

(B)

Statistikk

Del A: Sannsynlighets teori

Her var modellen og dens parameter
kjent, og vi regnet basert på dette.

Eks: Binomial \rightarrow parameter P
Normalf. " μ, σ

| Del B: Statistikk er parameterne for
modellen typisk ikke kjent. Vi skal:

- i) anslå (estimer) parameterne
- ii) ta stikkprøve til påstår om
parameterne

Vi skal også se på enkel linear
regresjon, for å studere sammenhengen
mellan to variabler.

Populasjon og utvalg

Populasjon:

Alle individer/objekter vi er interessert i.

Utvalg:

En delmengde av populasjonen, som vi innhenter data fra.

Viktige egenskaper:

- tilfeldig utvalg
- start utvalg



når utvalget
er tilfeldig
og start,
kan vi trekke
mer presise
konklusjoner
om helt
populasjonen

Tilfeldig:

All individ/objekt
i populasjonen har
same sannsynlighet
for å bli trukket ut

Datar-analyse: én variabel

Utdelgsstørrelse: n

Variabel: X

Datasett:

i	x_i
1	x_1
2	x_2
3	x_3
:	:
n	x_n

gjennomsnitt:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \dots + x_n}{n}$$

(utdelsesgjennomsnitt)

median:

når observasjonene ordnes i stigende rekkefølge er observasjon $(n+1)/2$ om n er oddetall, og gjennomsnittet av observasjon $n/2, n/2+1$ om n er partall

Standard avvik:

$$S_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

(utvalgs std. avvik)

$$\left(S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{kallas utvalgs-kovarians} \right)$$

kvartiler: finnes ved å ordne observasjonen i stigende relativ følge

øvre kvartil:

$$\frac{\text{øvre kvartil}}{\text{observasjoner}} = \frac{3}{4} \cdot (n+1)$$

nedre kvartil:

$$\frac{\text{neder kvartil}}{\text{observasjoner}} = \frac{1}{4} \cdot (n+1)$$

kvartil bredde

$$\frac{\text{kvartil bredde}}{\text{øvre kvartil} - \text{neder kvartil}}$$

Tolkning:

gjennom snitt / median: sentralmas

std. avvik / kvartil bredde: sprengningsmas

Beregning:

\bar{x} , s_x kan finnes ved hjelp av kalkulator

Eles:

7
13
9
12
14

datasett

legges inn på kalkulator:

- ① **C STAT** kommer statistikkmenne
- ② 7 $\Sigma +$
13 $\Sigma +$
9 $\Sigma +$
12 $\Sigma +$
14 $\Sigma +$
- ③ Brukes statistikkfn.
 $\bar{x} \rightarrow \bar{x} = 11$
 $s_x \rightarrow s_x \approx 2.915$

median / quartiler: finnes ved å ordne observasjonene i stigende rekkefølge
(ikke mulig på kalkulator)

Gratisk fra stikkjeg:

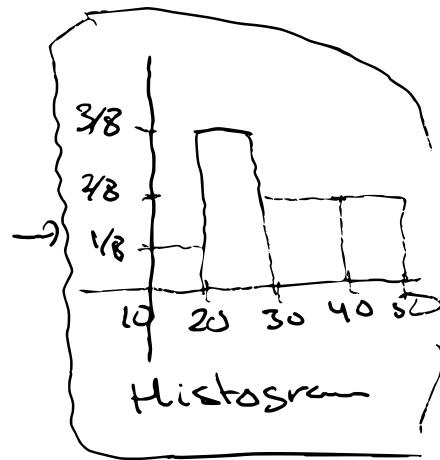
- histogram
- bokstaplett

} ber vte hvordan
man lager disse

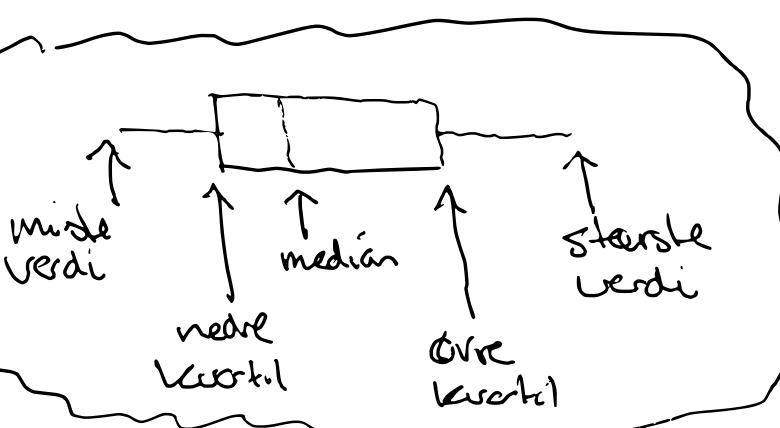
Histogram:

Frekvensstabell

	antall	frekvens
10-20	1	1/8
20-30	3	3/8
30-40	2	2/8
40-50	2	2/8
	8	



Bokstaplett:



Trekning fra en fordeling

Når vi skal bruke teori for å si noe populasjons ("alle") basert på et datasett fra et utvalg (noen "få" som er trukket ut), tenker vi ofte slik:

- i) Vi velger en modell for variablene vi ser på, for eksempel at X er binoisk fordelt eller normalfordelt med ulikerte parametere.
- ii) Dersom utvalgsstørrelsen er n , tenker vi at X_1, X_2, \dots, X_n er stokastiske variabler som følger samme fordelis som X . Om $X \sim N(\mu, \sigma^2)$ er normalfordelt, sier vi at X_1, \dots, X_n er trukket fra en normalfordeling med forventning μ og std. avvik σ .

iii) Vi skriver x_1, x_2, \dots, x_n for tallverdiene i datasettet (merk: sia bolsteuer) og kaller det de observerte verdiene. Vi tenker at x_i er en av de mulige verdiene X kan ha (med sannsynligheter som følges fra fordelingen).

iv) Dersom utvalget er tilfeldig, sier vi at x_1, \dots, x_n er washnigse treninger fra fordelingen til X . De er de stokastiske variablene x_1, \dots, x_n uavhengige.

Estimatorer: Stokastisk variabel
som bruges til at
estimere en
parameter

Navn: estimator for en parameter
af koeff. $\hat{\theta}$, men noen
estimatorer har egne navn

μ : Parameteren μ har
$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$
 som estimator. Det er vanlig
at skrive \bar{X} , ikke $\hat{\mu}$.

σ^2 : Parameteren σ^2 har
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$
 som estimator. Det er vanlig
at skrive S^2 , ikke $\hat{\sigma}^2$.

P: Parameteren p har

$$\hat{p} = \frac{x}{n}$$

Som estimator når X er binomial med parameter (n, p) .

Krav til estimatoren:

- i) Forventingsret: $E(\hat{\theta}) = \theta$
- ii) Var($\hat{\theta}$) mist null
- iii) Var($\hat{\theta}$) $\rightarrow 0$ når $n \rightarrow \infty$

Dette er oppfylt for \bar{x}, S^2, \hat{p} .

Nøg viktige formler:

Når x_1, \dots, x_n er uavhengige trekkninger fra $N(\mu, \sigma^2)$, så har vi:

$$E(\bar{x}) = \mu, \quad \text{Var}(\bar{x}) = \sigma^2/n$$

Hvis X er binomialt fordelt ned
parametre (n, p) , så er

$$E(\hat{p}) = p, \quad \text{Var}(\hat{p}) = \frac{pq}{n}$$

med $\hat{p} = \bar{x}/n$ ($q = 1 - p$).

Standard feil:

Standardfeilen til en estimatør
er standartavviket til estimatøren

$$SE(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$$

Punkt-estimat for en parameter

Når vi setter inn observerte verdier
fra et utvalg i en estimatør for θ ,
får vi et punkt-estimat for θ

$$\mu: \bar{x} \quad \sigma^2: s^2 \quad p: \bar{x}/n$$

Konfidensintervall:

Et konfidensintervall for parameteren θ på signifikansnivå α er et intervall $[A, B]$, med stokastiske grenser A og B , slik at

$$P(A \leq \theta \leq B) = 1 - \alpha$$

Det kallas gjerne et $(1 - \alpha) \cdot 100\%$ konfidensintervall for parameteren θ .

Eks: $\alpha = 0.05 \iff 95\%$ konfidens-intervall

$$P(A \leq \theta \leq B) = 95\%$$

i) Konfidensintervall for μ , σ ligent:

Anta: $\{x_1, \dots, x_n\}$ værdier, idet de
fordelte med $E(x_i) = \mu$ og
 $Var(x_i) = \sigma^2$ med σ ligent
og

\bar{x} tilnæmt normalfordelt
(x_i normalfordelt eller $n \geq 30$)
for hver i

Konfidensintervall:

$$\bar{x} \pm z_{\alpha/2} \cdot \underbrace{\frac{\sigma}{\sqrt{n}}, \text{du)}_{\begin{array}{l} \uparrow \\ \text{punkt-} \\ \text{estimat} \\ \text{for } \mu \end{array}} \quad \left\{ \begin{array}{l} A = \bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \\ B = \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \end{array} \right.$$

\uparrow

$SE(\bar{x})$
std fejl
t.l. estimatorer

\bar{x}

"Z-intervall for μ "

ii) Konfidensintervall for μ , σ ukjent

Anta: x_1, \dots, x_n uavhengige normalfordelte
med $E(x_i) = \mu$ og $Var(x_i) = \sigma^2$

Konfidensintervall:

$$\bar{x} \pm t_{\alpha/2}^{n-1} \cdot s/\sqrt{n}, \text{ dvs } \left\{ \begin{array}{l} A = \bar{x} - t_{\alpha/2}^{n-1} \cdot s/\sqrt{n} \\ B = \bar{x} + t_{\alpha/2}^{n-1} \cdot s/\sqrt{n} \end{array} \right.$$

↑
estimat
for std. feil
t.o.l \bar{x}

↑
punkt-
estimat
for μ

"T-intervall for μ "

iii) Konfidensintervall for σ^2

Anta: x_1, \dots, x_n uavhengige normalfordelte med $E(x_i) = \mu$ og $\text{Var}(x_i) = \sigma^2$

Kritiske verdier:

La Y være χ^2_{n-1} -fordelt (n-1 df.).

Vi skriver χ^2_α for den kritiske verdi defineret ved

$$P(Y > \chi^2_\alpha) = \alpha$$

for enhver α .

Merk: Vi kan ikke finne χ^2_α med kalkulator, via kravne tabell.
(Tabell E6 i boken)

Konfidensintervall:

$$\left[S^2 \cdot (n-1) / \chi^2_{\alpha/2}, S^2 \cdot (n-1) / \chi^2_{1-\alpha/2} \right]$$

Merk: χ^2_{n-1} -fordelingen er ikke symmetrisk.

iv) Konfidensintervall for p

Anta: $\{X$ binomial fordelt med parameter (n, p)
og

$\{n$ ster slik at X er tilnærmet normalfordelt

Konfidensintervall:

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \text{ dvs} \quad \begin{cases} A = \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \\ B = \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \end{cases}$$

↑
punkt-
estimat
 $\hat{p} = x/n$
for p

↑
estimat
for std. feiler
til \hat{p}

Hvor ster må n være?

$$n \hat{p} \hat{q} = n \hat{p}(1-\hat{p}) \geq 5$$

Hypotese test : Vi tar stilling til en påstård om en parameter θ .

Std oppsett:

- i) Utsl modell og hypoteser

H_0 : nullhypote

H_1 : alternativ hypote

Vi velger $H_1 = \text{det vi ønsker å vise}$, og $H_0 = \text{komplementet til } H_1$

- ii) Utsl testobservator og former
på forkostningsområde (nøyresidig,
venstresidig, tosidig)

- iii) Utsl signifikansnivå α = akseptert
sannsynlighet for type I - feil.

- iv) Finn kritisk verdi c bestem
forkostningsområde = verdier for
testobservator som gir at vi forskarter H_0

- v) Sørle inn data, regne ut testobservator
og ta stilling til næstandene.

To mulige
utfall:

Vi forkaster H_0
(" H_0 sann")

Vi beholder H_0
(" H_0 sann")

	H_0 sann	H_1 sann
Beholder H_0	✓	Type II feil
Forkaster H_0	Type I feil	✓

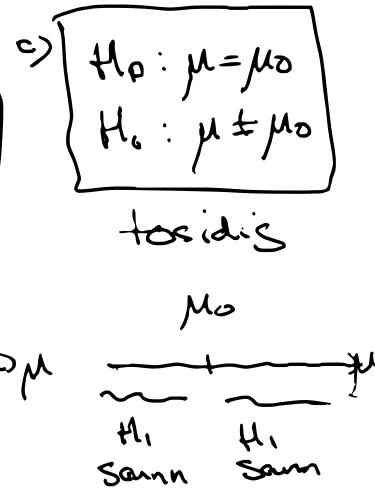
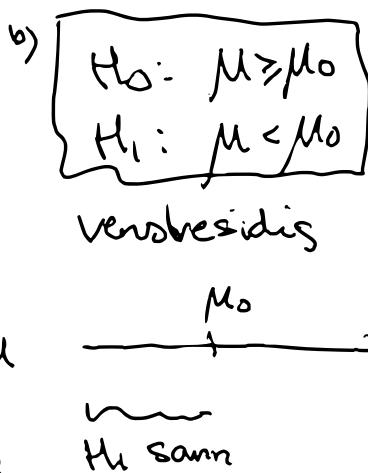
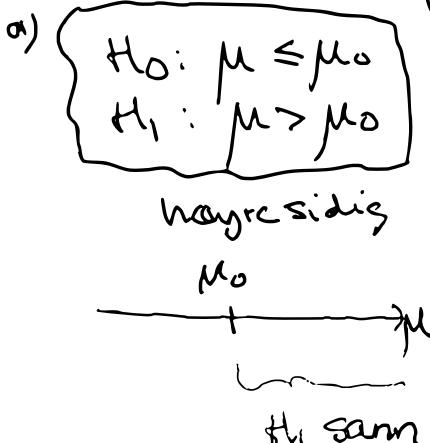
↑
Beslutning

Type I feil: "uskyldig dømt"
Type II feil: "skyldig gøn fri"

i) Hypotesetest for μ

Anta: x_1, \dots, x_n uavhensige normal-fordelte med $E(x_i) = \mu$, $\text{Var}(x_i) = \sigma^2$

Test typer:



Testobservator:

$$T = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} \quad \text{eller} \quad Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

(o vijent) : (o vijent)

T-fordelt
 $n-1$ df.
 når $\mu = \mu_0$
 "T-test"

std.
 normal-fordelt
 "Z-test"

Forkastningssområde:

a) Høyresidig

$$T > t_{\alpha}$$

eller

$$Z > Z_{\alpha}$$

b) Venstresidig

$$T < -t_{\alpha}$$

eller

$$Z < -Z_{\alpha}$$

c) Totsidig

$$|T| > t_{\alpha/2}$$

eller

$$|Z| > Z_{\alpha/2}$$

(kritiske t-verdier her n-1 df.)

Alternativer:

- Bruk \bar{x} i stedet for T/Z som testobservator (skriv ut forkastningsområdet)
- Bruk p-verdi:

Regn ut p-verdi og sammenlikn med signifikansnivå α i stedet for å finne forkastningsområde

p-verdi: sannsynlighet for en minst like elektrisk verdi for testobservator som den vi finner

Vi forkaster H_0 om $P < \alpha$
Vi beholder H_0 om $P \geq \alpha$

Styrkefunksjon og type II -feil:

Styrkefunksjon:

$$\delta(\mu) = P(\text{forkaster } H_0 | \mu)$$

sannsynlighet for type II -feil for
hver sitt verdi (sann verdi) av μ .

ii) Hypotestest for p = populasjonsandel

Anta: $\{ X \text{ binomialt fordelt med parameter } (n, p) \}$
og
 $\} n \text{ stor slik at } X \text{ er tilnærmet normalfordelt } (n\hat{p}(1-\hat{p}) \geq 5)$

Merk: vi må ikke forveksle p her (parameter for binomialt ford.) med p -verdi

Testobser verdi: $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \quad (\hat{p} = \frac{X}{n})$

Hypoteser:

- a) $H_1: p > p_0$ b) $p < p_0$ c) $p \neq p_0$

Førhetsgranside:

- a) $Z > z_\alpha$ b) $Z < z_\alpha$ c) $|Z| > z_{\alpha/2}$