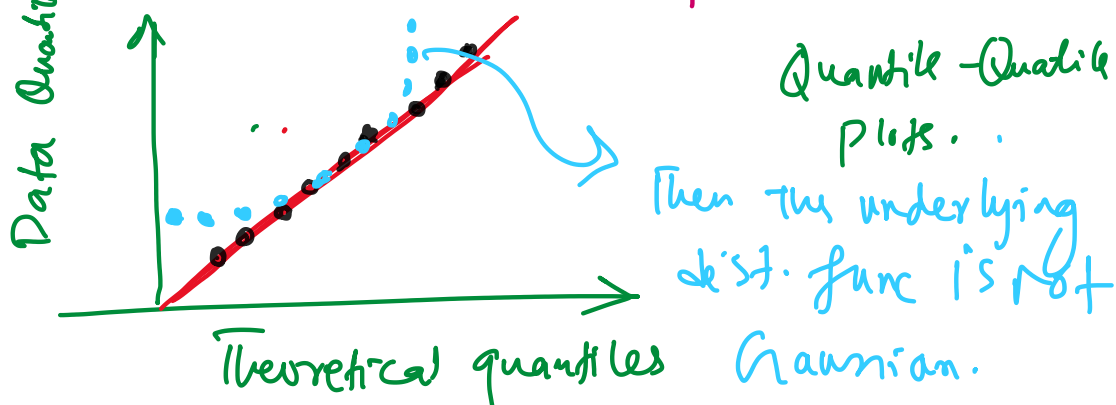


Monday, 7 February 2022 11:49

Sample: $\{x_1, x_2, \dots, x_n\}$

- underlying distribution?
- My guess is that this sample comes from the Gaussian distribution. \Rightarrow Hypothesis.



Data-Based Estimates of Descriptive Statistics

if you $P(x)$ probability dist.

Mean: $E(x) = \sum x_i P(x=x_i)$ & $E(x) = \int_{-\infty}^{+\infty} x P(x) dx.$

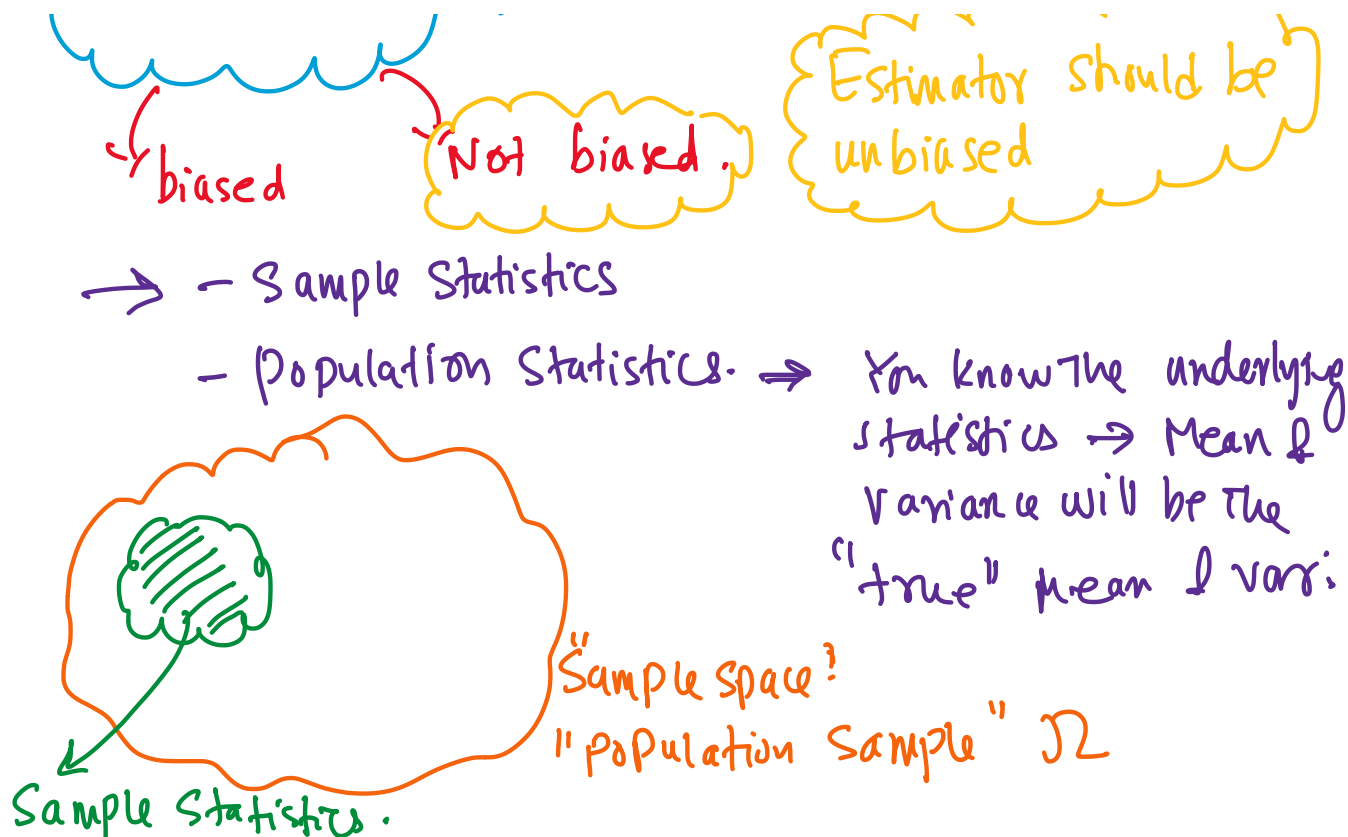
Var: $\text{Var}(x) = \sum_{x=i} E[(x-x_i)^2]$ & $E(x^2) = \int_{-\infty}^{+\infty} x^2 P(x) dx.$

$$\text{Var}(x) = E(x^2) - (E(x))^2 = E(x^2) - \mu^2.$$

Suppose now you don't know the underlying distribution function.

$\{x_1, x_2, \dots, x_n\}$

Estimators \rightarrow Statistics.



N Measurements. n_i : $i = 1, 2, \dots, N$ $\{x_i\}$

Sample Mean: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

Sample Standard deviation: $\bar{s} = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2}$

↓
 $N-1$

The reason for the $N-1$ term instead of the naively expected N is due to the fact that \bar{x} is also determined from the data. → Bessel's Correction.

True standard deviation is σ^2

Sample standard deviation is S^2 .

for a Gaussian dist. the underestimation varies from 20% for $N=2$ to 3% for $N=10$.

less than 1% for $N > 30$.

For the large sample size we reduce the 'underestimation' to less than 1%.

→ What do you mean by large Number?

→ It depends on your particular case and level of accuracy.

but generally the transition $N=10$ to $N \approx 100$.

"Massive" data set the transition may occur at N of the order of Million or even billion.

(\bar{X}, S)



estimators
from the sample.

(μ, σ)



Truth from the
population.

These estimators have a variance and a bias

Mean Squared Errors (MSE)

$$MSE = \sigma^2 + \text{bias}^2.$$

MSE = V + Bias

Bias: Expectation of the difference b/w the estimator and the truth value.

$$N \rightarrow \infty, V \rightarrow 0, \text{bias} \rightarrow 0$$

Consistent estimators

- Central Limit Theorem
- Law of Large Numbers.

(1) Law of large Numbers:

Let x_1, x_2, \dots be a sequence of independent random variables with a common distribution and $E(|x_i|) < \infty$ then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n x_i \rightarrow \mu = E[x_i]$$

as $n \rightarrow \infty$; for all $\epsilon > 0$

$$P(|\bar{X}_n - \mu| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Sample Mean \bar{X}_n gives a good approximation of the population mean $\mu = E(x_i)$ when N is large.

Sample Mean \rightarrow Random Quantity.
 Population Mean \rightarrow fixed Number - (Par).

(2) Central limit Theorem:

let x_1, x_2, \dots be a sequence of i.i.d. random variables with mean $\mu = E(x_i)$ and finite variance $\sigma^2 = E((x_1 - \mu)^2) > 0$ Then

$$P\left(\sqrt{n}(\bar{x}_n - \mu) < x\sigma\right) \rightarrow \bar{\Phi}(x)$$

$$\bar{\Phi}(x) = \int_{-\infty}^x \phi(t) dt \quad \text{and} \quad \phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}.$$

Distribution approaches the Gaussian dist for large sample size.

Next lecture:

- Examples
- Discuss some discrete I

Continuous distribution
functional.