

Statistics

① define problem: — hypothesis

② collect the data

③ organize the data

④ cleanse the data

⑤ perform operation / algorithm : find result

⑥ hypothesis testing

↳ distributions - (T) F, chi² success
failure

ML-ops - workflow

→ R

→ python

x	y
1	3
2	5
3	7
4	9
5	11

known data

$$y = 2x + 1$$

formula = model (lm)

predict() → (R)

unknown data

$$x = 1 \text{ to } x = 100$$

$y = \{ \dots \}$

distribution

def generate_values(x):
return 2x + 1

→ distribution function

python



Sci-kit learn



Linear Regression



model

Distribution ∘ function

→ going to generate listing

- The distribution of a statistical data set (or a population) is a listing or function showing all the possible values (or intervals) of the data and how often they occur
- When a distribution of categorical data is organized, you see the number or percentage of individuals in each group → factors → levels → unique values
- When a distribution of numerical data is organized
 - they're often ordered from smallest to largest, broken into reasonably sized groups (if appropriate)
 - then put into graphs and charts to examine the shape, center, and amount of variability in the data
- Types
 - Discrete probability distribution
 - Continuous probability distribution



Terminologies

- **Frequency distribution**
 - A frequency distribution is a table that displays the frequency of various outcomes in a sample. → unique values
- **Relative frequency distribution:**
 - A frequency distribution where each value has been divided (normalized) by a number of outcomes in a sample i.e. sample size.
- **Probability distribution:**
 - alias for *Relative frequency distribution*
 - indicates the way the total probability of 1 is distributed over all various possible outcomes



$x_1 = \{ \underline{10}, \underline{20}, \underline{10}, \underline{15}, \underline{18}, \underline{20}, \underline{18}, \underline{10}, \underline{15} \}$

frequency distribution

x	frequency
10	3
20	2
15	2
18	2
9	

Relative frequency distribution or probability distribution

x	frequency	relative frequency / probability
<u>10</u>	3	$p(10) = 3/9$
<u>20</u>	2	$p(20) = 2/9$
<u>15</u>	2	$p(15) = 2/9$
<u>18</u>	2	$p(18) = 2/9$
9		$= \frac{3+2+2+2}{9} = 1$

total outcomes

addition of all probabilities = 1

$$\alpha = \{ 5, 8, 9, 7, 3, 4, 8, 9, 10, 2, 3, 5, 7 \}$$

Discrete Probability Distribution



Discrete Probability Distribution (function)

- Is also called as probability mass function (PMF)
- The sum of all the individual probabilities must equal to 1

- Types

- ① ✓ Uniform distribution : same probabilities for all outcomes
- ② ✓ Binomial distribution : only two possible outcomes (n Bernoulli trials)
- ③ ✓ Negative Binomial distribution
- ④ ✓ Poisson distribution
- ⑤ ✓ Geometric distribution



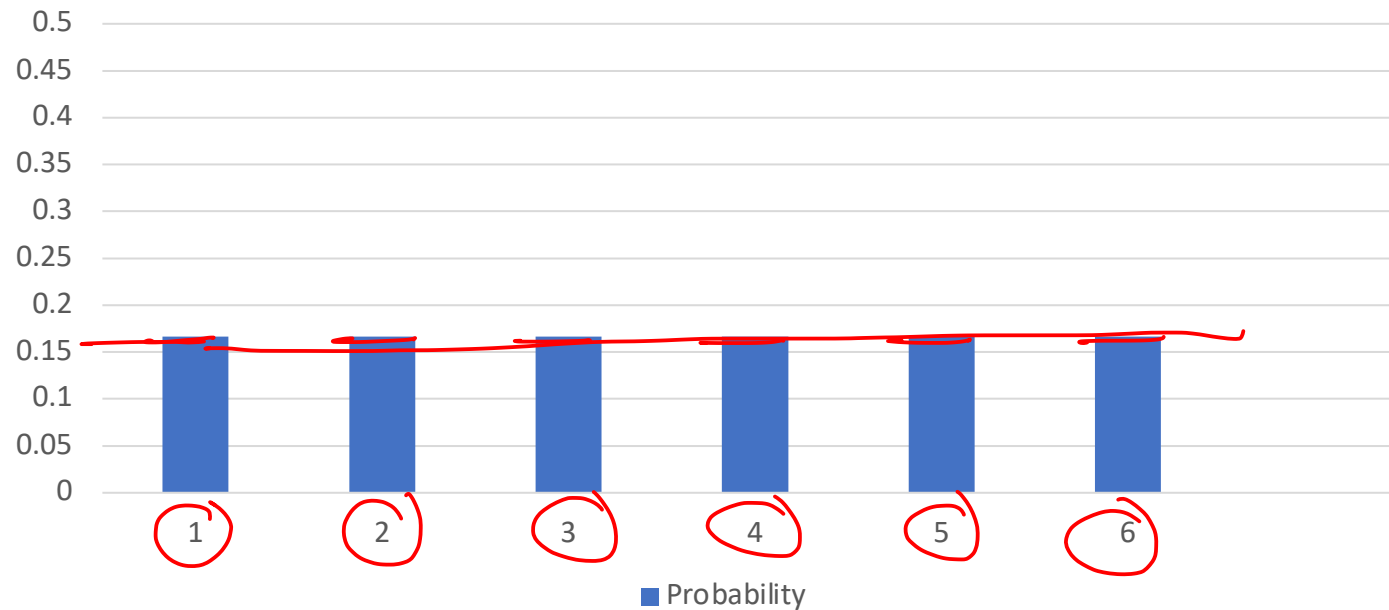
- ① Binomial : probability of x success in n trials
- ② Negative Binomial : find no of trials for constant no of success
- ③ Geometric : find no of trials needed for first success
- ④ Poisson : find no of success in n trials

Uniform distribution

- The probabilities of each outcome are evenly distributed across the sample space
- E.g. Rolling a fair die has 6 discrete equally probable outcomes

$$\begin{aligned} p(1) &= 1/6 \\ p(2) &= 1/6 \\ &\vdots \\ p(6) &= 1/6 \end{aligned}$$

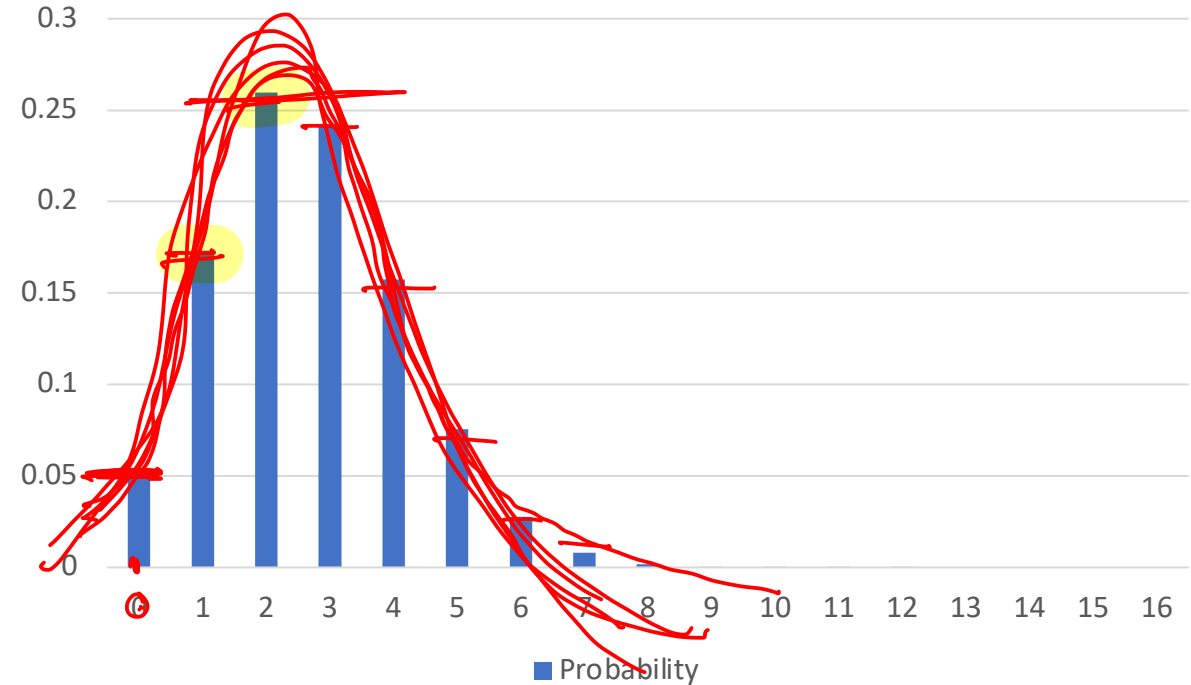
$$6/6 = 1$$



Binomial distribution

- Binomial means two discrete outcomes
- The outcomes of a trial are mutually exclusive
- Bernoulli Trial → tossing a coin
 - When a trial is conducted there will be only two possible outcomes
 - Success or Failure, true (false, 1 & 0, H & T)
- A series of n trials follow binary distribution when
 - The probability of success (p) is a constant
 - Trials are independent of one another
- E.g. Flipping a coin

(H) → not (T)



Binomial distribution function

- Probability mass function

- Generates probability of observing x successes in n trials
- Probability of single trial is constant (p)

2 - 10 trials

$$1 = 10 = 2 = 1$$

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

$$2 = 10 = 0 = 0$$

$$3 = 10 \Rightarrow 10 = 0$$

1000

success

→ head

$$p(x; n, p) = \binom{n}{x} p^x (1-p)^{(n-x)}$$

$$x \binom{n}{x}$$



Negative Binomial distribution (function)

2 success

- Negative Binomial Distribution is of number of trials needed to get a fixed number of successes
- Each trial results in one of the two possible outcomes
 - Success
 - Failure
- Probability of success = $P(\text{Success}) = p$ (constant)
- Probability of failure = $P(\text{Failure}) = 1 - p$
- X represents the trial number of the rth success
- Probability Mass Function

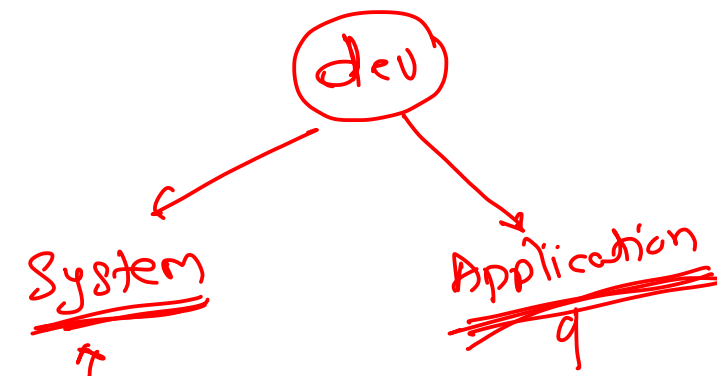
$r = 2$
 $n = 10$

$2 \text{ success} = \frac{2}{n} \binom{n}{r} \rightarrow \boxed{v} \quad 0.5$

str = abcdef
sub = (cd)
→ a' b' c' d'

$r = 2$
 $p(H) = 0.5$

$P(x) = \binom{x-1}{r-1} * (1-p)^{x-r} * p$



Geometric distribution

- Is a distribution of the number of trials needed to get the first success in the repeated Bernoulli trials
- Every trial results in one of two possible outcomes
 - Success
 - Failure
- $P(\text{Success}) = p$ (stays constant)
- $P(\text{Failure}) = 1 - p$
- X represents the number of trials needed to get the first success
- For the first success to occur on the x th trial:
 - The first $x - 1$ trials must be failures
 - The x th trial must be a success
- Probability mass function

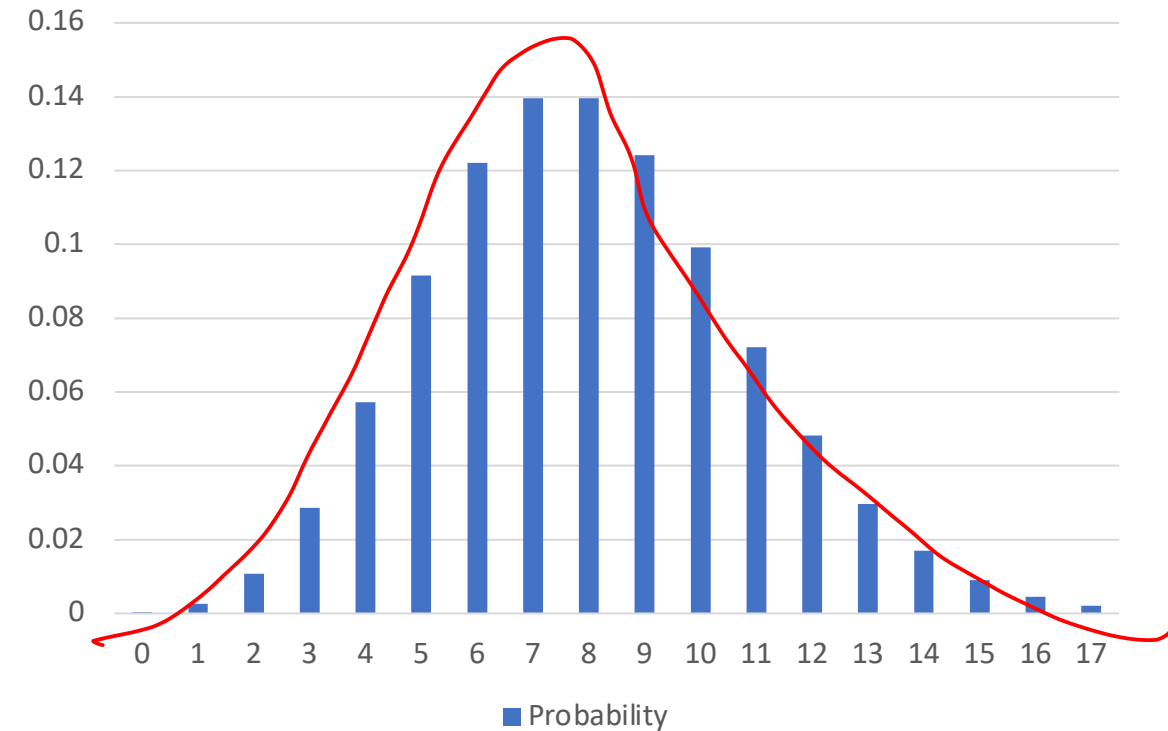
$$P(x) = (1 - p)^{x-1} * p$$



Poisson distribution

- A Binomial distribution considers number of successes out of n trial
- A Poisson distribution considers number of successes per unit of time over the course of many units
- Terminology
 - X: Variable
 - E: Expectation (Mean) $\Rightarrow E(X)$
 - μ : Expectation (Mean) $= E(X)$
 - $\lambda = \frac{\# \text{ occurrences}}{\text{interval}} = \mu$
- Probability Mass Function

$$p(x) = \frac{\lambda^x e^{-\lambda}}{x!}$$



interval

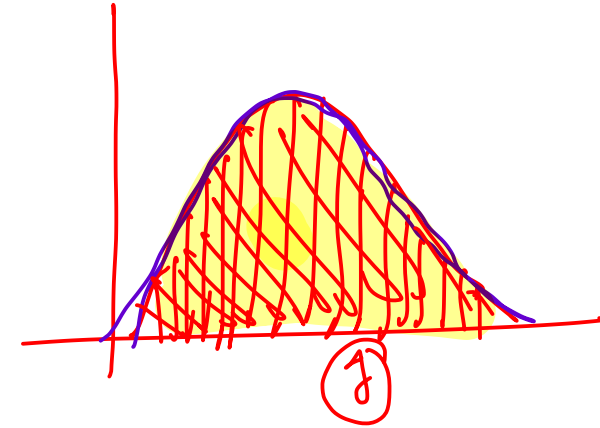
Continuous Probability Distribution



Continuous Probability Distribution (PDF)

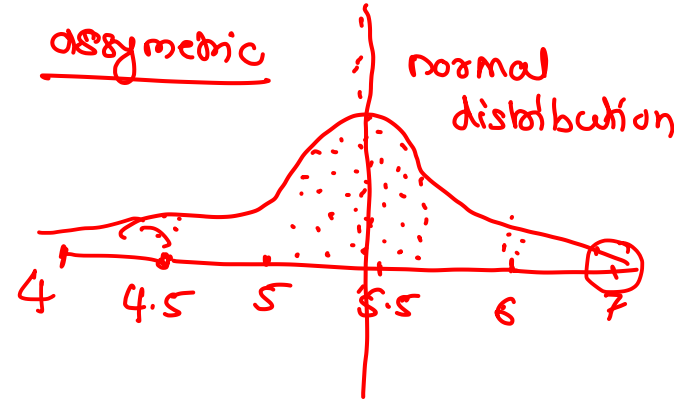
- Also called as probability density functions (PDF)
- Where the variable values are infinite \rightarrow interval
- The area under the probability curve equals 1 \rightarrow
- Types
 - ✓ Normal distribution
 - ✓ Student T distribution
 - ✓ Chi-Squared distribution
 - ✓ Exponential distribution
 - ✓ Logistic distribution

\rightarrow classification model
evaluation
(AUC)

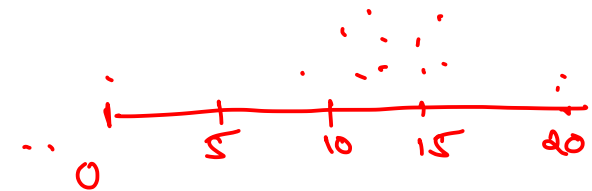


Normal Distribution

- Data sources tend to be around central value with no bias left or right
- Many real life data points follow normal distribution
- E.g.
 - Peoples height and weights
 - Population blood pressure
 - Test scores
 - Measurement errors
- Probability distribution function

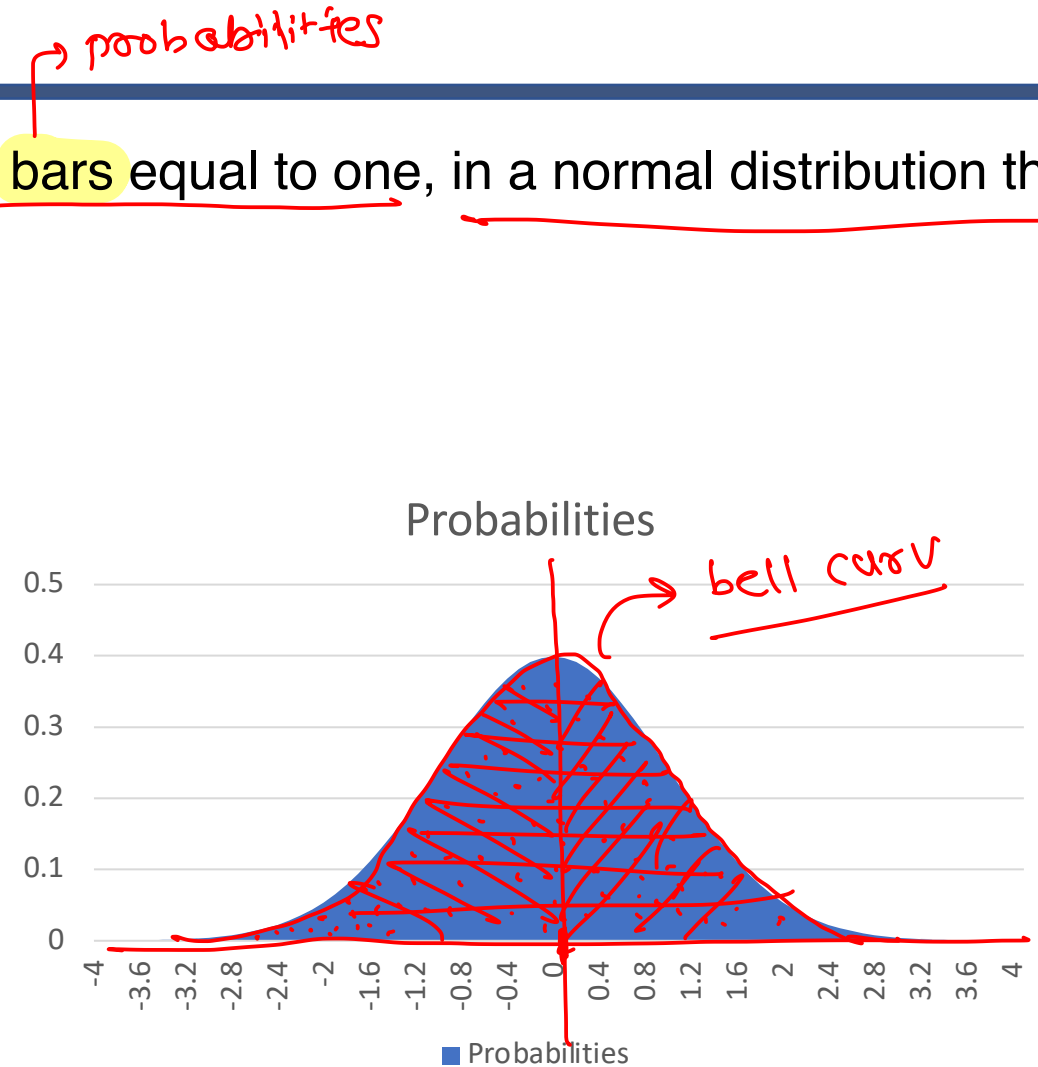
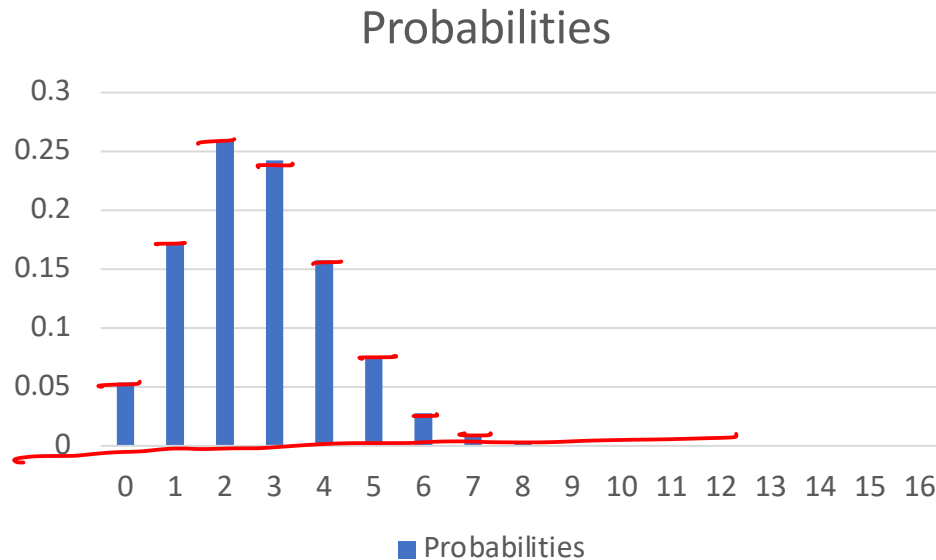


$$f(x) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} * e^{\frac{-(x-\mu)^2}{2\sigma^2}} \right)$$



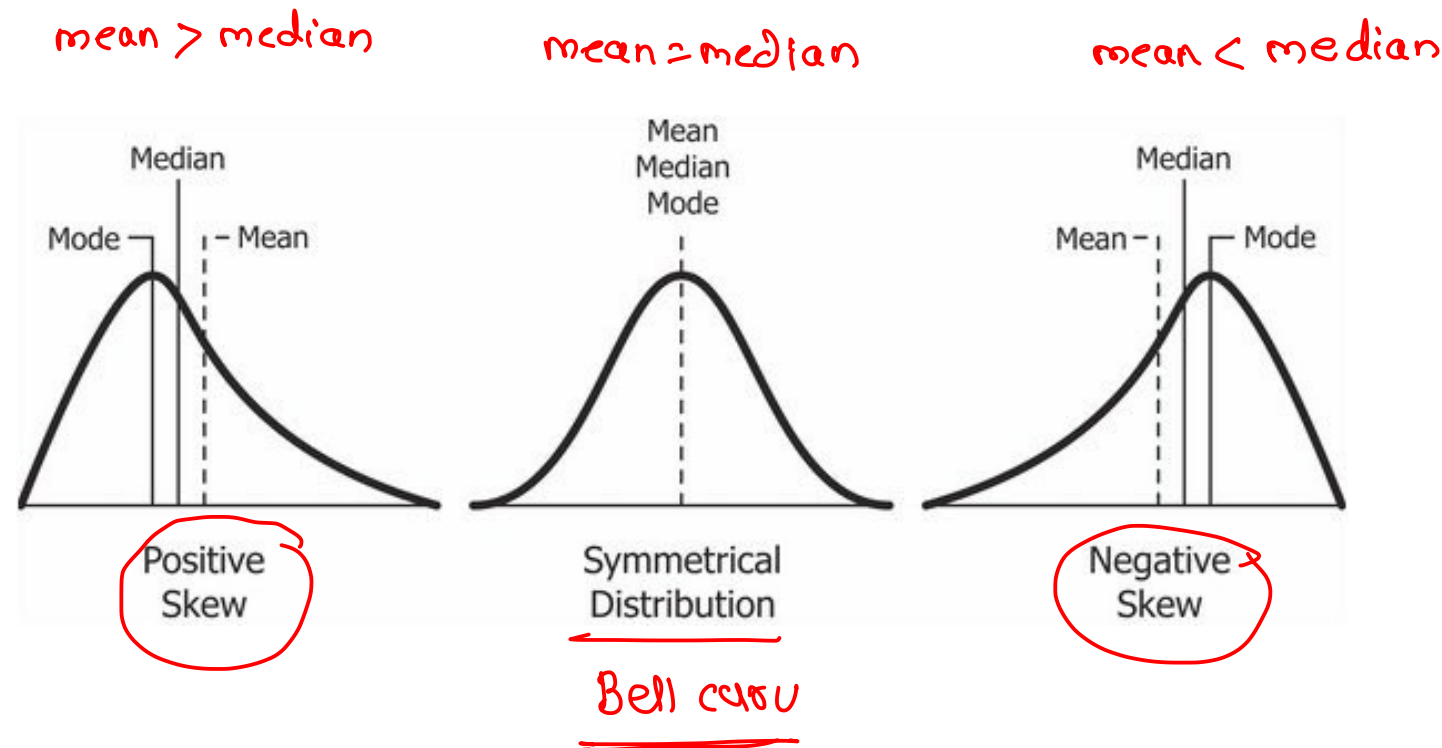
Normal Distribution

- Unlike discrete distributions, where the sum of all the **bars** equal to one, in a normal distribution the area under the curve equals to 1



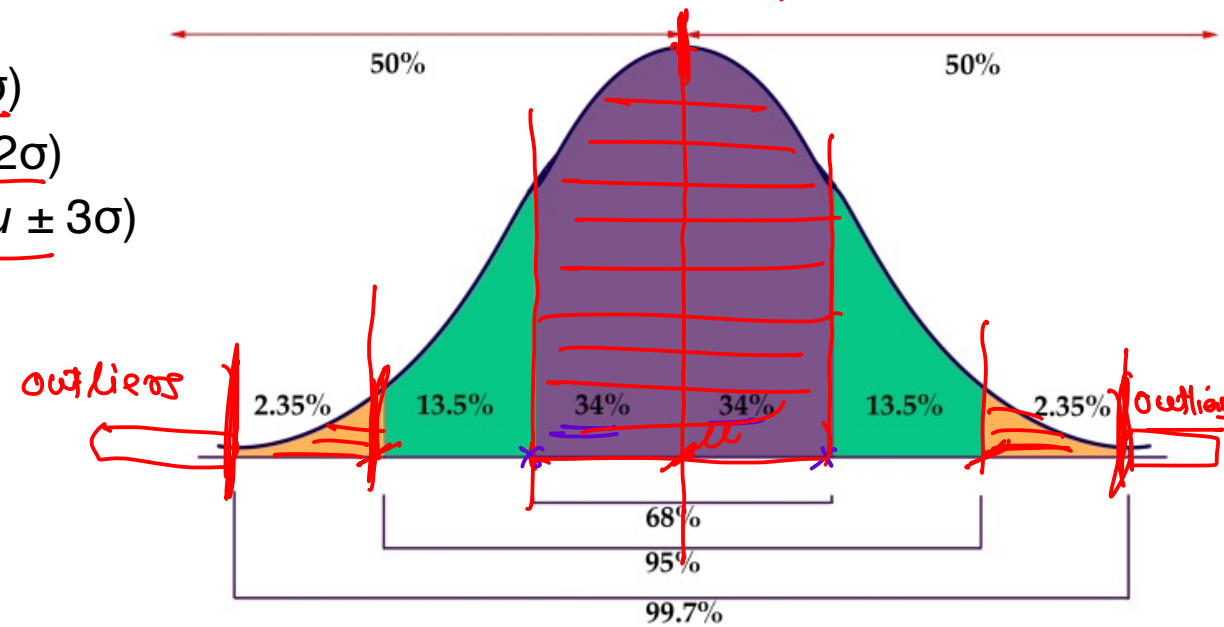
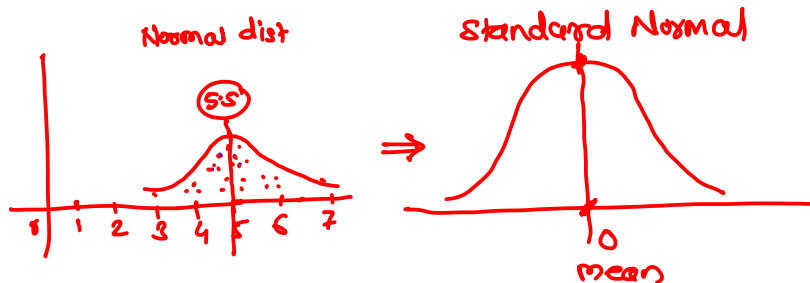
Normal Distribution

- Also called as Bell Curve or Gaussian Distribution
- Always symmetrical
- Asymmetrical curves display skew and are NOT normal



Standard Normal Distribution

- The standard normal distribution is a special case of the *normal distribution*
- It is the distribution that occurs when a normal random variable has a mean of zero and a standard deviation of one *[$\mu = 0, \sigma = 1$]*
- Empirical Rule *68-95-99.7*
 - Often used for forecasting
 - For a normal distribution, almost all data falls within three standard deviations (σ) of the mean (μ)
 - The empirical rule shows that
 - 68% falls within the first standard deviation ($\mu \pm \sigma$)
 - 95% within the first two standard deviations ($\mu \pm 2\sigma$)
 - 99.7% within the first three standard deviations ($\mu \pm 3\sigma$)



Standardizing Normal Distribution

- We can take a normal distribution and standardize it to a standard normal distribution
- If we can model our data as a normal distribution we can convert the values in the normal distribution to a standard normal distribution to calculate percentile
- To calculate z-score use formula:

z-value →
$$z = \frac{x - \mu}{\sigma}$$

x : observation
 μ : mean
 σ : std devian

- Using Z-table we can get the percentile of the value
- Since total area under the curve is 1, probabilities are bounded by 0 and 1



Central Limit Theorem

- The mean values from a group of samples will be normally distributed about the population mean, even if the population itself is not normally distributed
- That is, 95% of all sample means should fall within 2σ of the population mean
- To read more about CLT, please visit
 - https://en.wikipedia.org/wiki/Central_limit_theorem

