

1.3 Beskrivende statistikk

- sentralt mål: gjennomsnitt og median
- spredningsmål: standardavvik og kvartilbredde
- hensiktsmessig presentasjon av data

Histogram

- eksempel: datasett med alder på studenter:

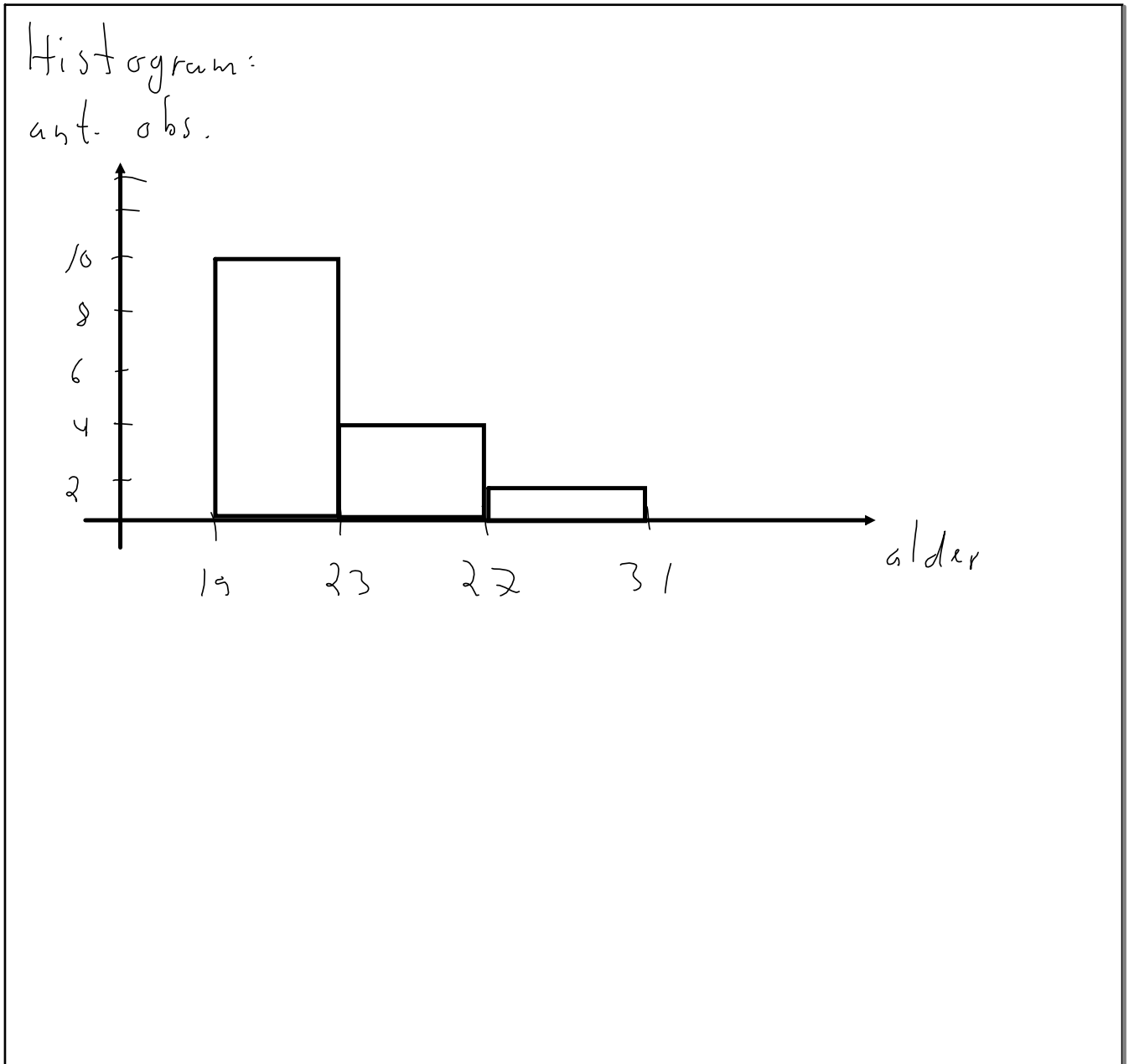
19, 20, 22, 25, 28 osv.

- deler datamaterialet inn i intervaller/
"klasser", f.eks.

$[19, 23)$, $[23, 27)$, $[27, 31)$

- klassebredden er 4
- utgangspunkt: frekvens tabell som viser antall i hver klasse:

Klasse	antall
$[19, 23)$	10
$[23, 27)$	4
$[27, 31)$	2



Sentralmål (mål på hvor "sentrum" for et datasett er)

To vanligste sentralmål:

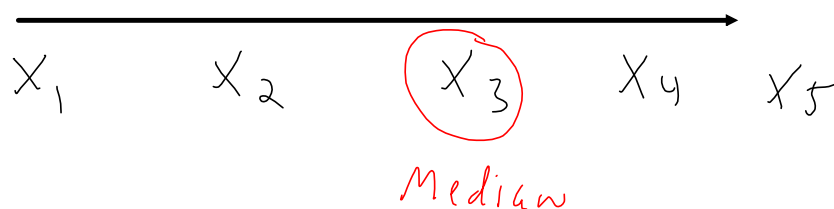
1) Gjennomsnittet av målingene X_1, X_2, \dots, X_n

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n} = \frac{\sum_{i=1}^n X_i}{n}$$

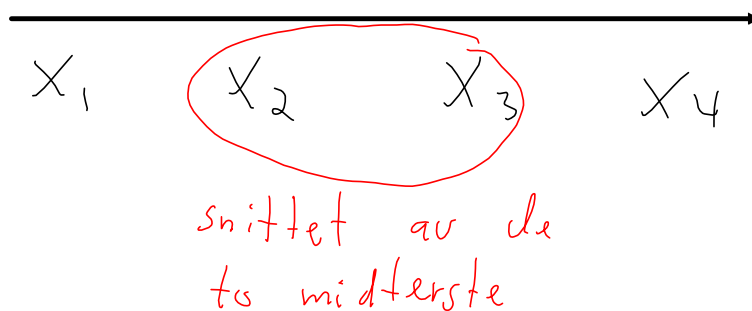
- ulempe: gjennomsnitt er følsomt for utliggere / ekstremverdier

2) Median: observasjonen i midten av datasettet, når dette er sortert i stigende rekkefølge

Oddesantall



Partall



- fordel: median er mindre følsom for utliggere

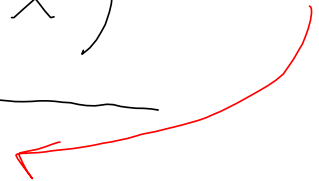
Spredningsmål (mål på hvor "spredt" et datasett er)

1) Standard avvik / varians

Varians for et datasett X_1, X_2, \dots, X_n
er definert som

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

WTF?



Standard avviket er da rota av
variansen :

$$S = \sqrt{\text{Varians}} = \sqrt{S^2}$$

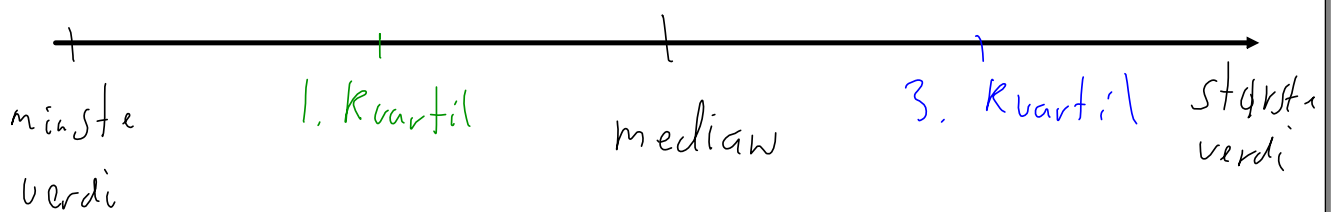
- alemp: følgsomt for utliggere

2) Kvartilbredde

En kvartil ("fjerdedel") for et datasett er en observasjon som deler datasettet (sortert i stigende rekkefølge) i fire like store deler:

$$Q_1 = \text{obs. nr. } \frac{n+1}{4}$$

$$Q_3 = \text{obs. nr. } 3 \cdot \left(\frac{n+1}{4}\right)$$

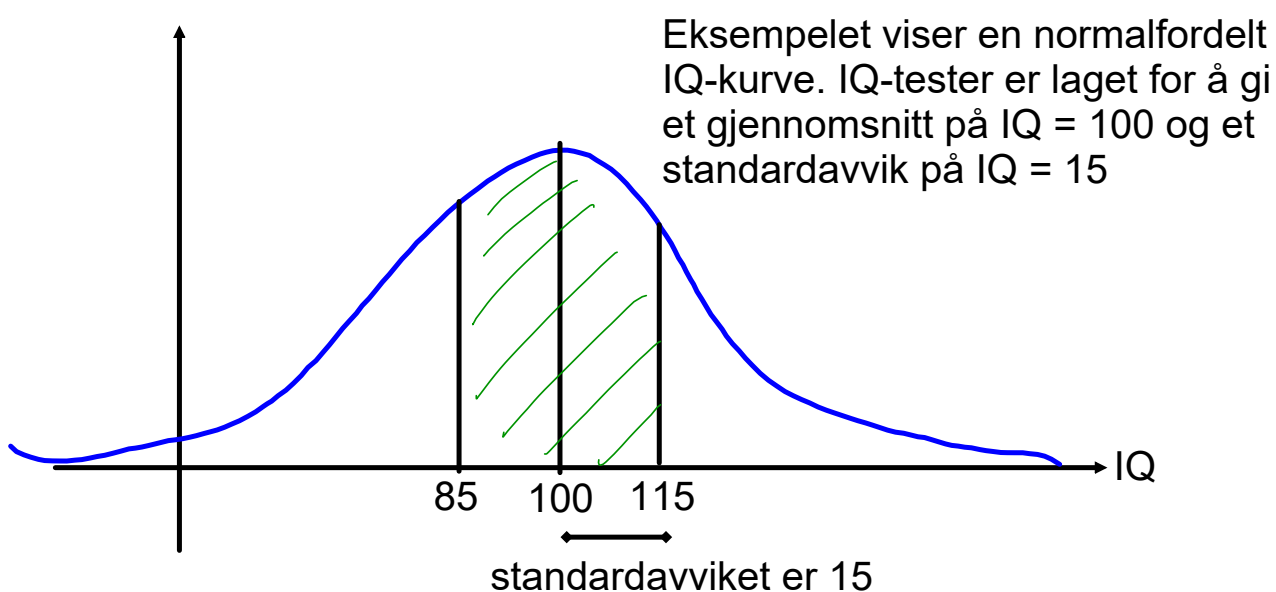


Spredningsmålet kvartilbredde er definert som

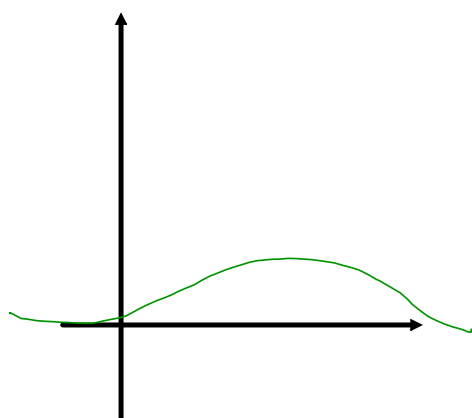
$$Q_3 - Q_1$$

- fordel: mindre følsomt for utliggere enn standardavvik

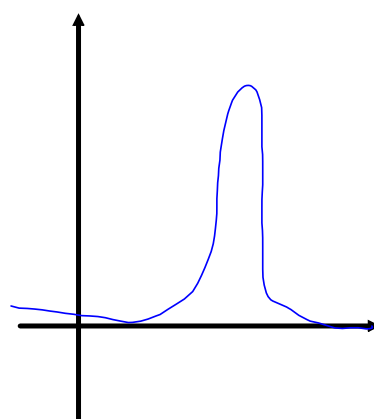
Figurer som illustrerer sentral- og spredningsmål



Standardavvikets betydning for utseendet til fordelingen



stort standardavvik
= stor "spredning"



lite standardavvik
= liten "spredning"

Oppsummering av nyttige Excel-formler for sentral- og spredningsmål og datavisualisering

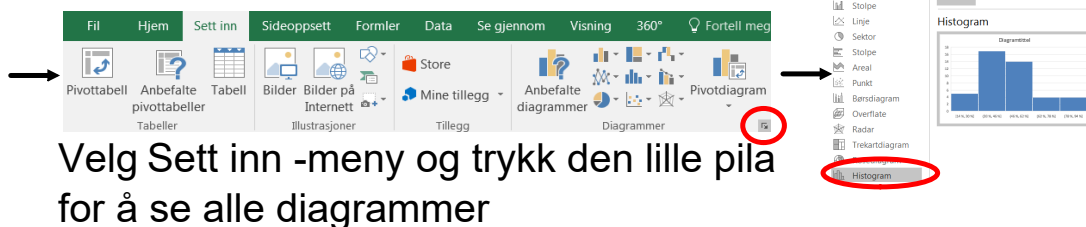
Anta at dataene ligger i én kolonne, i området A1-A100.

Gjennomsnitt	=GJENNOMSNITT(A1:A100)
Standardavvik	=STDAV.S(A1:A100)
Median	=MEDIAN(A1:A100)
Øvre kvartil Q3	=KVARTIL.EKS(A1:A100;3)
Nedre kvartil Q1	=KVARTIL.EKS(A1:A100;1)
Kvartilbredde	(ber Excel beregne Q3-Q1)

Forskjellen mellom "inkludert" og "ekskludert" kvartil er ikke noe vi skal bry oss om, men handler om hvordan man "veker" dersom kvartilene havner mellom forskjellige observasjoner

Sette inn histogram:

Merk data-området



Eksempel på beregning av median/kvartiler

Alderen til noen studenter er gjengitt i datasettene under:

Odde antall obs.

18, 19, 20, 21, 25

$$\text{Median: } M = \underline{\underline{20}}$$

$$Q_1 = \text{obs. nr. } \frac{(5+1)}{4} = \text{nr. } 1,5$$

18 Q_1 19 20 21 Q_3 25

1 1,5 2 3 4 5

Her er

$$Q_1 = \underline{\underline{18,5}}$$

$$Q_3 = \text{obs. nr. } \frac{3(5+1)}{4} = 4,5$$

Her er

$$Q_3 = \underline{\underline{23}}$$

Partall antall obs.

19, 20, 23, 25, 27, 31

$$\text{Median: } M = \frac{23+25}{2} = \underline{\underline{24}}$$

$$Q_1 = \text{obs. nr. } \frac{(6+1)}{2} = \text{nr. } 1,75$$

19 Q_1 20 23 25 27 Q_3 31

1 1,75 2 3 4 5 6

Her er

$$Q_1 = 19 + 0,75 \cdot (20 - 19) = \underline{\underline{19,75}}$$

$$Q_3 = \text{obs. nr. } \frac{3(6+1)}{4} = 5,35$$

Her er

$$Q_3 = 27 + 0,25 \cdot (31 - 27) = \underline{\underline{28}}$$