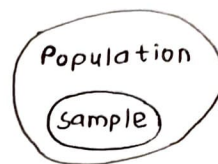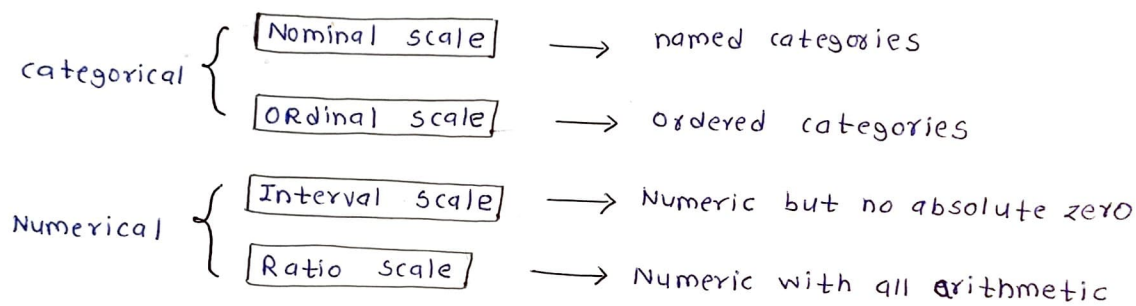# Statistics-1

## * Basics :

→ 2 branches : (a) Descriptive statistics
(b) Inferencial statistics

Population
Sample

→ Cross-sectional data : data collected at a point of time
Time-series data : data collected over a period of time

→ scales of measurement / data types :

categorical { 
| Nominal scale | ⟶ named categories
| ORdinal scale | ⟶ ordered categories

Numerical {
| Interval scale | ⟶ Numeric but no absolute zero
| Ratio scale | ⟶ Numeric with all arithmetic

## * Categorical Data :

→ Frequency distribution table / Relative frequency table
→ charts : (a) Pie chart
(b) Bar chart
(c) Pareto chart ( frequency sorted bar chart)
→ mode : Variable / Category with max. frequency
median : Variable / category of middle observation

→ continuous data category :
20-30 :  20 is included but 30 does not.
Upper limit is 30 and lower limit is 20.
10 is class width. (i.e. 30-20 =10)
25 is class mark. (i.e. $\frac{20+30}{2} = 25$)

## * Numericale Data:

→ Sample size : $n$     Sample mean : $\bar{x}$     Sample S.D. : $S$

    Population Size : $N$     Population mean : $\mu$     Population S.D. : $\sigma$

→ Mean:    $\dfrac{\Sigma f_i m_i}{\Sigma f_i}$   or   $\dfrac{\Sigma x_i}{n}$

→ Adding constant :   $y_i = x_i + c \Rightarrow \bar{y} = \bar{x} + c$

   multiplying with constant : $y_i = cx_i \Rightarrow \bar{y} = c\bar{x}$

---

→ median : middle value in ordered List

    (a) $n$ is odd : $\dfrac{n+1}{2}^{th}$ observation.

    (b) $n$ is even : average of $\dfrac{n}{2}^{th}$ and $\dfrac{n}{2}+1^{th}$ observation.

→ Adding constant :   New median = Old median + C

   multiplying with constant : New median = $c \times$ old median

---

→ mode : most frequently occuring value

→ Adding constant : New mode = Old mode + c

   multiplying with constant : New mode = $c \times$ old mode

---

→ Range : Diff. of largest & smallest value

---

→ Variance: Let $x_i - \bar{x}$ is deviation.

$$S^2 = \dfrac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \ldots + (x_n - \bar{x})^2}{n-1}$$

$$\sigma^2 = \dfrac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \ldots + (x_N - \mu)^2}{N}$$

→ Adding constant : New variance = old variance

   multiplying with constant : New variance = $c^2 \times$ old variance

→ Standard Deviation : Positive square root of variance

→ Adding constant : New S.D. = old S.D.

   multiplying with constant : New S.D. = C × old S.D.

---

→ measure of centrle tendency : Mode, median & mean

   measure of dispersion/spread : Variance & S.D.

---

→ Percentile :   x percentile means    $x\%$ of data are ≤ it and ;

                  $(100 - x)\%$ of data are ≥ it.

→ First / Lower Quartile ($25^{th}$ percentile)   $(Q_1)$

   Second / median Quartile ($50^{th}$ percentile) $(Q_2)$

   Third / upper Quartile ($75^{th}$ percentile) $(Q_3)$

* **Five-Number Summary :**

   1. minimum
   2. Q1
   3. Q2                      $IQR = Q_3 - Q_1$
   4. Q3
   5. maximum

* **Scatter Plot :**

   → x-axis : Explanetory / Independent variable
      y-axis : Responsive / Dependent variable

   → Gives insight about association

   → Covairance and corelation are measure of linear association.

## * Covariance :

→ Population's covariance = $\dfrac{\sum\limits_{i=1}^{N} (x_i - \bar{\mu}_x)(y_i - \bar{\mu}_y)}{N}$

Sample's covariance = $\dfrac{\sum\limits_{i=1}^{N} (x_i - \bar{x})(y_i - \bar{y})}{n-1}$

→ If covariance $> 0$ ⟹ Positive association

covariance $< 0$ ⟹ Negative association

→ Covariance has unit i.e. unit of $x_i$ × unit of $y_i$

## * Correlation :

→ Pearson correlation $\rho = \dfrac{\text{Covariance}}{\text{S.D. of } x_i \times \text{S.D. of } y_i}$

∴ $\rho = \dfrac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2}\,\sqrt{\sum(y_i - \bar{y})^2}}$        Range $(\rho)$ is $[-1, 1]$

$$\underset{\text{Strong}}{\overset{-1}{\bullet}} \quad \xleftarrow{\text{weak}} \overset{0}{\circ} \xrightarrow{\quad} \quad \underset{\text{strong}}{\overset{+1}{\bullet}}$$

→ Correlation coefficient does not have any units.

→ $\rho^2$ or $R^2$ gives goodness of straight line / curve fitting. $R^2 \in [0, 1]$

\* <u>Association b/w numeric and categorical variables</u> :

→ Point - Bi - Serial correlation coefficient $r_{ps}$

$$r_{ps} = \left( \frac{\bar{Y}_0 - \bar{Y}_1}{S_x} \right) \sqrt{P_0 P_1}$$

where ; $P_0$ = Proportion of categorical variable coded with 0
$P_1$ = Proportion of categorical variable coded with 1
$\bar{Y}_0$ = mean of categorical variable coded with 0
$\bar{Y}_1$ = mean of categorical variable coded with 1
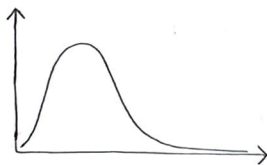$S_x$ = S.D. of numerical variable

\* <u>Infinite series</u> :

→ $\sum\limits_{n=0}^{\infty} ar^n = \dfrac{a}{1 - r}$   where $|r| < 1$

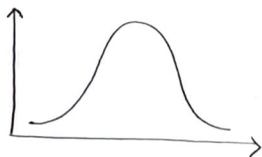→ $\sum\limits_{n=0}^{\infty} \left( \dfrac{x^n}{n!} \right) = e^x$
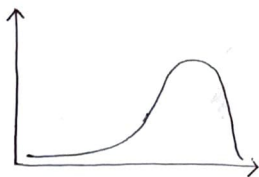
\* <u>Distributions</u> :

(a) Right skewed
(Positively skewed)



(b) Symmetric



(c) Left skewed
(Negatively skewed)

# * Permutation :

→ Definition: Ordered arrangement of all or $r$ objects from $n$.

$$nP_r = \frac{n!}{(n-r)!} \quad \text{(without repetition)}$$

$$= n^r \quad \text{(with repetation)}$$

→ Non-distinct permutation :

$$nP_n = \frac{n!}{P_1!P_2! \ldots P_K!} \quad \text{where} \quad P_i = \text{group of same objects}$$

→ Circular Permutation :

cw & ccw are considered different $= (n-1)!$

cw & ccw are considered same $= \dfrac{(n-1)!}{2}$

# * Combination :

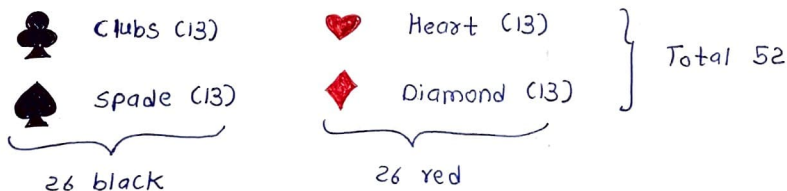→ Definition : selecting / choosing $r$ objects from $n$.

$$nC_r = \frac{nP_r}{r!} \quad \text{or} \quad \frac{n!}{r!(n-r)!}$$

→ Identity : $\quad nC_r = nC_{n-r}$

$$nC_r = n-1C_r + n-1C_{r-1}$$

$$\sum_{r=0}^{n} nC_r = 2^n$$

# * Playing cards :

♣ Clubs (13)          ♥ Heart (13)

♠ Spade (13)          ♦ Diamond (13)          } Total 52

‿‿‿‿‿                    ‿‿‿‿‿
26 black               26 red

# * Axiometic Probability :

→ P(E) is in accord with ;                    — $(E_1 \cap E_2 = \phi)$

(i)   $0 \leq P(E) \leq 1$

(ii)  $P(S) = 1$

(iii) For <u>mutually exclusive</u> / disjoint events $E_1, E_2, \dots E_n$

$$P(\cup E_i) = \Sigma P(E_i)$$

→ General properties :

(a)  $P(E^c) = 1 - P(E)$

(b)  $P(\phi) = 0$

(c)  $P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2)$

→ other approaches of probability are ;

• classical / Apriori / Theoretical approach

• Relative frequency / Aposteriori / Impirical approach

• Subjective approach

# * Independent Event :

→ Two event E & F are independent if and only if

(i)   $P(E \cap F) = P(E) \times P(F)$

→ Three event E, F, & G are independent if and only if

(i)   $P(E \cap F \cap G) = P(E) \times P(F) \times P(G)$

(ii)  $P(E \cap F) = P(E) \times P(F)$

(iii) $P(F \cap G) = P(F) \times P(G)$

(iv)  $P(G \cap E) = P(G) \times P(E)$

→ If E & F are independent, then following are also independent.

(i)   E and $F^c$

(ii)  $E^c$ and F

(iii) $E^c$ and $F^c$

# * Conditional Probability :

→ Probability of E conditioned on F :

$$P(E/F) = \frac{P(E \cap F)}{P(F)} \quad ; \quad P(F) > 0$$

→ If E and F are independent then $P(E/F) = P(E)$.

→ multiplication rule :

- $P(E \cap F) = P(E) \cdot P(F|E)$

- $P(E_1 \cap E_2 \cap E_3 \cap \ldots \cap E_n) = P(E_1) \cdot P(E_2|E_1) \cdot P(E_3|E_1 \cap E_2) \ldots P(E_n | E_1 \cap E_2 \ldots \cap E_{n-1})$

# * Total Probability :

→ Let's $F_1, F_2, \ldots, F_K$ are mutually exclusive and exhaustive then for any event E;

$$P(E) = \sum_{i=1}^{k} P(E/F_i) \cdot P(F_i)$$

# * Bay's Rule :

→ Let's $F_1, F_2, \ldots, F_K$ are mutually exclusive and exhaustive then for any event E;

$$P(F_i | E) = \frac{P(E|F_i) \cdot P(F_i)}{\sum_{i=1}^{k} P(E|F_i) \cdot P(F_i)}$$

# * Random variable :

→ Definition : The quantities of interest or real valued functions defined on the sample space are known as random variable.

→ Discrete random variable
Continuous random variable

* <u>Probability mass Function (PMF)</u>:

→ PMF for discrete random variable defined as;

$P(x_i) = P(X = x_i)$  (i.e. probability of occuring $x_i$)

→ General properties:

$P(x_i) \geq 0$

$\sum_{i=1}^{\infty} P(x_i) = 1$

* <u>Cumulative Distribution Function (CDF)</u>:

→ CDF for discrete random variable defined as;

$F(a) = P(x < a)$  (where $x_1 < x_2 < \ldots < x_n$)

→ For discrete random variable CDF is step function.

| | |
|---|---|
| $F(x_1)$ | 0 |
| $F(x_2)$ | $0 + P(x_1)$ |
| $F(x_3)$ | $0 + P(x_1) + P(x_2)$ |
| ⋮ | |
| $F(x_n)$ | $P(x_1) + P(x_2) + \ldots + P(x_{n-1})$ |
| $F(x_{n+1})$ | 1 |

* <u>Expectation of Random Variable</u>:

→ Expectation of Random variable $x_i$ where $i = 1$ to $n$ is;

$E(x) = \sum_{i=1}^{\infty} x_i P(x_i)$

→ It is 'Long run Average' value of random variable.

**\* Variance of Random Variable :**

→ Variance of Random variable $x_i$ where $i = 1$ to n is;

$$V(x) = E((x - \mu)^2) \quad \text{where } \mu = E(x)$$

$$\therefore V(x) = E(x^2) - E(x)^2$$

**\* standard Deviation of Random variable :**

→ It is positive square root of variance of random variable.

**\* Proposition Rules :**

→ 
$$E(ax + b) = a E(x) + b$$
$$V(ax + b) = a^2 V(x)$$
$$SD(ax + b) = a SD(x)$$

→ $E(x) + E(y) = E(x + y)$ is always true

→ $V(x + y) = V(x) + V(y)$ is only true when $x$ & $y$ are independent.

→ $SD(x + y) = \sqrt{V(x) + V(y)}$ is only true when $x$ & $y$ are independent.

**\* Bernoulli Random Variable :**

→ A random variable that takes either 0 or 1.

| X | 0 | 1 |
|---|---|---|
| $P(X = x_i)$ | $1 - p$ | $p$ |

$E(x) = p$

$V(x) = p(1 - p)$

**\* Uniformly Distributed Random variable :**

→ A random variable that takes values 1 to n.

| X | 1 | 2 | | n |
|---|---|---|---|---|
| $P(X = x_i)$ | $\frac{1}{n}$ | $\frac{1}{n}$ | …… | $\frac{1}{n}$ |

$E(x) = \dfrac{n + 1}{2}$

$V(x) = \dfrac{n^2 - 1}{12}$

# ✱ Hypergeometric Random Variable:

→ Let's there is two category in population of size N. If size of category 1 is $m$, then size of another category is $N-m$. Choosing a sample of size $n$ is having $i$ from category 1 then;

$$P(X=i) = \frac{\binom{m}{i}\binom{N-m}{n-i}}{\binom{N}{n}} \quad , \quad i = 0,1,2,\dots n$$

Here, $X$ is called hypergeometric random variable.

→ $E(X) = \dfrac{nm}{N}$

$$V(X) = \frac{nm}{N}\left[\frac{(n-1)(m-1)}{(N-1)} - \frac{nm}{N} + 1\right]$$
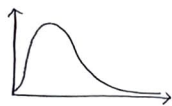
## ✱ Binomial Random Variable:

→ Independent and identically distributed (IID) Bernouli random variables is called Binomial Random variable because it's PMF is binomial.

→ Let $n$ Bernouli trails are performed with probability of success in each trial $p$. Let $X$ denotes the no. of successes in $n$ trials then;

$$P(X=i) = \binom{n}{i} p^i (1-p)^{n-i} \qquad (X \sim B(n,p))$$

## ✱ Graph of Binomial Distribution:

(i) $P < 0.5$ and $n$ small.



(ii) $P > 0.5$ and $n$ small.



(iii) $P = 0.5$



(iv) $n$ large.

Approches symmetry.

\* Expectation & Variance of Binomial Distribution :

$\rightarrow$ $E(X) = nP$

$V(X) = np(1-p)$

$\rightarrow$ $P = 1 - \dfrac{V(X)}{E(X)}$   (only for Bernoulli & Binomial Random variables.)

\* Straight Line Fit by Min(SSE) :

$\rightarrow$ Square Sum Error (SSE) $= \sum ( y_i - (mx_i + c))^2$

$\rightarrow$ $m = \dfrac{n\sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$

$c = \bar{y} - m\bar{x}$