## Measure of central Tendency
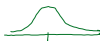
Measures of central tendency are statistical metrics that describe the center point or typical value of a dataset. They provide a single value that summarizes a set of data by identifying the central position within that dataset.

1) Mean or average     Ages = [24, 32, 12, 48, 16, 20]
2) Median
3) Mode

$\downarrow$

Center point
central position



### 1) Mean

Mean is the sum of all values divided by the number of values.

Population Mean ($\mu$)        Sample mean ($\bar{x}$)

Population (N)     $n \leq N$     Sample(n)

Sample n is a subset of population.

$$\mu = \sum_{i=1}^{N} \frac{X_i}{N}$$        $$\bar{x} = \sum_{i=1}^{n} \frac{X_i}{n}$$

Here $X$ is a random variable $\{N \rightarrow$ population size$\}$     $n \rightarrow$ is sample size

$X = \{5, 8, 12, 15, 20\}$
$N = 5$

$\mu = \dfrac{5+8+12+15+20}{5} = \dfrac{60}{5} = 12$

### Characteristics

1) Affected by extreme outliers.

2) Used for interval and ratio data

without outlier
$X = \{1, 2, 8, 4, 5\}$
$\mu = \dfrac{1+2+3+4+5}{5} = 3$

with outlier
$X = \{1, 2, 3, 4, 5, 100\}$
$\mu = \dfrac{1+2+3+4+5+100}{6} = \dfrac{115}{6} = 19.166$

### 2) Median

The median is the middle value in a dataset, where the values are arranged in ascending or descending order.

$X = \{1, 2, 3, 4, 5\}$       $X = \{3, 4, 1, 8, 2, 100\}$

The numero de elements is 5, 5 is odd    $= \{1, 2, 3, 4, 5, 100\}$

Median = 3        No of elements $= 6$
       6 is even
       Median $= \dfrac{3+4}{2} = 3.5$

### Characteristics

i) Not affected by extreme outliers.

ii) Used for ordinal interval and ratio data.

### 3) Mode

Definition: The mode is the value that appears most frequently in a dataset.

Dataset: 2,4,4,6,7,7,7,9

Mode = 7 (most frequent value)

Mode = 5, 6 $\{$ bimodal $\}$

*) Characteristics

1) Not effected by extreme values.
2) used for nominal, ordinal interval, and ratio data.

### Choosing the appropriate measure

1. Mean: Best used when data is symmetrically distributed without outliers. Provides a mathematical average, which is useful for further stadistial calculations.



1. Median: Best used when data is skewed or contains outliers. Provides the middle value, which better represents the center of a skewed dataset.



outlier

1. Mode: Best used for categorical data to identify the most common category. Also useful for identifying the most frequent value in ordinal, interval or ratio data.

### Real word application

Feacture engineering   $\Rightarrow$ EDA

for missing   Mode
Mode      $\uparrow$ {Nominal}
       + (ordinal)

mean
median

| | Age | Weight | Salary | Gender | Degree | + |
|---|---|---|---|---|---|---|
| 1 | 24 | 70 | 40k | M | BE | |
| 2 | 25 | 80 | 70k | F | - | |
| 3 | 27 | 95 | 45k | F | - | |
| 4 | 24 | - | 50k | M | PHD | |
| 5 | 32 | - | 60k | - | BE | |
| 6 | - | 60 | - | - | Master | |
| 7 | - | 65 | 55k | - | BSC | |
| 8 | 40 | 72 | - | M | BE | |
| + | | | | | | |

Handling the missing value