# ✳ Variance & Standard Deviation ⤵

| Variance (Population Data) | Variance (Sample Data) |
|---|---|
| $$\sigma^2 = \sum_{i=1}^{N} \frac{(x-\mu)^2}{N}$$ | $$\sigma^2 = \sum_{i=1}^{n} \frac{(x-\bar{x})^2}{n-1}$$ |
| $$\boxed{std = \sqrt{\sigma^2}}$$ | $$std = \sqrt{\sigma^2}$$ |

Bessel error

$$\langle d.o.f = n-1 \rangle$$

# ✳ Covariance ⤵

$$\left\langle Cov(x,y) = \sum_{i=1}^{n} \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1} \right\rangle$$

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

$$\langle Cov(x,x) = Var(x) \rangle$$

# Adv. of Covariance → Quantify the relationship b/w X & Y

# Disadv. of Covariance → Covariance does not have a specific limit value
↳ $Cov(X, Y) \Rightarrow -\infty$ to $+\infty$

\* Correlation → Pearson Corr. Coeff.
→ Spearman "

① Pearson corr. coeff. $\Rightarrow [-1$ to $+1]$ → limit

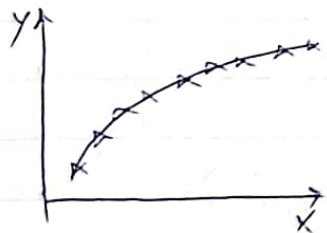↳ $$\rho_{x,y} = \frac{Cov(X, Y)}{\sigma_x \cdot \sigma_y}$$

↳ More the value to $+1$ → more +ve correlated X & Y

" " " " $-1$ → " $-ve$ " " "

Note:- Pearson corr. not suitable for non-linear data

② Spearman rank corr. ⇒

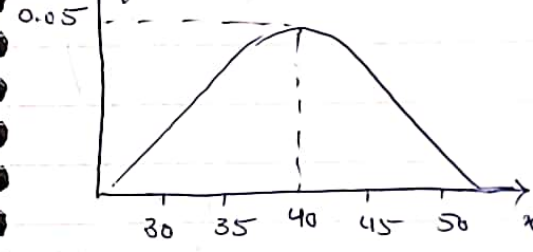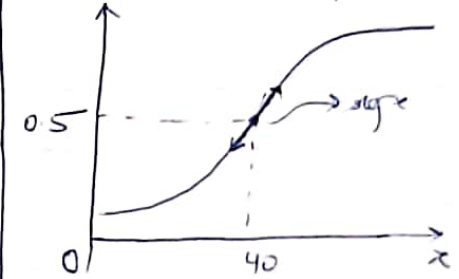$$r_s = \frac{Cov(R_{(x)}, R_{(y)})}{\sigma(R_{(x)}) * \sigma(R_{(y)})}$$

---

Ex:-

| x | y | R(x) | R(y) |
|---|---|------|------|
| 1 | 2 | 2 | 1 |
| 3 | 4 | 3 | 2 |
| 5 | 6 | 4 | 3 |
| 7 | 8 | 5 | 5 |
| 0 | 7 | 1 | 4 |

\* Probab. Density func"
↑ prob density.

Cummulative probab.



① probab. Mass func" (PMF) → used for discrete random variab.

② probab. Density func" (PDF) → used for contin. random varia

③ Cummu. Distributive func"

\* Probab. Density ⇒ Gradient of Cummulative density func"

## * p.d.f. properties

i) Non-negative; $f(x) \geq 0 \quad \forall x$

ii) The total area under the p.d.f. curve is equal to 1

$$\int_{-\infty}^{\infty} f(x) \, dx = 1$$
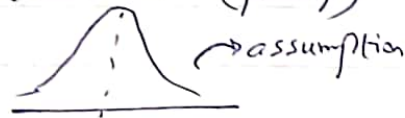
## # Types of Probab. Distribution → [p.d.f, p.m.f, c.d.f]

① **Bernoulli Distribution**
→ Outcomes are binary (p.m.f) like tossing of coins

② **Binomial Distribution** → p.m.f.

③ **Normal / Gaussian distrib.** → (p.d.f)


→ assumption

④ **Poisson Distribution** ⇒ (p.m.f)
⑤ **Log Normal Distrib.** ⇒ (p.d.f)
⑥ **Uniform Distrib.** (p.m.f)

---

## Example:-

Dataset → House price prediction Dataset

| Sie of House | No. of Rooms | Location | Floor | Sea side | Price |
|---|---|---|---|---|---|
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| Continuous random variable | Discrete | | Discrete | 0 & 1 | Continuos |
| pdf | p.m.f. | | p.m.f | p.m.f | pdf |

① **Bernoulli Distribution** < success (1) / fail (0)

↳ i) Discrete Random variable (p.m.f)

ii) Outcomes are Binary

↳ P(success) → (P)   so, P(lose) = (1-P)

NOTE:- $\boxed{P.m.f = p^k \times (1-p)^{1-k}}$

where $\{k = 0, 1\}$
       ↓     ↓
    failure  success

Also,

$$p.m.f = \begin{cases} q = 1-p & ; \text{ if } k = 0 \\ p & ; \text{ if } k = 1 \end{cases}$$

→ Mean of Bernoulli distrib: ⌐

$$\boxed{E(x) = \sum_{k=0}^{1} k \cdot P(k)}$$

→ Median of Bernoulli Distrib ⌐

$$Median = \begin{cases} 0 & \text{if } p < \frac{1}{2} \text{ (or) } q > \\ [0.1] & \text{if } p = \frac{1}{2} \text{ (or) } q = p \\ 1 & \text{if } p > \frac{1}{2} \text{ (or) } q < p \end{cases}$$

→ Mode of Bernoulli ⌐

$\langle p > q \rangle$ → $p$ will be mode
Else $q$ will be mode

→ Variance →

$$\boxed{\sigma^2 = pq = p(1-p)}$$
$$\sigma = \sqrt{pq} = \sqrt{p(1-p)}$$

---

② Binomial Distribution

→ Multiple sequence / cases of Bernoulli distribution
→ success
→ Failure [Binary]
→ discrete r.v.
→ perform for n-trials

Ex:- Tossing a coin 10 times → $\boxed{n = 10}$

↳ Notation → $B(n,p)$

parameters → $n \in \{0,1,2,---\}$ → No of trails
$p \in [0.1]$ → success probab.
for each trail
$q = 1-p$

support → $k \in \{0,1,2,---n\}$ → No. of success

* P.m.f for Binomial ⌐

→ $\langle P_r(k,n,p) = {}^nC_k \, p^k (1-p)^{n-k} \rangle$

for $k = 0,1,2,----n$ where

Binomial coeff. ← $\boxed{{}^nC_k = \frac{n!}{k!(n-k)!}}$

$\rightarrow$ Mean $\Rightarrow$ $n \cdot p$

$\rightarrow$ Variance $\Rightarrow$ $n \cdot p \cdot q$

$\rightarrow$ $\sigma$ (std) $\Rightarrow$ $\sqrt{n \cdot p \cdot q}$

③ Poisson Distribution ⌐

$\quad\quad\hookrightarrow$ Discrete r.v. (p.m.f)

$\quad\quad\hookrightarrow$ no. of events in fixed interval of time

Ex:- No. of people visiting hospital every hr
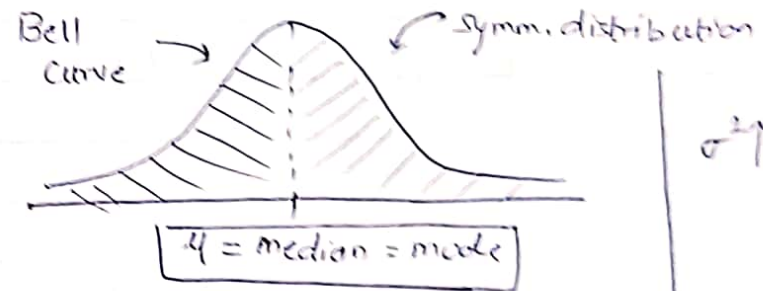
" " " banks " ,

$\Rightarrow \{\lambda \rightarrow$ Expected no. of events occurring at every time interval

$\rightarrow$ p.m.f for poisson $\Rightarrow$ $\boxed{p(x) = \dfrac{e^{-\lambda} \lambda^{x}}{x!}}$

$\rightarrow$ mean $= E(x) = \mu = \lambda * t$

$\quad\quad$ where $t \rightarrow$ time interval

④ Normal / Gaussian Distribution ⌐

$\quad\quad\hookrightarrow$ continuous r.v. (p.d.f)



Bell curve

Symm. distribution

$\mu = $ median $=$ mode

$\sigma \uparrow \Rightarrow$ spread $\uparrow$

Notation $\Rightarrow$ $N(\mu, \sigma^2)$

Parameters $\Rightarrow$ $\begin{cases} \mu \in R \longrightarrow \text{mean} \\ \sigma^2 \in R > 0 \longrightarrow \text{variance} \\ x \in R \end{cases}$
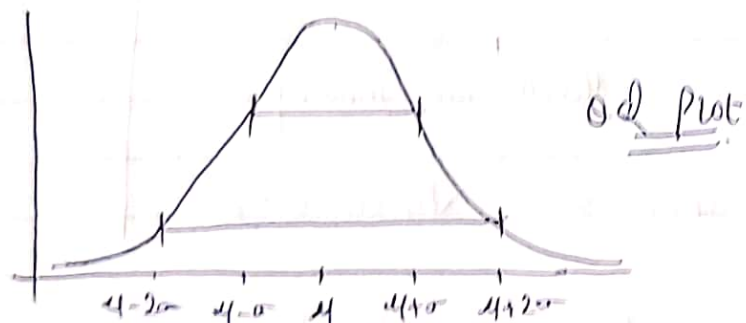
$\rightarrow$ $\boxed{\begin{array}{l} \text{p.d.f for Gaussian} \\[4pt] = \dfrac{1}{\sigma\sqrt{2\pi}} e^{*\left(\dfrac{x_i - \mu}{\sigma}\right)^2} \end{array}}$

$\rightarrow$ Mean $\Rightarrow$ $\mu = \sum\limits_{i=1}^{n} \dfrac{x_i}{n}$

→ Variance ⟹

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

* Emperical rule of Gaussian distrib.



0d Plot

$$\boxed{66 - 95 - 99.7} \text{ Rule}$$

$$\hookrightarrow P(\mu - \sigma \leq x \leq \sigma + \mu) \simeq 66\%$$

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

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

NOTE :- Std. Normal Distrib. is the one when we convert p.d.f to diff. p.d.f with $\boxed{\mu = 0 \text{ , } \sigma = 1}$

│ Can be converted through →

$$\boxed{Z\text{-score} = \frac{x_i - \mu}{\sigma}}$$

*** 
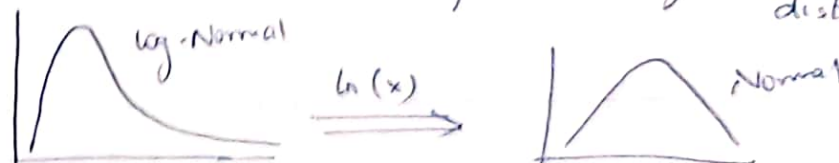( We do standardization just to bring every columns 
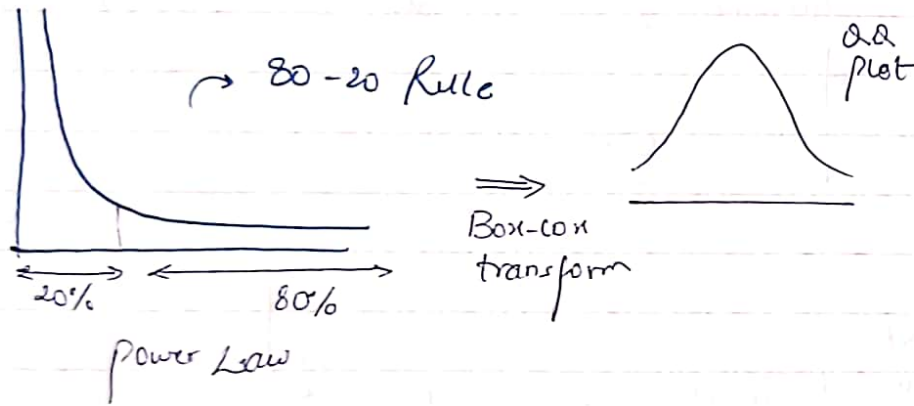to same unit of measurement )

* Log Normal Distribution
$$\hookrightarrow p.d.f \longrightarrow \text{random variable}$$

If r.v "x" is log normally distributed, the $\boxed{Y = \ln(x)}$ → Normal Distribution

Simi,

$$Y \longrightarrow \text{normally distributed, then}$$
$$X = \exp(Y) \longrightarrow \text{log-normally distributed}$$



Log-Normal    $\ln(x)$    Normal

Scanned with CamScanner

§ Power-law Distrib. ⌐



→ 80-20 Rule

QQ plot

⟹

Box-cox transform

20%      80%

Power Law

---

§ Inferential Statistics ⌐ { Hypothesis testing }

↳ conclusion or Inference



sample

Conclusion

↳ Hypothesis testing

population

---

☑ Hypothesis Testing Mechanism

① Null Hypothesis (Ho)
↳ assumption to begin with.

---

② Alternate Hypothesis (H₁)
↳ opp. of (Ho)

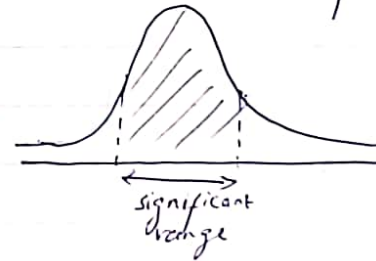③ Experiments ⟶ statistical Analysis
⟶ Correct Proof

④ Accept the (Ho) or Reject the (Ho)

# P-Value ⟹
use to accept or reject the null Hypothesis
↳ lies in significant value range →Accept Ho
if not       → Reject Ho



significant range

---

☀ Hypothesis Testing & Statistical Analysis ⌐

① Z-Test     ⎫→Average ⟹    Z table → Z score & P value
② t-Test     ⎭              t table
③ Chi-Square  ⟹ Categorical data
④ Annova      ⟹ Variance

NOTE :- For Z test → we need pop. std. deviation
& n > 30

So,

$$Z = \dfrac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$ for pop. data.

↳ Now, for case when we don't know pop. std. devi.
then we use t-test

$$t = \dfrac{\bar{x} - \mu}{S/\sqrt{n}}$$ S = sample std. deviation

→ we use dof = n-1

* Type 1 & type 2 Errors ⌐

Outcome 1 :- We reject the null hypo. when in reality
it is false → Good

Outcome 2:- We reject " " " " it is true
→ Type 1 Error

Outcome 3:- We retain the null hypo. when in reality
it is false → Type 2 Error

Outcome 4:- " " " " " it is true → Good
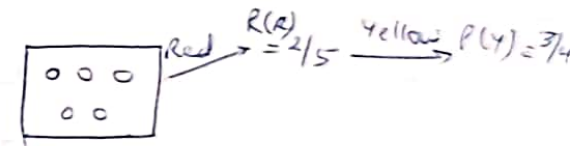
---

* Bayes Statistics (Bayes Theorem) ⌐

probab. ⟶ Independent Events
⟶ Dependent Events

① Independent Events | ② Dependent Events

Ex:- Rolling a dice | Red →$\frac{R(R)}{=2/5}$ Yellow $P(Y)=3/4$

$P(x) = \frac{1}{6}$ |

$$P(R \& Y) = P(R) \times P(Y/R)$$

conditional probab.

$= 2/5 \times 3/4 = 6/20$

↳ P(A and B) = P(B and A)

$P(A) \times P(B/A) = P(B) \times P(A/B)$

$$P(B/A) = \dfrac{P(B) \times P(A/B)}{P(A)}$$

$$P(A/B) = \dfrac{P(A) \times P(B/A)}{P(B)}$$

$A, B \rightarrow$ events

$\begin{bmatrix} P(A/B) \Rightarrow \text{prob. of } A \text{ given } B \text{ is true} \\ P(B/A) \Rightarrow \quad " \quad " \quad B \quad " \quad A \text{ is true} \\ \\ P(A), P(B) \Rightarrow \text{Independ. probab of } A \& B \end{bmatrix}$

* Confidence Interval �len

     Z-test $\Rightarrow$ point Estimate $\pm$ Margin Error

     High & Low             $\bar{x} \pm Z_{\alpha/2} \sigma/\sqrt{n}$
         C.I.

     t-test $\Rightarrow \bar{x} \pm t_{\alpha/2} \sigma/\sqrt{n}$

* Chi - Square Test �len

→ The chi sq. test for goodness of fit test claims about population proportions.

→ It is a non-parametric test i.e. performed on categorical [ordinal & nominal] data.

$$\boxed{X^2 = \Sigma \frac{(0 - E)^2}{E}}$$

     $0 \rightarrow$ observed
     $E \rightarrow$ Expected