

**Skriftlig tentamen** avseende kursen **Introduktion till statistik för statsvetare (ST131G)**

2013-06-05

\*\*\*\*\*

Skrivtid: 5 timmar.

Hjälpmedel: Miniräknare, samt vidhäftat formelblad.

Genomgång: 2013-06-19 klockan 14.00 i sal B307.

Obs! Datum, tid och lokal !!!

\*\*\*\*\*

Tentamen består av fem uppgifter vilka kan ge maximalt fyra poäng vardera, totalt tjugo poäng. Då en uppgift i sin tur består av två eller flera deluppgifter värderas dessa lika. För full poäng på en uppgift/deluppgift krävs att tydliga, fullständiga och välmotiverade lösningar samt svar inlämnas. Lycka till !!! / Peter Claësson

\*\*\*\*\*

1.

Ett urval av studenter i ett visst studentbostadsområde i närheten av ett svenskt lärosäte tillfrågades om hur länge var och en hade köat innan de fick tillgång till sin lägenhet.

Kötid (år)	0	1 - 2	3 - 4
Antal studenter	12	20	8

- Ange vad som ovan är variabel, vilken variabeltyp och vilken datanivå som föreligger samt vad som är frekvens. Motivera.
- Illustrera fördelningen ovan i ett lämpligt diagram samt motivera valet av diagram.
- Såväl medelvärde som median beräknades för fördelningen ovan, varvid man erhöll 1,95 respektive 1,80. Hur ska dessa två värden tolkas?
- Bland de utvalda studenterna visade det sig att 40% var män. Illustrera även denna fördelning i ett lämpligt diagram, samt motivera valet av diagram.

2.

De faktiska utgifterna av ett särskilt slag vid en stor statlig myndighet har sammanställts för ett antal år. I tabellen nedan ges även konsumentprisindex (KPI) för respektive år.

År	1990	2000	2005	2010
Utgifter (Mkr)	48	40	32	12
KPI	208	260	280	303

- Beräkna en indexserie i löpande priser (basår 2000) avseende de aktuella utgifterna.
- Beräkna därefter en motsvarande indexserie i fasta priser (2000 års prisnivå).
- Illustrera även de två ovan beräknade indexserierna i ett lämpligt diagram. Motivera.
- Hur stor är den procentuella utgiftsförändringen mellan år 1990 och år 2010?

3.

För ett antal år sedan gjordes en studie avseende priset på fritidshus. Inom ett visst avgränsat geografiskt område valdes slumpmässigt tio fritidshus, samtliga sålda under året ifråga. Förutom att försäljningspriset studerades för vart och ett av de utvalda fritidshusen noterades även avståndet till kusten. Avsikten var att för aktuella variabler genomföra en regressions- och korrelationsanalys, där man angav försäljningspriset i tusentals kronor och avståndet till kusten i kilometer.

Då en regressionslinje anpassades erhöles följande ekvation:  $\hat{y} = 1234,0 - 56,0x$  (avrundade värden).

- Efter nödvändiga beräkningar fick man att korrelationskoefficienten blev  $-0,7$ . Hur tolkas detta värde?
- Rita in linjen i ett lämpligt diagram samt skissera hur sambandet mellan variablerna bör ha sett ut.
- Tolka de erhållna koefficientvärdena i linjens ekvation, uttryckta i termer av de aktuella variablerna.
- Ge en prediktion med hjälp av linjen för ett fritidshus på en mils avstånd från kusten, samt kommentera.

4.

Antag att man inom en mycket stor svensk kommun önskar genomföra en omfattande urvalsundersökning bland samtliga invånare (18 år eller äldre), med syfte att ta reda på inställningen till ett planerat vägprojekt.

Precisera vad som här kan anses vara mål- respektive rampopulation. Redogör därefter för, samt ge även exempel på, nedanstående åtta statistiska begrepp/metoder - utifrån det givna antagandet ovan.

- |                       |                        |
|-----------------------|------------------------|
| a) Undertäckningsfel, | e) Systematiskt urval, |
| b) Övertäckningsfel,  | f) Stratifierat urval, |
| c) Mätfel,            | g) Kvoturval,          |
| d) Bortfallsfel,      | h) Urvalsfraktion.     |

5.

Redogör för de åtta statistiska begreppen/metoderna nedan, samt ge realistiska exempel på respektive.

- |                      |                        |
|----------------------|------------------------|
| a) Nominaldata,      | e) Nonsenskorrelation, |
| b) Diskret variabel, | f) Fastbasindex,       |
| c) Frekvenspolygon,  | g) Bias,               |
| d) Tratt-teknik,     | h) Konfidensintervall. |

## 1. Beskrivande statistik

### 1.1 Medelvärde, varians, standardavvikelse

Ett statistiskt material består av  $n$  observationer

$$x_1, x_2, \dots, x_n$$

Medelvärdet är

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum x}{n} \quad (1.1.1)$$

Variansen är

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1} = \frac{\sum x^2 - (\sum x)^2 / n}{n-1} \quad (1.1.2)$$

Standardavvikelsen är

$$s = \sqrt{s^2} \quad (1.1.3)$$

När materialet redovisas i en frekvenstabell, där värdet  $x_i$  förekommer med frekvensen  $f_i$ , är medelvärdet och variansen

$$\bar{x} = \frac{\sum f_i x_i}{n} \quad (1.1.4)$$

$$s^2 = \frac{\sum f_i (x_i - \bar{x})^2}{n-1} = \frac{\sum f_i x_i^2 - (\sum f_i x_i)^2 / n}{n-1} \quad (1.1.5)$$

Räkne regler

Om  $y = a + bx$ , där  $a$  och  $b$  är konstanter, är

$$\bar{y} = a + b\bar{x} \quad (1.1.6)$$

$$s_y^2 = b^2 s_x^2 \quad (1.1.7)$$

### 1.2 Regression, korrelation

Regressionslinjen är  $y = a + bx$ .

Regressionskoefficienten är

$$b = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{\sum xy - \sum x \sum y / n}{\sum x^2 - (\sum x)^2 / n} \quad (1.2.1)$$

$$a = \bar{y} - b\bar{x} \quad (1.2.2)$$

Korrelationskoefficienten är

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}} = \frac{\sum xy - \sum x \sum y / n}{\sqrt{[\sum x^2 - (\sum x)^2 / n][\sum y^2 - (\sum y)^2 / n]}} = b \frac{s_x}{s_y} \quad (1.2.3)$$

Residualvariansen är

$$s_e^2 = \frac{n-1}{n-2} s_y^2 (1-r^2) = \frac{1}{n-2} (\sum y^2 - a \sum y - b \sum xy) \quad (1.2.4)$$

### 1.3 Prisindex

Laspeyres index är

$$\frac{\sum p_t q_0}{\sum p_0 q_0} \cdot 100 \quad (1.3.1)$$

Paasches index är

$$\frac{\sum p_t q_t}{\sum p_0 q_t} \cdot 100 \quad (1.3.2)$$

Statistiska institutionen



Stockholms  
universitet

## Rättningsblad

**Datum:** 5/6-2013

**Sal:** Brunnsvikssalen

**Tenta:** Statistik för statsvetare

**Kurs:** Introduktion till statistik för statsvetare

**ANONYMKOD:**



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. 3.5	2.75	3	2.25	1					

POÄNG	BETYG	Lärarens sign.
12.5	D	

1.

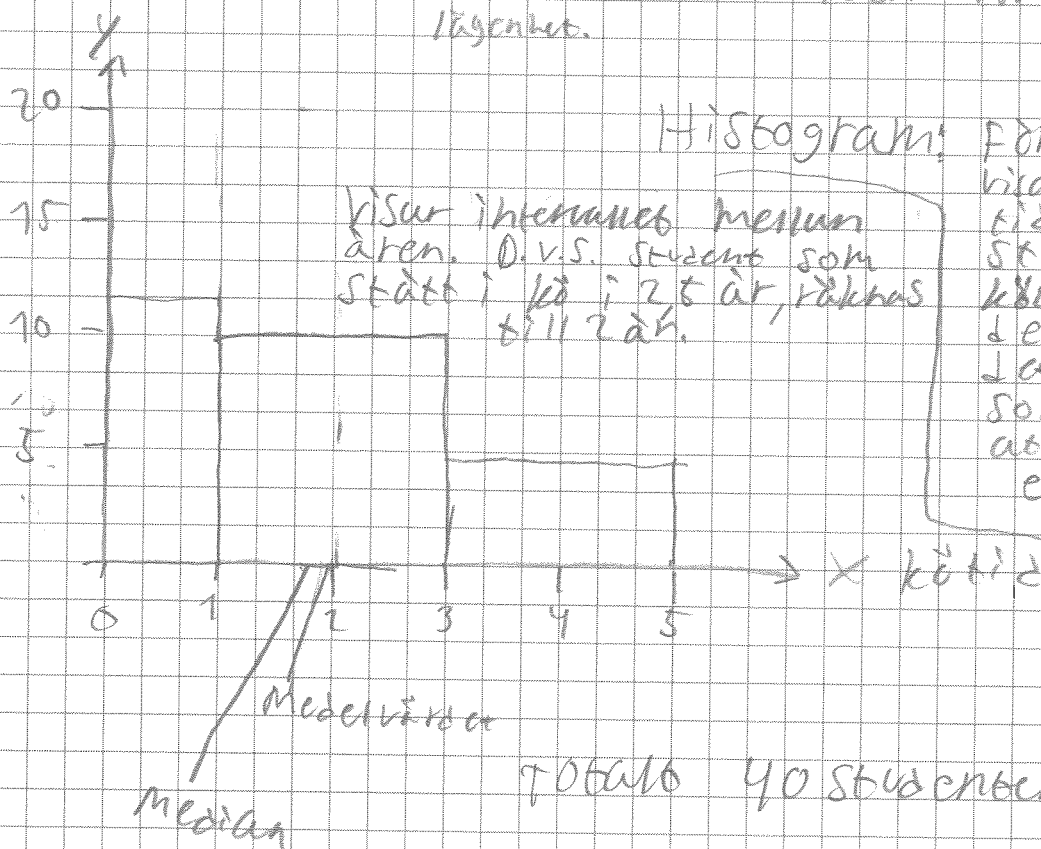
a)	kötid (år)	0	1-2	3-4
	Antal studenter	12	20	8

a) kötiden är variabeln. Den är en kvantitativ diskret variabel. Datatypen är kvantitativ eftersom kötiden har en absolut nollpunkt. Kötiden är bakgrundsvariabeln till att få en lägenhet. Studenterna är alltså beroende av kötiden.

0.75

b)

Antal studenter är frekvensen eftersom studenterna är beroende av kötiden för att få lägenhet.



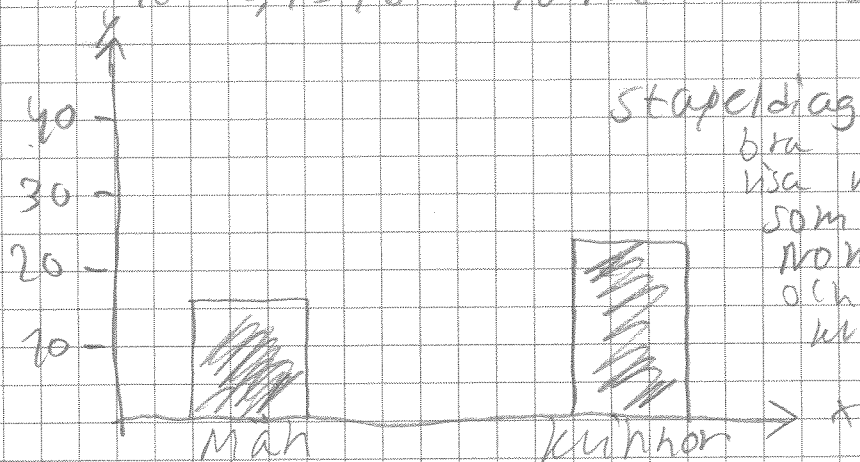
c medelvärden är tolkat som medelkötid. d.v.s. snitttiden har köat i 1.95 år. Det är beräknat genom att den totala summan har dividerats med antalet observationer.

Medianen är det mittre värde i observation. d.v.s. En sammansättning av alla studenter där det mittre värde väljs ut av de totalt 40 studenterna.

0.75

d Av de 40 studenterna var 40% män

$$40 \cdot 0,4 = 16 \quad 16 \text{ män och } 24 \text{ kvinnor.}$$



Stackediagram: Passar bra här man ska visa värden som köns som går under nominaldata och är en kvalitativ variabel.

1

Bas år 1990

2. År	1990	2000	2005	2010
utgifter	48	40	32	12
KPI	208	260	280	303

a)

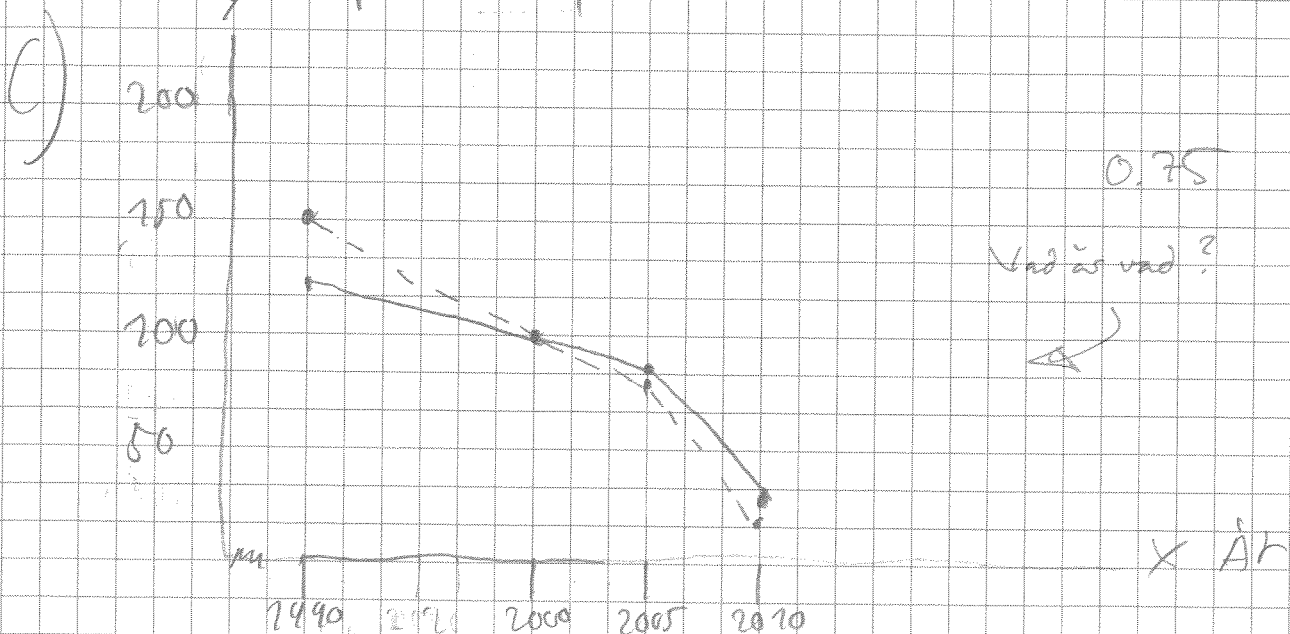
År	lönande	Hade basåret som källare i den Deflator lönande
1990	83	$208/260 \cdot 0,83 = 0,66 = 66$
2000	100	$260/260 \cdot 1,0 = 1,0 = 100$
2005	125	$280/260 \cdot 1,25 = 1,35 = 135$
2010	300	$303/260 \cdot 3 = 3,5 = 350$

a)

År	lönande	Deflator
1990	120	<del>150</del> 150
2000	100	100 = 100
2005	80	0,74 = 74
2010	30	0,26 = 26

1 + (1)

Index



d)

1 Löpande index är den procentuella skillnaden

400%

$$\frac{48}{12} = 4,000$$

1 Pass index är den procentuella skillnaden

68%

$$\frac{208}{303} = 0,68$$

0



3) 10 Prithus hus slumprättiga värden - avståndet till kusten studerades.

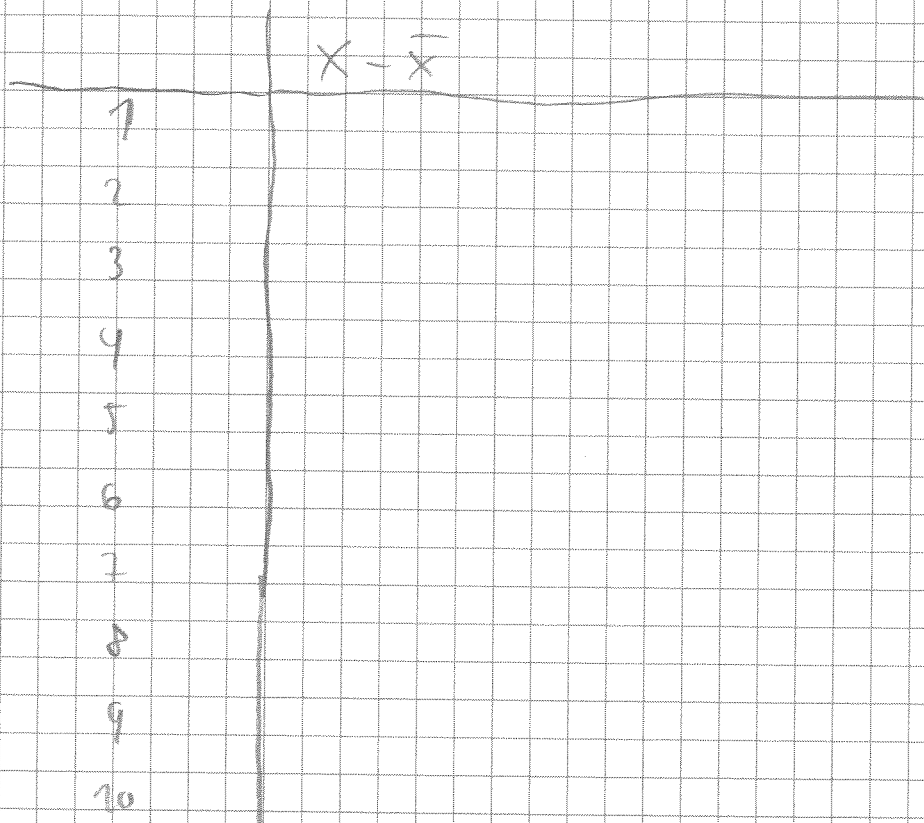
Avvik - hitta ett samband mellan pris och avstånd till kusten.

Läs av värdena i tabellen och rita ett spridningsdiagram.

Pris (tusen kronor)	Avstånd till kusten (km)
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10

Fördelningsgenomsnittet är 1234 tusen kronor

$$\bar{x} = 1234 \text{ tusen } \text{€}$$

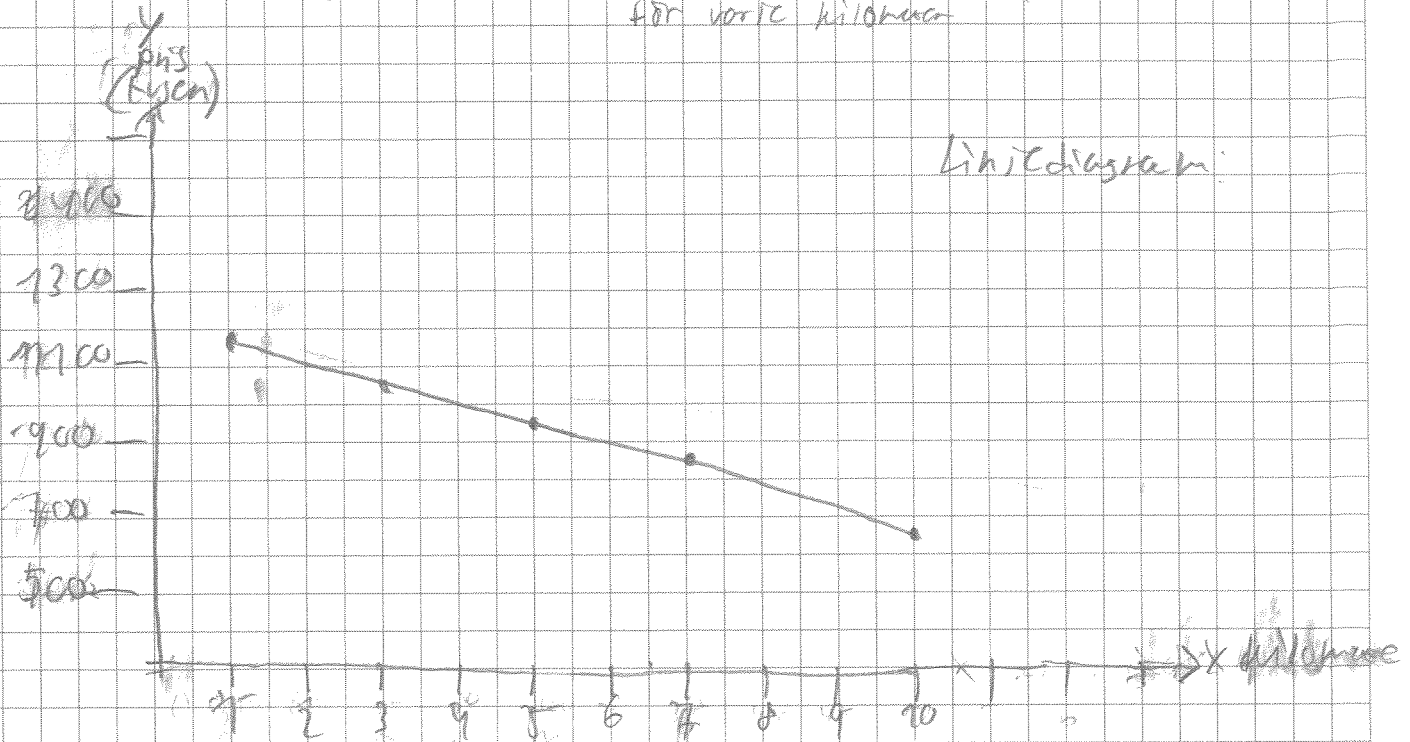


a)  $-0,7$  visar på ett lägt <sup>linjärt</sup> negativt samband d.v.s. att sambandet efter beräkningarna inte är jättelika. Detta kommer också påverka residualernas <sup>?</sup> avstånd till genomsnittslinjan.

0.5

b)

$1234,0 - 56,0 \cdot x$  visar hur mycket priset sjunker för varje kilometer



c)

$1234,0$  = start för husen d.v.s  $1234000$  kr  
 där  $-56000$  är att värdet minskar med  $56000$  kr  
 för varje  $x$  värde. Husen blir billigare ju längre  
 från kusten de ligger

0.5

d)

$$1234,0 - 56,0 \cdot 10 = 674 = 674000 \text{ kr}$$

och med detta kommer ett hus enligt dessa  
 beräkningar kunna vara gratis bara för  
 att det ligger för långt från kusten.

Det står andra värden som följer in så om vilken skillnad

Undertäcknings inte här ut till mål

4)

Svens kommun Undersökning  
vägvisning.Målpopulationen kan anses vara alla myndiga  
personer som bor eller är verksamma i kommunen  
d.v.s. de som även arbetar i kommunen.Rampopulationen kan vara alla invånare  
i kommunen. Eftersom det finns register över  
i vilket kommun människor är bosatta i och  
även var de bor. ???a) Undertäckningsfel = Är här man här ut till  
Rampopulationen men inte till målpopulationen.  
EX. Att man här ut till alla bosatta  
i kommunen men inte dem som arbetar  
i den. (0,5)b) Övertäckningsfel = Är här man här man här  
ut till den målpopulationen men i större  
utsträckning än vad som var tänkt. Man täcker  
ramen.  
Övertäckning. D.v.s. att man här ut till fler  
än vad det var tänkt från början.  
(0,5)EX. Att man här ut till fler än de  
som bor och arbetar i kommunen.Anledning till detta kan vara på grund  
av undersökningen som inte begränsar  
målet till målpopulationen. ?

c) Mätfel = Fel av de olika feltyperna.

Bortfall, Täckningsfel, Beredningsfel, mätfel, ...

Mätfel är det uppstått fel i mätningarna. Att någon svarar dåligt på frågor eller inte förstår frågorna.

Ex. I kommunen kanske man väljer att göra en undersökning där man gör fel i mätningarna. Prägar fel ~~med~~ <sup>kan</sup> såv människor och genom dessa får man inte övertensäkmaras svar i undersökningen.  
(0.5)

d) Bortfall<sup>z</sup> = Det uppstår bortfall i undersökningen, man får inte in antalet svar som man räknas med.

Ex. kommunen väljer att göra ett systematiskt urval där man väljer ut var tionde. 10, 20, 30, 40. Men där 10% av de som slumpmässigt har valts ut inte svarar.  
0.25

e) Systematiskt urval = Bundet slumpmässigt urval.

Ex: I kommunen anländer man sig av cbs register över alla kommunmedlemmar där man slumpmässigt väljer en sitta<sup>hur?</sup> ex. 4 då tar man 4, 14, 24, 34, ... och ber de svara på undersökningen.<sup>varför?</sup>  
0.25

f) Stratifierat urval = Delar upp i strata ex. kön där man slumpmässigt väljer ut värdet från varje stratum.

Ex: Delar upp kommunmedlemmar i låg, medel, hög inkomst. Därefter väljer man slumpmässigt ut värdena från varje strata.  
0.25

4

g)

kvoturval = icke slumpmässigt urval.

Undersökningsansvarige lämnas ganska fri  
att ange vilka som ska undersökas.

EX: Undersökningsansvarige står vid centrum  
Hilssalmans med 50 kolleger under en vecka  
De ska totalt fråga ut - 1000 personer  
d.v.s. 20 personer var. Kommunmedlemmarna  
är totalt bestående av 60% kvinnor &  
40% män. Därför ska 600 kvinnor &  
400 män frågas ut. Individerna får  
den undersökningsansvarige Hilssalmans med  
hans kolleger välja själva. o.s

h)

Urvalsfraction = Hur mycket det skämmer  
i jämförelse med en totalundersökning.

I fraktionen ska det kunna tas ett stickprov  
där det ska kunna jämföras med ett stickprov  
i en totalundersökning. Slumpmässigt bundet.

EX. Urvalsfractionen ska motsvara samma  
värden som skulle komma in i en  
total undersökning.

Urvalsfractionens resultat ska motsvara samma  
resultat som i en totalundersökning.

100 slumpmässiga valda kommunmedlemmar  
sär ska motsvara hela kommunens. 2)



5) a) Nominelldata = kvalitativa, används när man ska mäta via kön eller civilstånd dvs 2 värden som är kvalitativa, men där ... ???

Ex: könsfördelningen i en gymnasieklass.

0.25

b) Diskret variabel = En kvantitativ variabel (med absolut nollpunkt) Ex. på dessa kan vara antalet barn i en familj som då går under datatypen Kvotdata. Definition?

0.25

c) Frekvenspolygon = Flera frekvenser dvs där du har en frekvens som inte kan separeras från en annan. Hör ihop med kumulativ frekvens, där frekvenserna för olika variabelvärden sammansätts till total frekvens. ( )

0

d) Tratt-teknik = Där information sällas ut genom en visualiserad tratt. Man sällar i den totala informationen för att underlätta samsamställningsarbetet.

Ex: I den totalundersökningen i gymnasieskolan vill man använda sig av en tratt-teknik för att snabbt kunna få en bild av hur resultatet ser ut. ( )

0

e) nonsens korrelation: Ett samband som egentligen inte hör ihop.  
 Ex. Att man tv-tittande skulle visa sig leda till att människor får högre IQ.  
 Det är ett nonsenssamband eftersom detta är ju helt beroende på vad du väljer att tittar på när du kollar på tv.

(0.5)

f) Fast basindex = (När du räknar ut de fasta indexet som en deflator. Genom detta ser du bara inflationens påverkan av indexet. Ex när man räknar ut en indexserie genom att använda sig av KPI.  $\frac{1}{1}$ )

g) DIAS I — 0

h) konsistensmetoden I Intervall — 0