reg} = 107,4 + 4,849926435 (10,3 - 23,9) = 41,941000 \text{ R}$$

(32 91)

Standard fel:

$$s(A_{\text{res}}) = \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \left( \frac{\sum y_k^2 - n\bar{y}^2}{n-2} - \frac{b^2 \left(\sum x_k^2 - n\bar{x}^2\right)}{n-2} \right) =$$

$$= 16,9175913 \quad \text{förklara lite bättre hur du får detta}$$

$$\sqrt{16,9175913} = 4,051856525 \approx 4,05 \quad \mathbb{R}$$

~~10~~

svari standard fel är 4,05  
och den ska ha de genomsnittliga  
åldern är ca 41 år

$$b) \hat{A}_{\text{svot}} = \hat{r} \cdot \mu_x = \frac{\sum y_k}{\sum x_k} \cdot \mu_x = \frac{2148}{478} \cdot 0,3 = 16,785556$$

( $\approx 46$ )  $\mathbb{R}$

$$s(A_{\text{svot}}) = \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot \left( \frac{\sum (y_k - \hat{r}x_k)^2}{n-1} \right)$$

$$\begin{aligned} \sum (y_k - \hat{r}x_k)^2 &= \sum y_k^2 - 2\hat{r} \sum x_k y_k + \hat{r}^2 \sum x_k^2 = \\ &= 246304 - 2 \cdot \frac{2148}{478} \cdot 53315 + \left(\frac{2148}{478}\right)^2 \cdot 11832 \end{aligned}$$

$$= 6068,357777 \quad \mathbb{R}$$

(se också ta 61a)

$$\sqrt{CA_{kvot}} = \left(9 - \frac{20}{1132}\right) \cdot \frac{1}{20} \cdot \frac{6068,3572777}{20-1} =$$

$$= 15,6872706 \quad \mathcal{R}$$

$$\sqrt{15,6872706} = 3,960916 \approx 3,96 \quad \mathcal{R}$$

Så var den genomsnittliga åldern är 46 år och standardfelet är 3,96. ~~10~~

C) Svaret är inte jätte stort skillnad mellan de olika svartningarna och inte heller jätte stort skillnad i standardfelet. Däremot är kvotkvantitetens standardfelet lite mindre och därför mest lämpligt. Som man ser i diagrammet är det utspritt både mellan ålder och mellan längd. Detta ger en hög varians och således ett högt standardfelet. ~~2~~

g) Teleskopeffekten är när en respondent flyttar en undersökande hänvisning så att den inte passar in i det intervall som är i en undersökning. Det kan till exempel vara en undersökning som handlar om läkarbesök under det senaste tre månaderna. Respondenten kanske då säger att hen haft en läkarbesök de senaste tre månaderna när det egentligen var 4-5 månader sedan. Teleskopeffekten brukar innebära motsatsen, alltså att en respondent flyttar sin hänvisning så att den inte gör det. Exempel kan vara att respondenten säger att hen ej varit hos läkaren senaste 3 månaderna trots att hen varit det. Problemet blir att man får ett intervall som ej stämmer överens med mått och värden på variabeln. Detta är ett täckningsproblem. Man får övertäckning vid teleskopeffekten och undertäckning vid teleskopeffekten brukar. /k

b) De ska kodas från 0 och uppåt. Alla vet ej skall kodas 0. Man ska tänka på att ha så tydlig alternativ som möjligt.

/1- se lösningsskiss.

g) c) SSYK-klassifikation av yrken  
 SSUN-klassifikation av utbildning <sup>R</sup>  
 Klassifikationsstandarder gör det  
 väldigt lätt att dela in kvantitet  
 i grupper för att se dra en  
 slutsats. *Jämför resultat i olika undersök-  
 och i tidsserier* / 3

d) Snövrval är när man helt  
 slumpmässigt väljer ut en/ flera  
 personer. Det kan t ex vara genom  
 att står på stan, om frågor  
 om någon vill va med i en  
 undersökning. Det kan vara bra  
 om man vill få in många  
 svarsalternativ. Det är ett  
 sannolikhets vrval. Då alla människor  
 har teoretiskt samma chans att  
 komma med. *Inte nödvändigtvis kända  
 de är därför ej sannolikhets vrval* <sup>ej kända</sup>  
*De beskriver bekvämlighets- / kvot vrval* / 0

(8)

5)  $N=6$  element

$$U = \{2, 6, 8, 14, 18, 24\}$$

$$S_1 = \{2, 14\}$$

$$S_2 = \{6, 18\}$$

$$S_3 = \{8, 24\}$$

a)  $P(2) = \frac{1}{6}$   $P(6) = \frac{1}{6}$   $P(8) = \frac{1}{6}$   $P(14) = \frac{1}{6}$   $P(18) = \frac{1}{6}$   $P(24) = \frac{1}{6}$  ✓

Svara ja, ja alla urvalen har lika stor chans att bli valda ska vara större än noll ochända

b)

urval

$\bar{y}$

$\sigma_y^2$

2, 14

$$\frac{2+14}{2} = 8$$

$$\frac{(16-8)^2 \cdot \frac{1}{3}}{3} = 21,333$$

6, 18

$$\frac{6+18}{2} = 12$$

$$\frac{(24-12)^2 \cdot \frac{1}{3}}{3} = 48$$

8, 24

$$\frac{8+24}{2} = 16$$

$$\frac{(32-16)^2 \cdot \frac{1}{3}}{3} = 85,333$$

R

?

Svara  $\sigma_y^2 = 21,333$   
48  
85,333

3) stickprovsmätvärdena = 8, 12, 16



$$c) \quad v(\bar{Y}) = \underbrace{\left(\frac{k-h}{k-1}\right)}_R \cdot \underbrace{\frac{\sigma^2}{n}}_R = \left(\frac{6-2}{6-1}\right) \cdot \frac{1(7-1)}{6} = 0,055556 \quad \checkmark$$

d) Den som ger läggst varians är OSU. Detta eftersom urvalstypen är ett klusterurval *Rätt felaktigt* och klusterurval har en hög varians, *men fel slutsats eftersom fel tidigare* inom klustret. Listan för elementen är dessutom redan organiserade *Hur?* och man skall dra ett urval av kluster. Det är även organiserade så att variansen mellan varannan ska bli så liten som möjligt. */2*

$$I) \quad z_{0,05} \cdot D^* = z_{0,025} \cdot D$$

*D = standard felot*

$$1,6449 \sqrt{0,1599852968} = z_{0,025} \cdot D$$

$$D = 0,1599852968$$

$$n \geq \frac{N \cdot \sigma^2}{D^2 \cdot (k-1) + \sigma^2}$$

$$\bar{y} = \frac{\sum y_k}{N} = \frac{661,73}{290} = 2,281827586$$

$$\sigma^2 = \frac{\sum y_k^2}{N} - n \cdot \bar{y}^2$$

$$2,281827586$$

$$= \frac{14751,37}{290} - 290 \cdot 2,281827586^2 = 45,6005597$$

*Fel användning av formeln*      *Använd sig fr. av helt enkelt*

$$n \approx \underline{290.45,66005597}$$

$$0,1594852968^2 \cdot (290 - 1) + 45,6600597$$

$$n \approx 299,7865164 \approx 250$$

svart stikprov är bär vara på  
250 kommuner.

orimligt stort?

~~3~~

A