

MET4: Pensumoversikt

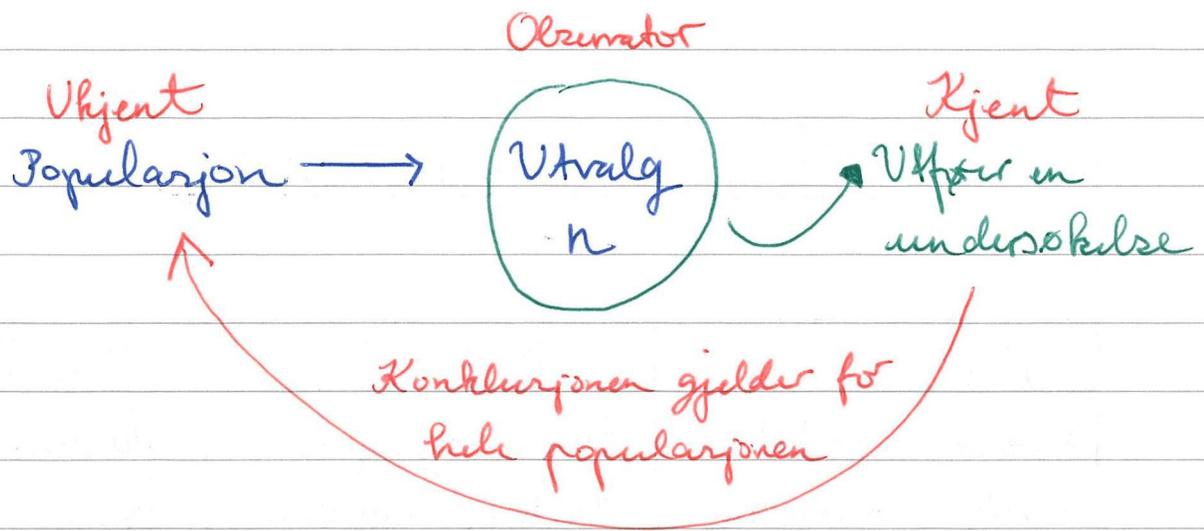
Del I	Del II	Del III
Beskrivende statistikk Utvalg og estimering Inferens om en og to populasjoner Enveis og toveis variansanalyse Fordingsfrie tester Kjipkvadrattester	Enkel regresjon Multipl regressjon Modellbygging Paneldataøkonomi OLS	Tidsserieanalyse Logistisk regresjon

DELI

I. Beskrivende statistikk

Populasjon: Alle enhetene som en problemstilling gjelder for

Utvalg: Den delmengden av populasjonen som blir undersøkt. Vi forutsetter tilfeldige utvalg: dvs at alle har samme sannsynlighet for å komme med i utvalget. Et tilfeldig utvalg sikrer representativitet og muligheter for å generalisere fra utvalget til populasjonen.



μ
 σ^2

Parameter: En tallstørrelse som beskriver en eller annen egenskap ved en populasjon, f.eks forventningen eller variansen. Dette er sanne, men ukjente størrelser.

\bar{X}
 S^2

Observator: En observator er en tallstørrelse som beskriver en egenskap ved utvalget, eks gjennomsnittet og utvalgsvarians.

Utvalgsplan:

- Enkel tilfeldig trekning
- Proporsjonal stratifisering
- Disproporsjonal stratifisering
- Klyngeutvelgelse

NB! Det er den absolutte størrelsen på utvalget som angjør presisjonen av estimatorene, ikke den relative.

- Beskrivende statistikk** : Gir informasjon om
- fordelinger \Rightarrow sentraltendens og variasjon
 - sammenhenger \Rightarrow korrelasjon og regresjon

Variablers målenivå : To hovedkategorier;

(1) **Kategori data** (inkl. rangeringsdata)

- Dikotome (dummy, indikator)
 - to kategorier 0/1
- Nominale variable
 - kategorier som ikke kan rangeres
- Ordinal data (rangeringsdata)

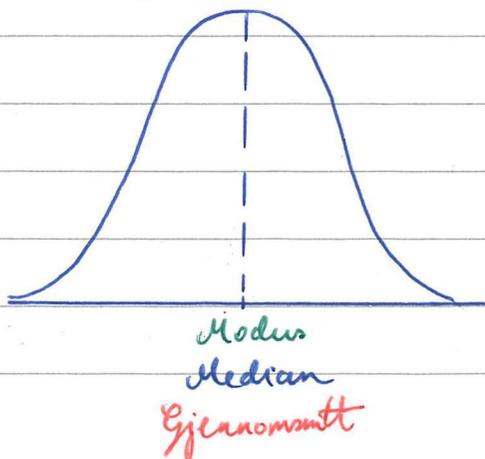
(2) **Målevariable**

- Intervall data
- Forholdsdata

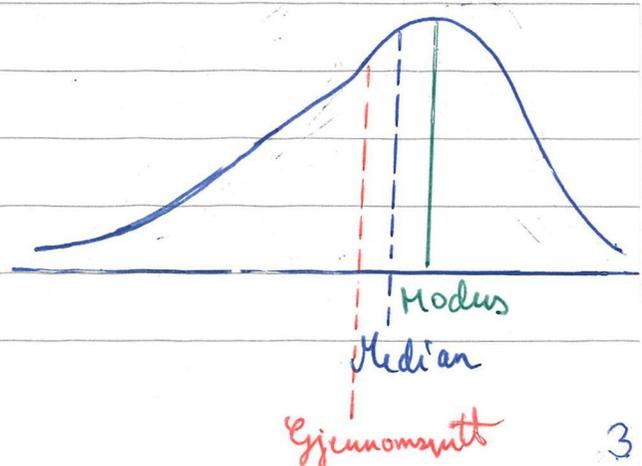
Sentraltendens :

- (1) Gjennomsnitt \bar{x} \rightarrow gjelder kun for målevariable
- (2) Median \rightarrow målevariable og ordinale variable
- (3) Typetall (modus) \rightarrow den vanligste verdien

Symmetrisk fordeling



Skewed fordeling



Spredning:

Variansbredde

Varians og standardavvik

↳ et mål på det gjennomsnittlige avviket fra gjennomsnittet

Varianskoeffisienten (CV): standardavviket delt på gjennomsnittet

- Variansen til populasjonen:

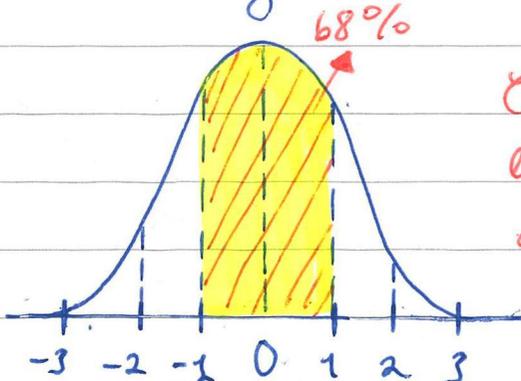
$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{N}$$

- Variansen til utvalget:

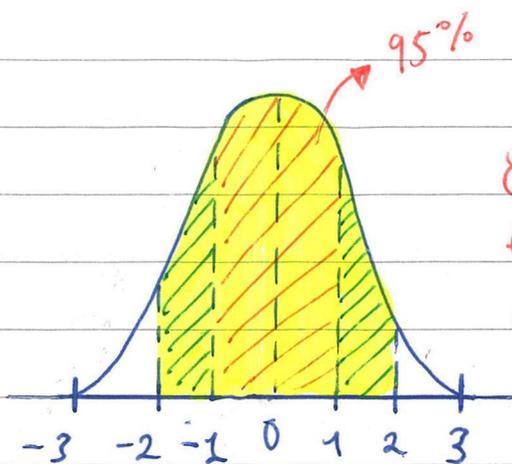
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

- Varianskoeffisientene: $\frac{\sigma}{\mu}$ og $\frac{s}{\bar{x}}$

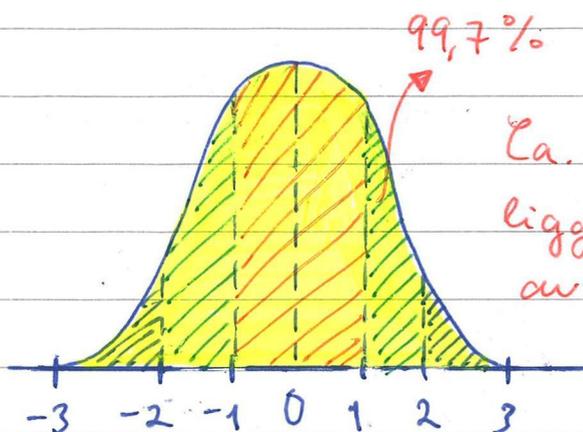
- Følging av standardavvik:



Ca. 68% av observasjonene ligger mindre enn 1 standardavvik fra gjennomsnittet



Ca. 95% av observasjonene ligger mindre enn 2 standardavvik fra gjennomsnittet.



Ca. 99,7% av observasjonene ligger mindre enn 3 standardavvik fra gjennomsnittet.

• **Perisentiler** : Angir den verdien som har p prosent av observasjonene under seg

• Merk: Median = 50% perentilen

• **Korrelasjon** :

Korrelasjonskoeffisienten r forteller oss hvor nær vi er en linear sammenheng

• $r \in [-1, 1]$

• Merk! $r = 0$ sier ikke at det ikke er en sammenheng, bare at sammenhengen ikke er linear!

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)S_x S_y} = \frac{\text{Cov}(X, Y)}{S_x S_y}$$

NB! Med $r=0$ kan det likevel være en sammenheng.

→ Ukorrelertitet medfører ikke uavhengighet

→ Men uavhengighet medfører ukorrelertitet

II. Utvalg og estimering

Stokastisk variabel: En funksjon (eller regel) som tilordner en tallverdi, x , til et hvert utfall ω , i en sannsynlighetsmodell.

Foventning: $E(X) = \sum x \cdot p(x)$

Varians: $Var(X) = E(X - E(X))^2$
 $\sigma^2 = E(X^2) - [E(X)]^2$

Standardavvik: $\sigma = \sqrt{\sigma^2} = \sqrt{Var(X)}$

Kovarians: $Cov(X, Y) = E[(X - E(X)) \cdot (Y - E(Y))]$
 $= E(X \cdot Y) - E(X) \cdot E(Y)$

Regneregler for foventning:

(1) $E(k) = k$

(2) $E(k + X) = E(k) + E(X) = k + E(X)$

(3) $E(kX) = k \cdot E(X)$

(4) $E(X + Y) = E(X) + E(Y)$

(5) $E(X \cdot Y) = E(X) \cdot E(Y)$

(6) $E(X \cdot Y) = E(X) \cdot E(Y) + Cov(X, Y)$

hvis uavhengige

hvis avhengige

Regneregler for varians:

(1) $\text{Var}(X) = E(X^2) - [E(X)]^2$

(2) $\text{Var}(k) = 0$

(3) $\text{Var}(k+X) = \text{Var}(k) + \text{Var}(X) = \text{Var}(X)$

(4) $\text{Var}(kX) = k^2 \cdot \text{Var}(X)$

(5) $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$ *hvis uavhengige*

(6) $\text{Var}(X-Y) = \text{Var}(X) + \text{Var}(Y)$ *hvis uavhengige*

NB!
avhengig { (7) $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{COV}(X,Y)$

(8) $\text{Var}(X-Y) = \text{Var}(X) + \text{Var}(Y) - 2\text{COV}(X,Y)$

Forventning og varians til et gjennomsnitt:

- $X_1, X_2, X_3, \dots, X_n$ uavhengige stokastiske variable
- utvalg: n observasjoner
- hver X har fordelingen $X \sim N(\mu, \sigma^2)$
- vi vil finne fordelingen til gjennomsnittet; \bar{X}

$$\bar{X} = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

$$E(\bar{X}) = E\left(\frac{1}{n} (X_1 + X_2 + \dots + X_n)\right)$$

$$= \frac{1}{n} E(X_1 + X_2 + \dots + X_n)$$

$$= \frac{1}{n} (\underbrace{E(X_1)}_{\mu} + \underbrace{E(X_2)}_{\mu} + \dots + \underbrace{E(X_n)}_{\mu})$$

$$= \frac{1}{n} (\mu + \mu + \dots + \mu)$$

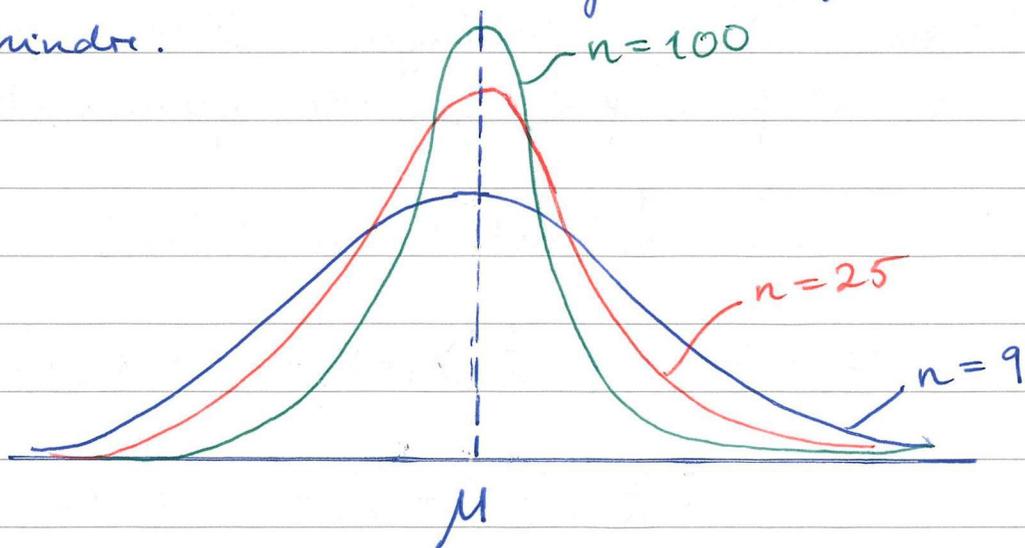
$$= \frac{1}{n} (\mu \cdot n) = \frac{\mu \cdot n}{n} = \underline{\underline{\mu}}$$

$$\begin{aligned}
 \text{Var}(\bar{X}) &= \text{Var}\left(\frac{1}{n}(X_1 + X_2 + \dots + X_n)\right) \\
 &= \left(\frac{1}{n}\right)^2 \cdot \text{Var}(X_1 + \dots + X_n) \\
 &= \left(\frac{1}{n}\right)^2 \cdot \underbrace{\text{Var}(X_1)}_{\sigma^2} + \underbrace{\text{Var}(X_2)}_{\sigma^2} + \dots + \underbrace{\text{Var}(X_n)}_{\sigma^2} \\
 &= \left(\frac{1}{n}\right)^2 \cdot n \cdot \sigma^2 = \frac{n \cdot \sigma^2}{n^2} = \underline{\underline{\frac{\sigma^2}{n}}}
 \end{aligned}$$

Konklusjon: $X \sim N(\mu, \sigma^2)$
 $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

$$\sigma(\bar{X}) = \sqrt{\text{Var}(\bar{X})} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$$

Merk: når n øker utvalgsstørrelsen, blir variansen mindre.



Variansen (spredningen mellom observasjonene) reduseres når n øker. Hølene trekkes seg sammen slik at de fleste observasjonene befinner seg i et intervall som sentrerer seg rundt gjennomsnittet i populasjonen (μ).

Centralgrenseteorem:

Et gjennomsnitt vil alltid nærme seg normalfordelingen når n har mange observasjoner, uansett hvilken fordeling X har.

$$S_n = X_1 + X_2 + \dots + X_n \quad X_i = \text{enkelt-observasjoner}$$

$$E(S_n) = E(X_1) + E(X_2) + \dots + E(X_n) = n \cdot E(X)$$

$$\begin{aligned} \text{Var}(S_n) &= \text{Var}(X_1 + X_2 + \dots + X_n) \\ &= \sum \text{Var}(X_i) + \sum \text{Cov}(X_i, X_j) \\ &= \sum \text{Var}(X_i) \quad \text{hvis uavhengige} \\ &= n \cdot \text{Var}(X) \quad \hookrightarrow \text{Cov}(X_i, X_j) = 0 \end{aligned}$$

$$\text{Når } n \rightarrow \infty \Rightarrow S_n \sim N(n \cdot \mu, n \cdot \sigma^2)$$

Forventning og varians til en andel:

- Binomisk forøksrekke: $I_j = 1 + 0 + 1 + \dots + 1_j$
- $E(I_j) = p \cdot 1 + (1-p) \cdot 0 = p$
- $\text{Var}(I_j) = E(I^2) - E(I)^2 = p - p^2 = p(1-p)$
- Vi er interessert i antall suksesser i et utvalg på n :

$$X_n = \sum_{j=1}^n I_j = I_1 + I_2 + \dots + I_n$$

$$\begin{aligned} E(X_n) &= E(I_1 + I_2 + \dots + I_n) \\ &= \underbrace{E(I_1)}_p + \dots + \underbrace{E(I_n)}_p = \underline{n \cdot p} \end{aligned}$$

$$\begin{aligned} \text{Var}(X_n) &= \text{Var}(I_1 + \dots + I_n) \\ &= \text{Var}(I_1) + \dots + \text{Var}(I_n) \\ &= p(1-p) + \dots + p(1-p) \\ &= \underline{\underline{np(1-p)}} \end{aligned}$$

$$X_n \sim \text{Bin}(n, p) \quad \text{der} \quad E(X_n) = np \\ \text{Var}(X_n) = np(1-p)$$

- Vi kan nå regne ut forventning og varians for andelen suksesser: $\frac{X_n}{n}$

$$E\left(\frac{X_n}{n}\right) = \frac{1}{n} E(X_n) = \frac{1}{n} \cdot np = \underline{\underline{p}}$$

$$\text{Var}\left(\frac{X_n}{n}\right) = \left(\frac{1}{n}\right)^2 \text{Var}(X_n) = \frac{1}{n^2} \cdot np(1-p) = \underline{\underline{\frac{p(1-p)}{n}}}$$

$$\text{Merk: } E(X_n) = np \quad E\left(\frac{X_n}{n}\right) = p$$

$$\text{Var}(X_n) = np(1-p) \quad \text{Var}\left(\frac{X_n}{n}\right) = \frac{p(1-p)}{n}$$

$$\frac{X_n}{n} \sim N\left(p, \frac{p(1-p)}{n}\right) \quad \text{når } n \text{ er stor}$$

$$\frac{X_n}{n} = \hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right)$$

$$\begin{aligned} E(\hat{p}) &= \mu = p \\ \text{Var}(\hat{p}) &= \sigma^2 = \frac{p(1-p)}{n} \end{aligned}$$

Inferens: Vi trekkes konklusjoner om populasjonens parametre basert på det vi observerer i utvalg.

Ukjente populasjonsparametre	Estimatorer
μ	\bar{X}
p	\hat{p}
σ^2	S^2

En god estimator er forventningsrett og konsistent:

Forventningsrett: $E(\hat{\theta}) = \theta$

Konsistent: $S^2 \rightarrow \sigma^2$ når $n \rightarrow \infty$

Merke En estimator kan være konsistent (lav varians) uten å være forventningsrett, og forventningsrett uten å være konsistent.

Viktig: at estimatoren er konsistent!

Hypotesetesting:
 $H_0: \mu = \dots$
 $H_A: \mu \neq \dots, \mu <, \mu > \dots$

P-verdi: den maksimale sannsynligheten for å observere det vi fikk, eller høyere enn det vi fikk, gitt at forutsetningene under H_0 stemmer.

$$\Rightarrow P_{H_0}(\bar{X} > \underline{\text{observert verdi}}) = P\text{-verdi}$$

Tosidig test: $2 \cdot P_{H_0}(\bar{X} > \text{observert verdi}) = P\text{-verdi}$

To typer feil:

(1) Forhastningsfeil: Sannsynligheten for å forkaste en hypotese som egentlig er sann: α . Type 1-feil

(2) Godkjenningsfeil: Sannsynligheten for å godta en hypotese som egentlig er feil. β . Type 2-feil.

	Naturens skjulte sannhet	
	H_0 er sann	H_A er sann
Vår beslutning: Beholde H_0	😊 Riktig konklusjon	😞 Gal konklusjon Type 2-feil: β
Forkaste H_0	😞 Gal konklusjon Type 1-feil: α	😊 Riktig konklusjon



Forhastningsfeil



Godkjenningsfeil

$$\text{Styrken til en test} = 1 - \beta$$

\Rightarrow sannsynligheten for at vi forkaster en nullhypotese som er feil

Inferens om én populasjon med ukjent standardavvik:

- I. Inferens om gjennomsnittet \bar{X}
- II. Inferens om standardavviket / variansen σ^2
- III. Inferens om en andel \hat{p}

I. Inferens om gjennomsnittet t -test

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \hat{\mu} \rightarrow \mu$$

$$E(\hat{\mu}) = E(\bar{X}) = \mu \quad \text{Forventningsrett}$$
$$\text{Var}(\hat{\mu}) = \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

Problem: vi kjenner ikke σ^2 . Må estimeres med utgangspunkt i observasjonene; med utvalgsvariansen.

$$S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2 \quad S^2 \rightarrow \sigma^2$$
$$E(S^2) = \sigma^2$$

$$\hat{\sigma}(\bar{X}) = S(\bar{X}) = \frac{S}{\sqrt{n}}$$

Når X -ene er $N \sim (\mu, \sigma^2)$ er

1. $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$ når $n > 30$

2. $T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$ når $n < 30$

\downarrow
 t -test: vi kjenner ikke variansen (utvalget er lite) 13

- Når σ er kjent (eller n stor) \Rightarrow bruk Z som testobservator.
- Når σ er ukjent (eller n liten) \Rightarrow bruk T som testobservator.

II. Inferens om standardavviket

- Når vi har ukjent standardavvik / varians kan det være aktuelt å teste hypoteser om σ^2
- Slike tester er basert på **kjiradfordelingen (χ^2)**
- Anta at $X_1, X_2, \dots, X_n \sim N(0, 1)$. Da er

$$Q = \sum_{i=1}^n X_i^2 \sim \chi_n^2$$

$$E(Q) = n$$

- Testobservator for variansen:

$$\frac{n-1}{\sigma^2} S^2 \sim \chi_{n-1}^2$$

III. Inferens om en andel

- X_n = antall suksesser i n forsøk
- Estimator for den sanne (ukjente) andelen p :

$$\hat{p} = \frac{X_n}{n} \sim N\left(p, \frac{p(1-p)}{n}\right)$$

Marit Helene Gladhaug

• Teststatistikk: $Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1)$ under H_0

• Konfidensintervall: $\hat{p} \pm k \cdot \hat{\sigma}(\hat{p})$

$\Rightarrow \hat{p} \pm k \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

$\Rightarrow \hat{p} \pm \underbrace{B}_{\text{feilmargin}} = \hat{p} \pm \underbrace{k \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}_{\text{Feilmargin}}$

$\Rightarrow n$ som en funksjon av feilmarginen/precisjonen:

$$n = \left(\frac{k \cdot \sqrt{\hat{p}(1-\hat{p})}}{B} \right)^2$$

Inferens om to populasjoner:

I. Sammenligning av to gjennomsnitt fra normalfordelte populasjoner

- Uavhengig utvalg med lik/ulik varians
- Matchede par

Toutvalgsmodellen: $n_1, n_2, X_1, X_2, \bar{X}_1, \bar{X}_2$

- Vi har to populasjoner med ukjente forventninger μ_1 og μ_2 .

Merk: $\bar{X}_1 \sim N(\mu_1, \frac{\sigma_1^2}{n_1})$ $\bar{X}_2 \sim N(\mu_2, \frac{\sigma_2^2}{n_2})$

- n_1 observasjoner fra populasjon 1
- n_2 observasjoner fra populasjon 2
- gjennomsnitt: \bar{X}_1 og \bar{X}_2
- Vi gjør inferens om $(\bar{X}_1 - \bar{X}_2)$:

$$E(\bar{X}_1 - \bar{X}_2) = \mu_1 - \mu_2$$

$$\begin{aligned} \text{Var}(\bar{X}_1 - \bar{X}_2) &= \text{Var}(\bar{X}_1) + \text{Var}(\bar{X}_2) && \text{antar} \\ &= \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} && \text{uavhengige} \end{aligned}$$

- Testobservator:

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

under H_0

Merk: standardavvikene σ_1^2 og σ_2^2 er ukjent.
 Bruk den empiriske variansen S^2 ; OBS! Må
 bruke t -tabellen og ikke normalfordelingstabellen.

- Vi kan anta at de to populasjonene har **lik varians** og estimere denne som

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$



"Pooled" varians

Vi setter S_p^2 inn for σ^2 og får en observator som er t -fordelt med $v = n_1 + n_2 - 2$ frihetsgrader.

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim t_{n_1 + n_2 - 2}$$

• Vi kan anta at de to populasjonene har **ulik varians** og setter de estimerte standardavvikene S_1 og S_2 for σ_1 og σ_2 og får testobservatoren

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \sim t_v$$

som er tilnærmet t -fordelt med antall frihetsgrader:

$$v = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \right)^2}{\frac{\left(\frac{S_1^2}{n_1} \right)^2}{(n_1 - 1)} + \frac{\left(\frac{S_2^2}{n_2} \right)^2}{(n_2 - 1)}}$$

Merk: Når vi har flere frihetsgrader får vi lettere forkastning.

II. Sammenligning av to varianser

- Når vi sammenligner to varianser ser vi på forholdet S_1^2/S_2^2 og ikke differansen slik som for gjennomsnittet.

$$\Rightarrow \text{Varianser} : \frac{S_1^2}{S_2^2}$$

$$\Rightarrow \text{Gjennomsnitt} : \bar{X}_1 - \bar{X}_2$$

- $H_0 : \sigma_1^2 = \sigma_2^2$
 $H_A : \sigma_1^2 \neq \sigma_2^2$

- Testobservatoren :

$$F = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} = \frac{S_1^2}{S_2^2} \sim F_{n_1-1, n_2-1}$$

III. Sammenligning av to andeler

- To populasjoner n_1 og n_2
- Observerer $\hat{p}_1 = \frac{X_1}{n_1}$ og $\hat{p}_2 = \frac{X_2}{n_2}$
- $X_1 \sim \text{Bin}(n_1, p)$ med $E(X_1) = np$
 $\text{Var}(X_1) = np(1-p)$
 $X_2 \sim \text{Bin}(n_2, p)$ med $E(X_2) = np$
 $\text{Var}(X_2) = np(1-p)$

$$\bullet \quad E\left(\frac{X_{n_1}}{n_1}\right) = p_1 \quad \text{Var}\left(\frac{X_{n_1}}{n_1}\right) = \frac{p(1-p)}{n_1}$$

\hat{p}_1 \hat{p}_1

$$\bullet E\left(\frac{\overset{\hat{p}_2}{X_{n_2}}}{n_2}\right) = p_2 \quad \text{Var}\left(\frac{\overset{\hat{p}_2}{X_{n_2}}}{n_2}\right) = \frac{p(1-p)}{n_2}$$

$$\bullet E(\hat{p}_1 - \hat{p}_2) = p_1 - p_2$$

$$\bullet \text{Var}(\hat{p}_1 - \hat{p}_2) = \left(\frac{1}{n_1} + \frac{1}{n_2}\right) p(1-p)$$

• Testobservator :

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right) p(1-p)}} \sim N(0, 1)$$

der $p = p_1 = p_2$

• ANOVA : Enveis variansanalyse

• Variansanalyse er en utvidelse av t-testvalgsmodeller til situasjoner der vi ønsker å sammenligne to eller flere populasjoner (grupper)

• Antar normalfordelt responsvariabel med **lik varians** for alle grupper

• H_0 : Alle gruppene har lik **forventning**
 H_A : Minst én er forskjellig

⇒ H_0 : $\mu_1 = \mu_2 = \mu_3 = \dots = \mu_n$
 H_A : Minst én forskjellig

- X_{ij} = utfallet for observasjon i tilhørende gruppe j
- $E(X_{ij}) = \mu_j$ $\text{Var}(X_{ij}) = \sigma^2$
- $\hat{\mu}_j = \bar{X}_j = \frac{1}{n_j} \sum X_{ij} \rightarrow$ gruppeforventning
- Samlet forventningsverdi som er lik for hver gruppe under H_0 :

$$\hat{\mu}_j = \bar{X} = \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{n_j} X_{ij}$$

- Variasjon mellom gruppene er

$$SST = \sum_{j=1}^k n_j (\bar{X}_j - \bar{X})^2$$

- Variasjon innad i gruppene er

$$SSE = \sum_{j=1}^k \sum_{i=1}^{n_j} (X_{ij} - \bar{X})^2$$

- Totalvariasjon : $SS(\text{Total}) = SST + SSE$

$$SS(\text{Total}) = \sum_{j=1}^k \sum_{i=1}^{n_j} (X_{ij} - \bar{X})^2$$

- Testobservator :

$$F = \frac{SST / (k-1)}{SSE / (n-k)} \sim F_{(k-1), (n-k)}$$

Mart Helene Gladhaug

Merk: Dessom variasjonen mellom gruppene (SST) er stor relativt til variasjonen innen gruppene (SSE) er det lite sannsynlig at de ulike gruppene har lik forventning

\Rightarrow SST \uparrow og SSE \downarrow : vi deler et stort tall på et lite tall \Rightarrow vi får en stor F-verdi
 \Rightarrow Stor F-verdi forkaster H_0 om lik forventning

Merk ANOVA-testen sier ikke hvilken gruppe som skiller seg ut. For å finne den avvikende gruppen kan vi bruke **Tukey-testen**.

- ANOVA: toveis variansanalyse
- $E(X_{ij}) = \mu_{ij} = \mu + \alpha_i + \beta_j$
- Antar felles varians for alle grupper
- $\hat{\mu} = \bar{\bar{X}} = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b X_{ij}$
- Testobservatorer:

$$F_A = \frac{SS(A) / (a-1)}{SSE / (a-1)(b-1)} \sim F_{(a-1), (a-1), (b-1)}$$

• $F_B = \frac{SS(B) / (b-1)}{SSE / (a-1)(b-1)} \sim F_{(b-1), (a-1), (b-1)}$

• **Fordelingsfrie tester** → når vi ikke har normalfordelte observasjoner

I. Wilcoxon's test for to utvalg

$$\mu_1 = \frac{n_1(n_1 + n_2 + 1)}{2}$$

$$\sigma = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$$

H_0 : Medianforskjellen er lik null

H_A : Medianforskjellen er ulik null

II. Wilcoxon's test for pairede observasjoner

$$\mu = \frac{n(n+1)}{4}$$

$$\sigma = \sqrt{\frac{n(n+1)(2n+1)}{24}}$$

III. Kruskal-Wallis-testen: Fordelingsfri enveis-
variansanalyse

IV. Friedmans-test: Fordelingsfri toveis
variansanalyse

• Kjikkvadrattester

I. Kjikkvadrattest for modelltilpassning

II. Kjikkvadrattest for uavhengighet

H_0 : kjennetegnene er uavhengige

H_A : kjennetegnene er ikke uavhengige

$$\chi^2 = \frac{(X_1 - \bar{X}_1)^2}{\bar{X}_1} + \frac{(X_2 - \bar{X}_2)^2}{\bar{X}_2} + \dots + \frac{(X_n - \bar{X}_n)^2}{\bar{X}_n}$$

DEL II

Enkel regresjon

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$

$$E(Y|X) = \beta_0 + \beta_1 X$$

$$\text{der } \hat{\beta}_1 = \frac{\text{Cov}(X_1, Y)}{\text{Var}(X)}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Fortschingene for OLS

(i) $E(\varepsilon) = 0$

(ii) β_e er konstant for alle X **homoskedastisitet**

(iii) E_i er uavhengig av E_j for alle i og j
Brudd på denne forutsetningen kalles autokorrelasjon

(iv) Dersom X -variablene ikke er deterministisk, men stokastisk, må den være uavhengig av E .
Avhengighet mellom X og E kalles endogenitet.

(v) Normalitet: E er normalfordelt.

Standardavvikene

S_e	Estimert standardavvik til feillemmet
$S(\hat{\beta}_2)$	Estimert standardavvik til stigningstallet
$S(\hat{Y})$	Estimert standardavvik til estimert forventet Y
$S(Y - \hat{Y})$	Estimert standardavvik til prediksjonsfeilen

(1) Hvorfor en koeffisient kan endre fortegn når en annen forklaringsvariabel inkluderes

Tenk deg at du har en åker, og du måler avlingen Y . Så ønsker du å forklare avlingen med nedbør i en lineær regresjon, og finner ut at koeffisienten blir negativ. Altså, dersom nedbøren øker med en enhet, blir avlingen beta enheter dårligere. Så inkluderer du temperatur som forklaringsvariabel, og får at begge er positive. Det kan skje dersom temperatur og nedbør er *negativt* korrelerte.

Dersom du bare har nedbør, kan vi tenke oss at mye nedbør ofte henger sammen med kaldt vær, og at det ikke er bra for avlingen. Men for en **gitt** temperatur, vil mye nedbør være bra. Derfor kan nedbørskoeffisienten bli positiv når vi kontrollerer for temperatur.

(2) Hvorfor en koeffisient kan gå fra å være signifikant til ikke signifikant (og omvendt) når en annen forklaringsvariabel inkluderes i modellen

En koeffisient er signifikant dersom den er stor i forhold til sitt standardavvik, sjekk formel på s. 35 i forelesning 11. den kan miste sin signifikans pga multikolaritet (andre faktor) eller få signifikans dersom den nye variabelen forklarer mye variasjon, slik at S_{ϵ} blir mindre (første faktor).

(1) White-test : tester heteroskedastisitet

$$H_0 : \sigma_1^2 = \sigma_2^2 \text{ for alle } i \text{ (homoskedastisitet)}$$

$$H_A : \sigma_1^2 \neq \sigma_2^2 \text{ for alle } i \text{ (heteroskedastisitet)}$$

Høy P-verdi \Rightarrow Behold $H_0 \Rightarrow$ homoskedastisitet 😊

(2) Kji-kvadrat-test : tester uavhengighet

$$H_0 : \text{kjennetegn er uavhengige}$$

$$H_A : \text{kjennetegn er avhengige}$$

Høy P-verdi \Rightarrow Behold $H_0 \Rightarrow$ uavhengighet 😊

- Heteroskedastisitet : estimatene er forventningsrette, men inferens er ikke gyldig.
- Autokorrelasjon : estimatene er forventningsrette, men inferens er ikke gyldig.
- Ikke normalfordelte feilledd : inferens er ikke gyldig i små utvalg.
- Endogenitet : forventningsskjeve estimatører