

# STATISTICS

- ① Population → Population mean }
  - median
  - mean
- ② Sample → Sample mean
- ③ Random variable ~~→ Differential Random variable~~
  - Discrete
  - Continuous

## # Population

Population mean  $\mu = \frac{1}{N} \sum_{i=1}^N x_i$

## # Sample

Sample mean  $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$

## # Random Variable

- Discrete  $\Rightarrow$  whole no., not a floating number  
Eg: No of bank a/c a person has,  
Population of a state
- Continuous  $\Rightarrow$  Within a range of values,  
we can have any value.  
whole no., floating nos.  
E.g: Height of a person

## # Gaussian Distribution or Normal Distribution

$$X \sim G.D(\mu, \sigma)$$

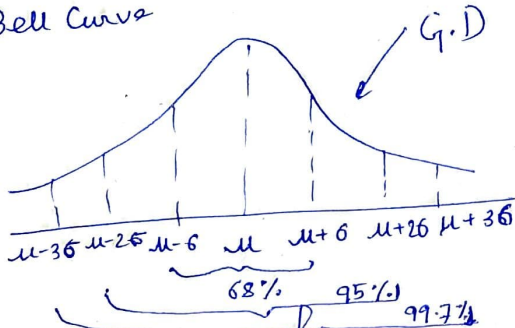
$\downarrow \quad \quad \downarrow$   
mean    s.d.

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{s.d} = \sigma = \sqrt{\text{Var}}$$

Bell Curve



Eg: Height of a population

Empirical formula in Gaussian Distribution

$$P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

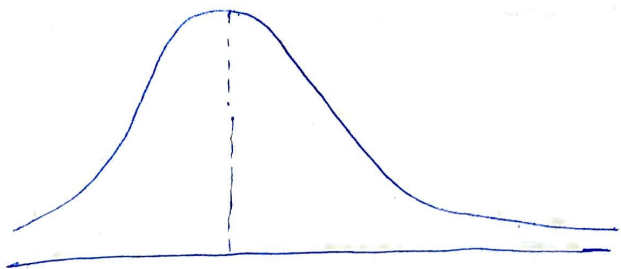
$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

## # Log Normal Distribution

$X \sim$  Log normal distribution  
if  $\ln(x)$  is normally distributed

i.e.  $\ln(x) \sim G.D(\mu, \sigma)$



Eg: Income of the people  
Product reviews

\* G.D  $\rightarrow$  Standard Normal Distribution (SND)  
( $\mu=0, \sigma=1$ )

$$= \frac{x_i - \mu}{\sigma}$$

# Covariance

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) * (y_i - \mu_y)$$

$x \uparrow \quad y \uparrow = \text{Cov} \rightarrow +ve$

$x \uparrow \quad y \downarrow = \text{Cov} \rightarrow -ve$

## # Mean, Median, Mode

## # Chebyshev's Inequality

$$\cancel{X \sim G.D} \quad X \approx G.D(\mu, \sigma)$$

$$P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

$$Y \not\sim G.D$$



$$P(\mu - k\sigma < x < \mu + k\sigma) > 1 - \frac{1}{k^2}$$

$$k = 2$$

$$P(\mu - 2\sigma < x < \mu + 2\sigma) > 1 - \frac{1}{2^2} = \frac{3}{4} = 75\%$$

## # Pearson Correlation Coefficient

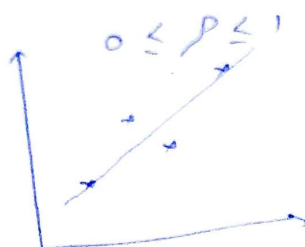
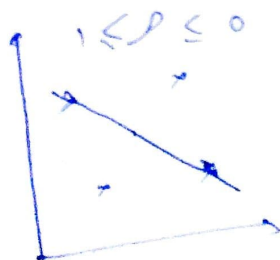
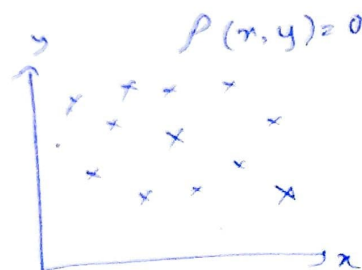
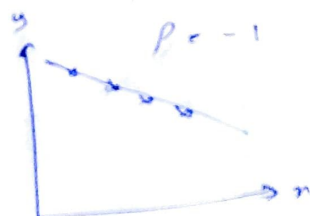
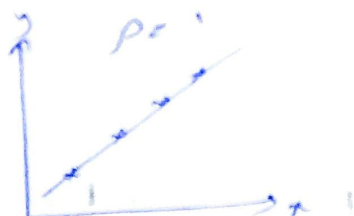
$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) * (y_i - \mu_y)$$

→ Gives the direction of relationship (x-y or x > y)

Pearson cc

$$P(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

$$P(x, y) = -1 \leq P \leq 1$$





## # Spearman's rank correlation

$$r = \rho_{r_g x, r_g y} = \frac{\text{cov}(r_g x, r_g y)}{\sigma_{r_g x} \sigma_{r_g y}}$$

$$r = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

$$r = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

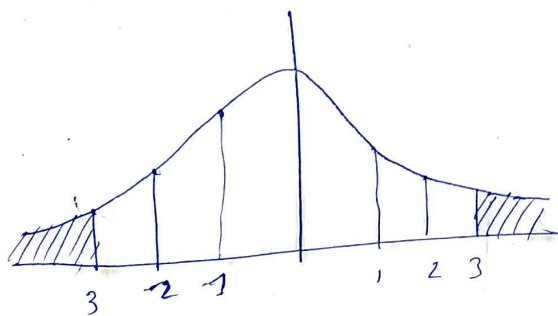
## # Z-score


$$\text{SND} \rightarrow \mu = 0 \quad \& \quad \sigma = 1$$

$$Z = \frac{x - \mu}{\sigma}$$

$$Z \leq 3 \rightarrow \text{Normal}$$

$$Z > 3 \rightarrow \text{outlier}$$

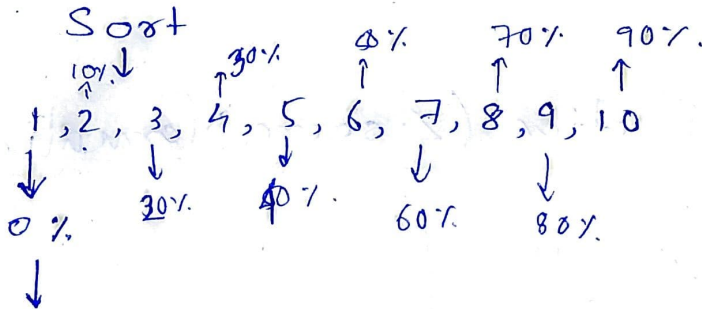


 ~~outlier~~  
Outliers

# # Inter Quantile Range (IQR)

## \* Percentiles

Let a set of numbers be 4, 9, 5, 6, 7, 1, 2, 8, 10



0 percentage of  
nos. lesser than 1

$$\boxed{IQR = 75\% - 25\%}$$

↓                      ↓

q3                      q1

~~lower bound =  $q1 * 1.5$~~

~~upper bound =  $q3 * 1.5$~~

$$\boxed{\begin{aligned} \text{lower bound value} &= q1 - (1.5 * IQR) \\ \text{upper bound value} &= q3 + (1.5 * IQR) \end{aligned}}$$

Any values away from lower bound & upper bound values ~~is~~ are outliers.

## # Normalization (min-max Normalization)

$$X_{\text{norm}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

\* Scale down your values between 0 and 1.

## # Standardization (Z-score Normalization)

$$Z = \frac{x - \mu}{\sigma}$$

\* Here all the features will be transformed in such a way that it will have the properties of a standard normal distribution with  $\mu = 0$  &  $\sigma = 1$

