

# **Descriptive Statistics using R**

## **Table of contents**

<b>1 The Dataset: Student Transcript</b>	<b>2</b>
<b>2 Objective</b>	<b>2</b>
<b>3 Dataset Overview</b>	<b>2</b>
<b>4 Loading the Dataset</b>	<b>3</b>
<b>5 Measures of Central Tendency - Mean &amp; Median</b>	<b>3</b>
<b>6 Measures of Central Tendency - Mode</b>	<b>4</b>
<b>7 Descriptive Statistics using DescTools</b>	<b>4</b>
<b>8 Measures of Spread - Range</b>	<b>5</b>
<b>9 Measures of Spread - IQR</b>	<b>6</b>
<b>10 Measures of Spread - Standard Deviation &amp; Variance</b>	<b>6</b>
<b>11 Measures of Spread - Mean Absolute Deviation (MAD)</b>	<b>7</b>
<b>12 Exericse - Major GPA vs. non-Major GPA</b>	<b>7</b>
<b>13 Solution - Major GPA vs. non-Major GPA</b>	<b>8</b>

## 1 The Dataset: Student Transcript



[Student Transcript](#)

## 2 Objective

Analyzing a student's transcript dataset to understand performance metrics.

## 3 Dataset Overview

- Columns: year, semester, course\_number, credits, letter\_grade, numerical\_value (GPA)
- Four academic years of data
- Grades for both major and non-major courses

## 4 Loading the Dataset

```
1 df = read.csv ("https://raw.githubusercontent.com/ahmedmoustafa/datasets/main/transcript/t
2 head(df)
```

year	semester	course_number	credits	letter_grade	numerical_value
Freshman	Fall	COMP100	3	A-	3.7
Freshman	Fall	COMP110	3	A	4.0
Freshman	Fall	HUMA181	3	A-	3.7
Freshman	Fall	SOCI181	3	A-	3.7
Freshman	Fall	ELEC181	3	B+	3.3
Freshman	Spring	COMP120	3	A	4.0

## 5 Measures of Central Tendency - Mean & Median

- **Mean:** The average of a set of numbers,  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

```
1 mean(df$numerical_value)
```

```
[1] 3.78
```

- **Weighted Mean:** The mean where some values contribute more than others,  $\bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$

```
1 sum(df$numerical_value * df$credits)/sum(df$credits)
```

```
[1] 3.790244
```

```
1 weighted.mean(df$numerical_value, w = df$credits)
```

```
[1] 3.790244
```

- **Trimmed Mean:** The mean after removing a specified number of the highest and lowest values.

```
1 mean(df$numerical_value, trim = 0.1)
```

```
[1] 3.8125
```

- **Median:** The middle value in a sorted list of numbers

```
1 median(df$numerical_value)
```

```
[1] 3.85
```

## 6 Measures of Central Tendency - Mode

- **Mode:** The value that appears most frequently in a set.
- As discussed before, mode is more appropriate for qualitative data values.
- So, let's compute mode for the `letter_grade`
- However, in R, there is no built-in function to compute the mode directly.
- Therefore, we need to install the `DescTools` package

```
1 if(!require(DescTools))
2   install.packages("DescTools",repos = "http://cran.us.r-project.org")
```

```
Loading required package: DescTools
```

- Now we can run the `Mode()` function from `DescTools` package

```
1 library(DescTools)
2 Mode(df$letter_grade)
```

```
[1] "A"
attr(,"freq")
[1] 20
```

## 7 Descriptive Statistics using DescTools

Measure	Function	Description
<b>Mode</b>	<code>Mode(data)</code>	Computes the mode. Returns multiple modes if they exist.
<b>Mean</b>	<code>Mean(data)</code>	Computes the arithmetic mean.
<b>Weighted Mean</b>	<code>WtdMean(data)</code>	Computes the weighted mean.
<b>Median</b>	<code>Median(data)</code>	Computes the median.
<b>Trimmed Mean</b>	<code>Mean(data, trim)</code>	Computes trimmed mean. <code>trim</code> is fraction (0 to 0.5) of observations to be trimmed.
<b>Standard Deviation</b>	<code>Std(data)</code>	Computes the sample standard deviation.
<b>Variance</b>	<code>Var(data)</code>	Computes the variance.
<b>Range</b>	<code>Range(data)</code>	Computes the range (difference between max and min).
<b>Interquartile Range</b>	<code>IQR(data)</code>	Computes the interquartile range.

## 8 Measures of Spread - Range

- **Range:** Difference between the largest and smallest values,

$$\text{Range} = x_{\max} - x_{\min}$$

```
1 max(df$numerical_value) - min(df$numerical_value)
```

```
[1] 0.7
```

```
1 range(df$numerical_value) # The base R (built-in)
```

```
[1] 3.3 4.0
```

```
1 Range(df$numerical_value) # From DescTools
```

```
[1] 0.7
attr(,"bounds")
[1] 3.3 4.0
```

## 9 Measures of Spread - IQR

- **Interquartile Range (IQR):** Difference between the first and third quartiles,

$$\text{IQR} = Q_3 - Q_1$$

– First Quartile:

- \* Also known as the lower quartile or the 25th percentile.
- \* It is the value below which 25% of the data falls. In other words, it cuts off the lowest 25% of the data.

– Third Quartile:

- \* Also known as the upper quartile or the 75th percentile.
- \* It is the value below which 75% of the data falls, meaning it cuts off the lowest 75% of data points.

```
1 quantile(df$numerical_value)
```

```
0% 25% 50% 75% 100%
3.30 3.70 3.85 4.00 4.00
```

```
1 quantiles = quantile(df$numerical_value)
2 quantiles[4] - quantiles[2]
```

```
75%
0.3
```

```
1 IQR(df$numerical_value)
```

```
[1] 0.3
```

## 10 Measures of Spread - Standard Deviation & Variance

- **Standard Deviation:** Measures the amount of variation or dispersion of a set of values,

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

```
1 sd(df$numerical_value)
```

```
[1] 0.2613574
```

- **Variance:**

$$var(x) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} = s^2$$

```
1 s = sd(df$numerical_value)
2 s^2
```

```
[1] 0.06830769
```

```
1 Var(df$numerical_value) # From DescTools
```

```
[1] 0.06830769
```

## 11 Measures of Spread - Mean Absolute Deviation (MAD)

- **Mean Absolute Deviation:** (MAD) a measure of dispersion representing the average distance of each data point from the mean

$$MAD = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

```
1 mean(abs(df$numerical_value - mean(df$numerical_value)))
```

```
[1] 0.22
```

- MAD is sensitive to outliers.

## 12 Exercise - Major GPA vs. non-Major GPA

- Using the provided dataset, compare the GPA (the `numerical_value` column) of the student in their major courses versus the non-major courses. For this dataset, Computer Science courses are the major courses, and their course numbers start with "COMP".
- Hint: You might find the `startsWith()` function in R useful to filter rows based on the course number.

## 13 Solution - Major GPA vs. non-Major GPA

- Major Courses

We can search for the rows with major courses using `startsWith()`

```
1 flag = startsWith(df$course_number, "COMP")
2 flag
```

```
[1] TRUE TRUE FALSE FALSE FALSE TRUE TRUE FALSE FALSE FALSE TRUE TRUE
[13] FALSE FALSE FALSE TRUE TRUE FALSE FALSE FALSE TRUE TRUE FALSE FALSE
[25] FALSE TRUE TRUE FALSE FALSE FALSE TRUE TRUE FALSE FALSE FALSE TRUE
[37] TRUE FALSE FALSE FALSE
```

- `TRUE` : a major course
- `FALSE` : a non-major course

```
1 major_courses = df[flag, ]
2 head(major_courses)
```

	year	semester	course_number	credits	letter_grade	numerical_value
1	Freshman	Fall	COMP100	3	A-	3.7
2	Freshman	Fall	COMP110	3	A	4.0
6	Freshman	Spring	COMP120	3	A	4.0
7	Freshman	Spring	COMP130	3	A-	3.7
11	Sophomore	Fall	COMP200	4	A	4.0
12	Sophomore	Fall	COMP210	4	A	4.0

```
1 major_gpa = median(major_courses$numerical_value)
2 major_gpa
```

```
[1] 4
```

- Non-Major Courses

To filter for non-major courses, we can just *negate* `flag` i.e.,

```
1 !flag
```

```
[1] FALSE FALSE TRUE TRUE TRUE FALSE FALSE TRUE TRUE TRUE TRUE FALSE FALSE  
[13] TRUE TRUE TRUE FALSE FALSE TRUE TRUE TRUE FALSE FALSE TRUE TRUE  
[25] TRUE FALSE FALSE TRUE TRUE TRUE FALSE FALSE TRUE TRUE TRUE FALSE  
[37] FALSE TRUE TRUE TRUE
```

- TRUE : a non-major course
- FALSE : a major course

```
1 nonmajor_courses = df[!flag, ]  
2 head(nonmajor_courses)
```

	year	semester	course_number	credits	letter_grade	numerical_value
3	Freshman	Fall	HUMA181	3	A-	3.7
4	Freshman	Fall	SOCI181	3	A-	3.7
5	Freshman	Fall	ELEC181	3	B+	3.3
8	Freshman	Spring	HUMA191	4	B+	3.3
9	Freshman	Spring	SOCI191	2	A-	3.7
10	Freshman	Spring	ELEC191	3	A-	3.7

```
1 nonmajor_gpa = median(nonmajor_courses$numerical_value)  
2 nonmajor_gpa
```

```
[1] 3.7
```

- Using the `summary()` function

```
1 summary(major_courses$numerical_value)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
3.3	3.7	4	3.8375	4	4

```
1 summary(nonmajor_courses$numerical_value)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
3.3	3.7	3.7	3.741667	4	4