# Machine Learning Course

**Rishav Das**
**Lead Data Scientist- Wipro**
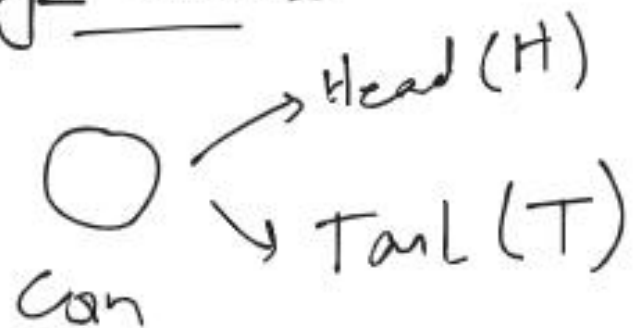
# Probability & Stats

**Rishav Das**
**Lead Data Scientist- Wipro**

# Probability stats



Head (H)

Coin → Tail (T)

Toss - coin - one

$$H = \frac{1}{2} = 0\cdot5$$
$$T = \frac{1}{2} = 0\cdot5$$ } → 1

$$HH = \frac{1}{4} = 0\cdot25$$
$$HT = 1/4 = 0\cdot25$$
$$TH = 1/4 = 0\cdot25$$
$$TT = \frac{1}{4} = 0\cdot25$$
$$\overline{\Sigma P(coin) = 1}$$

] Toss two coin's together

casino / Lottery

i) Ticket price → 6      → 6 Cr

     10 Cr → 60 Cr
         − 6
        54 Cr

Probability will be more

↓

500 ⟶ 1,50,000

Question 1: A bag consists of 3 red balls, 5 blue balls, and 8 green balls. A ball is selected at random. Find the probability of

1. Getting a red ball.

2. Getting a green ball.

3. Not getting a blue ball.

**Answer :** Total number of the balls = 3 + 5 + 8 = 16.

1. Let R be the event of getting a red ball. The number of favorable outcome = 3. The required probability is $P(R) = \frac{3}{16}$

2. Let G be the event of getting a green ball. The number of favorable outcome = 8. The required probability is $P(G) = \frac{8}{16} = \frac{1}{2}$

3. Let B be the event of getting a blue ball. The number of favorable outcome = 5. The required probability of getting blue ball = $\frac{5}{16}$. The probability of not getting a blue ball = $1 - P(B) = 1 - \frac{5}{16} = \frac{11}{16}$

Also, the event of not getting a blue ball is the same as getting a red or green ball.

$P(B') = P(R) + P(G) = \frac{3}{16} + \frac{8}{16} = \frac{11}{16}$.

Rule/Theorem

i) Summation of all probability is 1

2) All the Probability scores are between 0 - 1

ex: $70\% = \dfrac{70}{100} = 0.7$

3) $P(A') = 1 - P(A)$

Probability $\rightarrow$ $\dfrac{\text{No of success/failure}}{\text{No of event}}$

Success
of HH $= \dfrac{1}{4} = 0.25$

Failure
of HH $= 1 - 0.25 = 0.75$

# Permutation

- Permutation is the arrangement of items in which **order matters**

- Number of ways of **selection and arrangement of items** in which Order Matters

$$^n P_r = \frac{n!}{(n-r)!}$$

# Combination

- Combination is the selection of items in which **order does not matters** .

- Number of ways of **selection of items** in which Order does not Matters

$$^n C_r = \frac{n!}{r!\,(n-r)!}$$

**Examples**:

There are seven boys and three girls on a school tennis team. The coach must select four people from this group to participate in the county championship?

a. *How many* four-person teams can be formed from the group of ten students?

$$_{10}C_4 = 210$$

b. *In how many ways* can two boys and two girls be chosen to participate in the county championship?

$$_7C_2 \cdot {_3C_2} = 63$$

c. *What is the probability* that two boys and two girls are chosen for the team?

$$\frac{63}{210} \Rightarrow .3$$

$$\frac{_7C_2 \cdot {_3C_2}}{_{10}C_4}$$
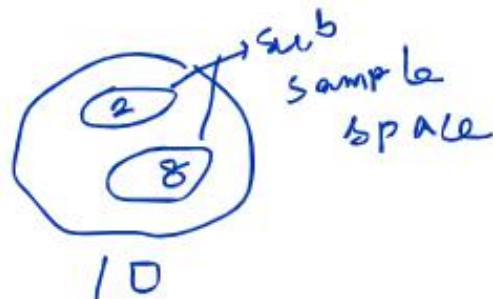
# Event

Win/loose

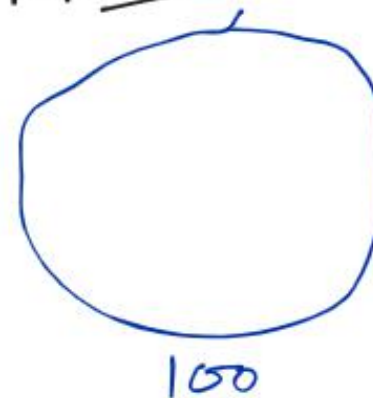i) Independent
ii) dependent

→ There is no relationship between Previous events and current events

→ There is a relationship between Previous event and current events
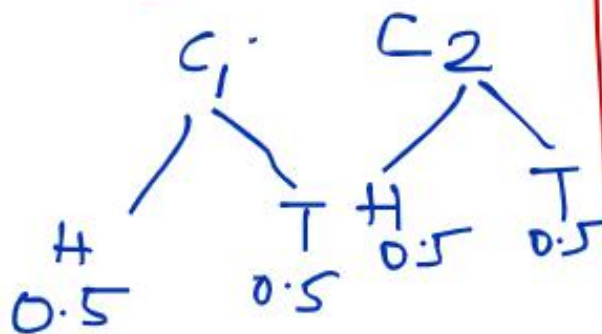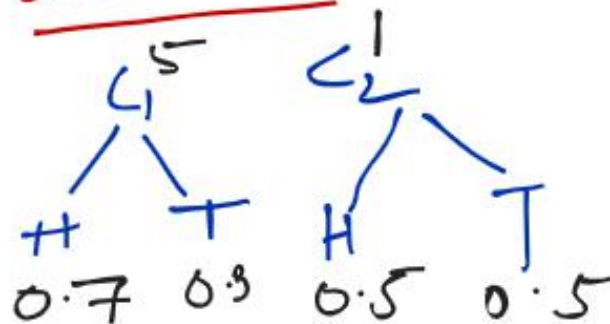
# Sample space



sub sample space

10

# Population



100

## Independent

$C_1$            $C_2$

H        T      H        T
0.5      0.5    0.5      0.5

## dependent

$C_1^5$            $C_2^1$

H     T      H     T
0.7   0.3    0.5   0.5

# Conditional probability

- To find the probability of the event **B** *given* the event **A**, we restrict our attention to the outcomes in **A**. We then find in what fraction of *those* outcomes **B** also occurred.

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

- Note: $P(A)$ cannot equal 0, since we know that **A** has occurred.

**Example:**

|  | Number of times students visited tutoring | | | |
|---|---|---|---|---|
|  | One or fewer times | Two to three times | Four or more times | Total |
| Full time student | 12 | 25 | 8 | 45 |
| Part time student | 2 | 5 | 6 | 13 |
| Total | 14 | 30 | 14 | 58 |

P( part time | visited four or more times) = $\frac{6}{14} \approx 0.43$

find          given

# Random Variables & Experiments

Random experiments

HH $\longrightarrow$ 2

HT $\searrow$ 1
TH
TT $\longrightarrow$ 0

Sample    Real
Space     Space

$P(H) = \{0, 1, 2\}$

HT   HH
TH

| $\frac{1}{3}$ | $\frac{2}{3}$ | $\frac{1}{3}$ | $= 0.33$ |

0.33

0.66

## Random Variable

Discrete RV

finite    Infinite
  − Whole no

ex:
  − no of marbles in
      a Jar
  − no of head/tail

continous Rv

− fraction    Infinite
− decimals    Values

ex: $10.00 - 10.30$.
   → Speed of car

# Discrete Random Variables

- Random variables (RVs) which may take on only a **countable** number of **distinct** values

  - E.g. the total number of tails X you get if you flip 100 coins

- X is a RV with arity $k$ if it can take on exactly one value out of $\{x_1, ..., x_k\}$

  - E.g. the possible values that X can take on are 0, 1, 2, ..., 100

$$P(HH) = \frac{1}{4}$$

$$P(HT, TH) = \frac{2}{4}$$

$$P(TT) = \frac{1}{4}$$

| $x_i$ | $P_i X_i$ $\overline{P(x_i)}$ | Result |
|---|---|---|
| 0 | $0 \cdot \frac{1}{4}$ | 0 |
| 1 | $1 \cdot \frac{2}{4}$ | $\frac{1}{2} = 0.5$ |
| 2 | $2 \cdot \frac{1}{4}_{\frac{1}{2}}$ | $\frac{1}{2} = 0.5$ |
| | | 1 |

# Continuous Random Variables

Probability density function (pdf) instead of probability mass function (pmf)

A pdf is any function $f(x)$ that describes the probability density in terms of the input variable $x$.

# Probability of Continuous RV

- Properties of pdf
    - $f(x) \geq 0, \forall x$
    - $\displaystyle\int_{-\infty}^{+\infty} f(x) = 1$

- Actual probability can be obtained by taking the integral of pdf
    - E.g. the probability of X being between 0 and 1 is

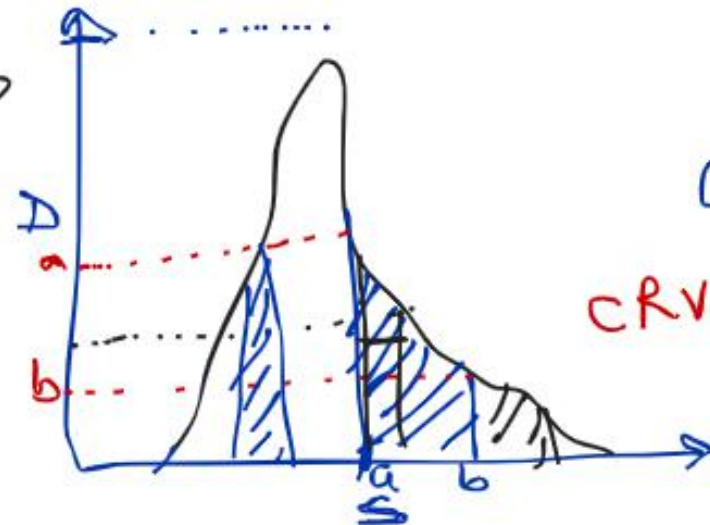$P(0 \leq X \leq 1) = \displaystyle\int_{0}^{1} f(x)dx$

PDF

computing $\underline{PDF}$

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x \cdot f(x) \cdot dx \begin{cases} \int_{-\infty}^{1} x \cdot f(x) \cdot dx \\ \int_{1}^{\infty} x \cdot f(x) \cdot dx \end{cases}$$

$$\int_{10 \cdot 00}^{10:30} x \cdot f(x) \, dx = 20$$

$$\int_{16:00:01}^{17:00:05} x f(x) \cdot dx$$



CRV

**Bayes' theorem**

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

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

ex

Bag-I
4 – W
6 – b

$E_1$

Bag II
4 – W
3 – B

$E_2$

1B
A

Find the Probability that it was Bag I

$$P(E_1) = \frac{1}{2} \quad P(E_2) = \frac{1}{2}$$

$$P(A|E_1) = \frac{6}{10} = \frac{3}{5} \qquad P(A|E_2) = \frac{3}{7}$$

$$P(E_1|A) = \frac{P(A|E_1) \times P(E_1)}{P(E_1)\,P(A|E_1) + P(E_2)\,P(A|E_2)} = \frac{\frac{3}{5} \times \frac{1}{2}}{\frac{1}{2}\frac{3}{5} + \frac{1}{2}\frac{3}{7}} = \frac{7}{12} = 58\%$$

Bag I Black ball

Picking the ball in $E_1$

Picking the ball in $E_2$

# Joint Probability Distribution

| | HHH | HHT | HTH | HTT | THH | THT | TTH | TTT |
|---|---|---|---|---|---|---|---|---|
| $E_X \to$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| $E_y \to$ | 3 | 2 | 2 | 1 | 2 | 1 | 1 | 0 |

$$X = \{0, 1\} \qquad y = \{0, 1, 2, 3\}$$

| | $y_0$ | $y_1$ | $y_2$ | $y_3$ |
|---|---|---|---|---|
| $x_0$ | $\phi$ | $\frac{1}{8}$ | $\frac{2}{8}$ | $\frac{1}{8}$ |
| $x_1$ | $\frac{1}{8}$ | $\frac{2}{8}$ | $\frac{1}{8}$ | $\phi$ |

# Central Limit Theorem

ex $P = \{ 1, 2, 5, 3, 4, 6 \}$

$S_1 = \{ 1, 1, 3, 6 \}$

$$\frac{1+1+3+6}{4} = \frac{11}{4}$$

$$\Rightarrow 2.7$$

$S_2 = \{ 3, 4, 5, 6 \}$

$$\frac{3+4+5+6}{4}$$

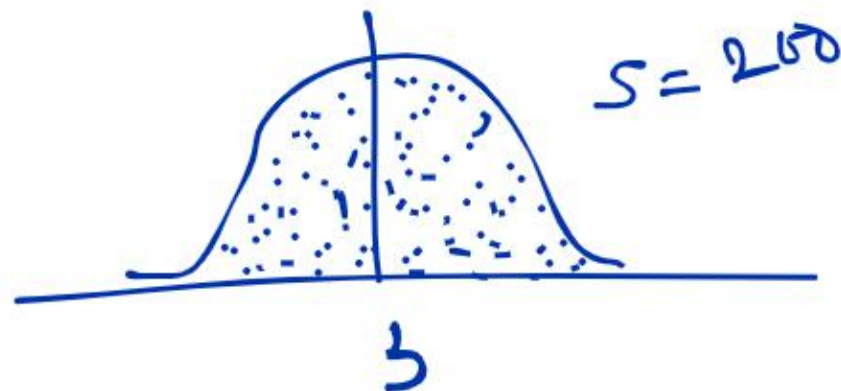$$\Rightarrow \frac{18}{4} = 4.4$$

$S_3 = \{ 1, 1, 6, 6 \}$

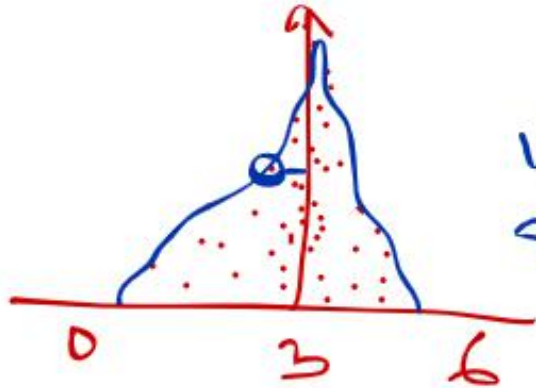$$\frac{1+1+6+6}{4}$$

$$= \frac{14}{4} = 3.5$$

$S_1 = 2.7$

$S_2 = 4.4$

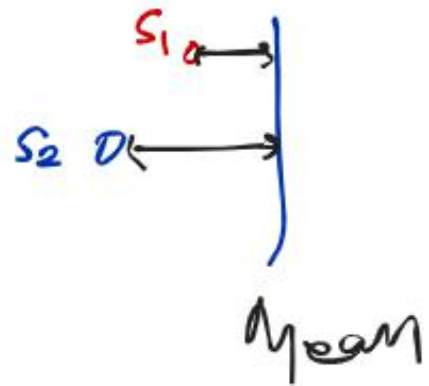$S_3 = 3.5$

$S = 250$

$\overline{x}$

# Probability distribution



$$\text{Mean} \rightarrow \mu = n \cdot p$$

$$\text{Variance} = \sigma^2 = npq$$

$$SD = \sqrt{\sigma^2} = \sqrt{npq}$$
$$\hookrightarrow \sigma$$

$n = $ no of trials

$P = $ success

$q = $ failures  [Probability]
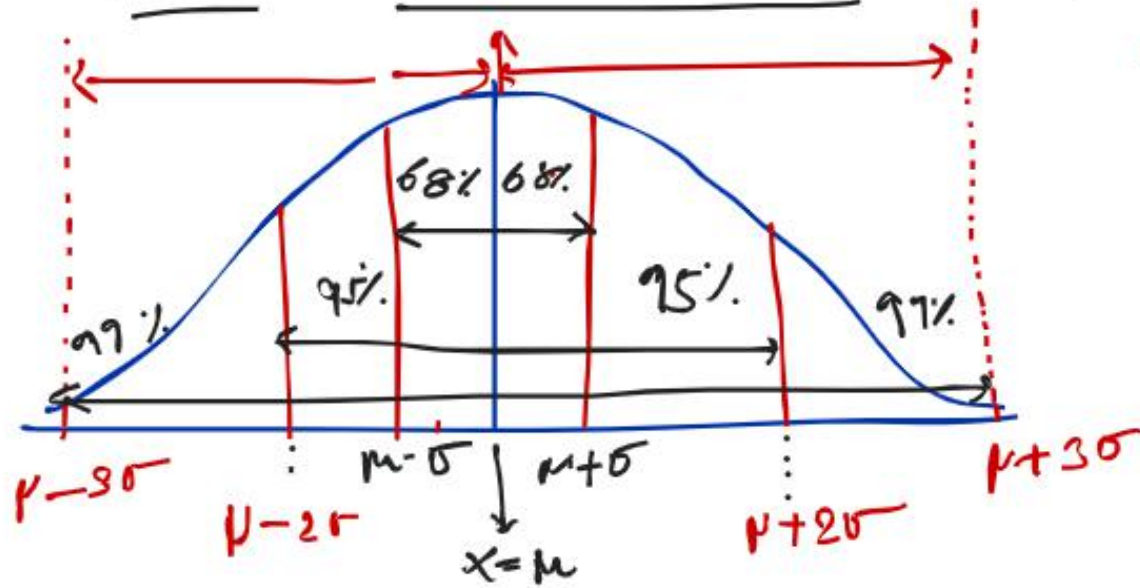
ex:

$n = 10$

$H = \frac{1}{2} \ (s)$

$T = \frac{1}{2} \ (f)$

$$\text{mean} = n \times p = 10 \times \frac{1}{2} = 5$$

$$\sigma^2 = n \times p \times q = 10 \times \frac{1}{2} \times \frac{1}{2}$$

$$= 2.5$$

$$\sigma = \sqrt{2.5} = 1.25$$

# Normal Distribution



Bell shaped curve

$\infty$

- Continous Probability distribution

$$Pdf = f(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$\int_0^{\infty} x \cdot f(x) \cdot dx$$

## standard Normal Distribution

$$z = \frac{X - \mu}{\sigma}$$

$$z = X$$

standard Normal varionce has $\mu = 0$ and $\sigma = 1$

# Poisson Distribution

Discrete R.V

$$P(x) = \frac{e^{-m} \, m^x}{x!}$$

$$= \frac{e^{np} \, (np)^x}{x!}$$

$m = np$ → should be very large/infinite

↳ should be very less

ex: casino / Lottery system

→ no of participant → $\alpha$

→ $P(win) \simeq 0.1\%$

150 cr

↳ 5cr

6 cr → 3 people

$$= \frac{3}{5cr}$$

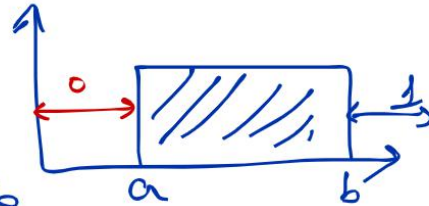| Difference between Variance and Standard Deviation ||
| --- | --- |
| **Variance** | **Standard Deviation** |
| It can simply be defined as the numerical value, which describes how variable the observations are. | It can simply be defined as the observations that get measured are measured through dispersion within a data set. |
| Variance is nothing but the average taken out of the squared deviations. | Standard Deviation is defined as the root of the mean square deviation |
| Variance is expressed in Squared units. | Standard deviation is expressed in the same units of the data available. |
| It is mathematically denoted as $(\sigma^2)$ | It is mathematically denoted as $(\sigma)$ |
| Variance is a perfect indicator of the individuals spread out in a group. | Standard deviation is the perfect indicator of the observations in a data set. |

# Uniform Distribution

In probability theory and statistics, the continuous uniform distribution or rectangular distribution is a family of symmetric probability distributions. The distribution describes an experiment where there is an arbitrary outcome that lies between certain bounds. The bounds are defined by the parameters, a and b, which are the minimum and maximum values.

Uniform Distribution

→ continous Rv

$$f(x) = \begin{cases} 0 & x < a \\ x-a/b-a & a \le x \le b \\ 1 & x \ge b \end{cases}$$

# Multinomial Distribution

Multinomial distribution, in statistics, a generalization of the binomial distribution, which admits only two values (such as success and failure), to more than two values. Like the binomial distribution, the multinomial distribution is a distribution function for discrete processes in which fixed probabilities prevail for each independently generated value.

## Multinomial distribution

Discute Rv | When events are Mutually exclusive

dice → 12 times | Find probability of each dice value occurring twice

1, 2, 3, 4, 5, 6

$$\Rightarrow \frac{n!}{n_1! \, n_2! \, n_3! \, n_4! \, n_5! \, n_6!} \cdot p_1^q \cdot p_2^2 \cdot p_3^2 \cdot p_4^2 \cdot p_5^2 \cdot p_6^2$$
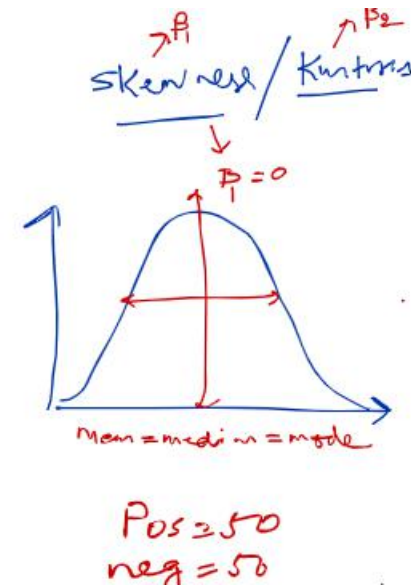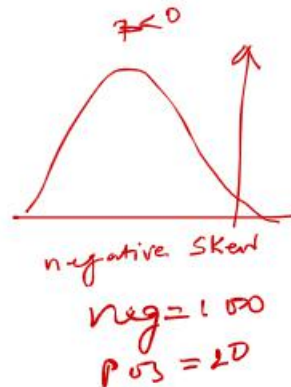
$$\Rightarrow \frac{12!}{2! \, 2! \, 2! \, 2! \, 2! \, 2!}$$

$$\left(\tfrac{1}{6}\right)^2 \cdot \left(\tfrac{1}{6}\right)^2 \cdot \left(\tfrac{1}{6}\right)^2 \cdot \left(\tfrac{1}{6}\right)^2 \cdot \left(\tfrac{1}{6}\right)^2 \cdot \left(\tfrac{1}{6}\right)^2$$
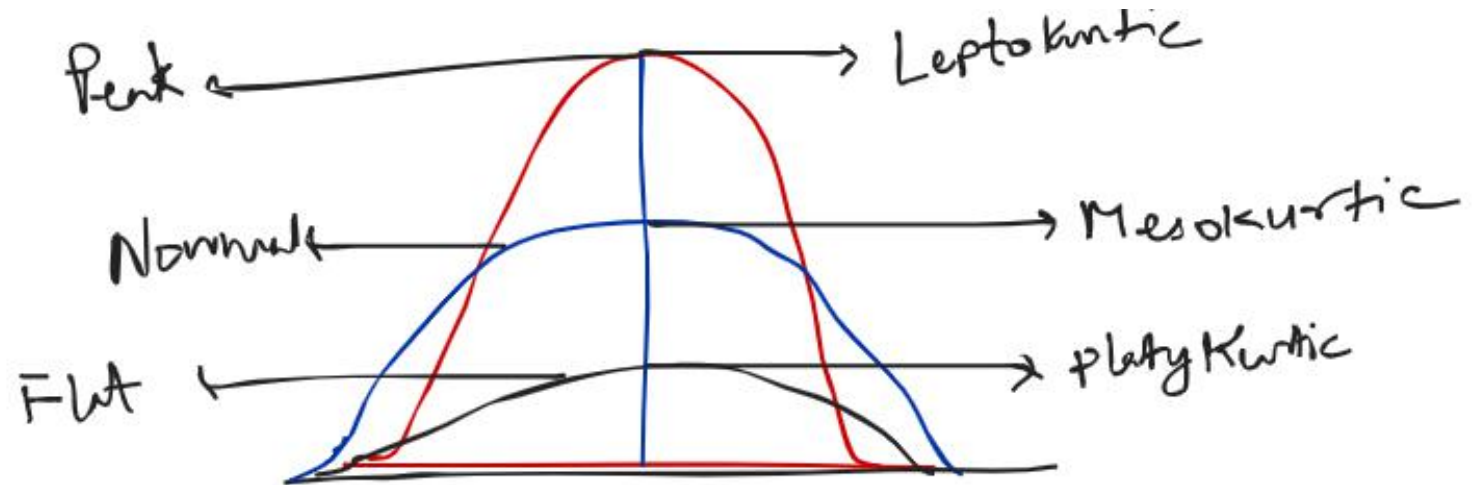
$$\Rightarrow 0.004$$

# Skewness

Skewness refers to a distortion or asymmetry that deviates from the symmetrical bell curve, or normal distribution, in a set of data. If the curve is shifted to the left or to the right, it is said to be skewed. Skewness can be quantified as a representation of the extent to which a given distribution varies from a normal distribution. A normal distribution has a skew of zero, while a lognormal distribution, for example, would exhibit some degree of right-skew.

# Kurtosis

kurtosis is a statistical measure that is used to describe distribution. Whereas skewness differentiates extreme values in one versus the other tail, kurtosis measures extreme values in either tail. Distributions with large kurtosis exhibit tail data exceeding the tails of the normal distribution (e.g., five or more standard deviations from the mean). Distributions with low kurtosis exhibit tail data that are generally less extreme than the tails of the normal distribution.

# Lines

*Lines*

**Every Straight lines** $y = mx + c$     $\longrightarrow$   $mx + c$

slope ↑ (over m)

constant (under c)

$m(\text{slope}) = \dfrac{\text{change in y value}}{\text{change in x value}}$
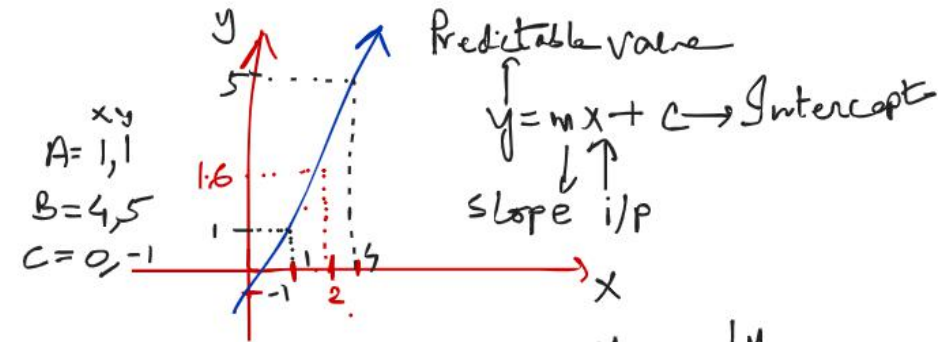
where m = 0 and x = 6, C= 1

then 0x6 +1 = 0+1 = 1
hence the constant helps to remove the
nullifcation from the straight line

thus we can say that when:
mx = 0

y = C,  hence C is referred as y intercept

Predictable value
$y = mx + c \longrightarrow$ Intercept
slope i/p

$A = 1,1$
$B = 4,5$
$C = 0,-1$

$m = \dfrac{5-1}{4-1} = \dfrac{4}{3} = 1.3$

$y = mx + c$
$= 1.3 \cdot x + (-1)$
$y = 1.3 x - 1$
$= 1.3 (2) - 1$
$= 2.6 - 1 = 1.6$

$m = \dfrac{\Delta y}{\Delta x} = \dfrac{dy}{dx}$

$m = 1$

$y = 1 \cdot x + c = x + c$
$y = 0 \cdot x + c = c$

# Introduction to Linear Regression

- The Pearson correlation measures the degree to which a set of data points form a straight line relationship.

- **Regression** is a statistical procedure that determines the equation for the straight line that best fits a specific set of data.

# Introduction to Linear Regression (cont.)

- Any straight line can be represented by an equation of the form $Y = bX + a$, where $b$ and $a$ are constants.

- The value of $b$ is called the slope constant and determines the direction and degree to which the line is tilted.

- The value of $a$ is called the Y-intercept and determines the point where the line crosses the Y-axis.

# Regression - Linear & MultiLinear



Dependent Variable → $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$

- Population Y intercept → $\beta_0$
- Population Slope Coefficient → $\beta_1$
- Independent Variable → $X_i$
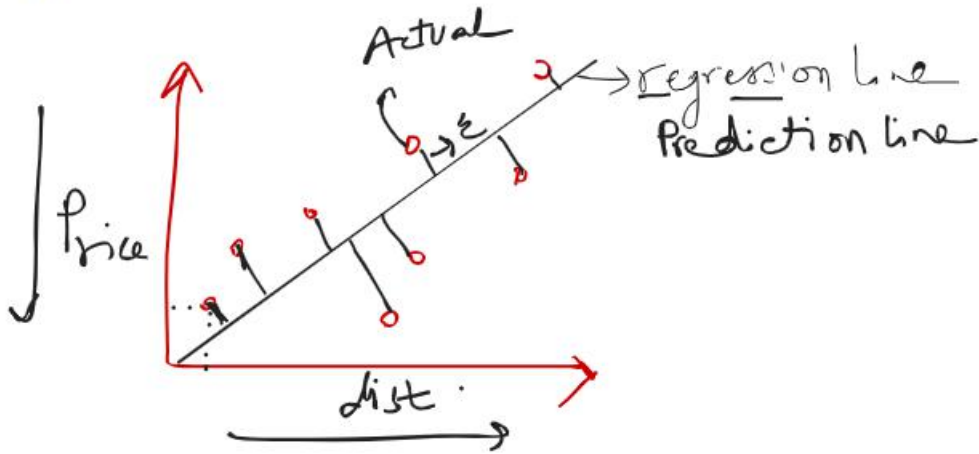- Random Error term → $\varepsilon_i$

Linear component: $\beta_0 + \beta_1 X_i$

Random Error component: $\varepsilon_i$

ex: House

1 → Loc
2 - sqft
3 - Medical
4 - School
5 - Market
6 - Airport
7 - Railway
8 - Quality

} Independent

Price → dependent

Actual

regression line
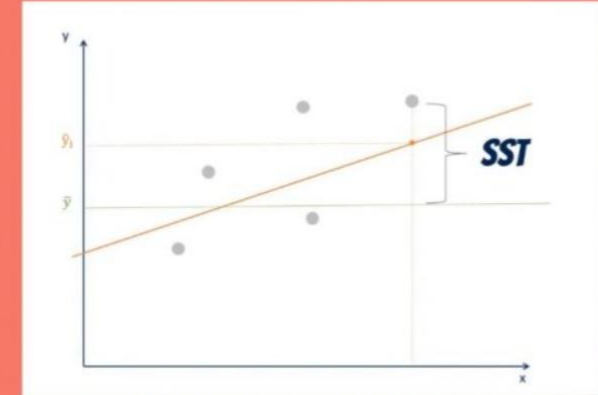Prediction line

Price

dist.

# Error Calculation

## Total Sum of Square

The total sum of squares , denoted TSS, is the squared differences between the observed dependent variable and its mean

$$\sum_{i=1}^{n}(y_i - \bar{y})^2$$

## Sum of Square Error - Regression

The second term is the sum of squares due to regression, or SSR. It is the sum of the differences between the predicted value and the mean of the dependent variable. Think of it as a measure that describes how well our line fits the data.

Measures the explained variablity by your line

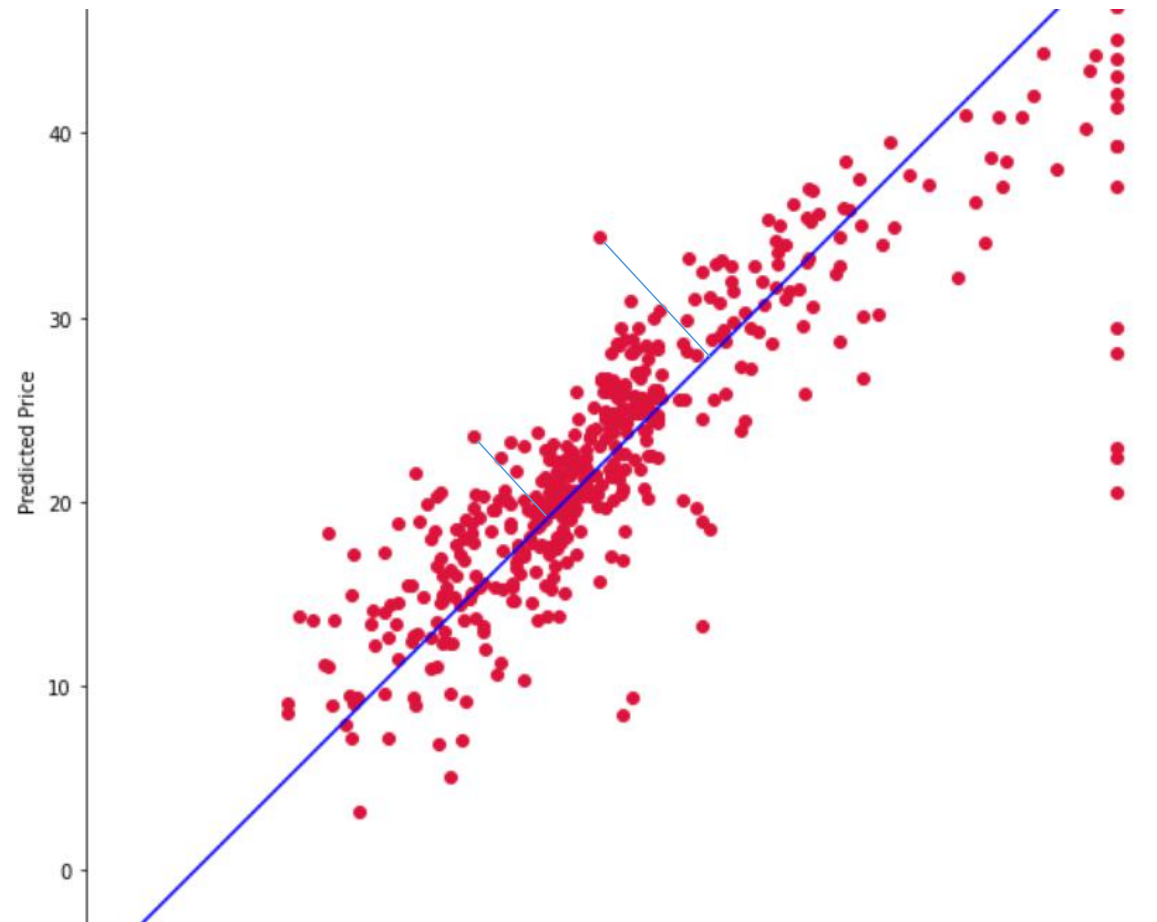$$\sum_{i=1}^{n}(\hat{y}_i - \bar{y})^2$$

# Error Calculation

**Sum of Squared Error - Prediction**

$$= \sum (y_i - \hat{y})^2 \quad \Big| \quad \begin{array}{l} \hat{y} = \text{Predicted Value} \\ y = \text{Actual Value} \end{array}$$

The last term is the sum of squares error, or SSE. The error is the difference between the observed value and the predicted value.

# Correlation Analysis - Using Rank

**Correlation formula**

| CRIM | Price | | $d1$ |
|------|-------|---|------|
| $2 - R_1$ | $40\ R_5$ | | $4$ |
| $10 - R_2$ | $35\ R_4$ | | $2$ |
| $20 - R_4$ | $25 - R_2$ | | $-2$ |
| $50 - R_5$ | $20 - R_1$ | | $-4$ |
| $15 - R_3$ | $30\ R_3$ | | $0$ |
| | | | $\overline{0}$ |

$$P = \frac{1 - 6\sum d^2}{n(n^2 - 1)}$$

$$P = 1 - \frac{6 \cdot (0)^2}{5(5^2 - 1)}$$

$$= \frac{1}{120} = 0.00$$

**Neg**

CRIM price

$\uparrow$

$\downarrow$

**Pos**

Age Price

$\downarrow$ $\uparrow$

$\uparrow$

$\downarrow$ $\downarrow$

# Normalization

$$X_{normal} = \frac{X - min}{max - min} = \frac{65 - 2.90}{100 - 2.90} = \frac{62.1}{97}$$

$$= 0.63$$

Range = 0 to 1

# Standardization

Standarization

$$X_{st} = \frac{X - \mu}{\sigma} = \frac{65 - 68.57}{28} = -0.12$$

$\mu = 0$
$\sigma = 1$

# Logistic Regression

$$\text{Logistic Regression}$$

$$y = \beta_0 + \beta_1 x + \xi$$

$$P = \frac{1}{1 + e^{-y}}$$
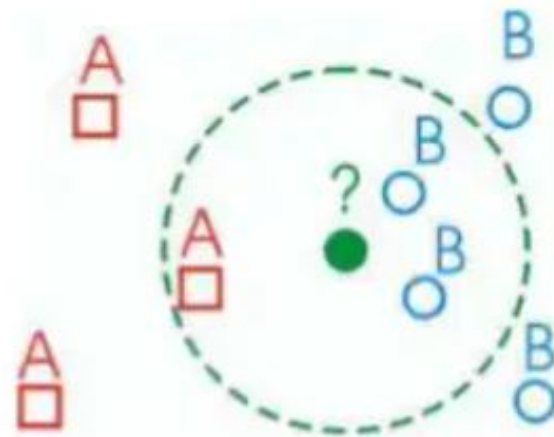
$$\frac{1}{1 + e^{-2}} = 0.88 \rightarrow 1$$

$$\frac{1}{1 + e^{-(-2)}} = 0.11 \rightarrow 0$$
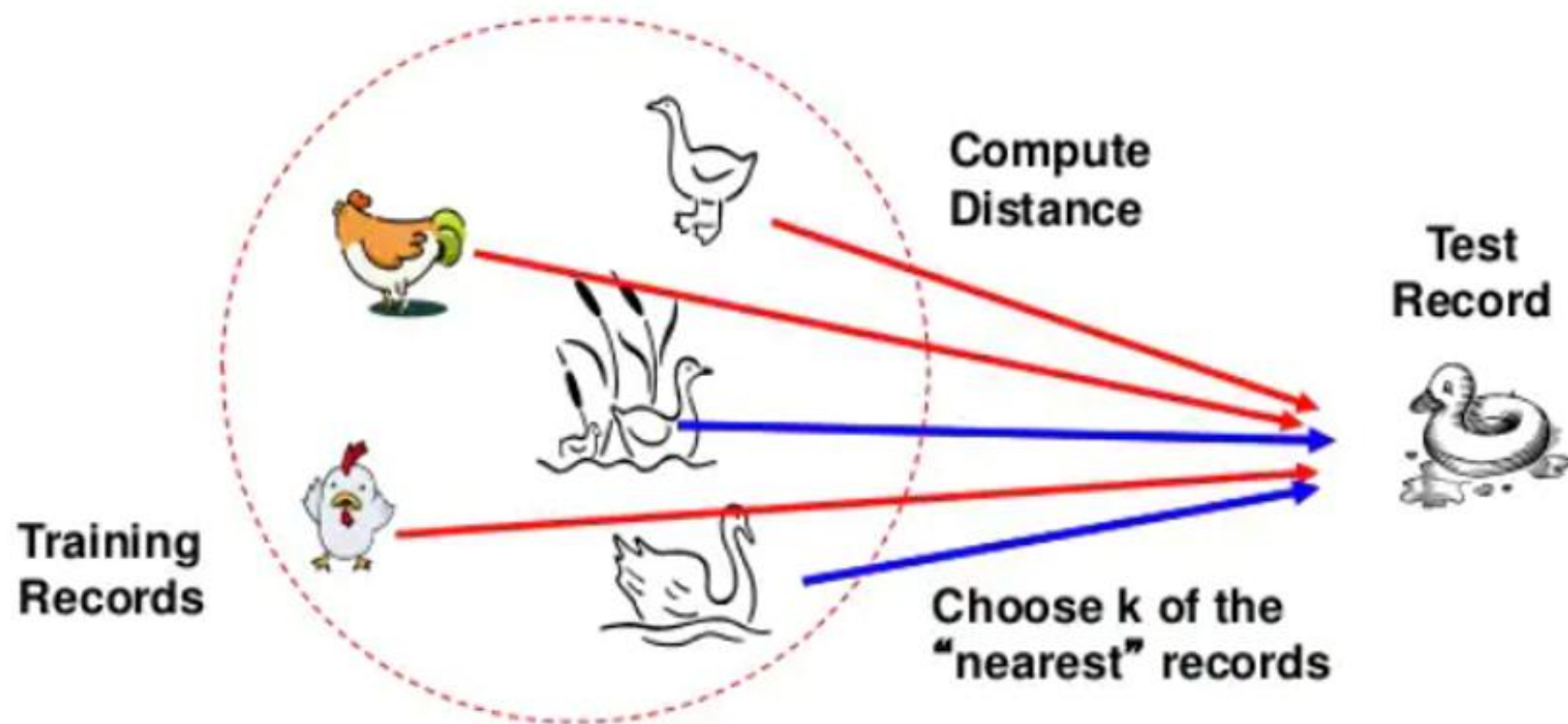
Range: $0 - 1$

# What is KNN?

- A powerful classification algorithm used in pattern recognition.

- K nearest neighbors stores all available cases and classifies new cases based on a *similarity measure*(e.g **distance function**)

- One of the top data mining algorithms used today.

- A non-parametric lazy learning algorithm (An Instance-based Learning method).

# KNN: Classification Approach

- An object (a new instance) is classified by a majority votes for its neighbor classes.

- The object is assigned to the most common class amongst its K nearest neighbors.(*measured by a distant function* )

# Distance Between Neighbors

- Calculate the distance between new example (E) and all examples in the training set.

- *Euclidean* distance between two examples.
  - $X = [x_1, x_2, x_3, .., x_n]$
  - $Y = [y_1, y_2, y_3, ..., y_n]$

  - The Euclidean distance between X and Y is defined as:

$$D(X,Y) = \sqrt{\sum_{i=1}^{n}(x_i - y_i)^2}$$

# Distance Between Neighbors

- Calculate the distance between new example (E) and all examples in the training set.

- *Euclidean* distance between two examples.
    - $X = [x_1, x_2, x_3, .., x_n]$
    - $Y = [y_1, y_2, y_3, ..., y_n]$

    - The Euclidean distance between *X* and *Y* is defined as:
    $$D(X,Y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

# Decision Tree

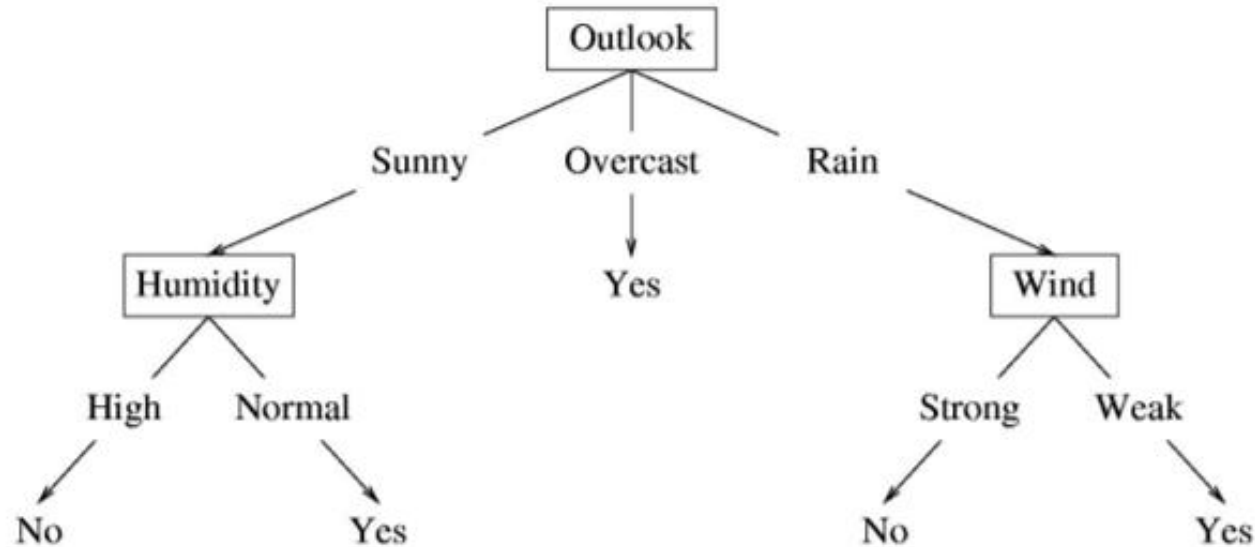## Sample Dataset (was Tennis Played?)

- Columns denote features $X_i$

- Rows denote labeled instances $\langle x_i, y_i \rangle$

- Class label denotes whether a tennis game was played

$\langle x_i, y_i \rangle$

| | Predictors | | | Response |
|---|---|---|---|---|
| Outlook | Temperature | Humidity | Wind | Class |
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

# Decision Tree

- A possible decision tree for the data:



- Each internal node: test one attribute $X_i$

- Each branch from a node: selects one value for $X_i$

- Each leaf node: predict $Y$

# Decision Tree

## Entropy

Entropy measures the impurity in the given dataset. In Physics and Mathematics, entropy is referred to as the randomness or uncertainty of a random variable X. In information theory, it refers to the impurity in a group of examples. **Information gain** is the decrease in entropy. Information gain computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values.

Entropy is represented by the following formula:-

$$Entropy = \sum_{i=1}^{C} -p_i * \log_2(p_i)$$

Here, **c** is the number of classes and **pi** is the probability associated with the ith class.

# Decision Tree

## Entropy

Entropy measures the impurity in the given dataset. In Physics and Mathematics, entropy is referred to as the randomness or uncertainty of a random variable X. In information theory, it refers to the impurity in a group of examples. **Information gain** is the decrease in entropy. Information gain computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values.

Entropy is represented by the following formula:-

$$Entropy = \sum_{i=1}^{C} -p_i * \log_2(p_i)$$

Here, **c** is the number of classes and **pi** is the probability associated with the ith class.

# Rough Work - calculation

## Entropy



$$\frac{3}{3+0} = \frac{3}{3}$$

$$N = \frac{0}{3}$$

$3y \mid 0N$

$3y \mid 3N$

$C_1$

$C_2$

$$-\left(\frac{3}{3}\right)\log_2\left(\frac{3}{3}\right) - \frac{0}{3}\log_2\left(\frac{0}{3}\right)$$

$$+(1)\log_2(1) - 0\log_2(0)$$

$$= 0$$

$$\underbrace{\qquad}_{1}$$

**FN:**

$$-\left(\frac{3}{6/2}\right)\cdot\log\left(\frac{3}{6/2}\right) - \left(\frac{3}{6/2}\right)\log_2\left(\frac{3}{6/2}\right)$$

$$= -\frac{1}{2}\log\left(\frac{1}{2}\right) - \left(\frac{1}{2}\right)\cdot\log\left(\frac{1}{2}\right)$$

$$= 1$$

**GI:**

$$1 - \left[(P_+)^2 + (P_-)^2\right]$$

$$1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2\right] \rightarrow 1 - \left[0.25 + 0.25\right]$$

$$1 - [0.5] = 0.5$$