

# Summary Statistics

## **Numerical**

- Centrality measure ( mean, median )
- Dispersion measure ( range, percentiles, variance , standard deviation )

## **Categorical**

- Total count
- Unique count
- Category Counts and proportions
- Per category statistics

Centrality  
Measure

**One number to represent entire set of values**

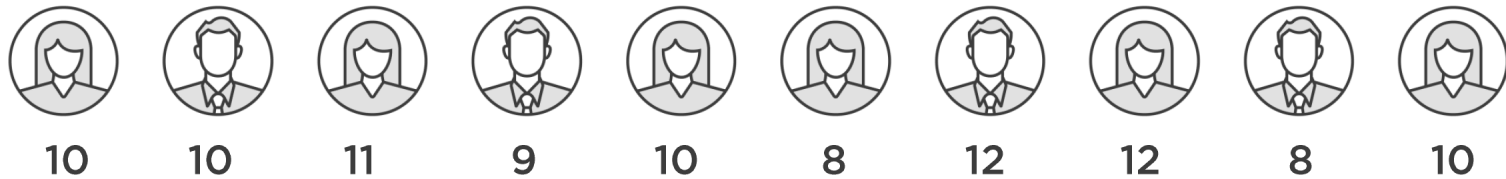
**Number central to the data**

**Central tendency**

Mean / Average

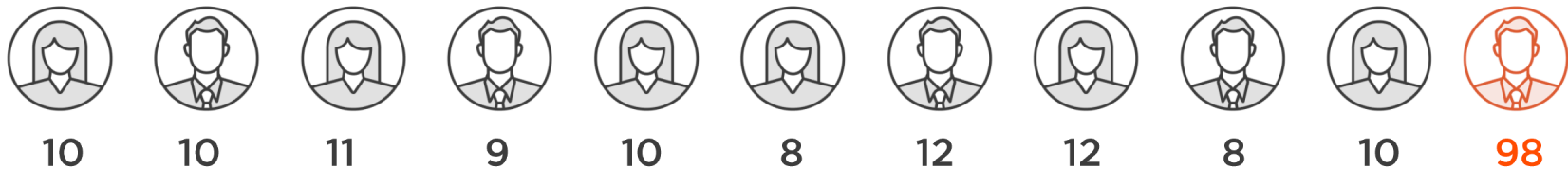
**Average behavior**

## Centrality Measure : Mean or Average



Mean age :  $\text{sum of ages} / \text{count} = 100 / 10 = 10$

Problem : Affected by extreme values



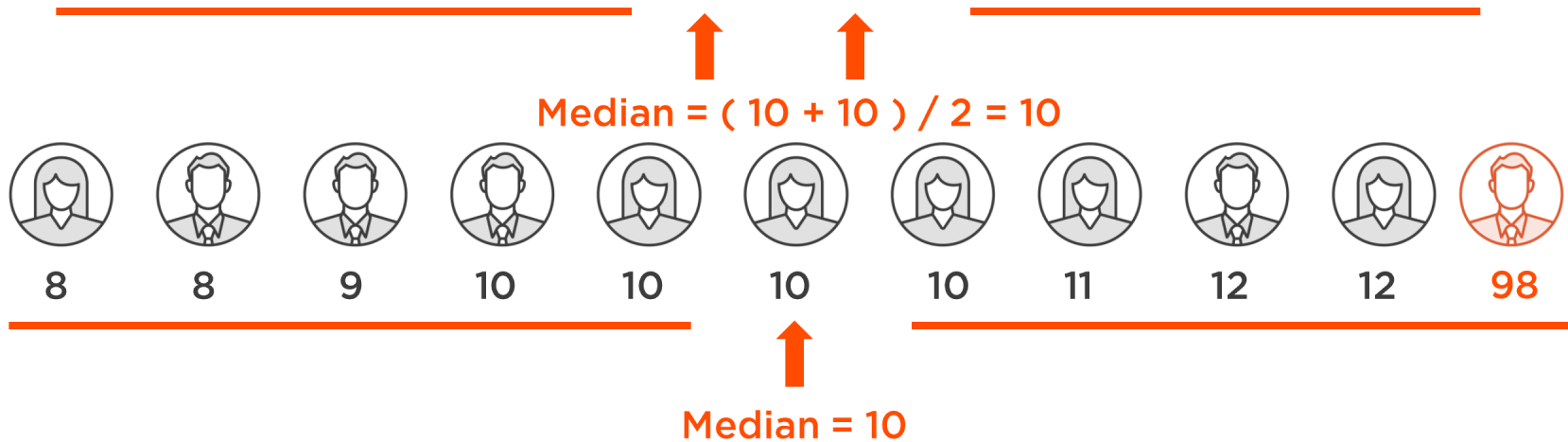
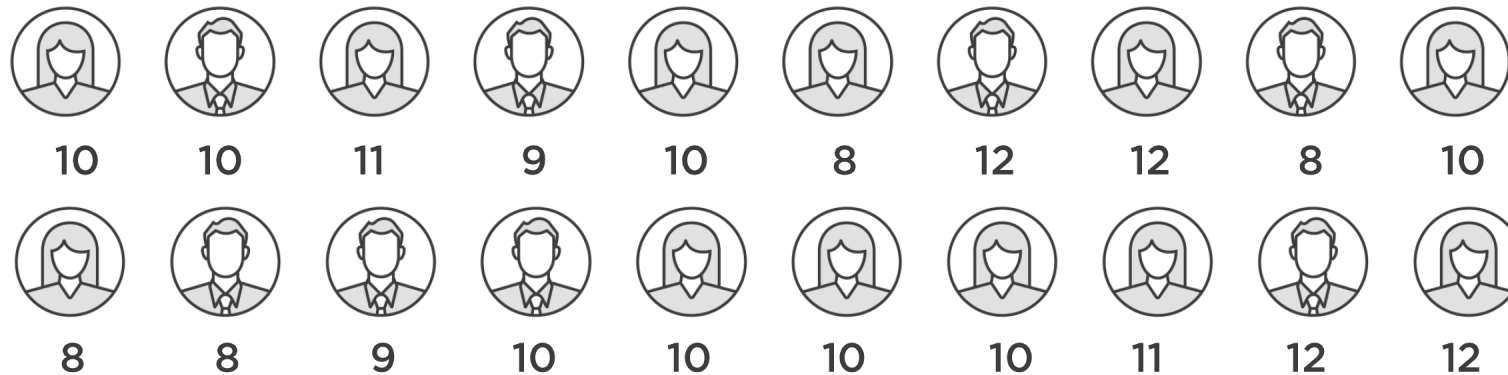
Mean age :  $\text{sum of ages} / \text{count} = 198 / 11 = 18$

A diagram illustrating the relationship between the Median and the Middle value in the sorted list. A vertical orange line is positioned between the two text labels. The word "Median" is on the left, and the phrase "Middle value in the sorted list" is on the right. The entire diagram is enclosed in a black border.

Median

**Middle value in the sorted list**

## Centrality Measure : Median



---

Spread /  
Dispersion  
Measure

How spread out values are from central  
value

Variability

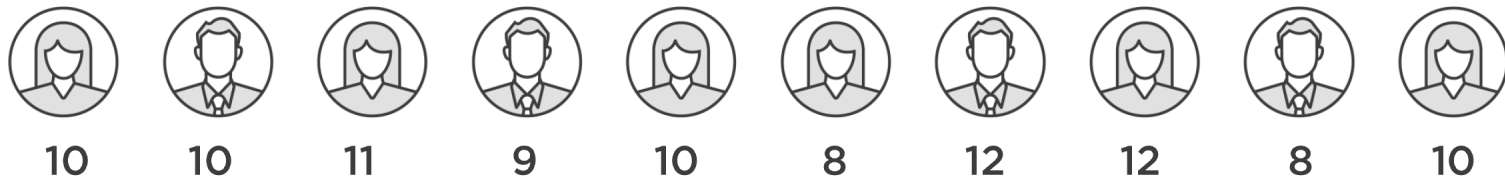


Range

**Difference between maximum and  
minimum**

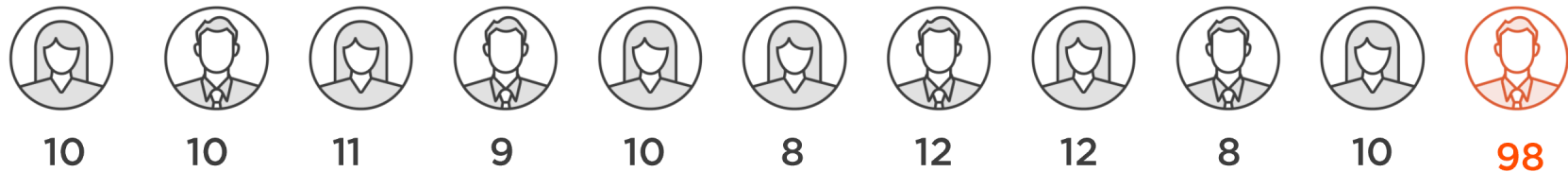


## Spread : Range



Age range :  $\text{max} - \text{min} = 12 - 8 = 4$

Problem : Affected by extreme values



Age range :  $\text{max} - \text{min} = 98 - 8 = 90$

# Percentiles

x percentile is y means x% of values are below y

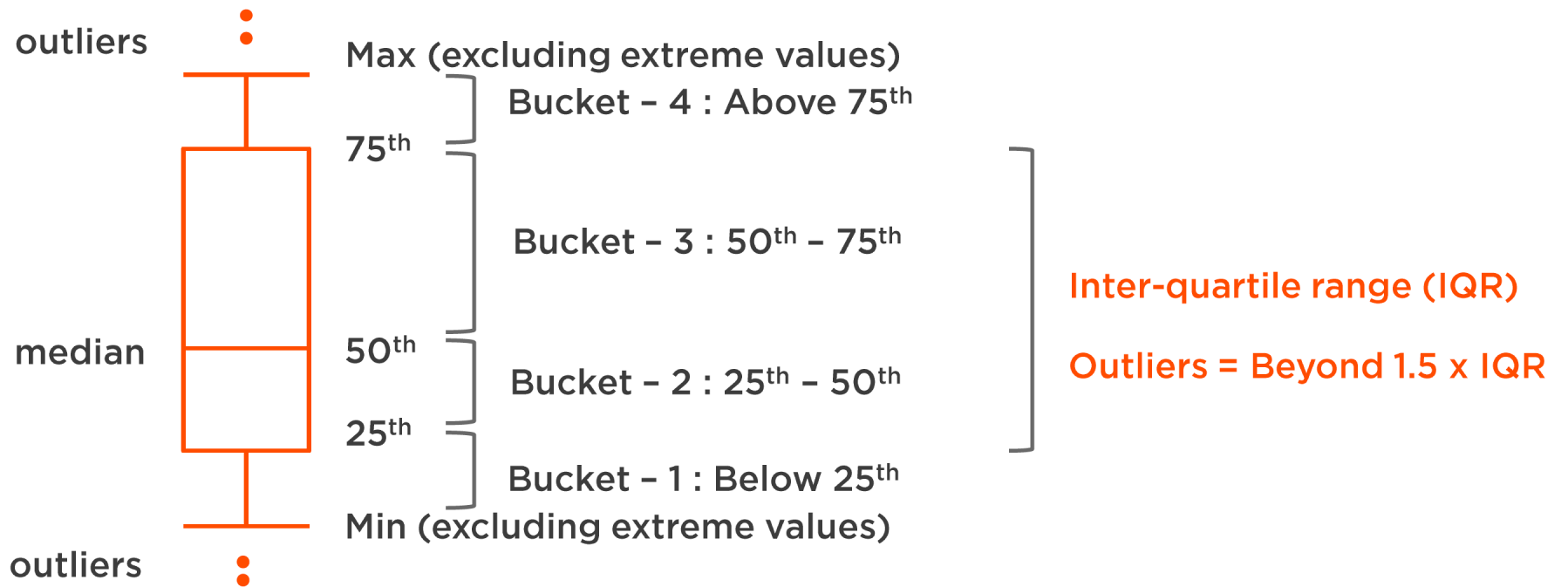
50 percentile is 10 means 50% of values are below 10

25<sup>th</sup>, 50<sup>th</sup>, 75<sup>th</sup>

- Bucket - 1 : Below 25<sup>th</sup>
- Bucket - 2 : 25<sup>th</sup> - 50<sup>th</sup>
- Bucket - 3 : 50<sup>th</sup> - 75<sup>th</sup>
- Bucket - 4 : above 75<sup>th</sup>

Quartiles

# Box-Whisker Plot



# Variance

Measure of variability

How far each value in list from mean value

Small variance = less spread

High variance = large spread

$$\text{Variance} = \frac{\text{sum}((\text{value} - \text{mean})^2)}{\text{count}}$$

Affected by extreme values

Unit is not clear

# Standard Deviation

**Standard deviation =  $\sqrt{\textit{variance}}$**

**Unit is same as that of the feature**

**Low standard deviation = less spread**

**High standard deviation = large spread**

Overview  
(Concepts)

**Exploratory data analysis**

- Distributions
- Grouping
- Crosstabs
- Pivots

# Distributions

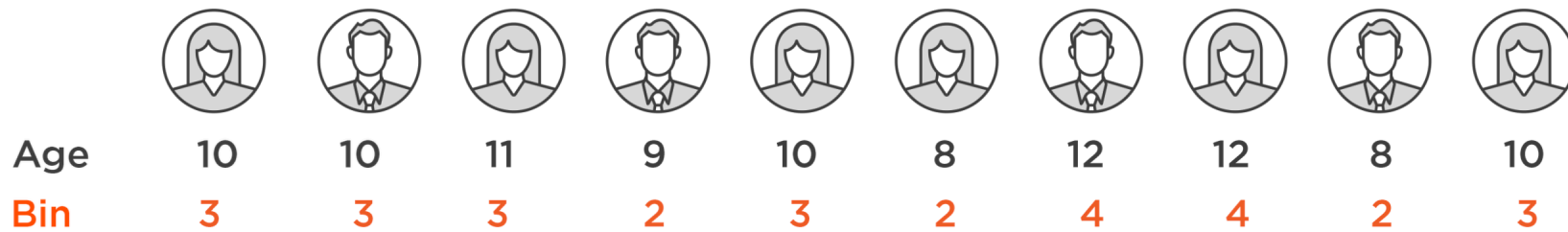
## Univariate

- Histogram
- Kernel Density Estimation (KDE) plot

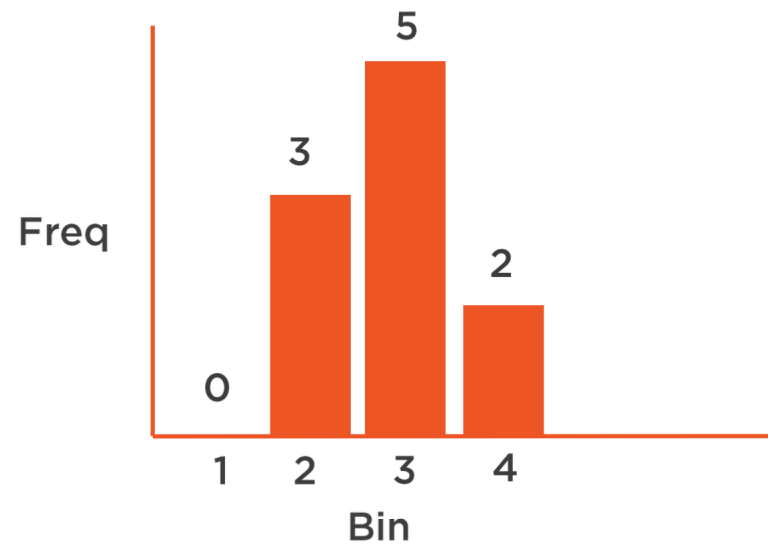
## Bivariate

- Scatter plot

# Histogram



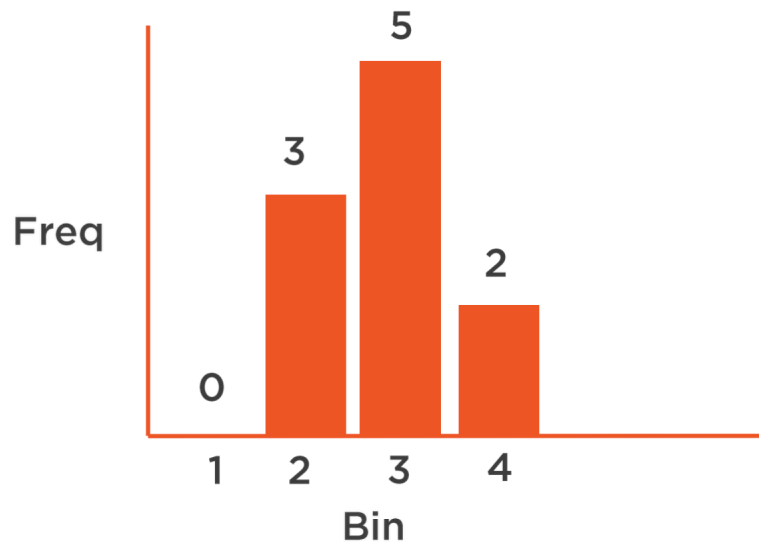
Bucket (bin) number	Bucket (bins)	Frequency
1	6-8	0
2	8-10	3
3	10-12	5
4	12-14	2





# Kernel Density Estimation (KDE) Plot

Histogram



KDE Plot

