

En population har en parameter vars värde är okänt. Vi vill (hitta) gissa värdet utan att kolla hela populationen.

Ex1: Coca Cola skall lansera den nya colan CCX som utmanare till Pepsi max. Hur stor del av populationen föredrar CCX framför Pepsi max?
(population = alla i Sverige/USA) värdet, parameter p = proportionen som föredrar CCX).

Ex2: Pigeas frallor har normalfördelad vikt med parametrar μ o σ^2 . Vad är dessa?
(kvalitetskontroll).

Hur skall vi gissa p, μ, σ^2 ?

1) vi samlar in data (stickprov)

2) Data sammanfattas i en teststatistiska

Fördelningen av teststatistiken är avgörande för hur bra vår gissning blir.

Ex1 (forts): CC frågar 1000 personer. Låt

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

Våra data är $x_1, x_2, \dots, x_{1000}$.

F7 (2)

För att gissa p tar vi $\frac{x_1 + x_2 + \dots + x_{1000}}{1000} = \bar{x}$

dvs andelen svarande som föredrog CCZ.

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

Ex 2: Påsen väger 500 frallor. Vikterna kallas x_1, x_2, \dots, x_{500} (data).

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

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

OBS!!! Det finns ett extremt viktigt före/etter (mätning) perspektiv.

1/ När modellen sätts upp (före mätning) är vikterna i ex 2 slumpvariabler ($\bar{x}_k \sim N(\mu, \sigma^2)$).

2/ Efter att mätningen genomförts är det numeriska data (x_k) .

Ex 1: Innan mätning genomförts har vi

$\bar{X} = \frac{\bar{x}_1 + \dots + \bar{x}_{1000}}{1000}$ som är en slumpvariabel.

Efter mätningen genomsnittet har vi

F7 (3)

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

Om väntevärdet på teststatistiken = värdet på sökta parametern så är teststatistiken Väntevärdesriktig (VVR) (unbiased)

Ex 1:
$$E[\bar{x}] = E\left[\frac{x_1 + \dots + x_{1000}}{1000}\right] = \frac{E[x_1] + \dots + E[x_{1000}]}{1000}$$

$$= E[x_1] = \rho \quad \text{så VVR!}$$

Ex 2: Även
$$E[s^2(\bar{x})] = E\left[\frac{1}{499} \sum_{k=1}^{500} (x_k - \bar{x})^2\right] = \sigma^2$$

så VVR ↑
lång räkning.

Låt oss fokusera på
$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k.$$

1/ Som ovan gäller att
$$E[\bar{x}] = E[x_1] = \mu$$

2/ vi har 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

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}.$$

OBS! Om $\sigma_{\bar{x}}$ är "stor" har vi stor spridning. Därmed är sann. stor att vår gissning missar målet rejält.

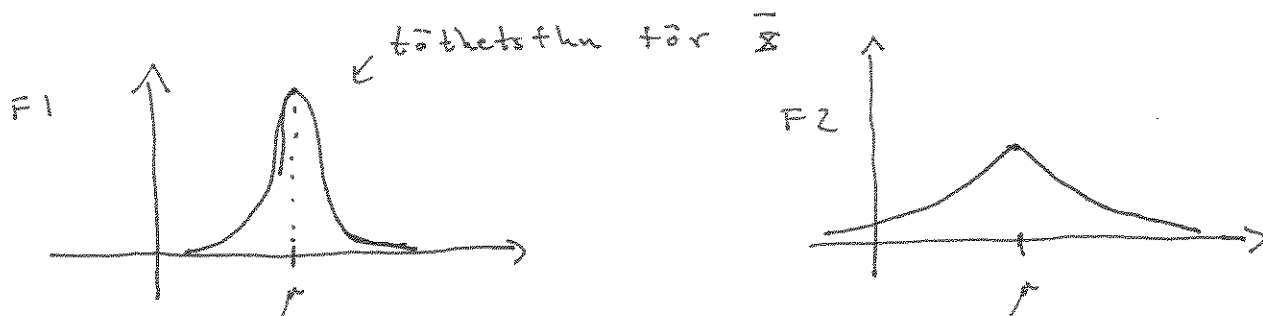
Sats 1: Om $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n$ är $N(\mu, \sigma^2)$
 (och oberoende) $\Rightarrow \bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

Sats 2: Om $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n$ har någon fördelning
 (vilken som helst!!!!) med $\mu = E[\bar{X}_i], \sigma^2 = \text{Var}(\bar{X}_i)$

gäller att $\bar{X} \approx N(\mu, \frac{\sigma^2}{n})$.
 ↑
 approximativt

Ann: Sats 2 kallas för CLT (Central Limit Theorem) och är en ~~de~~ statistikers viktigaste resultat. Om n är stort blir datans ursprung (nästan) irrelevant!

Varför är fördelningen hos \bar{X} viktig?



I F1 är det troligt att vi hamnar nära μ
 I F2 " " " " " " " långt från μ .