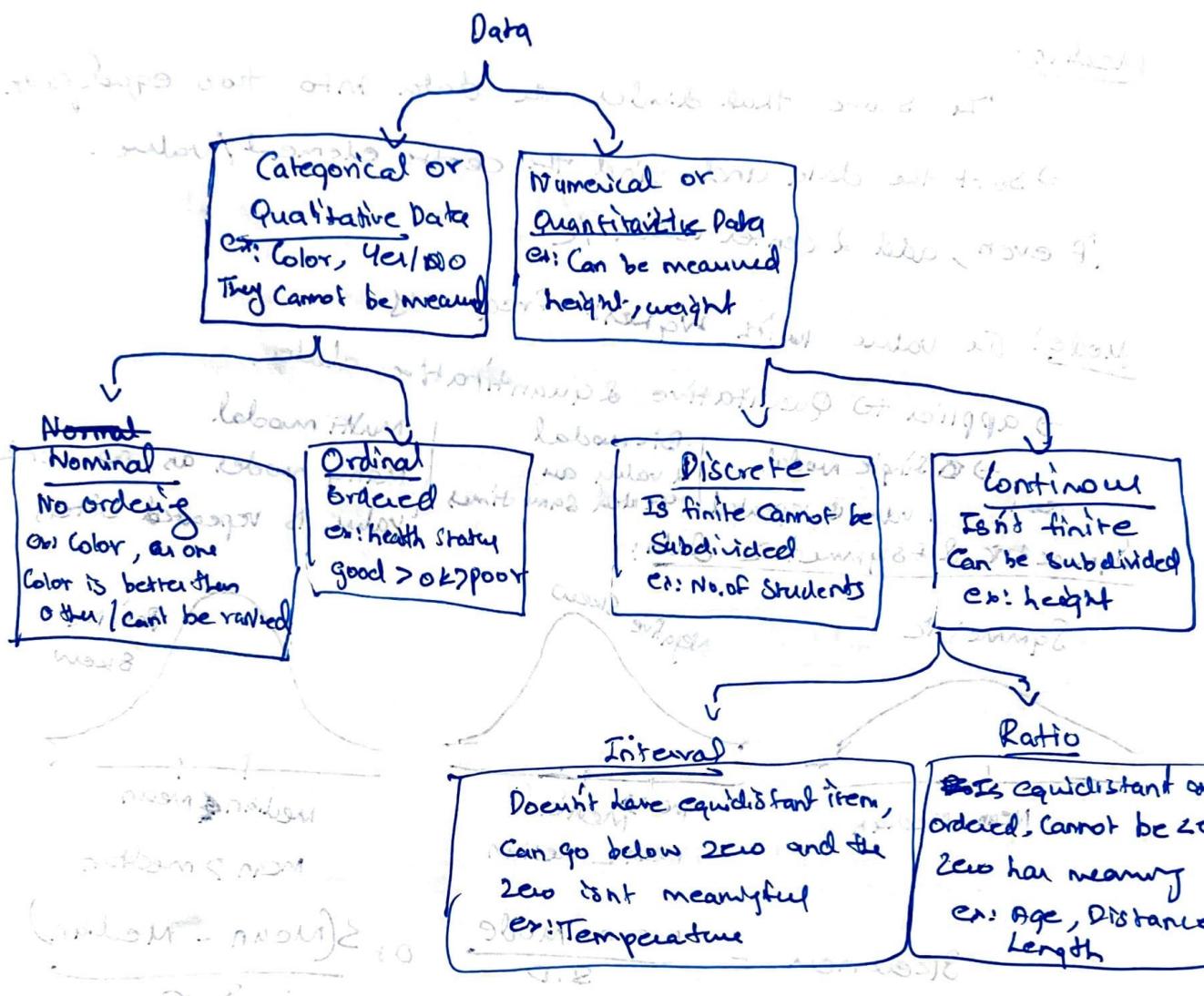


→ Statistics is the study of collection, analysis, interpretation, presentation and organization of data.

→ Types of data



Measure of Central Tendency:

Used to describe a set of ~~scores~~ with a single number that describes the performance of the group → center of data

- 3 used measures for central tendency

- Mean
- Median

Mode

- Mean: The arithmetic average $\bar{Y} = \frac{\sum Y}{N}$ Y is the value, N is the count of values

$$\text{ex: } 1, 2, 3, 4, 5 \quad \bar{Y} = \frac{1+2+3+4+5}{5} = 3$$

If presented in group for

→ A subset of population

\bar{Y} for Sample, μ for population

2	1	2	3
f	2	1	3

$$\bar{Y} = \frac{\sum f Y}{N} = \frac{1 \times 2 + 2 \times 1 + 3 \times 3}{6} = 2.166$$

→ entire data/group

→ Mean is the most stable and popular as it includes all data

→ Affected by extreme scores

→ Sum of each score's distance from mean is zero

→ Use when we want to measure standard deviation, coefficient of variation and skewness, and stability is desired.

Median:

The score that divides the data into two equal parts.

→ Sort the data and find the centre element/value.

If even, add 2 center value / 2

Mode: The value with highest frequency.

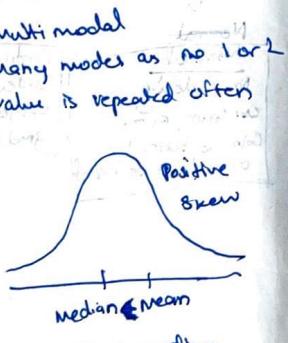
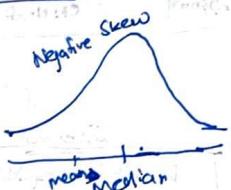
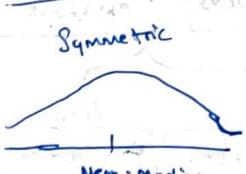
→ applies to Qualitative & quantitative data.

→ Single modal

→ Bi-modal 2 values are repeated sometimes

→ Multi-modal Many modes as no 1 or 2 values are repeated often

Symmetric & Asymmetric Data:



$$\text{Skewness} = \frac{\text{Mean} - \text{Mode}}{\text{S.D.}} \quad D = \frac{3(\text{Mean} - \text{Median})}{\text{S.D.}}$$

SD = Standard deviation

Measures of Variability

- Range
- Variance
- Standard deviation
- Variability is the distance measure of the difference between scores
- It describes the spread of scores or dist. of score from mean
- It captures how well the data is distributed and how well an individual score represents the distribution.

Range: Distance between smallest value to largest

$$= \text{Max} - \text{Min value} = \text{Range}$$

Standard Deviation: It represents how distributed the data is from mean

most common & imp measure of variability

→ Measures the standard or avg distance from mean

Ex: take the data 1, 3, 7, 8, 9

The mean is 6

$$\text{Deviation} = x_i - \bar{x} = (x_i - 6)$$

$$\text{Standard Deviation} = \sigma = \sqrt{\text{Variance}} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

$$= \sqrt{(1-6)^2 + (3-6)^2 + (7-6)^2 + (8-6)^2 + (9-6)^2}$$

$$= \sqrt{(1-6)^2 + (5-6)^2 + (7-6)^2 + (8-6)^2 + (9-6)^2} = \sqrt{25+1+1+1+9} = \sqrt{\frac{40}{5}} = \sqrt{8}$$

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N} \quad \text{for population}$$

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} \quad \text{for sample}$$

In a Sample
or below S
in below Z

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} \quad \text{Sum of squared differences}$$

Five point Summary: Represent data in 5 values

• Minimum

• Q1

• Median

• Q3

• Maximum

Interquartile range = Q3 - Q1

Q1 = IQR + (0.25 * LQ)

Q3 = IQR + (0.75 * UQ)

LQ = Median - (IQR * 1.5)

UQ = Median + (IQR * 1.5)

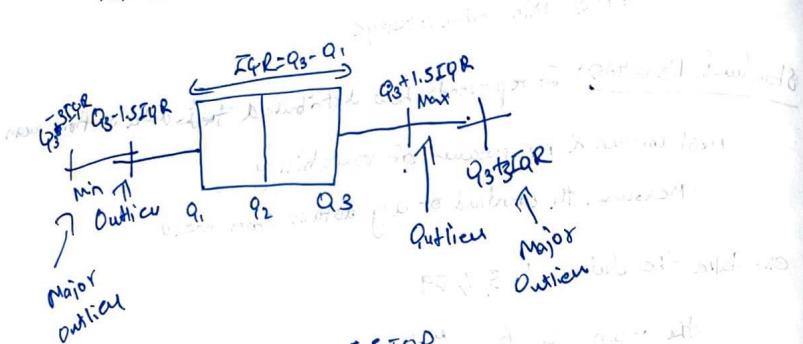
IQR = Q3 - Q1

Median = (Q1 + Q3) / 2

Outlier = Any value less than LQ - (IQR * 1.5)

or greater than UQ + (IQR * 1.5)

Box and Whisker Plot



$$\text{Lower limit} = Q_1 - 1.5 \text{ IQR}$$

Upper limit $Q_3 + 1.5 \text{ IQR}$. Values above or below these limits are potential Outliers

Probability:

→ Random experiment is used to describe an action whose outcome is not known in advance.

ex: tossing a coin

→ Sample space of a random exp. is a set S with all possible outcomes

→ Event is a subset of a sample space

Classical → Prob of an event = $\frac{\text{No. of favourable outcomes}}{\text{Total}}$ ex: $P(\text{Head}) = \frac{1}{2}$

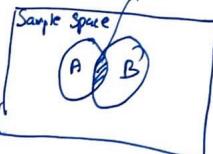
Empirical → $\frac{\text{No. of times the event occurs}}{\text{Total}}$ ex: on 1000 coin tosses 450 times Head

$$P(S) = \frac{450}{1000}$$

$$0 \leq P(S) \leq 1$$

$$P(E_1) \& P(E_2) \text{ such that } P(E_1 \cap E_2) = 0 \text{ then } P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

$$P(A \cap B)$$



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

2 Events are mutually exclusive if they don't have anything in common \Rightarrow Don't intersect
 $\hookrightarrow P(A \cap B) = 0$ (cannot happen at same time.)

→ An independent event is one in which one event doesn't affect the other
 → If two events happen with replacement, it's independent

$$P(A^c) \text{ compliment} = n(A^c) = n(S) - n(A)$$

\uparrow
sample space



Entire space

$$n(A \cap B^c) = n(A) - n(A \cap B)$$

$$n(A^c \cap B) = n(B) - n(A \cap B)$$

(a) Probability of event A or B not happening

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

If A & B are mutually exclusive

$$P(A \cup B) = P(A) + P(B)$$

For independent events $P(A \cap B) = P(A) \cdot P(B)$

Conditional probability:

Prob of A given B has happened
 \hookrightarrow conditioning event

$$P(A|B)$$

\hookrightarrow Happened $B \neq \emptyset$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Multiplication Rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \Rightarrow P(A \cap B) = P(A|B) \cdot P(B) \quad B \neq \emptyset$$

\hookrightarrow If A and B are independent events, then $P(A|B) = P(A)$

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2)P(A_4|A_1 \cap A_2 \cap A_3)$$

If A and B are independent events, then $P(A|B) = P(A)$

$$\therefore P(A|B) = \frac{P(A \cap B)}{P(B)} \Rightarrow P(A) = \frac{P(A \cap B)}{P(B)}$$

$P(A \cap B) = P(A) \cdot P(B) \rightarrow$ If this is true, the events are independent.

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(ANB) - P(BNC) + P(ANBNC)$$

Law of total prob

If in a sample space with events A_1, A_2, A_3
 $A_1 \cup A_2 \cup A_3 = S$ then it is called exhaustive events

If $A_1 \cap A_2 \cap A_3 = \emptyset$ (It is mutually exclusive all the events don't intersect) $P(A_1 \cap A_2 \cap A_3) = 0$

$$P(B) = P(B/A_1) \cdot P(A_1) + \dots + P(B/A_k) \cdot P(A_k)$$

such that A_i and B are disjoint for $i > k$

$$\sum_{i=1}^k P(B/A_i)P(A_i)$$

The prob of event B happening given A_1 has happened & given A_2 has happened & given A_3 has happened

Bayes theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Naive Bayes

→ Efficient ML algo for Solving classification & regression problem

→ Assumes conditional independence, Bayes Theorem does not.

↳ This means the relationship b/w all input features are independent.

→ Bayes Theorem for Single variable input.

$$P(A|B)$$

For more input w/ independence

→ It uses Bayes Theorem to calculate the Posterior prob

for 2 events like Yes or No

Ex: With given Weather data find out if play can happen

play on sunny day

frequency table

Weather	Yes	No
Overcast	5	0
Rainy	2	2
Sunny	3	2
Total	10	4

To find out, $P(\text{Yes} | \text{Sunny})$ play on sunny day

$$P(\text{Yes} | \text{Sunny}) = \frac{P(\text{Sunny} | \text{Yes}) \cdot P(\text{Yes})}{P(\text{Sunny})}$$

$$P(\text{Sunny} | \text{Yes}) = 3/10 = 0.3 \quad \text{and } P(\text{Yes}) = 10/14 = 0.71$$

$$P(\text{Sunny}) = (3+2)/14 = 0.35$$

$$P(\text{Sunny} | \text{No}) = 2/4 = 0.5 \quad \text{and } P(\text{No}) = 4/14 = 0.29$$

$$P(\text{No} | \text{Sunny}) = \frac{P(\text{Sunny} | \text{No}) \cdot P(\text{No})}{P(\text{Sunny})} = 0.29$$

$$P(\text{Sunny} | \text{No}) = 2/4 = 0.5 \quad \text{and } P(\text{No}) = 4/14 = 0.29$$

$$(0.29 \times 0.5) / 0.35 = 0.41$$

→ If we consider $b = 0.41$ at 0.41 is maximum

$$P(\text{Sunny} | \text{Yes}) > P(\text{Sunny} | \text{No})$$

→ Player can play on sunny day

→ $b = 0.41$ is maximum

→ To play on sunny day: Observe history

→ Relation $(P(\text{Yes} | \text{Sunny})) / (P(\text{Yes} | \text{Rainy}))$ is increasing
 because $(P(\text{Yes} | \text{Sunny})) = 3/10$ & $(P(\text{Yes} | \text{Rainy})) = 2/10$

→ A random variable assumes numerical values associated with random outcome of an experiment.
Its domain is the sample space and range is set of real numbers.

- Discrete random variable: Takes finite or countably infinite no. of values
 - Continuous random variable: Takes numerical value in an interval or collection of intervals.
- Continuous random var are generated in exp. when things are 'measured' and not 'counted'
e.g. measure of time, weight, volume, length

Discrete example: No of defective pieces, no. of steps,

If sum of all prob in a set is ≤ 1 , It is pdf (each entry is ≥ 0)
~~but~~

x	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
$P_1(x)$	0.25	0.25	0.25	0.25																
$P_2(x)$	0.4	-0.1	0.3	0.4																

$p(x) = \text{likelihood that random variable takes the value } x$

PDF: Expected value: mean of a random variable

$$\text{PDF: } E(x) = \mu = \sum x p(x)$$

Variance: Summarize the variability in the values of a random variable.

$$\text{Variance: } \text{Var}(x) = \sigma^2 = \sum (x - \mu)^2 p(x)$$

$$\text{Shortest formula: } E(x^2) - [E(x)]^2 = \sum x^2 p(x) - \mu^2$$

Standard deviation: positive square root of variance

Properties of Random Variables: $E(ax) = aE(x)$ a is const.

$$\text{Expected value: } E(x+b) = E(x) + b \quad b \text{ is const.}$$

$$E(ax+b) = aE(x) + b$$

Variance of discrete RV: $\sigma^2 = E[(x-\mu)^2] = \sum (x-\mu)^2 p(x)$
Standard deviation of RV: $\sigma = \sqrt{\sum (x-\mu)^2 p(x)}$

Prop of Variance: $V(a) = a^2 \cdot V(x)$ if a is const. $V(a) = 0$

$$\text{Prop of } V(x+a) = V(x)$$

$$\text{Prop of } \sigma_{x+y} = \sqrt{a^2 \cdot \sigma_x^2 + \sigma_y^2} = \sqrt{a^2 + 1} \cdot \sigma_x$$

Random Variable
and Probability distribution / Probability density function

Discrete vs Continuous
 $p(x)$ vs $f(x)$

Validation: i) $0 \leq P(x) \leq 1$ $0 \leq f(x) \leq 1$ Should satisfy
ii) $\sum P(x) = 1$ $\int f(x) dx = 1$ Satisfy

Probability distribution / Probability density function

example prob:

x	1	2	3	4	5	6
$P(x)$	0.05	0.10	0.35	0.40	0.10	

Computed

$$a) E(x)$$

$$b) V(x)$$

$$c) \sigma$$

d) $V(x)$ only shortcut formula?

$$x \quad P(x) \quad xP(x) \quad (x-\mu)^2 \cdot \text{const}^2 \cdot \text{const}^2 p(x) = \text{const}^2 \cdot \sigma^2$$

$$1 \quad 0.05 \quad 0.05 \quad 29.76 \quad 1.48 \quad 0.05$$

$$2 \quad 0.1 \quad 0.2 \quad 19.80 \quad 1.98 \quad 0.4$$

$$4 \quad 0.35 \quad 1.4 \quad 6.00 \quad 2.1 \quad 5.6$$

$$8 \quad 0.4 \quad 3.2 \quad 2.40 \quad 0.96 \quad 25.6$$

$$16 \quad 0.1 \quad 1.6 \quad 9.12 \quad 25.6 \quad 1.96$$

$$\frac{1}{1} \quad \frac{1}{1} \quad \frac{15.645}{15.645} \quad \frac{57.25}{57.25}$$

$$\text{PPF: } E(x) = \mu = \sum x p(x) = 6.45$$

a) $E(x) = \mu = \text{Exp}(\mu) = 6.45$
b) $V(x) = E((x-\mu)^2) p(x) = 15.645 - \sigma^2$
c) $\sigma = \sqrt{V(x)} = \sqrt{15.645} = 3.96$
d) $V(x)$ using shortcut = $E(x^2) - [E(x)]^2 = E(x^2)p(x) - \mu^2$

If the problem or PDF is given as \rightarrow Probability Density function.

$$f(x) = \begin{cases} Cx^2, & -1 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

it is a continuous distribution over the range is given

Imp formula: $\int_a^b f(x) dx = F(b) - F(a)$

① $\int x^n dx = \frac{x^{n+1}}{n+1}$ when $n \neq -1$ (Note: note about lower limit)

② $\int e^{ax} dx = \frac{e^{ax}}{a}$

③ $\int u dv = u v - \int v du$

④ $\int \frac{1}{1+x^2} dx = \tan^{-1} x$

problem:

x be a random variable with PDF

$$\text{then } f(x) = \begin{cases} cx^2, & -1 \leq x \leq 1 \\ 0, & \text{Otherwise} \end{cases}$$

a) find constant c

b) find $E(x)$ and $V(x)$

c) find $P(x \geq 1/2)$

a) as the limit is $-1 \leq x \leq 1$ for $f(x) \geq 0$

$$\int f(x) dx = 1$$

$$\int c x^2 dx = 1$$

$$\int c x^2 dx = 1 \rightarrow \text{new formula } ①$$

$$c \left[\frac{x^3}{3} \right]_1 = 1 \rightarrow \text{new formula } ②$$

$$c \left[\frac{1}{3} - \frac{(-1)}{3} \right] = 1 \rightarrow \text{new formula } ③$$

$$c \cdot \frac{2}{3} = 1 \rightarrow c = \frac{3}{2}$$

$$\Rightarrow f(x) = \begin{cases} \frac{3}{2}x^2, & -1 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

b)

$$E(x) = \int x f(x) dx \rightarrow \text{Not } \leq \text{ but } \leq$$

$$\int x \cdot \frac{3}{2}x^2 dx = \frac{3}{2} \int x^3 dx$$

$$= \int x \left[\frac{3}{2}x^2 \right] dx = \left[\frac{3}{2}x^3 \right]_1 = \frac{3}{2}(1)^3 = \frac{3}{2}$$

$$\int x \cdot \frac{3}{2}x^2 dx = \frac{3}{2} \int x^3 dx = \frac{3}{2} \left[\frac{x^4}{4} \right]_1 = \frac{3}{2} \cdot \frac{1}{4} = \frac{3}{8}$$

$$= \left[\frac{3}{2} - \frac{3}{2}(-1)^3 \right] = 0$$

$$E(x) = 0 \rightarrow \text{Indicates randomizing factor}$$

$$\text{Q3: } V(x) = E(x^2) - [E(x)]^2 = \int x^2 f(x) dx - 0$$

c) $P(x \geq 1/2) \rightarrow$ change interval

$$= \int_{1/2}^1 f(x) dx = \int_{1/2}^1 \frac{3}{2}x^2 dx$$

$$= \left[\frac{3}{2} \cdot \frac{x^3}{3} \right]_{1/2}^1 = \left[\frac{3}{2} \left(\frac{1}{8} - \frac{3}{2} \left(\frac{1}{8} \right)^2 \right) \right]$$

$$= \left[\frac{1}{16} - \frac{3}{16} \right] = \frac{3}{16}$$

$$= \left[\frac{3}{2} \cdot \frac{x^3}{3} \right]_{1/2}^1 = \left[\frac{3}{2} \left(\frac{1}{8} - \frac{3}{2} \left(\frac{1}{8} \right)^2 \right) \right]$$

$$= \left[\frac{3}{2} \cdot \frac{x^3}{3} \right]_{1/2}^1 = \left[\frac{3}{2} \left(\frac{1}{8} - \frac{3}{2} \left(\frac{1}{8} \right)^2 \right) \right]$$

Properties of PDF \rightarrow Probability Density Function

1) $f(x) \geq 0$

2) $\int_{-\infty}^{\infty} f(x) dx = 1$

3) $P(a \leq x \leq b) = \int_a^b f(x) dx$ area under $f(x)$ from a to b
 $= \int_a^b f(x) dx$ for any a and b

Cumulative Density Function

$$CDF = F(x) = P(X \leq x) = \int_{-\infty}^x f(x) dx$$

PDF \Rightarrow CDF
 $f(x) = \frac{dF(x)}{dx}$

Continuous Random Variable formulae

Same as PDF but we \int instead of \sum

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

$$\sigma^2 = V(x) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

$$SD = \sigma$$

Joint probability distribution:

Discrete
Now uses x_1, y_1 2 random variables instead of 1.

2 condition if $X = \{x_1, x_2, \dots, x_n\}$ & $Y = \{y_1, y_2, \dots, y_m\}$

① $0 \leq p_{ij} \leq 1$

② $\sum_{i=1}^n \sum_{j=1}^m p_{ij} = 1$ sum of 2 random variable ≥ 1

x	1	2	3
4	0.1	0.1	0.2
5	0.2	0.1	0.3
Total	0.3	0.2	0.5

Marginal PDF of $x = \begin{pmatrix} 1 & 2 & 3 \end{pmatrix}$ Marginal PDF of $y = \begin{pmatrix} 4 & 5 \end{pmatrix}$

Continued:

~~① $0 \leq f(x, y) \leq 1$~~

① $0 \leq f(x, y) \leq 1$

② $\int \int f(x, y) dxdy = 1$

Mean for RG using \int and also \sum and \int from 0 to ∞

Bernoulli Distribution:
Discrete random var

Probability mass function

$$P(x) = \begin{cases} p^x q^{1-x}, x=0,1 \\ 0, \text{ otherwise} \end{cases}$$

mean $\mu = \sum x P(x) = 0.p(0) + 1.p(1)$

variance $\sigma^2 = \sum (x - \mu)^2 P(x)$

$$\sigma^2 = pq \rightarrow \text{derivation not needed}$$

$$\sigma^2 = \sum x^2 P(x) - \mu^2 = \sum x^2 p - \mu^2$$

Binomial Distribution: If n times event

$$p(x) = \begin{cases} n C r p^r q^{n-r}, r=0,1,2,\dots,n \\ 0, \text{ otherwise} \end{cases}$$

To calculate prob of an event

happening in n attempts ex:

prob of getting a in dice roll = $1/6$

prob of not getting a = $5/6$

Let no of attempt $n = 10$

$$P(x=10) = {}^{10}C_{10} p^10 q^{10} = \left(\frac{1}{6}\right)^{10}$$

It is discrete prob dist

$$\mu = np \quad \sigma^2 = npq$$

$\rightarrow n$ is finite
 \rightarrow trials are independent of each other

\rightarrow prob of success p is constant for each trial

\rightarrow each trial results in two mutually exclusive events success & failure.

Poisson

$$P(n) = \begin{cases} e^{-\lambda} \frac{\lambda^n}{n!}, & n=0, 1, 2, \dots \\ 0, & \text{otherwise} \end{cases}$$

$\lambda = np$ or trials
= prob(poisson)

→ It is used for rare events when p prob is very small and trials n is very large or

→ Binomial dist. Can be approximated by poisson distribution
 $n \rightarrow \infty, p \rightarrow 0$ such that $\lambda = np$ is constant

$$\boxed{\mu \Rightarrow \sigma^2 \Rightarrow \text{standard normal distribution}}$$

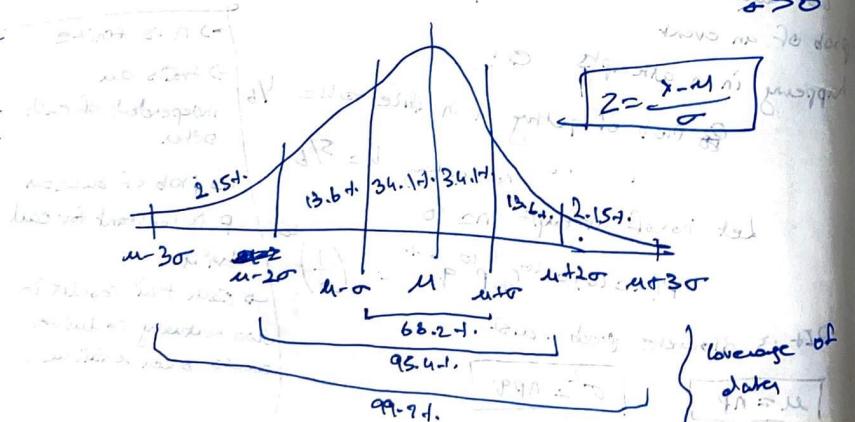
$\lambda = np$

Normal Distribution

Information such as heights, weights (length, size, profit, price etc) quantity large data with mean & variance follows normal distribution.

→ A continuous random variable X which assumes all possible values in the entire real space $-\infty < x < \infty$ is said to follow normal dist. with parameters μ & σ^2 .

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}, -\infty < x < \infty$$



Properties:

1) Normal curve is bell shaped and symmetric about mean

2) Mean = Mode = Median

3) Total area under normal curve = 1

4) Normal curve approaches but never touches x-axis as it extends

further away from the mean.

Normal Approximation to Binomial

If $n > 30$ & $np, nq > 5$ then use normal approximation

$$Z = \frac{x-\mu}{\sigma} = \frac{x-np}{\sqrt{npq}}$$

approx. & $\sigma = \sqrt{npq}$

Ex: mean to calculate probability at $(58 + 59)$ cases sigma?

Ex: for $n=100, q=0.5, p=0.5$

$$P(X \leq 12) = P(1.5 \leq Z \leq 12.5) = P(-2.68 \leq Z \leq -2.37)$$

and probabilities follow as multiplying bin two times to find prob using 2 table

$$= 0.00368$$

0.00521

$$= P(-2.37) - P(-2.68) = 0.00889 - 0.00368 = 0.00521$$

$$\text{for } P(X \leq 12) = P(X \leq 12.5)$$

Sampling: (i) 20.05 A.C. & want a representative sample from a population

ex: selecting marks of 10 students from a class of 40

Random Sampling: Randomly select from population

Stratified Sampling: Group the population and then select

Sample / obtain a representative sample

ex: $\begin{matrix} A & C & C \\ A & B & B \\ A & C \end{matrix} \rightarrow \begin{matrix} B & B & B \\ B & B & B \\ C & C & C \end{matrix} \rightarrow \begin{matrix} B & B & B \\ C & C & C \end{matrix}$

populations

$$\frac{1}{3} = \frac{1}{3}$$

Sampling error can occur due to a biased selection

Sampling variability is the variance in the sample data, larger sample size informs less about the true parameter.

The sample size leaves the variability.

$$\text{Sampling 2 formula} = \frac{\bar{x}-\mu}{\sigma/\sqrt{n}}$$

where \bar{x} is the mean of sample & n is the size of sample.

$$\text{When } \frac{n}{N} > 0.05 \quad \frac{\text{sample size}}{\text{pop. size}} > 0.05$$

use finite correction factor

$\text{use finite correction factor}$

No of samples of size n that can be drawn from a population of size N is
with replacement: n^N
without replacement: $\binom{N}{n}$

Central Limit theorem:
→ Regardless of the shape of population distribution, for large sample sizes ($n \geq 30$), the sampling distribution of mean is normally distributed.

→ If $n \leq 30$ but the population is normally distributed, then the sampling dist. of mean is also normally distributed.

→ The sample mean can be analyzed using Z score

$$Z = \frac{\bar{x}-\mu}{\sigma/\sqrt{n}} \quad (Z \sim N(0, 1))$$

→ If the population is finite & $N/n \geq 0.05$ (i.e.) more than 5% of the population is used in sample, use finite correction factor if $N/n > 0.05$

$$Z = \frac{\bar{x}-\mu}{\sigma/\sqrt{n}} \quad \text{if } N/n > 0.05$$

Sampling Proportion:

The group within a population

e.g. 10 red balls in a pack of 50

$$\hat{P} = \frac{r}{n} = \frac{10}{50}$$

• Expected value of proportion $E(\hat{P}) = P$

• SD of sample proportion $Sd(\hat{P}) = \sqrt{\frac{P(1-P)}{n}}$

2 formula for sample proportion

$$Z = \frac{\hat{P}-P}{\sqrt{\frac{P(1-P)}{n}}}$$

Confidence interval:

→ It is better compared to point estimate, used to estimate where the μ of a population lies.

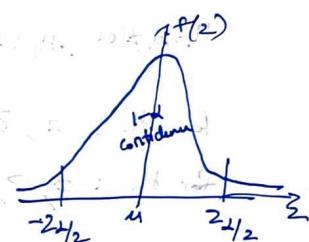
→ ~~the confidence interval~~ ~~is the range of values which we are 95% sure contains the true population mean.~~

the confidence interval is given by:

$$\bar{x} - 2\alpha/2 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + 2\alpha/2 \frac{\sigma}{\sqrt{n}}$$

α is the area under the normal curve

outside the confidence interval area



Ex: To find 95% confidence interval take α as 0.05 because

$$\alpha = 0.05 \quad 1-\alpha = 0.95 \quad \alpha/2 = 0.025 \quad \alpha = 5\% = \frac{5}{100} = 0.05$$

Substitute $\alpha/2 = 0.025$ and get 2 by finding Z value for 0.025

Look up Z table and get Z value at 0.025 for left tail

Z score for 90.1 = 1.65 92.1 = 1.75

and at right tail 95.1 = 1.96 and a midpoint of 0.05

99.1 = 2.58

→ If a finite population is given use finite correction factor $\sqrt{\frac{N-n}{n-1}}$

Population proportion: (used to estimate the proportion)

1) Confidence interval for population proportion

$$\hat{P} - Z_{\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}} \leq \mu \leq \hat{P} + Z_{\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$$

Hypothesis testing

- Hypothesis is a statement or a claim
- Hypothesis to be validated that claim based on data
- H_0 : Null hypothesis
↳ no significant change or difference

H_0 or H_a

H_a : Alternative hypothesis

H_1 or H_a

↳ opposite of null hypothesis

$H_0: \mu = 100$, $H_a: \mu \neq 100$ example of null hypothesis

$H_1: \mu \neq 100$, $H_1: \mu > \mu_2$ ex of alternative hypothesis

for H_0 : $=$, \geq or \leq is accepted

for H_1 : \neq , $>$, $<$ is used

Test:

When $\mu_1 < \mu_2$ \Rightarrow one-tailed test, or left tailed test

$\mu_1 > \mu_2 \Rightarrow$ one-tailed, or right tailed test

$\mu_1 \neq \mu_2 \Rightarrow$ two tailed test

→ Critical region is the area where the hypothesis is rejected even if the hypothesis is true. It is beyond the confidence

interval area or extreme area.

$\alpha = P(\text{Rejecting } H_0 \text{ while } H_0 \text{ is true})$

$\beta = P(\text{Accepting } H_0 \text{ while } H_0 \text{ is false})$

α is type 1 error

β is type 2 error

	H_0 is true	H_0 is false
H_0	✓	✗
H_a	✗	✓

2 test is used when

- Mean of single population (μ)
- Diff bet. 2 population means ($\mu_1 - \mu_2$)
- proportion of single population (p)
- Diff. bet. 2 proportions ($p_1 - p_2$)

→ Need to know the standard deviation for 2 test:

→ Step 1: $n \geq 30$, σ should be known

→ Step 2: $\bar{x} \pm z_{\alpha/2} \sigma / \sqrt{n}$

→ Applying the formula at 95% \Rightarrow $z = 1.96$

Hypothesis testing formula.

	Sample size	formula
One mean	Large	$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$
	Small	$t = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$
Two mean	Large	$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$
	Small	$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$
One proportion	Large	$Z = \frac{p - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$
Two proportion	Large	$Z = \frac{p_1 - p_2}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$
Student's paired t-test	Small	$t = \frac{\bar{d} - d_0}{s_d / \sqrt{n}}$

Critical value	Level of significance
$Z_{\alpha/2}$	1.96, 2.45, 5.91, 10.93
$Z_{\alpha/2} = 2.58$	$Z_{\alpha/2} = 1.96$
$Z_{\alpha/2} = 1.645$	$Z_{\alpha/2} = 1.645$
$Z_{\alpha} = 2.33$	$Z_{\alpha} = 1.645$
$Z_{\alpha} = 1.28$	$Z_{\alpha} = 1.28$

P-value can be used to determine the evidence of rejection or not

(a) $\text{rejecting H}_0 \rightarrow \text{not H}_0$

Give 2, find the value in 2 table

for ex: if $Z = 1.65$, we get p-value on $1 - 95.05\% = 4.95\%$

If p
 p value $< 1\%$ \Rightarrow Overwhelming evidence to support alt. hypothesis
 p value bet $1\% - 2.5\%$ \Rightarrow Strong evidence to support alt. hypothesis
 p value bet $5\% - 10\%$ \Rightarrow weak evidence
 p value $>$ above 10% \Rightarrow no evidence to support alt. hypothesis

planned prior to analysis

t-test		planned + post sample	
\rightarrow when Sample is Small < 30		then formula $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$	
		s_{sample}	
		s_{sample}	

the obtained value is t value and should be checked

in the t-table sample size

Consider the degrees of freedom $n-1$

Chi-square test:

- Independence test

- Goodness of fit

\rightarrow used only for frequencies

H_0 : Independent \Rightarrow Observed = expected

H_1 : Not-Independent \Rightarrow Observed \neq expected

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}, k = r \times c \text{ of cells}$$

short hand: $\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$

ex: Accident data by weekday

data for independence bet. weekday & accident

	O	E	$O-E$	$(O-E)^2/E$	Sum
Mon.	184	150	34	7.71	
Tues.	148	150	-2	0.03	
Wed.	145	150	-5	0.17	
Thurs.	153	150	3	0.06	
Fri.	150	150	0	0	
Sat.	184	150	34	0.11	
Sun.	116	150	-34	7.71	
	<u>1050</u>	<u>1050</u>		<u>15.27</u>	
	<u>$\bar{O} = 1050/7 = 150$</u>	<u>$E = 1050/7 = 150$</u>		<u>$\chi^2 = \sum \frac{(O-E)^2}{E}$</u>	

taking S.I. LOS at $\alpha = 0.05$ \Rightarrow degrees of freedom $= 7-1 = 6$

we get critical $\chi^2_{\text{critical}} = 12.6$

as $15.27 > 12.6$ we reject H_0 and accept H_1

effect of most variables on accident \rightarrow they are dependent on day of week

Observed			Expected		
	Cancer	No Cancer	Total		
Smoked	400	800	1200		
Non-smoked	800	500	1300		
Total	1200	1300	2500		
	c_1	c_2	n		

calculated as $\frac{R_{ij} \times C_j}{n}$

now total column total $= \frac{R_{ij} \times C_j}{n}$

$\chi^2_{\text{cal}} = \frac{(400 - 326.67)^2}{326.67} + \frac{(800 - 373.33)^2}{373.33} + \frac{(1200 - 426.67)^2}{426.67} + \frac{(1300 - 373.33)^2}{373.33}$

for n independent variables $\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$ when a, b, c, d are the entries

shortcut $\chi^2 = \frac{n}{r+s} \frac{(ad-bc)^2}{(a+b)(c+d)(a+c)(b+d)}$

$$\text{fisher's exact probability} \Rightarrow P = \frac{1}{n!} \frac{n_1!}{r_1!} \frac{n_2!}{r_2!} \frac{c_1!}{c_1!} \frac{c_2!}{c_2!}$$

$$\text{Yates' correct} = \frac{(10-6.1-0.5)^2}{E_{1,1}}$$

Variance based test:

Fishers F-test: Ratio of variances

$$\text{Single variance: } H_0: \sigma_1^2 = \sigma_2^2 \text{ or } H_1: \sigma_1^2 > \sigma_2^2$$

$$H_1: \sigma_1^2 > \sigma_2^2 \quad X = \frac{(n_1-n_2)S_1^2}{S_2^2}$$

$$H_0: \sigma_1^2 \leq \sigma_2^2 \quad X = \frac{(n_2-n_1)S_2^2}{S_1^2}$$

$$\text{Ratio of 2 varianc} \quad H_0: \sigma_1^2 = \sigma_2^2 \text{ or } H_1: \sigma_1^2 > \sigma_2^2$$

$$H_1: \sigma_1^2 > \sigma_2^2 \quad F = \frac{S_1^2}{S_2^2} \sim F(n_1-1, n_2-1)$$

Ratio of 2 variance test formulae \Rightarrow the 2nd + 3rd probab
of getting 2nd sample

$$F = \frac{S_1^2}{S_2^2} \sim f(n_1-1, n_2-1)$$

If F-value is less than or equal to 3.52 then H₀ is accepted.

~~critical value~~ \Rightarrow critical is obtained from F-table

with degrees of freedom of 1st & 2nd sample.

1 st Sample	2 nd Sample
n ₁	n ₂

Correlation analysis

\Rightarrow Statistically check for the existence & extent of relationship between 2 variables

covariance:

$$\text{Covariance}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

If it is -ve, it means $x \downarrow y$, i.e. increase the one decreases the other.

If +ve, $x \uparrow y$

If no relation, = 0

\Rightarrow Covariance only shows how 2 variables change together not the dependency between them.

correlation coefficient:

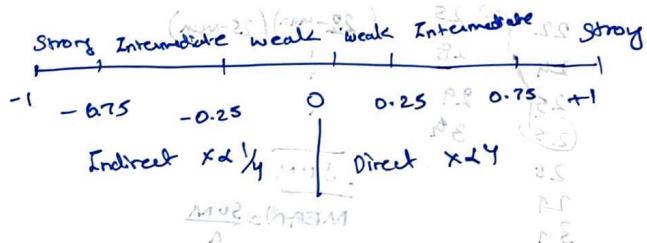
Value and the test result

\Rightarrow Can find the strength of relationship by eliminating the effect of units.

Unit of the value by dividing with SD.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

for r value



Regression analysis:

Linear regression creates a best fit line to predict y value.

$$\text{Equation: } y = \beta_0 + \beta_1 x$$

(Find values of both variables of y & x to substitute in the below formula)

Given values of y & x substitute in the below formula and solve the two eqns to get β_0 & β_1 .

$$\Sigma y = n\beta_0 + \beta_1 \Sigma x \quad \text{--- (1) Methodical}$$

$$\Sigma y = \beta_0 \Sigma x + \beta_1 \Sigma x^2 \quad \text{--- (2) Methodical}$$

Maximum likelihood Estimation:

$$\text{for } p = \frac{k}{n}$$

for binomial distribution $= \frac{\bar{x}}{n}$

for poisson distribution $= \bar{x}$

for normal distribution $\bar{x} = \frac{\sum x_i}{n} = \bar{x}$ and $\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$

for uniform distribution $= x_i \text{ for } n \text{ when } x_i \text{ is min value & } x_n \text{ is max value in the interval}$