

TENTAMEN I INTRODUKTION TILL STATISTIK FÖR
STATSVETARE
2017-04-27

Skrivtid: 10.00-15.00

Godkända hjälpmedel: Miniräknare.

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

Uppgift 1 (20 poäng)

Det totala antalet anmälningar som inkommit till Skolinspektionen under åren 2012-2016 gällande brister i skolhälsan visas i tabellen nedan.

År	2012	2013	2014	2015	2016
Antal anmälningar	8	10	23	13	11

- Vad har variabeln "Antal anmälningar" för variabeltyp och datanivå?
- Räkna ut medelvärdet och medianen för antal anmälningar under den angivna tidsperioden. Förklara eventuella skillnader mellan de två lägesmått.
- För att räkna ut variansen för antal anmälningar under den angivna tidsperioden använder vi formeln $\sigma^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}$ istället för formeln $S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$. Förklara varför.
- Skapa en indexserie av variabeln "Antal anmälningar" med 2013 som basår. Förklara hur indexserien tolkas.

Uppgift 2 (20 poäng)

Vattenverket i en kommun vill undersöka halten av järn i dricksvattnet. Man drar ett slumpmässigt urval om 40 hushåll och beräknar för dessa en genomsnittlig järnhalt på 0.16 mg/l (milligram per liter) och en standardavvikelse på 0.09 mg/l.

- Beräkna ett 90%-igt konfidensintervall för den genomsnittliga järnhalten i kommunens dricksvatten.
- Förklara hur ett 90%-igt konfidensintervall ska tolkas.
- Om man vill att den totala bredden på det 90%-iga konfidensintervallet ska vara som mest 0.03, hur stort stickprov måste man då minst ta?

Uppgift 3 (20 poäng)

Vi går tillbaka till uppgift 2 och använder återigen insamlade data från de 40 hushållen. Livsmedelsverkets riktlinjer är att dricksvattnet är tjänligt utan anmärkning om järnhalten understiger 0.2 mg/l.

a) Utför ett hypotestest för att undersöka om kommunens dricksvatten är tjänligt utan anmärkning. Använd signifikansnivån $\alpha = 0.01$.

I ett hypotestest anges oftast p-värdet. Antag att p-värdet för ett test blev 0.02.

b) Förklara vad som menas med ett p-värde och hur vi använder det för att besluta om vi förkastar eller inte förkastar en nollhypotes för olika signifikansnivåer α .

Uppgift 4 (20 poäng)

En småföretagare gjorde följande sammanställning över ålder och sjukfrånvarodagar under ett år för sina 6 anställda.

Namn	A-son	B-son	C-son	D-son	E-son	F-holm
Ålder	20	31	35	44	52	60
Sjukfrånvaro (dagar)	16	7	18	23	20	28

a) Anpassa regressionslinjen $\hat{y} = bx + a$ d.v.s. skatta a och b då $y =$ antal sjukfrånvarodagar.

b) Tolka den skattade koefficienten b i termer av ålder och sjukfrånvaro.

c) Kan man göra en meningsfull tolkning av koefficienten a ? Förklara varför eller varför inte.

d) Korrelationskoefficienten $r = 0.74$. Förklara begreppet korrelation generellt (med hjälp av en eller flera ritade bilder) samt ge en tolkning av det givna värdet på korrelationskoefficienten i detta fall.

Uppgift 5 (20 poäng)

a) Förklara begreppet determinationskoefficient

b) Förklara skillnaden mellan en diskret och en kontinuerlig fördelning. Rita ett exempel på vardera fördelning.

c) Förklara skillnaden mellan korrelation och kausalitet.

d) En tidsseriemodell kan tänkas bestå av 4 komponenter. Ange dessa fyra komponenter och förklara vad respektive komponent innebär.

e) Förklara begreppet "bias i en skattning" och vad det kan bero på.

Lycka till!

FORMELBLAD

DESKRIPTIV STATISTIK

Ett urval består av n stycken observationer.

Medelvärde:

$$\bar{x} = \frac{\sum x_i}{n}$$

Medelvärde från frekvenstabell:

$$\bar{x} = \frac{\sum f_i x_i}{n}$$

Vägt medelvärde:

$$\bar{x} = \frac{\sum_{i=1}^l w_i \bar{x}_i}{\sum_{i=1}^l w_i}$$

Kvartiler:

$$q_1 = \frac{n+1}{4} \qquad q_3 = \frac{3(n+1)}{4}$$

Varians:

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = \frac{\sum x_i^2 - n\bar{x}^2}{n-1}$$

Varians från frekvenstabell:

$$s^2 = \frac{\sum f_i (x_i - \bar{x})^2}{n-1} = \frac{\sum f_i x_i^2 - n\bar{x}^2}{n-1}$$

VÄNTEVÄRDE OCH VARIANS

Väntevärde för X :

$$\mu = E(X) = \sum x f(x)$$

Varians för X :

$$\sigma^2 = V(X) = \sum [x - E(X)]^2 f(x) = \sum x^2 f(x) - [E(X)]^2$$

SAMPLINGFÖRDELNINGAR OCH CENTRALA GRÄNSVÄRDESSATSEN

Om populationen är normalfördelad med väntevärde μ och varians σ^2 dvs $x \sim N(\mu, \sigma^2)$

så är $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$

Om populationen har en annan fördelning (vilken som helst) med väntevärde μ och varians σ^2

så är $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$ om $n \geq 30$ enligt centrala gränsvärdeessatsen

STATISTISK INFERENS (normalfördelade data eller $n \geq 30$)

Konfidensintervall:

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad \text{alternativt} \quad \bar{x} \pm z_{\alpha/2} \frac{S}{\sqrt{n}}$$

$$P \pm z_{\alpha/2} \sqrt{\frac{P(1-P)}{n}}$$

Antal observationer för en given bredd $2B$ på konfidensintervallet:

$$n = \frac{Z_{\alpha/2}^2 \sigma^2}{B^2}$$

Hypotesprövning:

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \quad \text{alternativt} \quad Z = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

$$Z = \frac{P' - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}}$$

Tabellvärden för standardiserad normalfördelning:

α	Z_{α}
0.005	2.58
0.01	2.33
0.025	1.96
0.05	1.64
0.1	1.28

REGRESSION

Skattning av regressionslinjen $\hat{y} = bx + a$

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$a = \bar{y} - b\bar{x}$$

$$\hat{y} = bx + a$$

Residualvarians:

$$s_e^2 = \frac{\sum e^2}{n-2} = \frac{\sum (y - \hat{y})^2}{n-2}$$

Residualspridning:

$$s_e = \sqrt{\frac{\sum e^2}{n-2}} = \sqrt{\frac{\sum (y - \hat{y})^2}{n-2}}$$

Korrelation:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{[n \sum x_i^2 - (\sum x_i)^2] [n \sum y_i^2 - (\sum y_i)^2]}}$$

Determinationskoefficient (enkel linjär regression):

$$R^2 = r^2$$

INDEX

$$I_t = \frac{I_t}{I_b} 100$$

$$I_t = \frac{I_t}{I_b} 100$$

Laspeyres index:

Paasches index:

$$P_t^L = \frac{\sum p_t q_t}{\sum p_0 q_0} \cdot 100$$

$$P_t^P = \frac{\sum p_t q_t}{\sum p_0 q_t} \cdot 100$$



Stockholms
universitet

Statistiska institutionen

Rättningsblad

Datum: 27/4-2017

Sal: Brunnsvikssalen *Vaglevikssalen*

Tenta: Statistik för statsvetare

Kurs: Introduktion till statistik för statsvetare

ANONYMKOD:

SFS-0013

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					6
Lär.ant. 20	17	17	19	18,5					

POÄNG 91,5	BETYG A	Lärarens sign. JF
---------------	------------	----------------------

1. a) Kvantitativ, diskret variabel, då den är siffermässig och endast antar heltal.

Datanivån är en kvotnivå. Detta då avståndet mellan värdena på skalan går att mäta, samt att det finns en given nollpunkt. 4 anmälningar är dubbelt så mycket som 2 anmälningar osv.

5

b) Medelvärdet: Det genomsnittliga värdet i datamaterialet.

$$\bar{x} = \frac{\sum x_i}{n} = \frac{8 + 10 + 23 + 13 + 11}{5} = 13$$

Median: Det mittersta värdet i ett storleksordnat datamaterial.

8, 10, 11, 13, 23

Pågrund av värdet 23 vilket är något extremt i relation till övriga observationerna blir medelvärdet snedvridet och således lägre än medianen som inte alls tar hänsyn till extremer åt något håll.

5

c) σ^2 = Verkliga variansen i hela populationen (Parameter).

s^2 = Variansen i ett urval. (variabel).

Då vi här utgår från data över det totala antalet anmälningar, räknar vi därför ut variansen för en totalundersökning σ^2 , och inte för en urvalsundersökning, s^2 .

d)	t (År)	X (antal anmälningar)	Index med 2013 som basår
	2012	8	80
	2013	10	100
	2014	23	230
	2015	13	130
	2016	11	110

$$I_t = \frac{X_t}{X_b} \cdot 100$$

$$I_{2012} = \frac{8}{10} \cdot 100 = 80$$

$$I_{2014} = \frac{23}{10} \cdot 100 = 230$$

$$I_{2015} = \frac{13}{10} \cdot 100 = 130$$

$$I_{2016} = \frac{11}{10} \cdot 100 = 110$$

Indexserien visar med hur mycket antalet anmälningarna skiljer sig för åren 2012, 2014, 2015 och 2016 i relation till anmälningarna vid basåret 2013. Man kan t.ex. utläsa att det kom in 10% fler anmälningar år 2016 i relation till antalet 2013.

2.

$$n = 40 \quad (40 \text{ hushåll})$$

$$\bar{x} = 0,16 \quad (\text{mg/liter})$$

$$s = 0,09 \quad (\text{mg/liter})$$

a) 90%-igt konfidensintervall.

$$\text{Ger } \alpha = 0,1$$

Vi tar $\frac{\alpha}{2}$ och tar ut motsvarande värde på Z_{α} ur tabellen över standardiserad normalfördelning.

$$\frac{0,1}{2} = 0,05.$$

$$\alpha: 0,05 \quad \text{ger } Z_{\alpha}: 1,64$$

Sätter in värdena i formeln: $\bar{x} \pm Z_{\alpha/2} \frac{s}{\sqrt{n}}$

$$0,16 \pm 1,64 \cdot \frac{0,09}{\sqrt{40}}$$

$$0,16 \pm 0,023$$

Ger intervallet =

$$(0,137; 0,183)$$

b) Den verkliga populationsparametern, i detta fall verkliga medelvärdet, μ , befinner sig med 90% sannolikhet i intervallet (0,137 - 0,183.)

c) Den totala bredden på konfidensintervall = $2B$.

B är således halva bredden.

Da $2B$ är givet till $0,03$, blir $B = \frac{0,03}{2} = 0,015$

Vi utgår från samma Z_{α} och S men använder oss nu av formeln för att välja ut n vid den givna bredden.

$$n = \frac{Z_{\alpha/2}^2 \cdot \left(\frac{S}{n}\right)^2}{B^2}$$

Da vi g har standard avvikelsen för hela populationen tar vi standardavvikelsen för vårt urval, S .

$$n = \frac{1,64^2 \cdot 0,09^2}{0,015^2} = 96,8256, \text{ avrundar uppåt } \approx 97$$

Tolkning: Desto större stickprov, alltså desto större antal n , desto högre precision i intervall! Tar vi 97 st i vårt urval istället för 40 st, minskar längden på intervall. Det blir lättare att se mer precist var det verkliga värdet befinner sig. För en total bredd på $0,03$ måste vi ha med minst 97 observationer i urvalet.

13

3. a) Nullhypotes: $H_0: \mu_0 = 0,2$

Alternativhypotes: $H_A: \mu < 0,2$ - vänstersidigt test

Sign $\alpha = 0,01$ ger enligt tabellen $z_\alpha = 2,33$

Använder följande formel för att räkna ut vår teststatistika:

$$Z = \frac{\bar{X} - \mu_0}{S / \sqrt{n}}$$

Använd tidigare angivna värden för \bar{X} , n och S .

$$\bar{X} = 0,16$$

$$n = 40$$

$$S = 0,09$$

$$Z = \frac{(0,16 - 0,2)}{(0,09 / \sqrt{40})} = \frac{-0,04}{0,0142302499} = -2,81$$

Då $-2,81$ är mindre än $-2,33$ i detta vänstersidiga test så kan vi förkasta nullhypotesen. Dricksvattnet är tjäntigt utan anmärkningar på signifikansnivån $\alpha = 0,01$.

K 10

b) P-värdet säger oss vad sannolikheten är att få vårt värde på Z_{obs} eller något ännu mer extremt. I detta fall är sannolikheten 2% då p -värdet = 0,02.

Är p -värdet större än signifikansnivån bör nullhypotesen förkastas, är värdet mindre bör nullhypotesen inte förkastas utan accepteras.

p -värdet används som komplement då olika hypoteser av samma datamaterial kan få olika resultat beroende på vald signifikansnivå.

4. a)	Alder	z	Sjukfrånvaro dagar		$x_i \cdot y_i$
	x_i		y_i	y_i^2	
	20	400	16	256	320
	31	961	7	49	.
	35	1225	18	.	.
	44	1936	23	.	.
	52	2704	20	.	.
	60	3600	28	.	.

$$\sum x_i = 242 \quad \sum x_i^2 = 10826 \quad \sum y_i = 112 \quad \sum y_i^2 = 2392 \quad \sum x_i \cdot y_i = 4899$$

$$n = 6$$

$$b = \frac{n \cdot (\sum x_i \cdot y_i) - (\sum x_i) \cdot (\sum y_i)}{n \cdot (\sum x_i^2) - (\sum x_i)^2}$$

$$b = \frac{6 \cdot 4899 - 242 \cdot 112}{6 \cdot 10826 - 242^2} = \frac{2290}{6392} = 0,358 \approx 0,36$$

Fortsättning 4a) →

$$a = \bar{y} - b \cdot \bar{x}$$

$$a = \frac{112}{6} - 0,36 \cdot \frac{292}{6} = \underline{4,15} \quad \mathcal{R}$$

$$\hat{y} = bx + a \quad \text{— Regressionslinje}$$

$$\hat{y} = 0,36x + 4,15$$

 \mathcal{R}

5

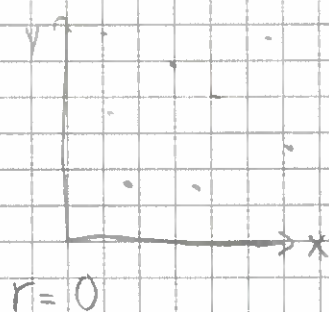
- b) b = riktningskoefficient. Bestämmer om ett samband är positivt eller negativt, samt hur stor lutningen är. b säger i detta fall att desto äldre en person är desto högre blir sjukfrånvaron, då sambandet är positivt. Lutningen är dock relativt liten vilket innebär att det inte är någon drastisk skillnad i sjukfrånvaro mellan åldrarna. Hade b varit = 1, hade det inneburit att +1 år i ålder leder till +1 sjukdag, mer/år. Är man 20 år hade sjukfrånvaron/år blivit 29,15 dagar, 21 år inneburit 29,51 dagar osv.

5

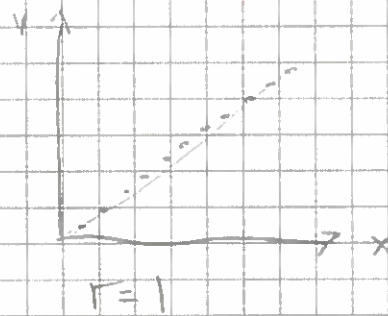
c) $a = \text{intercept}$. Bestämmer var linjen skär y-axeln.
 Negativt värde = under origo, positivt värde = över origo.
 a används för att veta var linjen ska ritas in. Är
 $a = 4,15$ ska linjen skära y-axeln vid $y = +4,15$.
 Man måste dock tänka kritiskt kring denna koefficient.
 Om vi hypotetiskt har åldern $x = 0$ hade vi utifrån vår
 formel [$y = 0,36 \cdot 0 + 4,15$] fått $y = 4,15$. Men någon
 som är 0 år kan omöjligt ha 4 sjukfrånvarodagar.
 Formeln måste sättas in i rätt sammanhang för att
 kunna göra en meningsfull tolkning.

5

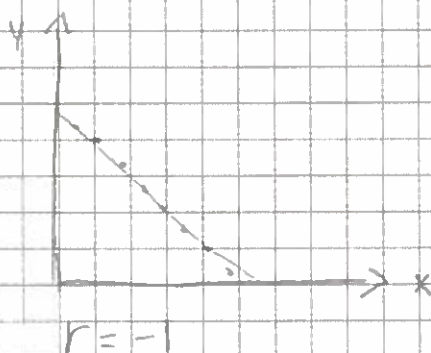
d) korrelation = hur starkt det ^{linjära} sambandet är.



Inget samband



Perfekt positivt samband



Perfekt negativt samband.

Då $r = 0,74$ innebär detta
 ett ganska starkt positivt
 samband mellan x och y .

4

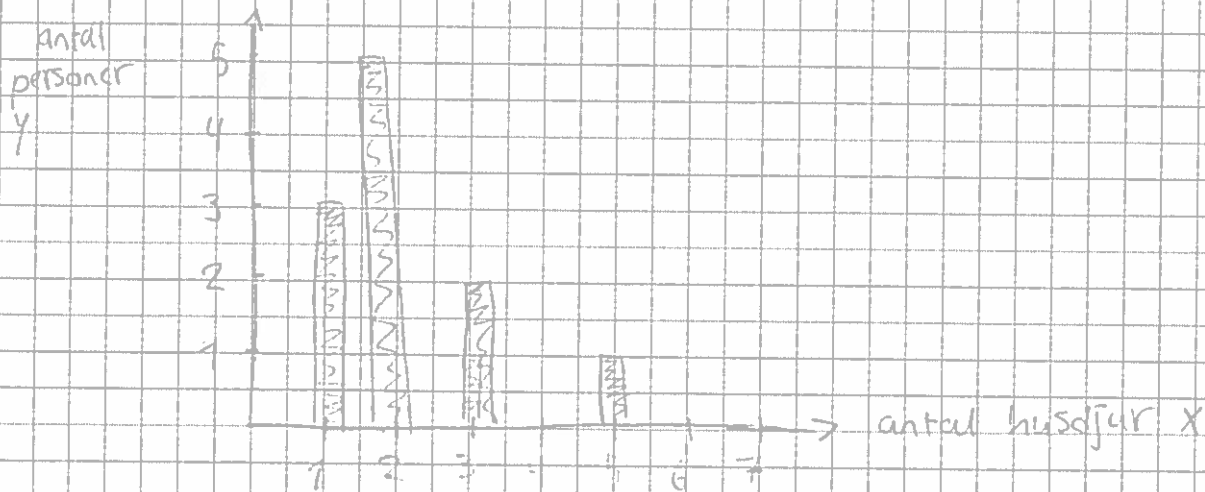
5. a) Determinationskoefficient = Förklaringsgrad. I hur stor grad förändringen y beror på förändringen i x i den skattade regressionslinjen. Skrivs som R^2 , alltså korrelationen r^2 . Kan anta värden mellan 0-1.

Är $R^2 = 0,9$ innebär detta förändringen i y till 90% beror på förändringen i x . Övriga 10% kan bero på naturlig variation och andra förklaringsvariabler.

4

b) Diskreta fördelningar:

Variabeln kan endast anta vissa värden, ofta heltal. Kan t.ex. illustreras genom ett stapeldiagram, där y står för frekvensen (antal personer) och x står för antal husdjur/per person.



5. b) - fortsättning →

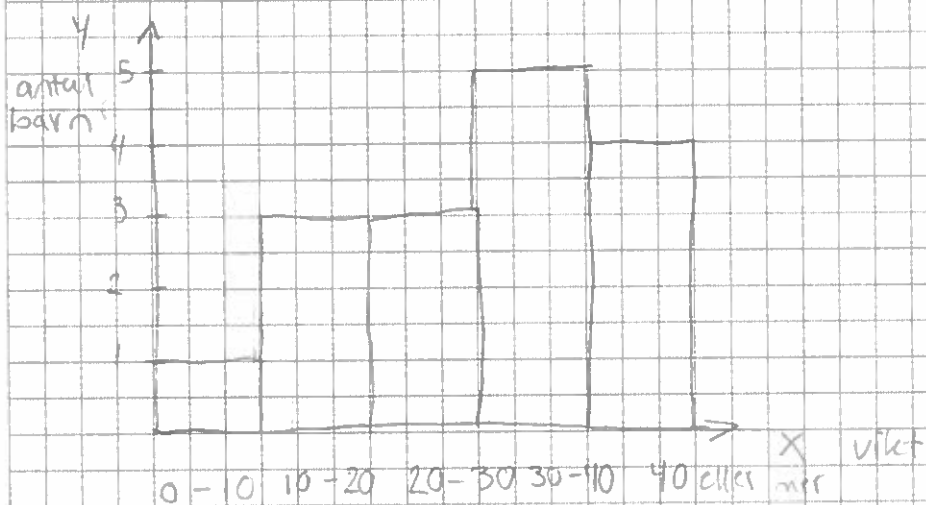
Kontinuerliga fördelningar =

Variabeln kan anta alla tänkbara värden i ett givet

intervall. Kan illustreras genom ett histogram är y

stär för frekvensen (antal barn) och x för vikt

hos barnen.



4

c) Korrelation står för hurvida ett samband finns eller inte och hur starkt detta samband i så fall är. Kausalitet är ett orsakssamband där förändringen i en variabel direkt leder till förändring i en annan variabel, alltså orsak-verkan. Kausalitet sker i en riktning: A leder till B men B leder inte till A. T.ex kan bra väder leda till ökad glassförsäljning, med detta innebär inte att ökad glassförsäljning leder till bra väder.

Det kan finnas en korrelation mellan två variabler utan att den ene orsakas av den andra. →

5 c) - fortsättning →

T.ex kan det finnas en korrelation mellan drunkningsolyckor och glassförsäljning. Men ökat antal drunkningsolyckor orsakar inte ökad glassförsäljning. Istället orsakas dessa båda mer troligt av bra väder, vilket blir de kausala sambanden.

4

d) Tidsseriemodell med 4 komponenter:

1. Trend = Tidsserien följer en trend, alltså de nya värdena baseras på hur utvecklingen hittills sett ut.
2. Slump = värdena i tidsserien beror på slumpen, kan ej kontrolleras, naturlig variation.
3. ~~Omvarter~~ konjunkturer = värdena påverkas av yttre förhållanden. Konjunkturer, ekonomiska, sociala förhållanden osv.

25

e) I totalfelet ingår både urvalsfelet och systematiska fel. Bias är ett annat ord för de systematiska felet. Alltså fel som inte beror på skillnader i skattade värden för urvalet och verkliga värden i populationen orsakat av naturlig variation (urvalsfel).

Bias kan bero på en rad feltyper.

Täckningsfel = Fel målpopulation täcks in i rimen.

Övertäckning då för många inklusive icke-relevanta täcks in. Undertäckning då individer/ objekt från

målparen inte alls täcks in i ramparen.

Bortfallsfel = Fel då undersökningen saknar data/svar från individer/objekt som ingår i urvalet.

Måtfel = Fel på t.ex mätinstrumentet. En trasig väg ger fel resultat.

Bearbetningsfel = Fel under bearbetningsprocessen. Då t.ex data ska kodas och göras om till digitalt, eller räknas om till önskat mått.

4