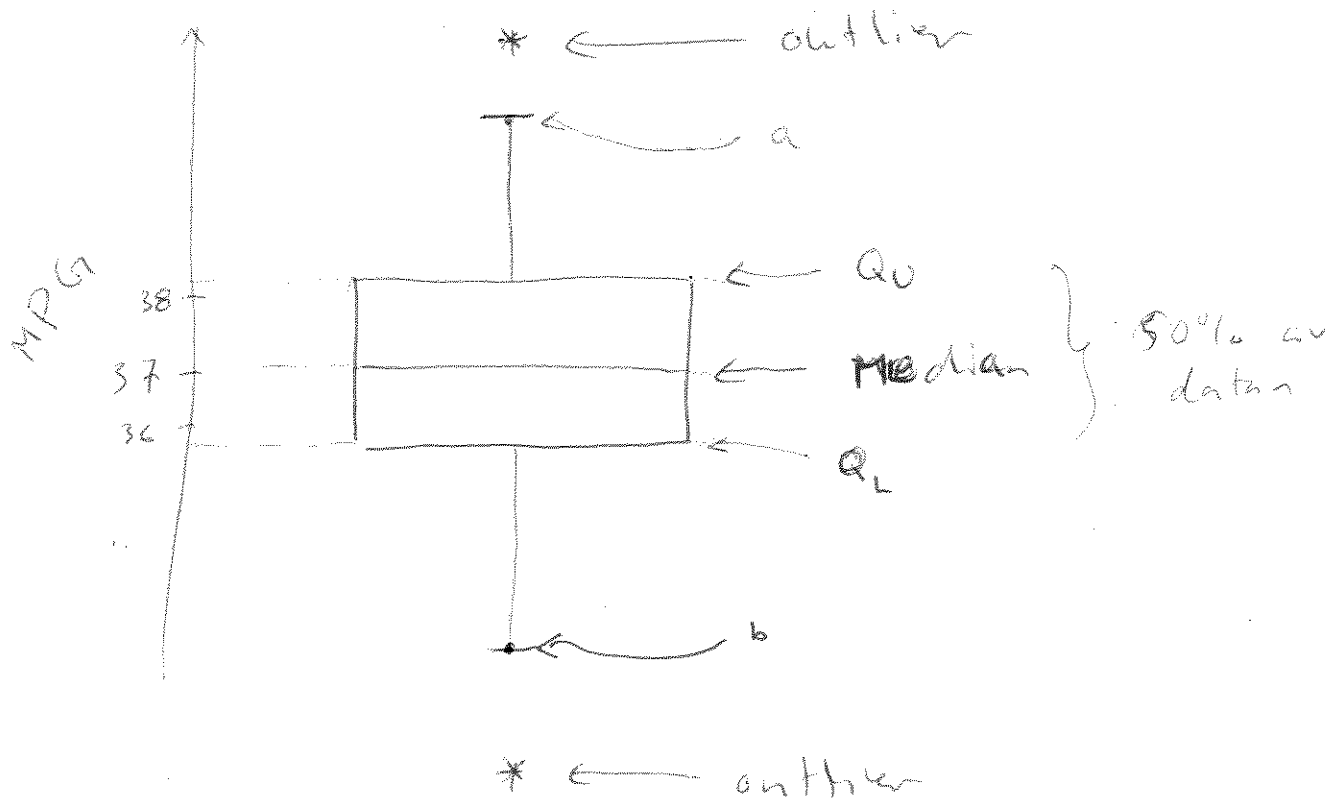


Boxplots (låddiagram)

(kap 2.8)

Boxplots är ett väldigt bra sätt att illustrera kontinuerlig data. En boxplot ~~kan~~ kan se ut på följande sätt (Figur 31)



Punkten a är den datapunkt (observationen) med högst värde som är mindre än $Q_U + 1.5 \times IQR = Q_U + 1.5(Q_U - Q_L)$

Punkten b är den datapunkt med lägst värde som är större än $Q_L - 1.5 \times IQR = Q_L - 1.5(Q_U - Q_L)$

PLOTTAR!

Outliers är ^{data} punkter som ligger väldigt långt ifrån resten av datasetet. I en boxplot ~~är~~ är outliers "utanför" a och b och betecknas med en stjärna eller liknande.

Skewhet

Ibland kan det vara viktigt att känna igen ett dataset som inte är symmetriskt. Dataset som inte är symmetriska kallas för skewa. I ett skett dataset är de vanligare med extrema observationer åt ena hållet än åt det andra hållet.

PLOT!

Ett tecken på att ett dataset är skett är att medianen och medelvärdet ~~är~~ ligger en bit ifrån varandra. PLOT!

En boxplot av ett skett dataset kännetecknas av att den ena ~~halvan~~ av lådan är större än den andra
"delan" "morrhärr"

och att motsvarande ~~del~~ är större än det andra.

PLOT!

Exempel på typiska skett dataset är lönenivåer, befolkning-
mängd.

Punktskattningar

(kap 6)

En population har ett parameter vars värde är okänt. Vi vill gissa värdet utan att undersöka hela populationen. Hur ska vi göra gissningarna och vad har gissningarna (slumpmässiga!) för egenskaper?

Ex 1: Marknadsundersökning. Coca cola skall lansera den nya colan CCX som utmanare till Pepsi max. Hur stor del av populationen föredrar CCX före Pepsi max?

(population = alla i Sverige/USA/världen)

parameter: p = proportionen som föredrar CCX)

Ex 2: Kvalitetskontroll. Pögens frallor har normalfördelad vikt med okända parametrar μ och σ^2 . Vad är värdet på dessa?

Hur ska vi gissa parametrarna p , μ och σ^2 ?

1) Vi samlar in data (stikprov)

2) Data sammanfattas med en skattning

Fördelningen för skattningen är avgörande för hur bra vår gissning blir. Det är skattningen som är vår gissning!

Ex 1 (forts): Coca Cola frågar 1000 personer. Låt

$$X_k = \begin{cases} 1 & \text{om person } k \text{ föredrar CCX} \\ 0 & \text{— — — — — Pepsi max} \end{cases}$$

Vår data är $x_1, x_2, \dots, x_{1000}$. För att gissa
värdet på μ så tar vi $\bar{x} = \frac{x_1 + \dots + x_{1000}}{1000}$
dvs andelen som föredrog CCX i undersökningen.

Är $\mu = \bar{x}$? Nej!!! Är μ nära \bar{x} ? Förhoppningsvis

Ex 2 (forts): Pigen säger 500 av sina frallor. Vikterna

~~blev~~ blev x_1, x_2, \dots, x_{500} (data).

$\bar{x} = \frac{x_1 + x_2 + \dots + x_{500}}{500}$ är skattaren för μ

$s^2 = \frac{1}{499} \sum_{k=1}^{500} (x_k - \bar{x})^2$ är skattaren för σ^2

OBS! Det finns ett väldigt viktigt före/efter
perspektiv.

- 1) När modellen sätts upp (före mätning)
är observationerna slumpvariabler (X_k)
- 2) Efter att vi har observerat vårt stick-
prov (efter mätning) så är observationerna
numeriska värden (x_k).

Ex 1 (forts.): Innan mätningarna är gjorda så

har vi $\bar{X} = \frac{X_1 + X_2 + \dots + X_{1000}}{1000}$ som är

en slumpvariabel. Efter mätningarna har vi

$\bar{x} = \frac{x_1 + x_2 + \dots + x_{1000}}{1000}$ som är ett numeriskt värde.

Skattningarna \bar{x}, \bar{x}, s^2 i exemplena ovan kallas för punktskattningar eftersom vi ~~skä~~ får ett värde. Senare kommer vi prata om intervallskattningar.

I "före" perspektivet så har vi slumpvariabler \bar{X} och S^2 , dessa kallas för skattare medan de numeriska värdena \bar{x} och s^2 i "efter" perspektiv kallas för skattningar.

Väntevärdesriktighet

En skattare kallas för väntevärdesriktig om väntevärdet ~~för~~^{för} skattaren är lika med värdet på parametern.

Ex 1 (forts.): $E(\bar{X}) = E\left(\frac{X_1 + X_2 + \dots + X_{1000}}{1000}\right) =$

$$= \frac{E(X_1) + E(X_2) + \dots + E(X_{1000})}{1000} = \frac{P + P + P + \dots + P}{1000} =$$

$$= \frac{1000P}{1000} = P$$

Ex 2 (forts). Även här har vi

$$E(\bar{X}) = E\left(\frac{X_1 + \dots + X_{500}}{500}\right) = \frac{E(X_1) + \dots + E(X_{500})}{500} = \frac{\mu + \mu + \dots + \mu}{500} = \frac{500\mu}{500} = \mu$$

Vi har också

$$E(S^2) = E\left(\frac{1}{499} \sum_{k=1}^{500} (X_k - \bar{X})^2\right) = \sigma^2$$

lång väntning

Låt oss fokusera på vår skattare för väntevärdet,

$$\text{dvs } \bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$$

1) Som ovan gäller att $E(\bar{X}) = \mu$

2) Vi har även att $\text{Var}(\bar{X}) = \frac{\text{Var}(X_1)}{n} = \frac{\sigma^2}{n}$

Med notationen $\sigma_{\bar{X}}^2 = \text{Var}(\bar{X})$ har vi att $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$

Observera att om $\sigma_{\bar{X}}^2$ är "stor" så har vi stor spridning. Därmed är sannolikheten stor att vår skattning \bar{X} missar målet rejält, dvs det är stor sannolikhet att \bar{X} faktiskt inte är nära μ .

Observera att om n är stort, dvs om vårt stickprov är stort, så är $\frac{\sigma^2}{n} = \sigma_{\bar{X}}^2$ litet så \bar{X} är nära μ med större sannolikhet!

PLOT!