# Sample properties

Often times, we do not know the underlying distribution (and thus the pdf) of a given dataset/population and therefore we cannot determine population properties.

$$E(X) = \mu = \int_{-\infty}^{\infty} x \cdot \boxed{f(x)} \, dx$$

If $f(x)$ is unknown, we can't compute $\mu, \sigma^2, \sigma$, etc.

An <u>estimator</u> $(\hat{\theta})$ is an expression used to estimate a statistical quantity $\theta$, <u>eg</u>: population mean or variance.

In statistical terms, one can define the <u>bias</u> of an

estimator.

$$B(\hat{\theta}) = E(\hat{\theta}) - \theta$$

bias of the estimator

expected value of the estimator

true value of the estimand.

$$B(\hat{\theta}) = 0$$
$$\left( \text{if } E(\hat{\theta}) = \theta \right)$$

$\hat{\theta}$ is called an UNBIASED ESTIMATOR.

$$B(\hat{\theta}) \neq 0$$
$$\left( \text{if } E(\hat{\theta}) \neq \theta \right)$$

$\hat{\theta}$ is called a BIASED ESTIMATOR.

(1) Sample mean $\longrightarrow$ unbiased estimator for the population

Given $n$ measurements constituting a sample, the sample mean is the arithmetic average of the values obtained for the random variable.

$n$ values: $x_1, x_2, \cdots, x_n$

$$\boxed{\overline{x} = \frac{1}{n}\left(\sum_{i=1}^{n} x_i\right)}$$

$$\overline{X} = \boxed{\mu} = \int_{-\infty}^{\infty} x\, f(x)\, dx.$$

true value of the population mean (estimand)

$$E(\overline{x}) = E\left(\frac{1}{n}\sum_{i=1}^{n} x_i\right)$$

$$= \frac{1}{n}\sum_{i=1}^{n} \boxed{E(x_i)} = \frac{1}{n}\sum_{i=1}^{n}\mu = \frac{1}{n}\times n\mu = \boxed{\mu}$$

each measurement is a random variable with the same underlying distribution

$$B(\bar{x}) = E(\bar{x}) - \mu = \mu - \mu = \boxed{0}$$

(2) **Sample variance**: is used as an estimator for the population variance.

$$\rightarrow s^2 = \frac{1}{\boxed{n-1}} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

sample mean.

is an unbiased estimator of $\sigma^2$ (population variance).

The factor of $(n-1)$ appears as the number of degrees of freedom <u>after</u> defining the sample mean is $(n-1)$.

$$\longrightarrow \tilde{s}^2 = \frac{1}{n} \sum_{i=1}^{n} (z_i - \bar{z})^2 \quad : \text{uncorrected sample variance}$$

(biased estimator)

③ <u>Sample standard deviation</u> is the square root of the sample variance.

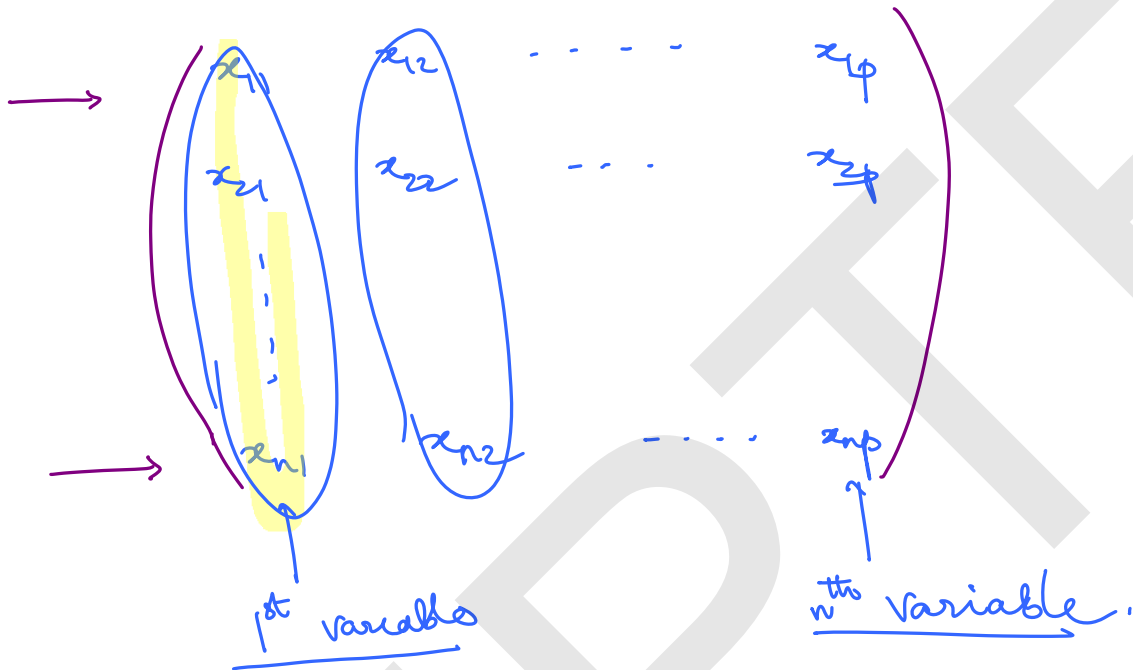$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2} \qquad \tilde{s} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2}$$

(stdev.p)

(stdev)
in Excel

④ <u>Sample covariance</u>

For two random variables X and Y, if there are n datapoints $(x_1, y_1)$, $(x_2, y_2)$, ...., $(x_n, y_n)$

$$q_{xy} = \frac{1}{n-1} \sum_{i=1}^{w} (x_i - \bar{x})(y_i - \bar{y}) \quad \longleftarrow$$

$1^{st}$ datapoint $\longrightarrow$

$n^{th}$ datapoint $\longrightarrow$

$$\begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}$$

$1^{st}$ variable

$n^{th}$ variable.

$$\vec{\bar{x}} = \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_p \end{pmatrix} = \left( \frac{1}{n} \sum_{i=1}^{n} x_{i1} \quad \frac{1}{n} \sum_{i=1}^{n} x_{i2} \cdots \frac{1}{n} \sum_{i=1}^{n} x_{ip} \right) = \frac{1}{n} \sum_{i=1}^{n} \begin{pmatrix} x_{i1} & x_{i2} \cdots \end{pmatrix}$$

$$\boxed{\vec{\bar{x}} = \frac{1}{n} \sum_{i=1}^{n} \vec{x}_i}$$

$$\boxed{\left(\overline{S_{jk}}\right) = S_{kj}} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{ij} - \overline{x}_j)(x_{ik} - \overline{x}_k)$$

$$j, k = 1, 2, 3, \cdots, \boxed{p.}$$

$$S = \begin{pmatrix} S_{11} & S_{12} & \cdots & S_{1p} \\ S_{21} & S_{22} & \cdots & S_{2p} \\ \vdots & & & \\ S_{p1} & S_{p2} & \cdots & S_{pp} \end{pmatrix}$$

Sample covariance matrix

$\longrightarrow$ diagonal entries are the sample variances

symmetric matrix since $S_{jk} = S_{kj}$.

$$S = \begin{pmatrix} \frac{1}{n-1}\sum_{i=1}^{n}(x_{i1}-\bar{x}_1)^2 & \frac{1}{n-1}\sum_{i=1}^{n}(x_{i1}-\bar{x}_1)(x_{i1}-\bar{x}_2) & \cdots & \frac{1}{n-1}\sum_{i=1}^{n}(x_{i1}-\bar{x}_1)(x_{ip}-\bar{x}_p) \\ & \vdots & & \\ \frac{1}{n-1}\sum_{i=1}^{n}(x_{ip}-\bar{x}_p)(x_{i1}-\bar{x}_1) & \cdots & & \frac{1}{n-1}\sum_{i=1}^{n}(x_{ip}-\bar{x}_p)^2 \end{pmatrix}_{p \times p}$$

$$= \frac{1}{n-1}\sum_{i=1}^{n} \begin{pmatrix} (x_{i1}-\bar{x}_1) \\ (x_{i2}-\bar{x}_2) \\ \vdots \\ (x_{ip}-\bar{x}_p) \end{pmatrix} \underbrace{\begin{pmatrix} (x_{i1}-\bar{x}_1) & (x_{i2}-\bar{x}_2) & \cdots & (x_{ip}-\bar{x}_p) \end{pmatrix}}_{1 \times p}$$

$p \times 1$

$$\boxed{S} = \boxed{\frac{1}{n-1}\sum_{i=1}^{n}(\vec{x}_i - \bar{\vec{x}})^T \cdot (\vec{x}_i - \bar{\vec{x}})} \longrightarrow \text{Sample covariance matrix.}$$