



Stockholms universitet

OBS! Läs noga igenom anvisningarna i tentamen, t.ex. hur du ska skriva svaren. Det är ditt ansvar som student att följa de anvisningar som ges.

NOTE! Read the examination instructions carefully, e.g. how to write the answers. It is your responsibility as a student to follow the given instructions.

Skriv din anonymiseringskod och dagens datum på allt material du lämnar in.
(Enter your anonymization code and today's date on all submitted materials)

Anonymiseringskod (Anonymization code)	3	1	1	-	0	0	0	6	-	D	A	S
Datum (Date YYYY-MM-DD)	2021-10-29						Plats nr. (Seat No.)	32				

Kurs/Kurskod (Course/Course code)	ST306G
Kursmoment (Course component)	Survey sampling

Fylls i av tentamensvärd (To be filled in by invigilator)

Direkt i skrivning: (kryss)		Svarsblankett: (kryss)		Lösa svarsblad: (antal)	6
--------------------------------	--	---------------------------	--	----------------------------	---

Lämnat in blankt: (kryss)		Dator: (kryss)	
------------------------------	--	-------------------	--

Inlämningstid: 18:00 Signatur tentamensvärd: DR

Fylls i av lärare/examinator (To be filled in by teacher/examinator)

Betyg:	A	Poäng:	48
--------	---	--------	----

Signatur rättande lärare/examinator: DR

Regler i skrivsalen

- Följ tentamensvärds anvisningar.
- Väskor och ytterkläder ska placeras på anvisad plats.
- Placera ID-handling väl synlig på bordet framför dig.
- Ingen student får lämna skrivsalen under de första 30 minuterna.
- Endast en student i taget får besöka toaletten. Vid toalettbesök skriv ditt namn och klockslag på avsedd lista. Efter toalettbesöket ska du åter ange klockslag på listan.
- Elektronisk utrustning som mobiltelefon eller Smartwatch ska vara avstängd och placerad på anvisad plats.
- Under tentamen gäller tystnad – det är förbjudet att prata, eller på annat sätt kommunicera, med andra studenter under pågående tentamen.
- Innan tentamenshandlingarna lämnas in; skriv sidnummer, anonymiseringskod och datum på alla inlämnade papper.

Om något är oklart – fråga gärna tentamensvärden. Lycka till!

Rules in the examination hall

- Follow the invigilator's instructions.
- Bags and outerwear must be placed at the designated place.
- Place your ID document clearly visible on the table in front of you.
- No student may leave the examination hall for the first 30 minutes.
- Only one student at a time may visit the toilet. Before visiting the toilet, write your name and time on the intended list. After the toilet visit, enter the time on the list again.
- Electronic equipment such as a mobile phone or Smartwatch must be switched off and placed at the designated place.
- During the exam, silence applies – you are not allowed to talk, or otherwise communicate, with other students during the exam.
- Before submitting the examination documents; remember to write the page number, anonymization code, and date on all papers.

Please do not hesitate to ask the invigilator if anything is unclear. Good luck!

Uppg.nr.:
(Task no.)

Lärares
kommentar:
(Teacher's
note)

Poäng:
(Points)



a. Because yield is a variable of farms, like income is for persons, and we want to estimate total yield in country X of runner beans and other crops, our target population is all active farms in country X. This assuming only farms grow crops. ^{good point} If one were to be philosophical one could argue that the target population is the runner beans that have been grown in country X and that their weight is the variable of interest. But I think the first definition is much more practical and realistic.

Uppg.nr.: (Task no.)

6.

Lärarens kommentar: (Teacher's note)

2

b. Primary sampling units are the farms in the frame for method 2. Secondary sampling units are the plots.

1

c. For method 1, the observation units are the farms that were sampled from the frame.

For method 2, the observation units are the plots of land sampled at the farms. Here the weight of the crops would be a variable of the plot.

2

Poäng: (Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



Datum: (Date YYYY-MM-DD)	2021-10-29	Kurs/Kurskod: (Course/Course code)	ST3069	Sidnr.: (Page no.)	
Anonymiseringskod (Anonymization code)	3 1 1 - 0		0 0 6 - D A S		3

<p>d./ The first method would probably be subject to non-response, with farms with low yield <u>not bothering</u> to answer the questionnaire because they got worse things to deal with. This could introduce non-response bias.</p> <p>But answers could be pretty <u>precise</u> as farms must keep track of their yields at a very exact scale.</p> <p><i>Yes true. But the survey request may have come before they know</i></p> <p>The second method could be more imprecise as only <u>2</u> plots are sampled at each farm, making the estimate potentially unstable. Also the persons sent may not be as <u>experienced</u> as the farmers themselves at measuring yield, which may cause measurement bias.</p> <p>The second method could be seen as an <u>2-stage cluster sample</u> with only two SSV sampled, which is not great! I would probably choose method 1 and trust the farmers precision and write that they <u>gain no subsidies</u> by <u>underreporting</u> yield in the survey so as to minimize lying. I would adjust for non-response best I can after. And design the survey with non-response in mind.</p> <p>Turk! →</p>	<p>Uppg.nr.: (Task no.)</p> <p>60</p> <p>Lärarens kommentar: (Teacher's note)</p> <p>Poäng: (Points)</p> <p>5</p>
--	---

e. I would use weighting class estimator and use the indicator variable as the class determiner. I think this would be preferable to just using the regular HT estimator as we are bound to have some non-response. I would assume MAR and if it would turn out to be so the weighting class estimator would probably be an improvement over HT-estimator.

Uppg.nr.:
(Task no.)

6r

Lärens
kommentar:
(Teacher's
note)

2

Poäng:
(Points)

12



Datum: (Date YYYY-MM-DD)	2021-10-29	Kurs/Kurskod: (Course/Course code)	ST306G	Sidnr.: (Page no.)	4
Anonymiseringskod (Anonymization code)	3 1 1 - 0 0 0 6 - D A S				

7. A stratified simple random sample.
 Each section is a stratum and two SRS are drawn independently.
 Moreover it has a disproportional allocation.

Uppg.nr.: (Task no.) 7.
 Lärarens kommentar: (Teacher's note)

b. $y = \begin{cases} 1 & \text{if woman} \\ 0 & \text{otherwise} \end{cases}$ $\bar{y} = \hat{p}$ $N = 2500$

Survey association →

$$\hat{t}_{\text{survey}} = \sum y_i \cdot \frac{N_{\text{survey}}}{N_{\text{survey}}} \rightarrow 80 \cdot \frac{2000}{200} = \underline{800}$$

$$\hat{p}_{\text{survey}} = \frac{\sum y_i}{N_{\text{survey}}} \rightarrow \frac{80}{200} = \underline{0.4} \quad R$$

the cramer society →

$$\hat{t}_{\text{cramer}} = \sum y_i \frac{N_{\text{cramer}}}{N_{\text{cramer}}} \rightarrow 40 \cdot \frac{500}{200} = \underline{100}$$

$$\hat{p}_{\text{cramer}} = \frac{\sum y_i}{N_{\text{cramer}}} \rightarrow \frac{40}{200} = \underline{0.2} \quad R$$

The whole statistical society →

$$\hat{t}_{\text{str}} = \sum \hat{t} \rightarrow 800 + 100 = \underline{900}$$

$$\hat{p}_{\text{str}} = \frac{\hat{t}_{\text{str}}}{N} \rightarrow \frac{900}{2500} = \underline{0.36} \quad R$$

$$\hat{V}(\hat{p}_{\text{str}}) = \frac{1}{N^2} \sum N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_h^2}{n_h}, \quad S_h^2 \approx \hat{p}_h(1 - \hat{p}_h)$$

$$\hat{V}(\hat{p}_{\text{str}}) = \left(2000^2 \left(1 - \frac{200}{2000}\right) \frac{0.4(1-0.4)}{200}\right) + \left(500^2 \left(1 - \frac{200}{500}\right) \frac{0.2(1-0.2)}{200}\right) = \underline{0.0007104}$$

correct
 same mis calculation?

TURN! →

c. yes it is fair to say that Survey association has a greater proportion of women than the Cramer society.

Uppg.nr.:
(Task no.)

7.

Lärarens
kommentar:
(Teacher's
note)

If we look at the estimates we see great difference $\rightarrow \hat{p}_{\text{Survey}} = 0.4$

$$\hat{p}_{\text{Cramer}} = 0.2$$

lets find $\widehat{SE}(\hat{p}_{\text{diff}})$

$$\sqrt{\left(1 - \frac{200}{2000}\right) \frac{0.4(1-0.4)}{200} + \left(1 - \frac{200}{500}\right) \frac{0.2(1-0.2)}{200}} = 0.0394 \text{ R}$$

2

$$2 \cdot 0.0394 = 0.0788$$

The estimated differences CI does not contain 0 \rightarrow strong indication of our conclusion! *Good*

d. Domain estimation \rightarrow

$$\bar{y}_d = \frac{\sum y_i}{n} - \frac{\sum u_i}{n} \quad \text{where } x = \begin{cases} 1 & \text{if new} \\ 0 & \text{otherwise} \end{cases} \quad u = \begin{cases} y & \text{if new} \\ 0 & \text{otherwise} \end{cases}$$

mean new members $\rightarrow \frac{28000 + 35000 + 36000}{3} = 33000 \text{ R}$

variance $\rightarrow \widehat{V}(\bar{y}_d) = \left(1 - \frac{n}{N}\right) \frac{s_y^2 d}{nd}$ $s_y^2 d = \frac{20'000}{2} = 10'000$

$$\widehat{V}(\bar{y}_d) = \left(1 - \frac{200}{2000}\right) \frac{10'000}{3} = 3000 \text{ R}$$

\uparrow
because the survey association

Poäng:
(Points)

12



$y = 1$ Interested $\bar{y} = p, p = 0,3$
 0 not interested

margin of error = 0,02

a. sample size for 95% confidence interval and $e = 0,02$. Ignoring FPC \rightarrow

$$n = \frac{1,96^2 \cdot S_y^2}{e^2}, S_y^2 \approx p(1-p) = 0,3$$

$$n = \frac{1,96^2 \cdot 0,3(1-0,3)}{0,02^2}, n = 2016,84 \rightarrow n = \underline{2017}$$

b. max $S_y^2 \rightarrow \max p(1-p)$, where $p \in (0,1) \rightarrow p = 0,5 \rightarrow S_y^2 = 0,5(1-0,5) = \underline{0,25}$

Safe sample size \rightarrow

$$n = \frac{1,96^2 \cdot 0,25}{0,02^2}, n = \underline{2401}$$

c. Assumptions: MCAR \rightarrow The response rate and response propensity & does not depend on y, x or sampling design. The response set is seen as a SRS subsample of the sample set. We can therefore ignore non-response. This assumption is often not very realistic!

The sample mean is the estimated population mean $\rightarrow \bar{y}_s = \hat{p} = \underline{0,3}$ = estimate for all customers

I also assume that respondents tell the truth! No non-sampling errors!
 \uparrow
measurement error!

Turn! \rightarrow

for 50% response

do/ Possible reasons could be choice of mode, (email and telephone surveys can have

lower response rate compared to in-person), poor reputation

of agency (i.e. not a strong authority) and bad presentation of the survey and

question (for example no personalization with the name of the sample unit in the email or greeting on the phone). Likelihood can play a part in response rate.

Uppg.nr.:
(Task no.)

8.

Lärarens
kommentar:
(Teacher's
note)

3

Poäng:
(Points)

10



a. First estimator \rightarrow HT estimator (unbiased)

$$\hat{\tau}_y = \sum y_i \frac{N}{n} \rightarrow \bar{y} \cdot N \quad N=1000$$

$$\hat{\tau}_y = 5,38 \cdot \frac{1000}{10} = 538 \quad n=10$$

$$\hat{V}(\hat{\tau}_y) = N^2 \left(1 - \frac{n}{N}\right) \frac{S_y^2}{n}$$

$$\hat{V}(\hat{\tau}_y) = 1000^2 \left(1 - \frac{10}{1000}\right) \frac{0,08}{10} = 79200$$

Second estimator \rightarrow Regression estimator (approx. unbiased)

$$\hat{\tau}_{reg} = \hat{\tau}_y + \hat{\beta}_1 (t_x - \hat{\tau}_x), \quad \hat{\beta}_1 = \frac{S_{yx}}{S_x^2} = \frac{\sum (y_k - \bar{y})(x_k - \bar{x})}{\sum (x_k - \bar{x})^2}$$

$$\hat{\beta}_1 = \frac{0,13}{0,54} = 0,2407$$

$$t_x = 447, \quad \hat{\tau}_x = 4,47 \cdot \frac{1000}{10} = 447$$

$$\hat{\tau}_{reg} = 538 + (0,2407)(447 - 447) = 538$$

$$\hat{V}(\hat{\tau}_{reg}) = N^2 \left(1 - \frac{n}{N}\right) \frac{S_y^2}{n} (1 - R^2), \quad R^2 = \left(\frac{S_{yx}}{S_y S_x}\right)^2$$

$$\hat{V}(\hat{\tau}_{reg}) = 1000^2 \left(1 - \frac{10}{1000}\right) \frac{0,08}{10} (1 - 0,216^2) = 7550,4844$$

b. Looking at the graph, the y values from the new sample looks about the same, so the HT estimator's estimate would not change much. The x values are vastly different although! We seem to have gotten a "bad sample" if considering the x values, compared to the first sample when the estimator was exactly right (for x). So the regression correction term would in this case be different (by quite a lot)! This would result in a much different estimate for $\hat{\tau}_y$ when using regression estimator. I would rank HT estimator better than regression estimator in this case.

b₀ cont.

Uppg.nr.:
(Task no.)

9.

I would rank HT better in this case as we now have two samples of y values giving similar HT estimates so we are getting good information that $\hat{E}y$ is around 530 or so. To choose the regression estimator in the second case would be to let the "bad sample" of X have too much weight in my opinion!

Lärarens
kommentar:
(Teacher's
note)

Not sure I would agree, but I accept your line of argument.

Poäng:
(Points)

6