

# Exploring Data With R

Abhishek Kumar  
ItsAbhishekKumar.com  
@MeAbhishekKumar



**pluralsight**  
hardcore dev and IT training



# Outline

Overall  
structure

Continuous  
data

Categorical  
data

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Types of Data

## Categorical data



Colors



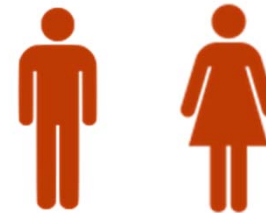
Gender

Use factor

## Continuous data



Mileage



Height, Weight, Age

Use numeric / integer

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Overall Structure

Number of  
observations

Number of  
features

Data  
types

Sample  
data

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Dataset

## Iris dataset

50 samples



**Iris-setosa**

50 samples



**Iris-virginica**

50 samples



**Iris-versicolor**

**Features :** sepal length, sepal width, petal length, petal width

**Available in **datasets** package**

Images: [http://en.wikipedia.org/wiki/Iris\\_flower\\_data\\_set](http://en.wikipedia.org/wiki/Iris_flower_data_set)

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Analysis of Continuous Data

Central  
tendency

Spread  
or dispersion

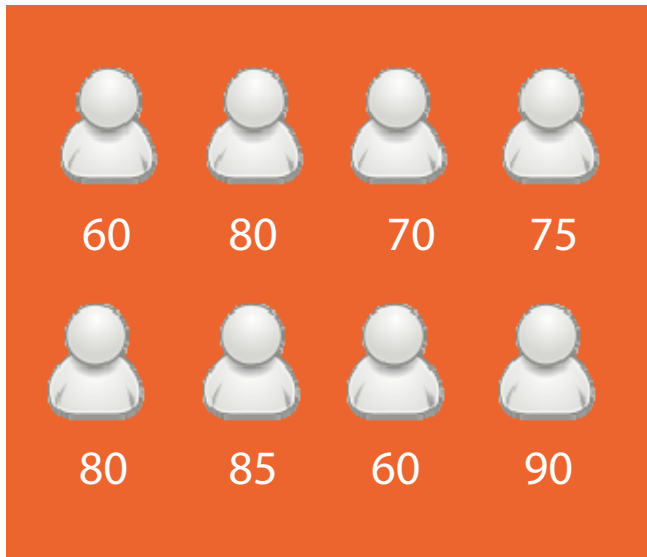
**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Central Tendency

## Mean (Average)



Set A

$$\text{Mean} = \frac{\text{Sum of all values}}{\text{Number of values}}$$

$$\text{Mean} = 75$$

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Central Tendency

## Median



Set A

{ 60, 60, 70, 75, 80, 80, 85, 90 }

Three orange arrows point to the 75, 80, and 80 values in the sorted set. Two arrows point down to 75 and 80, and one arrow points up to 80.

Median = **77.5**

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.



# Central Tendency

## Why Not Sufficient?



Set A

Mean = **75**

Median = **77.5**



Set B

Mean = **75**

Median = **77.5**

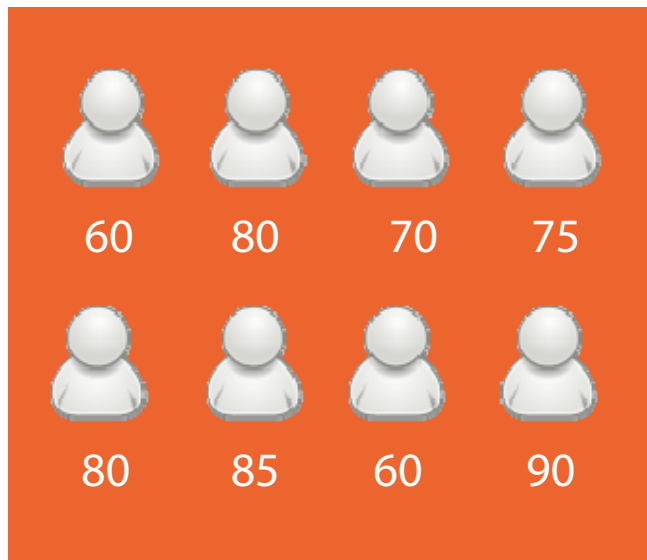
**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Spread

## Range



Set A

Range = **maximum - minimum**

Range = **90 - 60 = 30**

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Spread

## Range



Set A

Mean = **75**

Median = **77.5**

Range = **90 - 60 = 30**



Set B

Mean = **75**

Median = **77.5**

Range = **100 - 25 = 75**

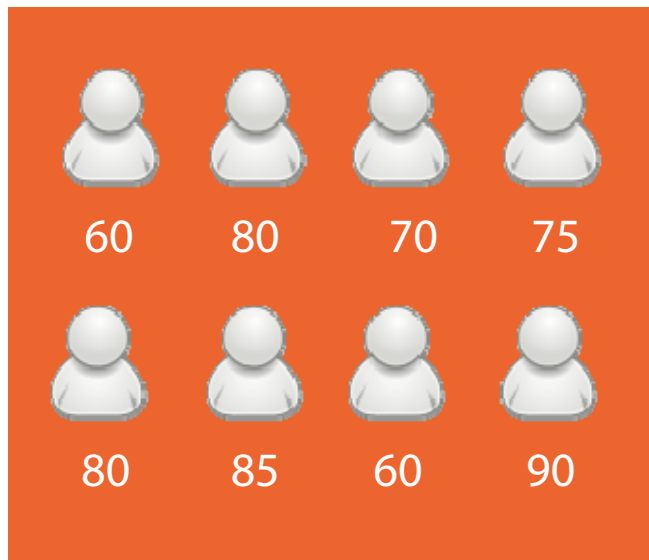
**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

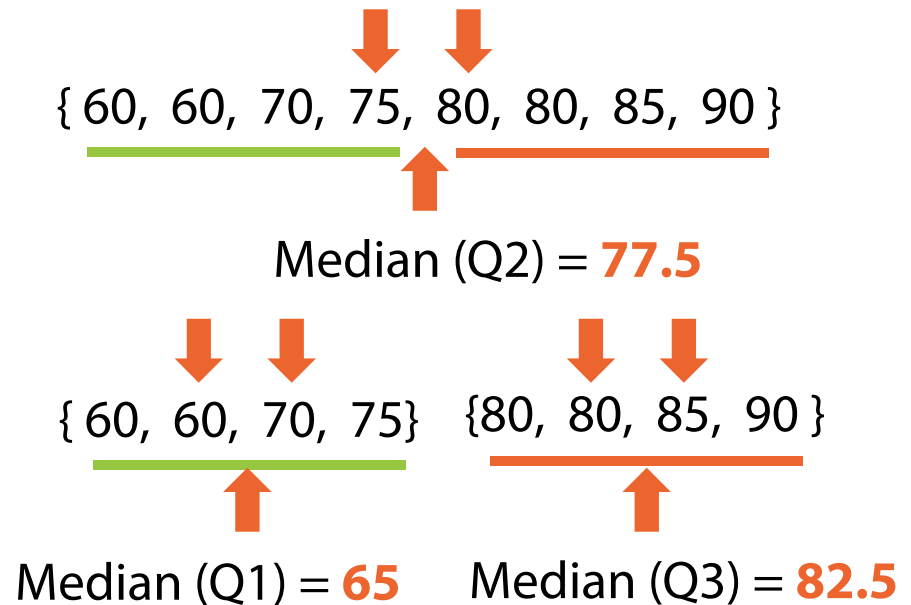
Note: Warning will not appear  
during Slide Show view.

# Spread

## Quartiles



Set A



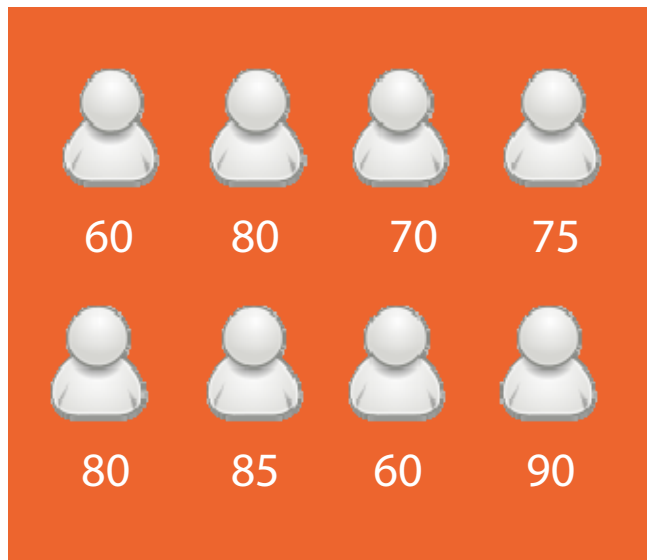
**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

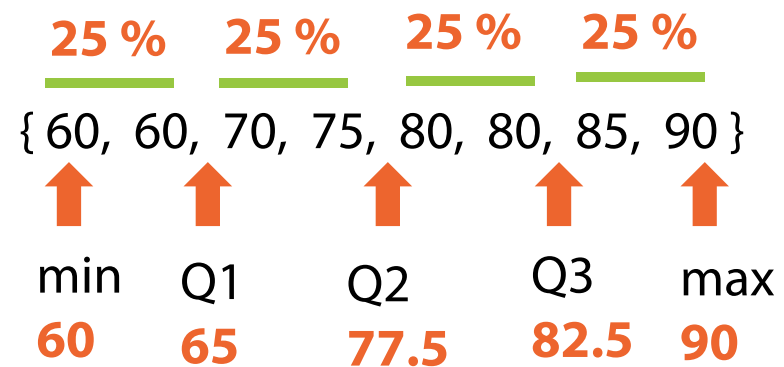
Note: Warning will not appear  
during Slide Show view.

# Spread

## Quartiles



Set A



Five point summary (**min, Q1, Q2, Q3, max**)

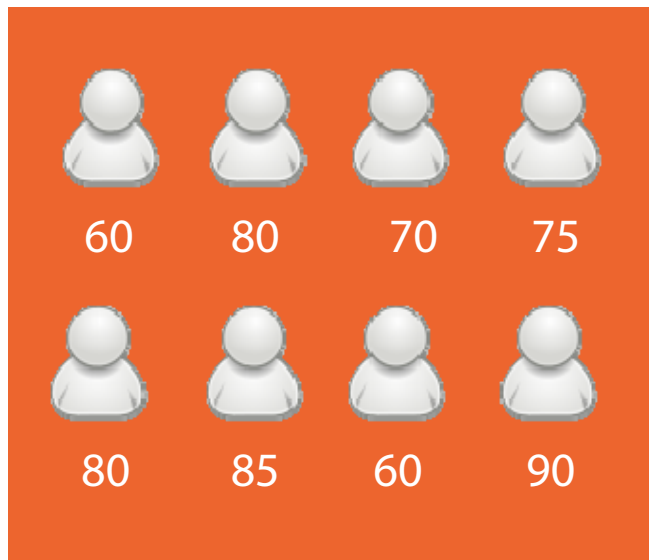
**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

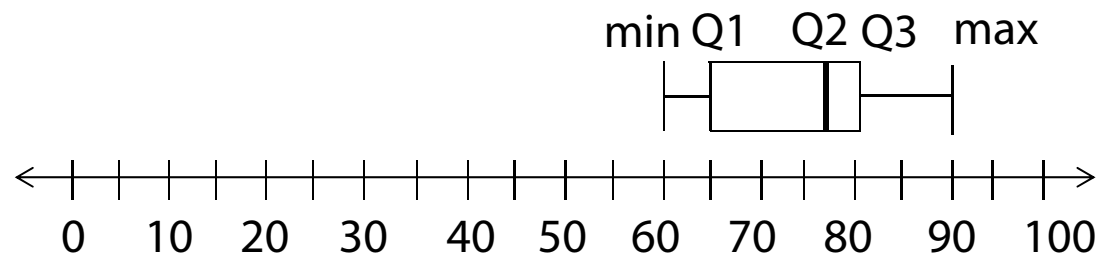
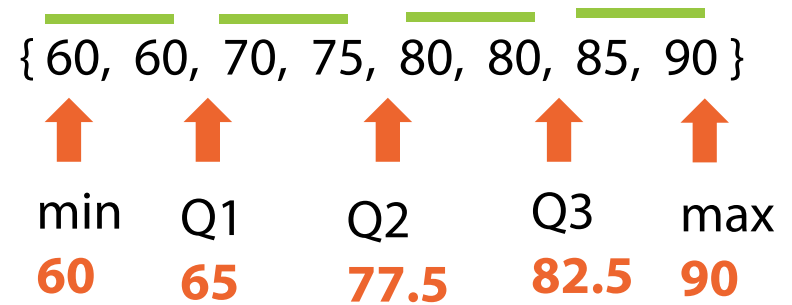
Note: Warning will not appear  
during Slide Show view.

# Spread

## Box Plot (Box – Whisker Plot)



Set A



**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Spread

## Box Plot

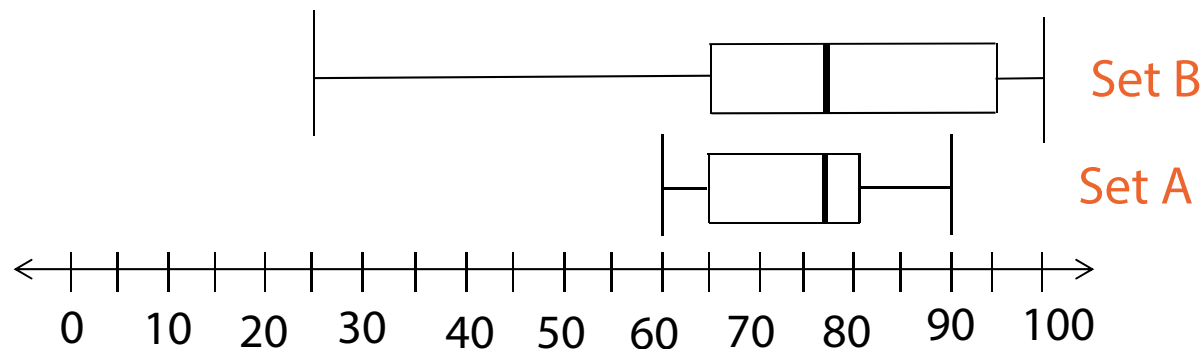


Set A

Set B

min	Q1	Q2	Q3	max
60	65	77.5	82.5	90

min	Q1	Q2	Q3	max
25	65	77.5	95	100

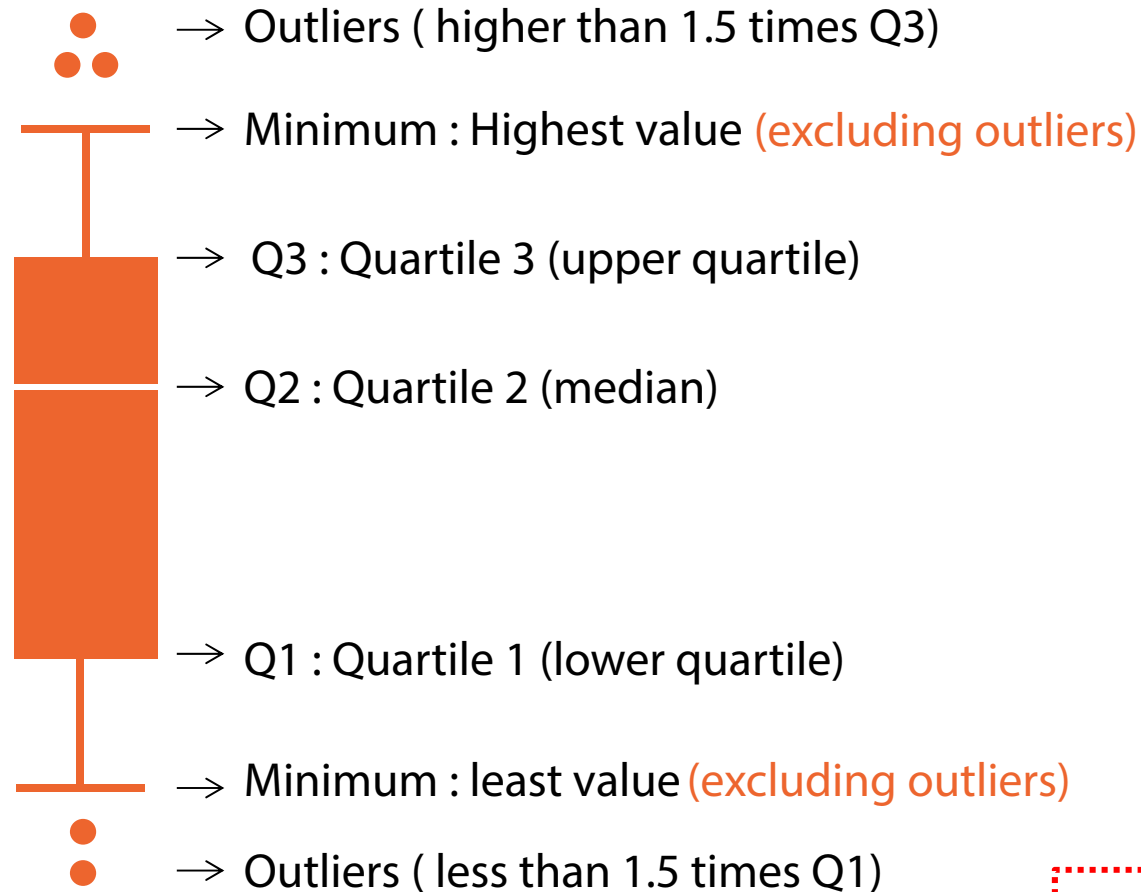


**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)  
Note: Warning will not appear  
during Slide Show view.

# Spread

## Box Plot



**Do Not Place Anything  
in This Space**

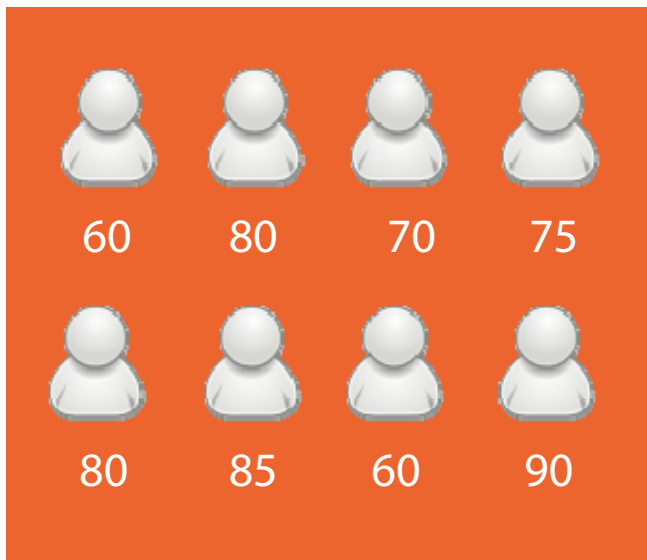
(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.



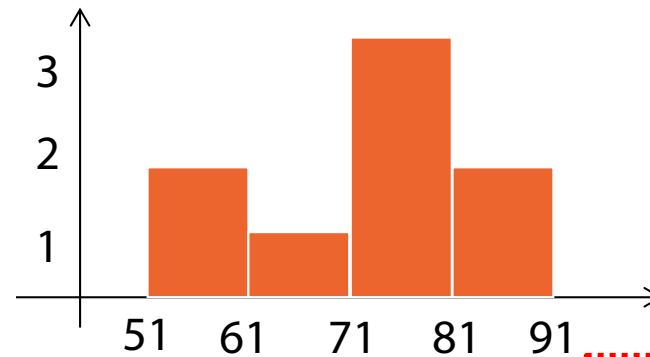
# Spread

## Histogram



Set A

Range	Count
51-61	2
61-71	1
71-81	3
81-91	2



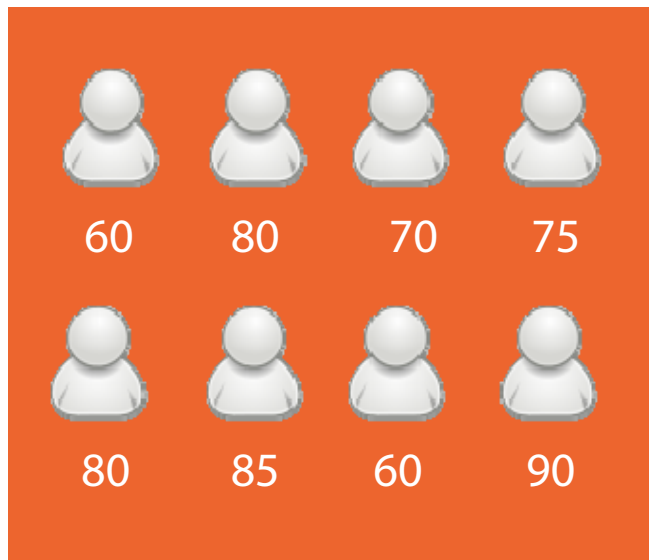
Histogram

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)  
Note: Warning will not appear  
during Slide Show view.

# Spread

## Variance & Standard deviation



Set A

{ 60, 80, 70, 75, 80, 85, 60, 90}      Mean = **75**

{ -15, 5, -5, 0, 5, 10, -15, 15 }

{ 225, 25, 25, 0, 25, 100, 225, 225 }

850

$850 / 8 = 106.25$

Variance = **106.25**

$\text{Sqrt}(106.25) = \sim 10.30$

Std. dev =  **$\sim 10.30$**

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Spread

## Variance & Standard Deviation



Mean = **75** Set A

Median = **77.5**

Std. deviation = **~10.3**

Variance = **106.25**



Mean = **75** Set B

Median = **77.5**

Std. deviation = **~20.9**

Variance = **437.5**

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Analysis of Categorical Data

Frequency  
distribution

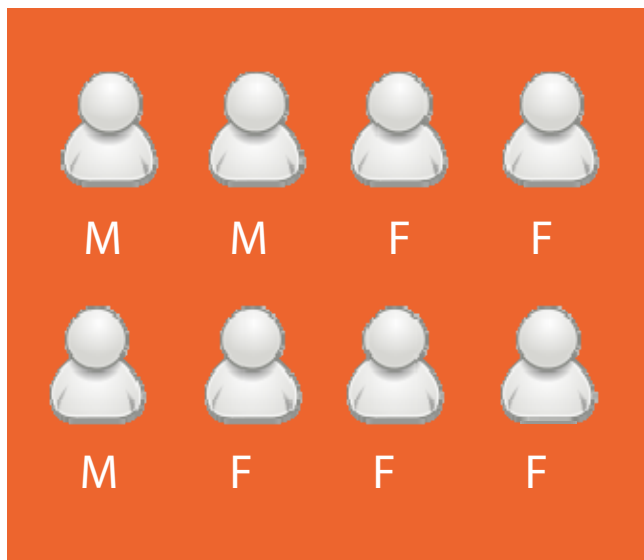
Category  
statistics

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

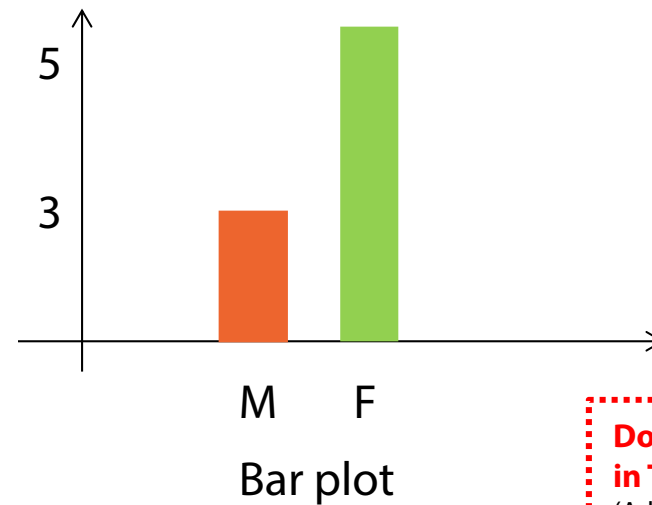
Note: Warning will not appear  
during Slide Show view.

# Frequency Distribution



Set A

Category	Count	Proportion
Male	3	$3/8 = 0.375$
Female	5	$5/8 = 0.625$



**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Category Statistics



Set A

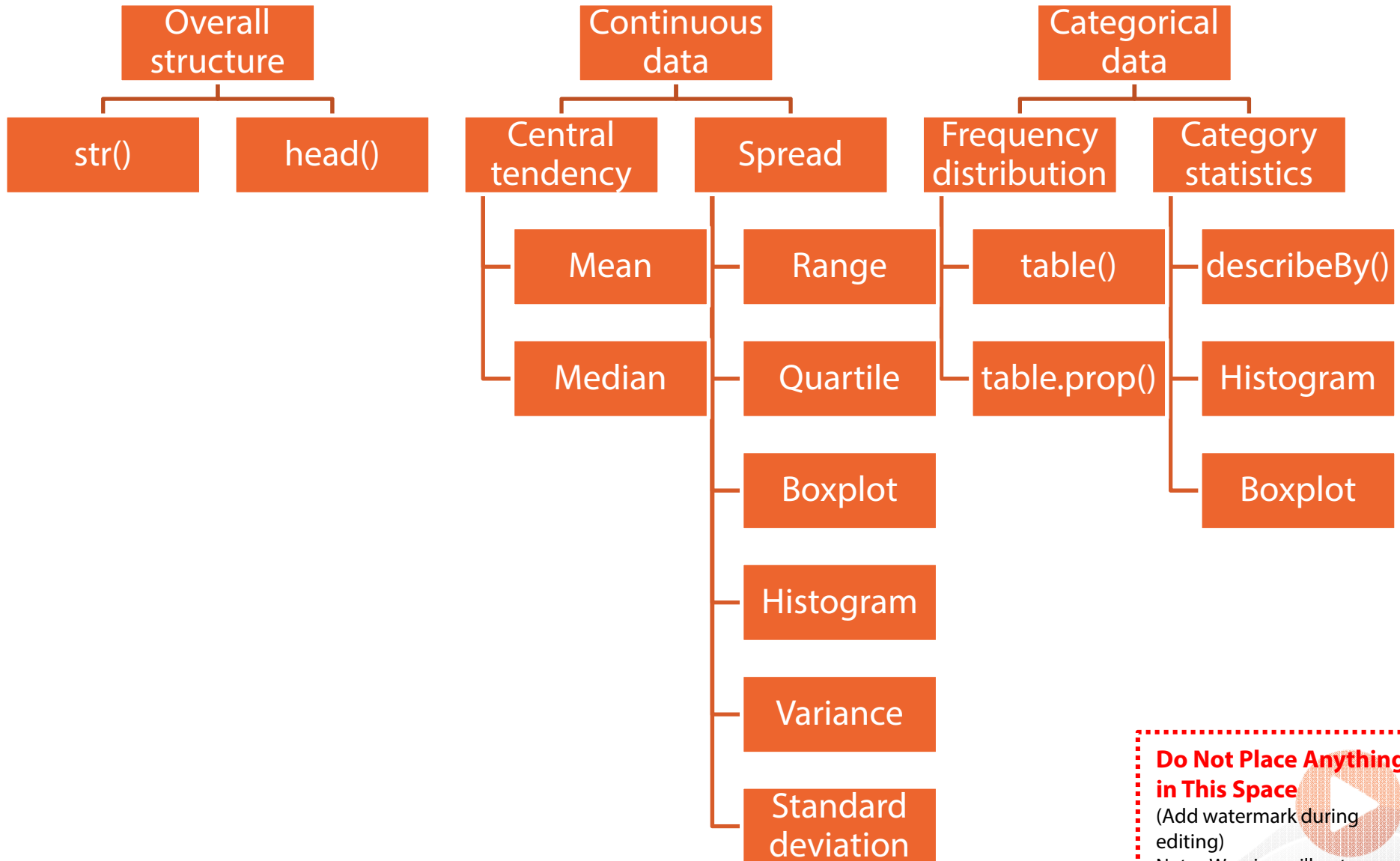
Category	Values	Mean
Male	{ 60, 80, 80 }	~ 73.3
Female	{ 70, 75, 85, 60, 90 }	76

**Do Not Place Anything  
in This Space**

(Add watermark during  
editing)

Note: Warning will not appear  
during Slide Show view.

# Summary



**Do Not Place Anything in This Space**

(Add watermark during editing)

Note: Warning will not appear during Slide Show view.