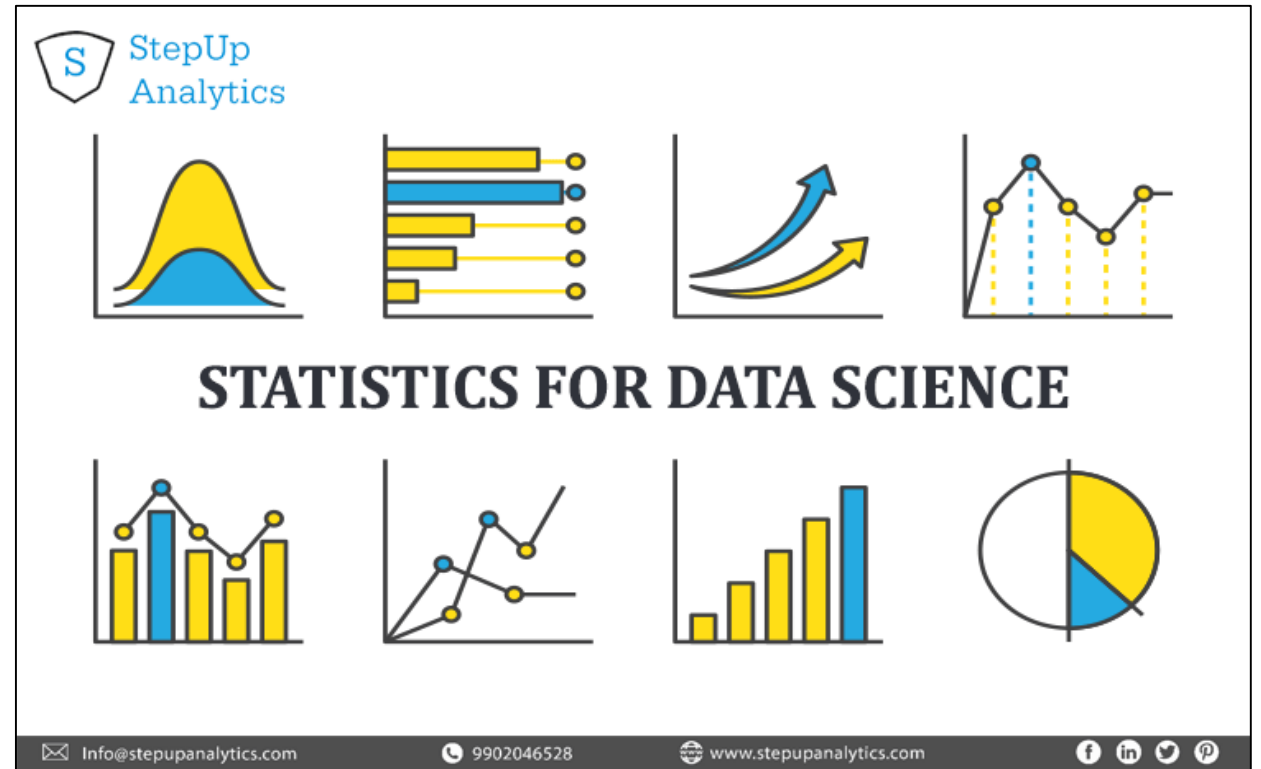# DS501 STATISTICAL AND MATHEMATICAL METHODS FOR DATA SCIENCE

**Dr. Muhammad Wasim**

PhD, MS, M.Phil, M.Sc, MCS

**Certified Data Analyst [KARACHI.AI]**

**Lecture Week 05**

➤ Statistics for Data Science
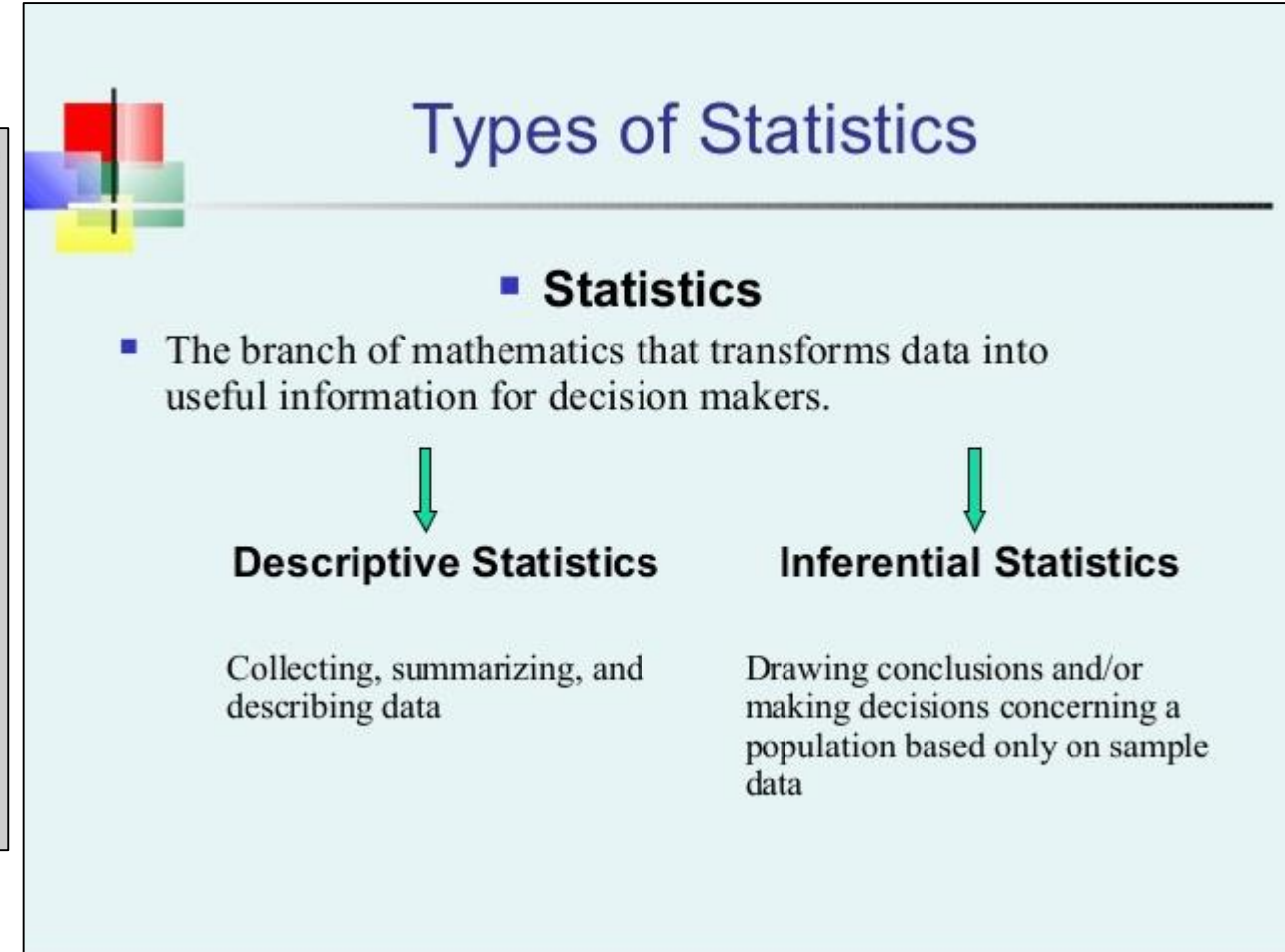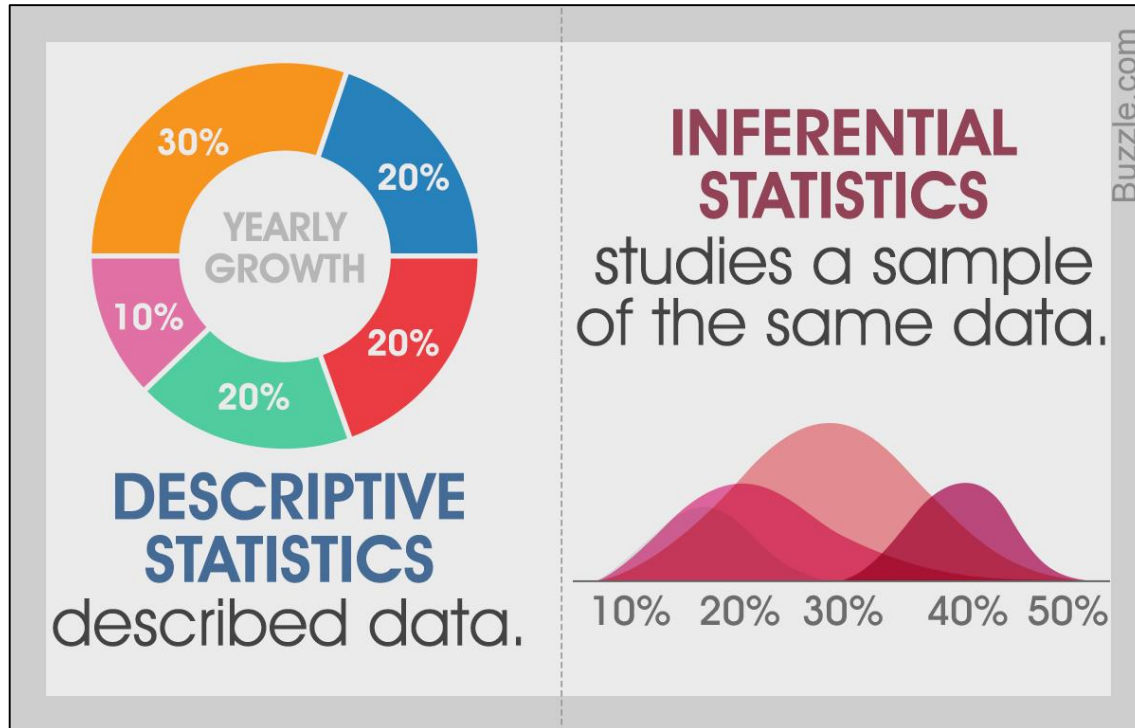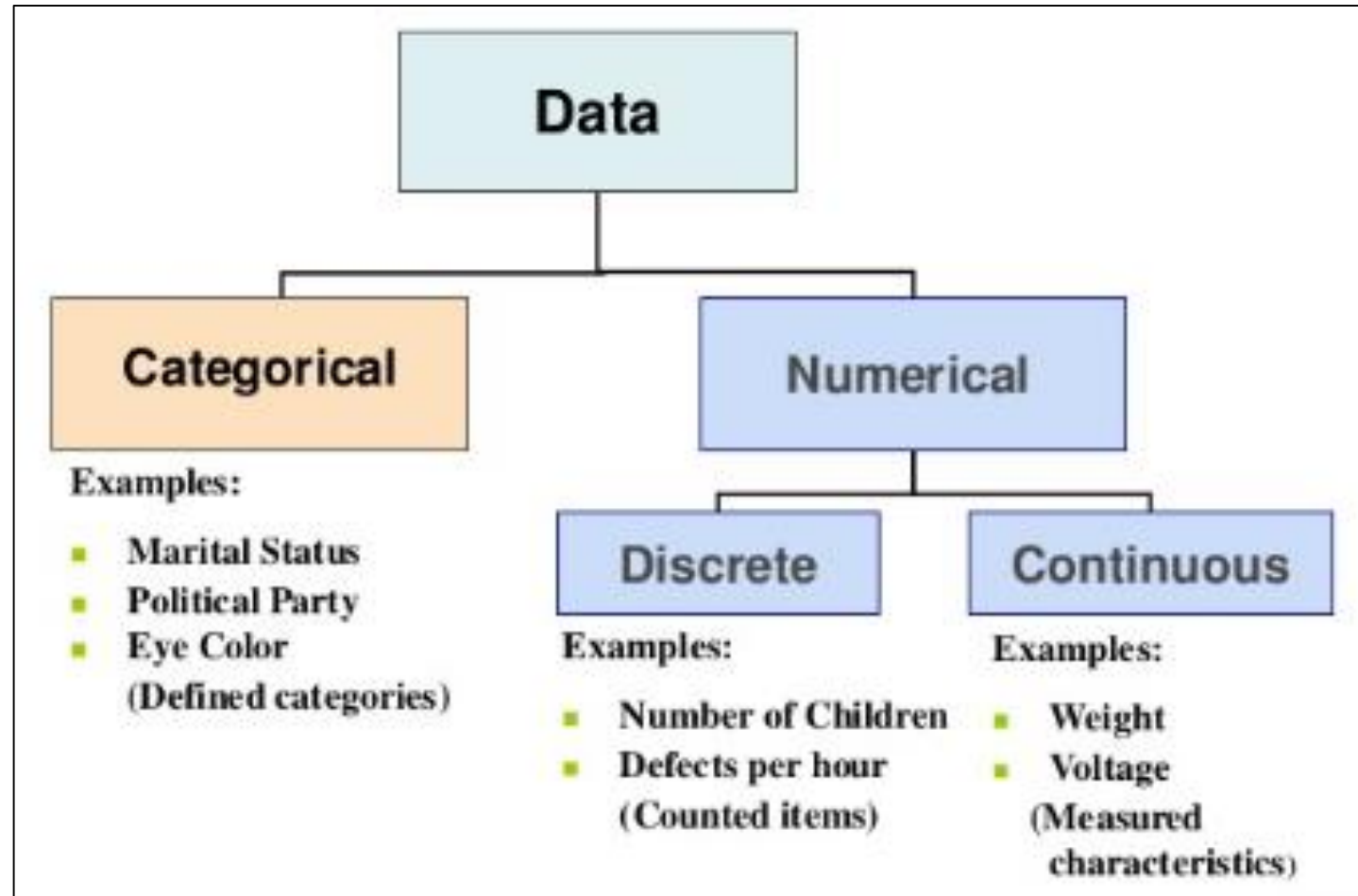
# STATISTICS FOR DATA SCIENCE

# STATISTICS FOR DATA SCIENCE

# STATISTICS FOR DATA SCIENCE

# STATISTICS FOR DATA SCIENCE

**Measures Of The Spread:**

Just like the measure of center, we also have measures of the spread, which comprises of the following measures:

**Range:** It is the given measure of how spread apart the values in a data set are.

**Inter Quartile Range (IQR):** It is the measure of variability, based on dividing a data set into quartiles.

**Variance:** It describes how much a random variable differs from its expected value. It entails computing squares of deviations.

- ***Deviation*** *is the difference between each element from the mean.*
- ***Population Variance*** *is the average of squared deviations*
- ***Sample Variance*** *is the average of squared differences from the mean*

**Standard Deviation:** It is the measure of the dispersion of a set of data from its mean.

# STATISTICS FOR DATA SCIENCE

1. Which class (A or B) is more consistent. How to study the variables by plotting a histogram

2. Write a python code for the solution with data visualizations.

Resource Link:

https://www.datacamp.com/tracks/statistics-fundamentals-with-python

| A | B |
|---|---|
| 56 | 63 |
| 82 | 89 |
| 68 | 65 |
| 67 | 75 |
| 59 | 21 |
| 95 | 80 |
| 2 | 82 |
| 64 | 71 |
| 42 | 14 |
| 93 | 73 |
| 56 | 15 |
| 79 | 21 |
| 48 | 10 |
| 98 | 59 |
| 48 | 24 |
| 33 | 39 |
| 53 | 93 |
| 44 | 11 |
| 60 | 96 |
| 77 | 35 |

# STATISTICS FOR DATA SCIENCE

Step 1: Import data for computation

```
1    >set.seed(1)
2    #Generate random numbers and store it in a variable called data
3    >data = runif(20,1,10)
```

Step 2: Calculate Mean for the data

```
1    #Calculate Mean
2    >mean = mean(data)
3    >print(mean)
4
5    [1] 5.996504
```

Step 3: Calculate the Median for the data

```
1    #Calculate Median
2    >median = median(data)
3    >print(median)
4
5    [1] 6.408853
```

# STATISTICS FOR DATA SCIENCE

Step 4: Calculate Mode for the data

```
1    #Create a function for calculating Mode
2    >mode <- function(x) { >ux <- unique(x) >ux[which.max(tabulate(match(x, ux)))]
3    }
4    >result <- mode(data) >print(data)
5
6    [1] 3.389578 4.349115 6.155680 9.173870 2.815137 9.085507 9.502077 6.947180 6.66
7    [10] 1.556076 2.853771 2.589011 7.183206 4.456933 7.928573 5.479293 7.458567 9.9
8    [19] 4.420317 7.997007
9
L0   >cat("mode= {}", result)
L1
L2   mode= {} 3.389578
```

Step 5: Calculate Variance & Std Deviation for the data

```
1    #Calculate Variance and std Deviation
2    >variance = var(data)
3    >standardDeviation = sqrt(var(data))
4    >print(standardDeviation)
5
6    [1] 2.575061
```

# STATISTICS FOR DATA SCIENCE

Step 6: Plot a Histogram

```
1    #Plot Histogram
2    >hist(data, bins=10, range= c(0,10), edgecolor='black')
```

The Histogram is used to display the frequency of data points:

**Histogram of data**