## Probability & Statistics for Machince Learning & Data Science

# Week 1 : Introduction to Probability and Probability Distribution

The probability of occurance of an event within a sample space of choices is:

$$P(E) = \frac{no.\ of\ Event}{no.\ of\ sample\ space}$$

The complement probability (probability of not occuring) is:

**Complement Rule**

$$P(E`) = 1 - P(E)$$

\# Disjoint Events A & B $\rightsquigarrow P(A \cap B) = 0$

\# Joint Events A & B $\rightsquigarrow 0 < P(A \cap B) \leq 1$

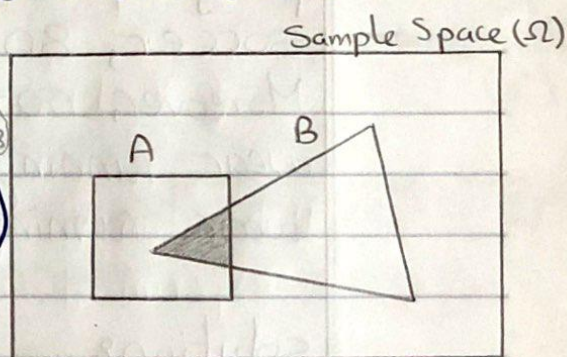Sample Space $(\Omega)$

### Sum Rule

(Added once with A & once with B)

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

\# OR $\rightarrow \cup$           \# AND $\rightarrow \cap$



$\rightarrow$ Independence:

An event is independent of the other, when its occurence doesn't affect of the occurence of the other.

\# If events X & Y are independent, then:

$$P(X \cap Y) = P(X) \cdot P(Y)$$

# General Product Rule

$$P(A \cap B) = P(A) \cdot P(B \mid A)$$

conditional Probability

\# The probability of B given that A is happening; means if A happens or occurs what is the probability of B. This is:

conditional

$$P(B \mid A) = \frac{\text{no. of B exist in A occurence}}{\text{no. of A}}$$

\# if B & A are independent $P(B \mid A) = P(B)$

$$P(A \cap B) = P(A) \cdot P(B)$$

## Example:

In a school of 100 students, 40 of them play soccer. Among the students who play soccer, 80% of them wear running shoes. Moreover, 50% of those of don't play soccer wear running shoes too. How many student wear running shoes?

## Solution:

S → Play Soccer          R → wear running shoes

$P(S) = 0.4$

$$P(R) = P(S \cap R) + P(S' \cap R)$$
$$= P(S) \cdot P(R \mid S)$$
$$+ P(S') \cdot P(R \mid S')$$
$$= (0.4) \cdot (0.8) + (0.6) \cdot (0.5)$$
$$P(R) = 0.62$$

Tree diagram:
- $S$ : $P(S) = 0.4$ → $R$  $P(R \mid S) = 0.8$
- → $R'$  $P(R' \mid S) = 0.2$
- $S'$ : $P(S') = 0.6$ → $R$  $P(R \mid S') = 0.5$
- → $R'$  $P(R' \mid S') = 0.5$

$$P(R) = \frac{\text{no. of R}}{\text{no. of students}} \longrightarrow \text{no. of R} = \underline{62} \text{ Students}$$
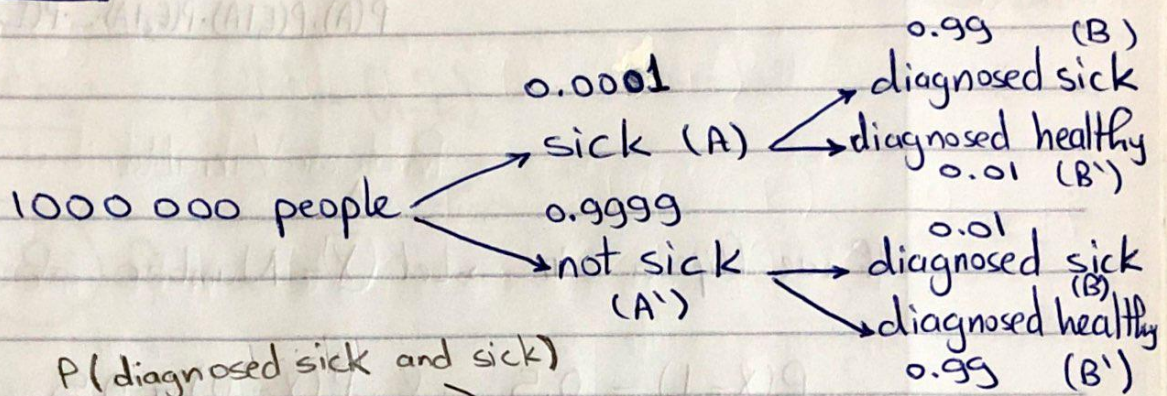
# From general product rule: $P(A|B) = \dfrac{P(A \cap B)}{P(B)}$

{Bayes Theorem}: $P(A|B) = \dfrac{P(A) \cdot P(B|A)}{P(A) \cdot P(B|A) + P(A') P(B|A')}$

Example:

    In a population of 1000000 people, there is a rare disease that gets 1 in each 10000 people. Moreover, the diagnostic test is 99% efficient. Find the probability of being sick given known that diagnosed sick.

solution: A → Sick            B → diagnosed sick

Tree diagram:
1000 000 people →
- 0.0001 sick (A) → diagnosed sick 0.99 (B), diagnosed healthy 0.01 (B')
- 0.9999 not sick (A') → diagnosed sick 0.01 (B), diagnosed healthy 0.99 (B')

P(diagnosed sick and sick)

$$P(\text{sick} \mid \text{diagnosed sick}) = \frac{(0.0001) \cdot (0.99)}{(0.0001) \cdot (0.99) + (0.9999) \cdot (0.01)}$$

P(diagnosed sick)

$$= 0.0098 = 0.98\%$$

another solution:

no. of diagnosed sick = $(0.99 * \overset{\text{sick}}{100} = 99 \text{ people})$

$+ (0.01 * 999900 = 9999 \text{ people})$

$= 10098 \overset{\text{healthy}}{\text{people}}$

no. of sick from those who are diagnosed sick = 99 people

$$P(\text{sick} \mid \text{diagnosed sick}) = \frac{\text{no. sick from diagnosed sick}}{\text{no. of diagnosed sick}}$$

$$= \frac{99}{10098} = 0.0098 = 0.98\%$$

⇒ In the previous example:

A → Prior: 100 out of 1000000 people are sick

E → Event: Diagnostic Test is 99% efficient

P(A|E) → Posterior: P(sick | diagnosed sick) = 0.0098

⇒ Naive Assumption:

It's assuming that events (E) being considered building the model are happening independently. This can ease math a lot, eventhough they are dependent in many cases.

Naive Assumption →

$$P(A|E_1 \& E_2 \& \ldots \& E_n) = \frac{P(A) \cdot P(E_1|A) \cdot P(E_2|A) \cdot \ldots \cdot P(E_n|A)}{P(A) \cdot P(E_1|A) \cdot P(E_2|A) \cdot \ldots \cdot P(E_n|A) + P(A^c) \cdot P(E_1|A^c) \cdot P(E_2|A^c) \cdot \ldots \cdot P(E_n|A^c)}$$

Random Variable
↑
If we flip a coin → Let $X$ = Number of heads ⟨ $X = 1$
  $X = 0$

$$P(X=1) = 0.5 \quad \& \quad P(X=0) = 0.5$$

\# Random Variables allow you to model the whole experiment at once.

| Discret Random Variables | Continuous Random Variables |
|---|---|
| ((Finite Number of Values)) | ((Infinite Number of Values)) |
| ((Can take only Countable Number of values)) | ((Take values on an interval)) |

Not Precise

can be put in a list

⇒ Variables ⟨
  → Deterministic: take fixed outcomes
    $x = 2$ , $P(x) = x^2$
  → Random: take uncertain outcomes
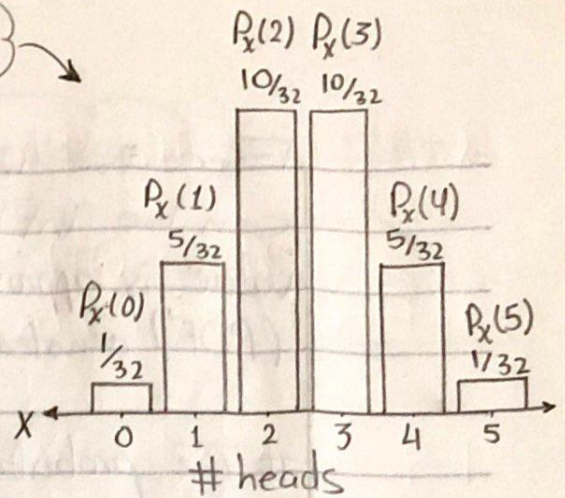    $X$ = number of defective item in a shipment

$X \rightarrow$ Number of heads in
5 coin tosses

⊛Probability Mass Function (PMF):



$P_x(2)$ $P_x(3)$
$10/32$ $10/32$
$P_x(1)$ $5/32$
$P_x(4)$ $5/32$
$P_x(0)$ $1/32$
$P_x(5)$ $1/32$

X-axis: # heads, values 0 1 2 3 4 5

$X = 0, 1, 2, 3, 4, 5$ & $P_x(x) = P(X=x)$

$P_x(x) \geq 0$ & $\sum_x P_x(x) = 1$

\# There are 10 ways to have 2 heads in 5 coin tosses:

Number of ways you can order 5 coins

$$\frac{5!}{2! \, (5-2)!} = 10 = C_2^5 = \binom{5}{2} \Rightarrow \text{Binomial Coefficient (combination)}$$

Number of heads $\rightarrow 2!$ $(5-2)! \rightarrow$ Number of not heads

Property: $\binom{n}{k} = \binom{n}{n-k}$  counts all combinations

All possible orders   Probability of seeing x heads

Binomial: $P_X(x) = \binom{5}{x} p^x (1-p)^{5-x}$, $\rightarrow$ Probability of seeing $5-x$ not heads

$x = 0, 1, 2, 3, 4, 5$

Event $X = x$ : $x$ is heads in 5 tosses
$\longrightarrow$ X follows a binomial distribution
$\longrightarrow$ $X \sim$ Binomial $(5 \cdot p)$

Number of Flips $\nearrow$   $\searrow$ P(H)

Bernoulli:
$\rightarrow$ Success: Occurance of the preferable event (X)
    $P$
$\rightarrow$ Failure: Not occurance of the peferable event (X')

Discrete   |   Continuous

\#Sum of heights equals 1.   \# Area under curve equals 1.

$$\sum_x P_x(x) = 1$$   $$\int_x P_x(x) \, dx = 1$$

Because in continuous random variables ↑values number of
can be infinite, and the probability of an exact
value is approximately zero. The probability density function
(PDF) denoted as $f_X(x)$ uses intervals of $x$ instead.

\# The probability density function $f_X(a \leq x \leq b) = \int_a^b p_X(x) \, dx$

\# The commulative distribution function $f_X(x \leq a) = \int_0^a p_X(x) \, dx$
$0 \leq CDF \leq 1$ and denoted as $\boxed{F_X(x)}$.
↝ It is a curve that starts from zero to 1 at the
end; where zero and one are the heights.

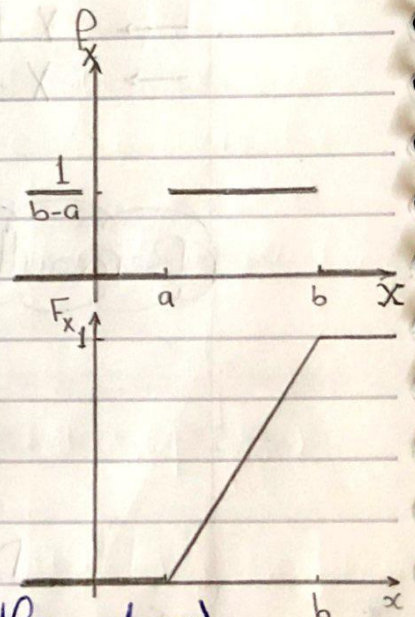↝ $F_X(x)$ can never decrease, it only increases till 1.

## Uniform Distribution

A continuous random variable can be modelled with it,
if all possible values lie in an interval have the same
frequency of occurance. Its parameters are:
⟶ a : beginning of the interval
⟶ b : end of the interval

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & x \notin (a,b) \end{cases}$$

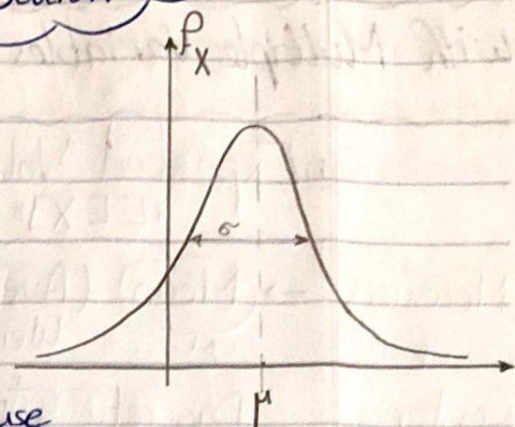$$F_X(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & a \leq x < b \\ 1 & b \leq x \end{cases}$$

\# mean ⟶ $\mu$ (average) or (center of the values)
\# standard deviation ⟶ $\sigma$ (measures the spread of values)

$$\overparen{\text{Normal (Gaussian)}\atop\text{Distribution}}$$

# It is a distribution that is symmetric and takes the bell-shaped.

# The closest function to this curve is $e^{-\frac{x^2}{2}}$. But because it doesn't fit well, we use this formula instead:

$$P_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \rightsquigarrow X \sim \mathcal{N}(\mu, \sigma^2)$$
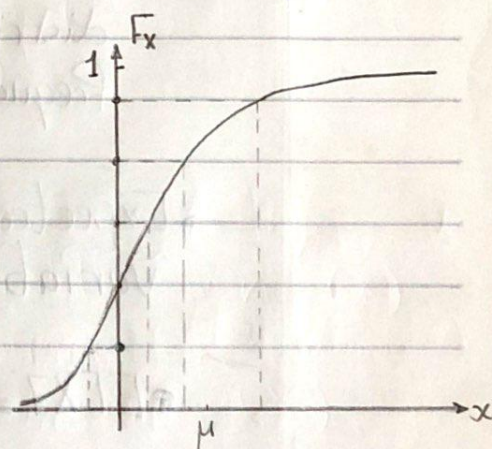
# Standardization: The way of transforming any normal distribution to a standard one.

→ we calculate $Z = \frac{x-\mu}{\sigma}$, then the new $\mu = 0$ & $\sigma = 1$

$$P_X(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$$

# Most of the natural phenomena can be represented using normal distribution.

→ To sample from distributions the y-axis of the CDF ($F_X$) gets divided into the number of samples, then these values gets translated into the sample values.

# erf(x) = $\frac{1}{\sqrt{\pi}}\int_0^x e^{-t^2} dt$

# erf → error function

# erfinv is from scipy.special

## Code

### Gaussian
$$y = F(x) = \frac{1}{2}\left[1 + \text{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right)\right]$$

$$x = F^{-1}(y) = \left[\sigma\sqrt{2} \cdot \text{erf}^{-1}(2y-1)\right] + \mu$$

### Binomial
$$y = F(x) = \sum_{k=0}^{\lfloor x \rfloor} \binom{n}{k} p^k (1-p)^{n-k}$$

$$0 \le x \le n$$

$$x = F^{-1}(y) = \text{scipy.stats.binom.ppf}(y, n, p)$$