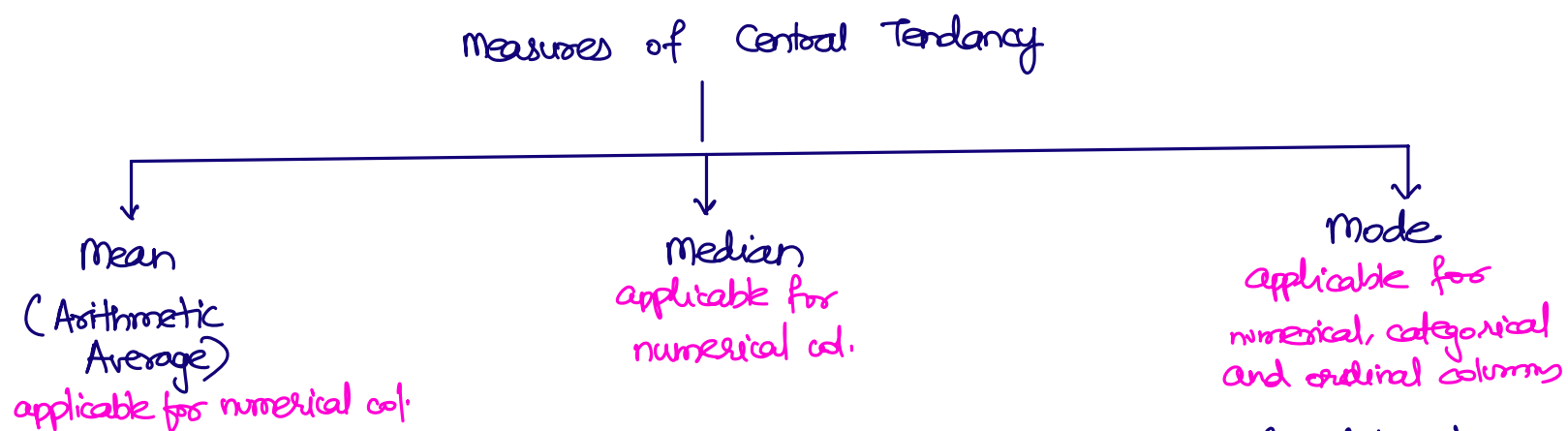## Descriptive Statistics

## Measures of Central Tendancy:

These are statistical tools used to summarize the data based on the central point of the data.

Measures of Central Tendancy

Mean
(Arithmetic Average)
applicable for numerical col.

Median
applicable for numerical col.

Mode
applicable for numerical, categorical and ordinal columns

ⓐ Mean: The sum of all data points divided by number of data pts.

$$mean = \frac{\sum_{i=1}^{n} x_i}{n}$$

$n \dots$ no. of records

eg: $[2, 3, 4, 5]$

$$mean = \frac{2 + 3 + 4 + 5}{4} = 3.5$$

ⓑ Median: Its the middle value of the dataset when its ordered from smallest to largest.

Calc. median

① Order the data in ascending order.

② Identify the middle position

- if no. of records is <u>odd</u>, the median value is at $\left(\frac{n+1}{2}\right)^{th}$ position.

eg: $[3,1,4,1,5]$
$\rightarrow [1,1,3,4,5]$

Here,
the middle position will be $\frac{5+1}{2} = 3^{rd}$ position $\Rightarrow \underline{\underline{3}}$

- if no. of records is even, the median value is the avg of $\left(\frac{n}{2}\right)^{th}$ position value and $\left(\frac{n}{2}+1\right)^{th}$ position value.

eg: $[3,1,4,1]$
$\rightarrow [1,1,3,4]$

Here,
middle position $\left(\frac{4}{2}\right) = 2^{nd}$ & $\left(\frac{4}{2}+1\right) = 3^{rd}$

$$median = \frac{1+3}{2} = \underline{\underline{2}}$$

③ Mode $\rightarrow$ most frequent data. / data with highest frequency (count)

<u>steps:</u>
① Count the frequency of each unique value.
② Identify the value with highest frequency.

eg $[1,2,2,3,3,3,4]$

$1 \rightarrow 1$      $4 \rightarrow 1$          $\therefore mode = \underline{\underline{3}}$
$2 \rightarrow 2$
$3 \rightarrow 3$

**Note:** if all unique data has same frequency, in that case, sort the data in ascending order and return the first value
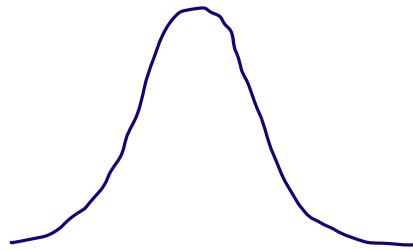
eg: [1,2,3,4] $\longrightarrow$ $\frac{1}{=}$

[1,2,2,3,3] $\rightarrow$ 1 → 1
2 → 2 }
3 → 2 } $\longrightarrow$ sort $\rightarrow$ 2//

## Data Distributions

This concept is applicable only for numerical columns

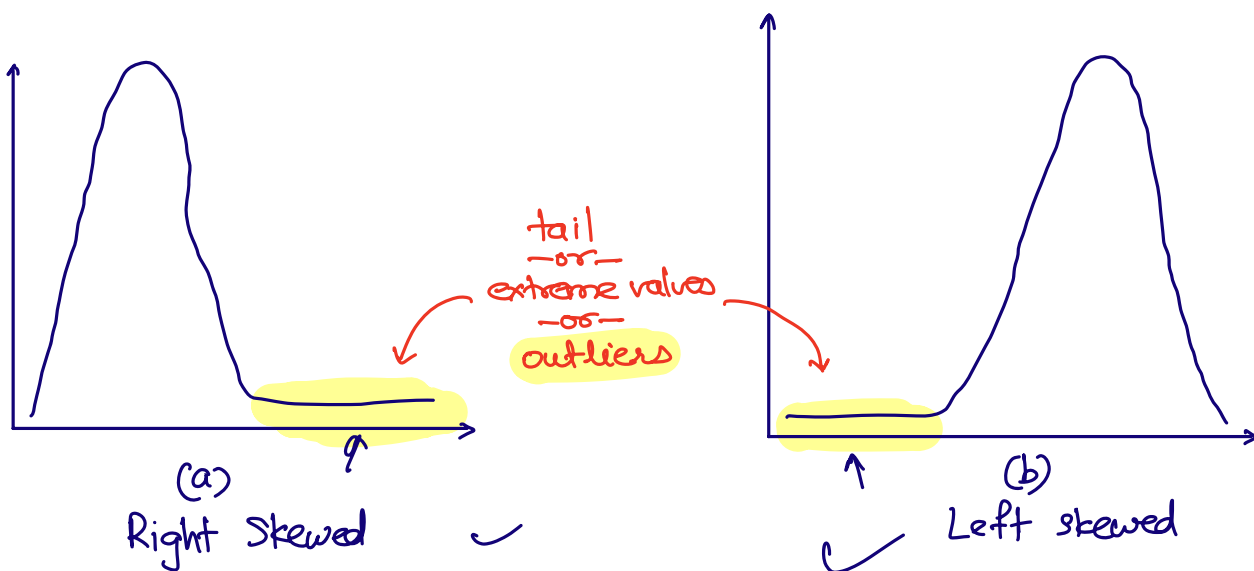ⓐ Normal Distribution | Gaussian Distribution | Bell curve

GENERALIZED
DATA



if data column forms a bell curve, chances are data is generalized.
(the data may represent the population)

ⓑ Skewed distribution.



tail
—or—
extreme values
—or—
outliers

(a)
Right Skewed ✓

(b)
Left skewed ✓

# Quartiles

Quartiles are the values that divide the dataset into 4 equal parts.

25%     50%     75%     100%

$Q_1$     $Q_2$     $Q_3$     $Q_4$

First     Second     Third     maximum
Quartile  Quartile   Quartile
          (median)

Consider a dataset $D = \{12, 15, 7, 23, 10, 8, 14, 18, 5, 20\}$

To calc quartiles

**step 1:** Arrange the data in ascending order

$$\{ 5, 7, 8, 10, 12, 14, 15, 18, 20, 23 \}$$

1  2  3  4  5  6  7  8  9  10

**Step 2:** Identify the median

median = 13 ⟹ $Q_2$

**Step 3:** Calc first quartile ($Q_1$)

$Q_1$ is the median of the lower half of the dataset

$$Q_1 = 8$$

**Step 4:** Calc third quartile ($Q_3$)

$Q_3$ is the median of the upper half of the dataset

$$Q_3 = 18$$

**Step 5:** Calc fourth quartile ($Q_4$)

$Q_4$ is the max value

$$Q_4 = 23$$

# Outliers

Outliers are those extreme values that affect the general tone of the domain.

To identify & remove an outlier, you can use Tukey's Method
-or-
1.5 IQR rule.

| Tukey's Method | 1.5 IQR rule |
| --- | --- |

$\longrightarrow$ Inter Quartile Range

$$IQR = Q3 - Q1$$

## Algo:

① Calc IQR

$$IQR = Q3 - Q1$$

② Calc the valid range of the given column.

$$lowerRange = Q1 - (1.5 * IQR)$$

$$upperRange = Q3 + (1.5 * IQR)$$

outliers                                          outliers

$-\infty$ ————|————————————|————————————|———— $+\infty$

lowerR                    O                    upperR