

6/11 En liten oppsummering for vi starter med
kursets høydepunkt: **REGRESJON.**

Oppsummering av hovedlinjene i kurset så langt:

Deskriptiv statistikk (Ch 1)

Kategorielle data
(2 eller flere kategorier ^{grupper})

tabeller: Frekvenser (antall)
%

Kontinuerlige data

↓
histogram
(boksplott)

↙

Symmetrisk

\bar{X} , SD

Basert på sjumennsritt

↘

skjævt

median, Q_1 , Q_3

Basert på rangeringer

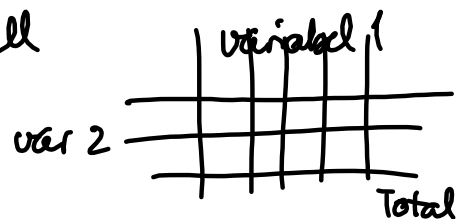
Bivariate sammenhenger (se neste side)

Neste: Regresjon

Bivariate sammenhenger: Sammenhenger mellom to ulike variable

kategorisk mot kategorisk

Krysstabell



Ch 2.6

Hypotesetest for ikke Pearson
 H_0 : Ingen sammenheng mellom de to variablene
 Ingen forskjell på gruppene
 Ch 8.2 og Ch 9

kategorisk mot kontinuerlig

2 kat
2 grupper

(histogram)
boksplott

s. 97 Ch 1.3 "side by side boxplot"

Symmetrisk



\bar{x}_1, SD_1 \bar{x}_2, SD_2

Skjeiv



median₁ median₂
Q₁ Q₃ Q₁ Q₃

Hypotesetest for
to-utvalgst-test

Ch 7.2

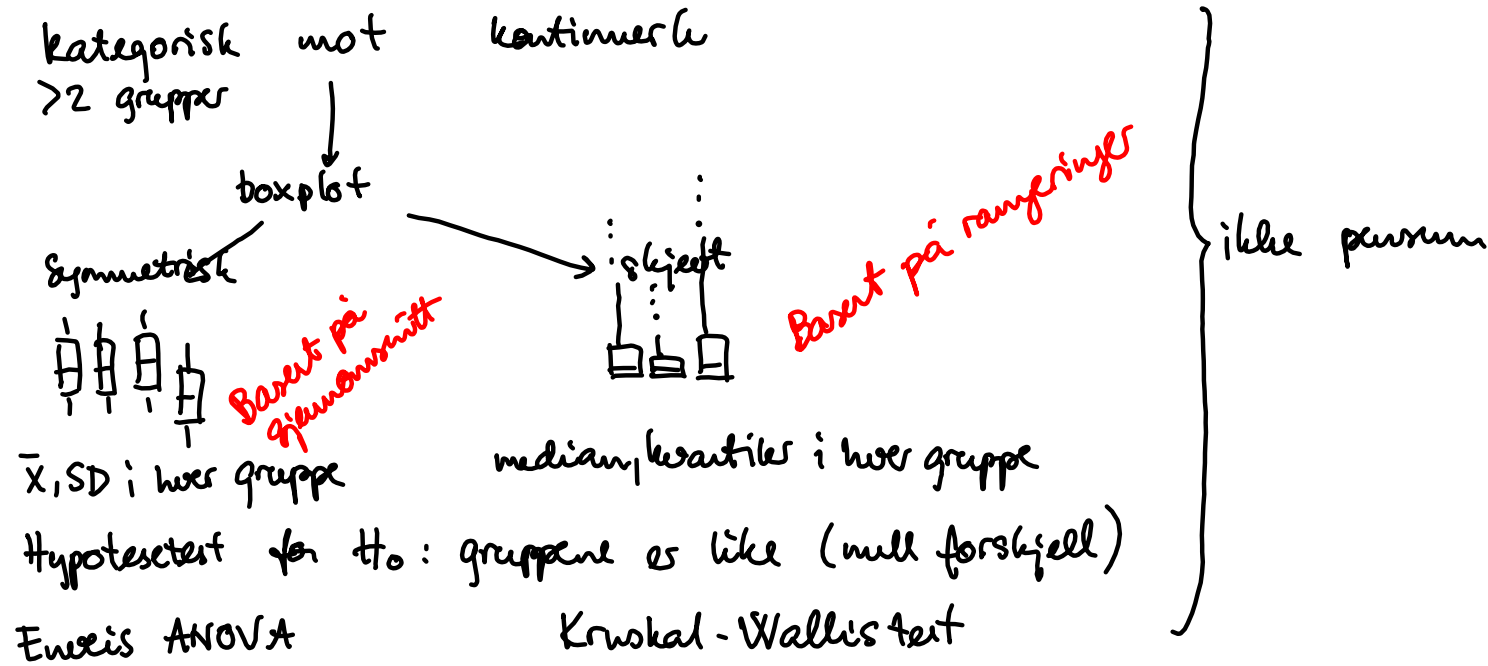
Basert på gjennomsnitt

H_0 : ingen forskjell på gruppene (ingen sammenheng mellom variablene)

Wilcoxon Rank sum test (Mann-Whitney-testen)

Ch 15.1

Basert på rangeringer av observasjoner



CLT etc : Teorien
 Bygges opp omkring
 normalfordeling og μ_1, μ_2

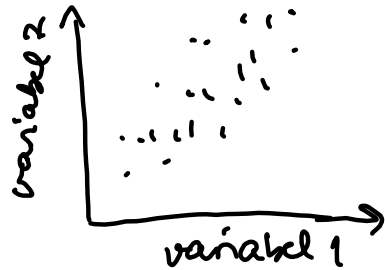
Parametriske tester
 Parametriske metoder

Bæret på rangeringer,
 ikke på normalfordeling og parametere
 Kaller de for

Ikke-parametriske tester
 Ikke-parametriske metoder.

Kontinuerlig mot kontinuerlig

Scatter-plot

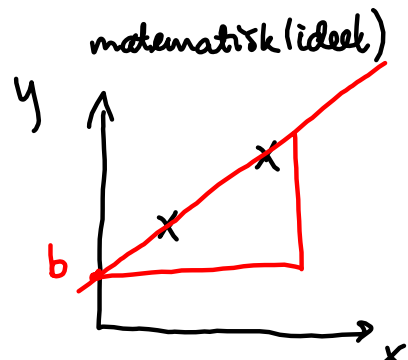


Oppsummerer sammenhengen mellom de to variablene med korrelasjonskoeffisienten, r (Ch 2)

r er et tall mellom -1 og 1 sul.no

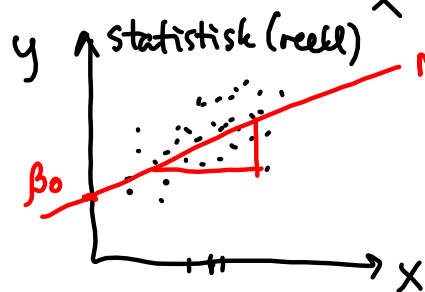
HM: Hva med det vi lærte om at

$$y = ax + b$$



Stigningstallet a : Hvor mange enheter y endres seg når x endres seg en enhet

PUGG



Må estimere denne linja

$$y_i = \beta_0 + \beta_1 \cdot x_i + \epsilon_i$$

Regressjonskoeffisienten β_1 sier hvor mange enheter vi forventer at y skal endre seg når x endres seg med en enhet

$$\epsilon_i \sim N(0, \sigma) \text{ uavhengig hva } x \text{ er}$$

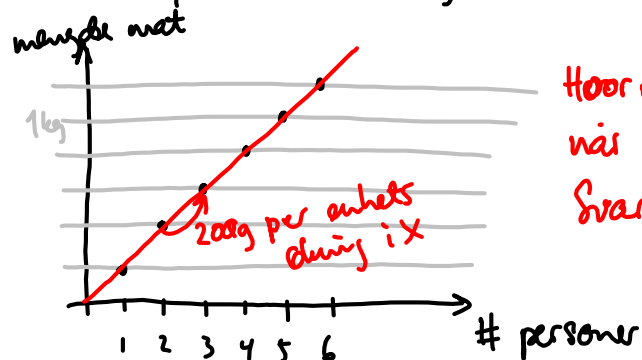
Linear regresjon : introduksjon

funksjon $y = f(x)$: Det er en sammenheng mellom y og x

$y = \text{mengden mat (i gram) man lager til middag} = f \left(\begin{matrix} \# \text{ personer} & \text{aktivitet} \\ \text{hvor sulten} & \text{anledning ...} \end{matrix} \right)$

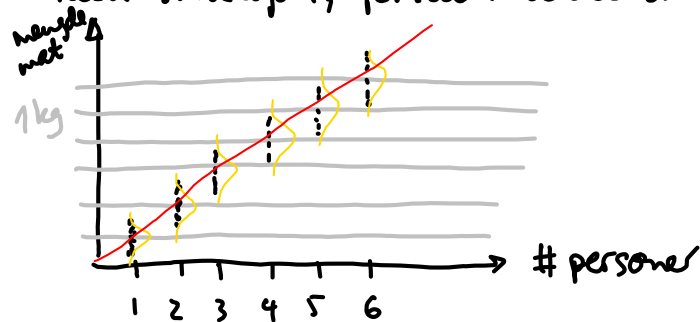
Linear funksjon (oppskrift)

f.eks. kjøtt : ca 200g per person



Hvor mye øker y for
når x øker med én enhet?
Svar : 200g

Reell situasjon, fortsatt linear smk :



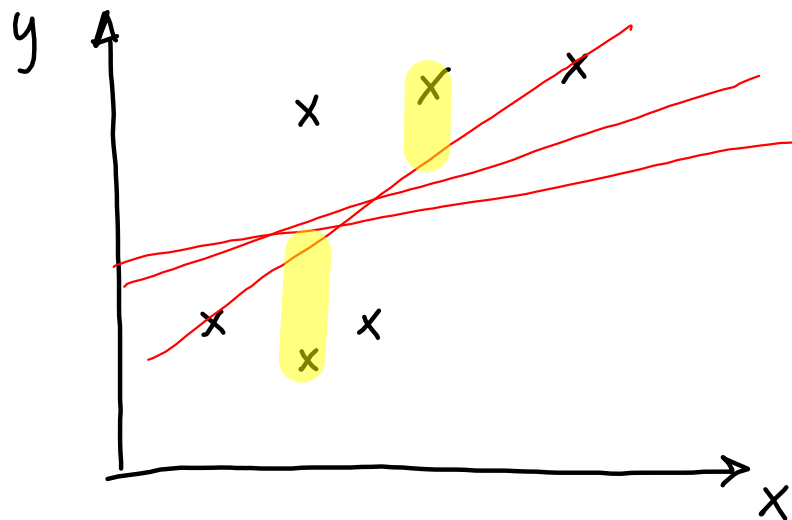
$$y = \beta_0 + \beta_1 \cdot x + \varepsilon$$

Hva vi forventer

$\varepsilon \sim N(0, \sigma)$ lik σ uansett hvor x er.

variasjonen

Estimering av regresjonslinja via minste kvadraters metode Method of least squares



$$\text{Hvis } y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

så må vi finne de β_0 og β_1 som gjør at linja "passer best" til observasjonene.

I minste kvadraters metode betyr det å finne de β_0 og β_1 som minimerer

den totale ^{sumerte} kvadratastanden til linja:

$$\sum_i \left(\overset{\text{avstand fra linja}}{y_i - (\beta_0 + \beta_1 x_i)} \right)^2$$

de β_0 og β_1 som minimerer denne kvadratastanden, kalles vi minste kvadraters estimater.

H0: Ingen smh mellom x og y er det samme som $\beta_1 = 0$