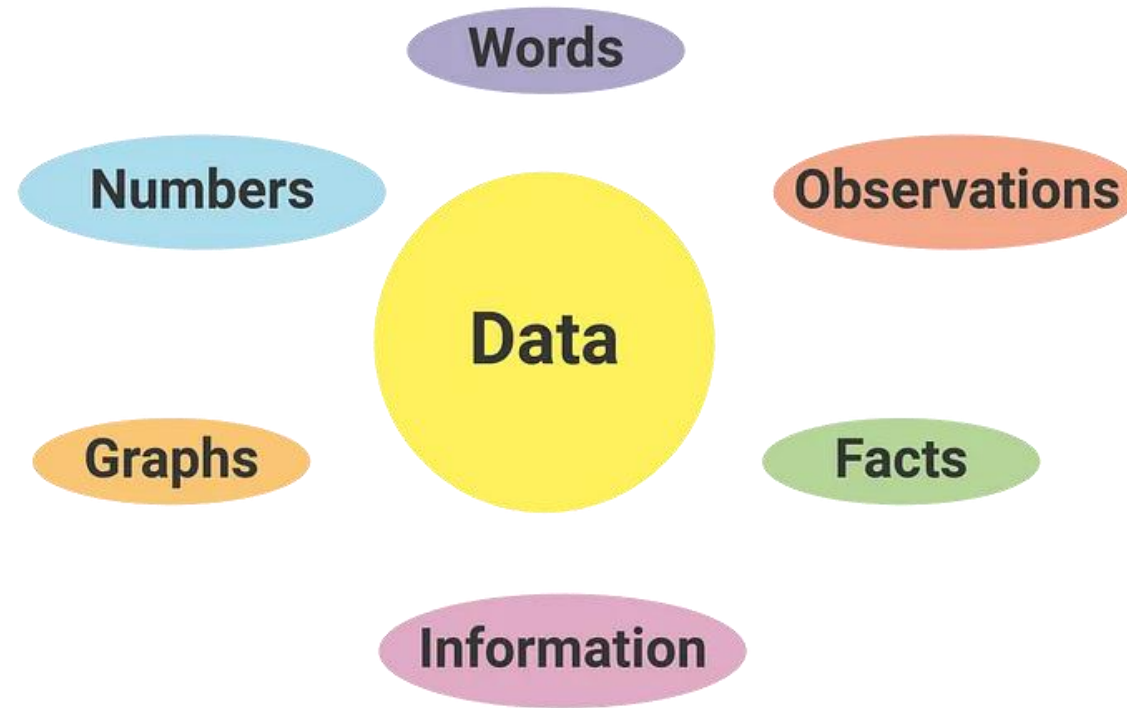


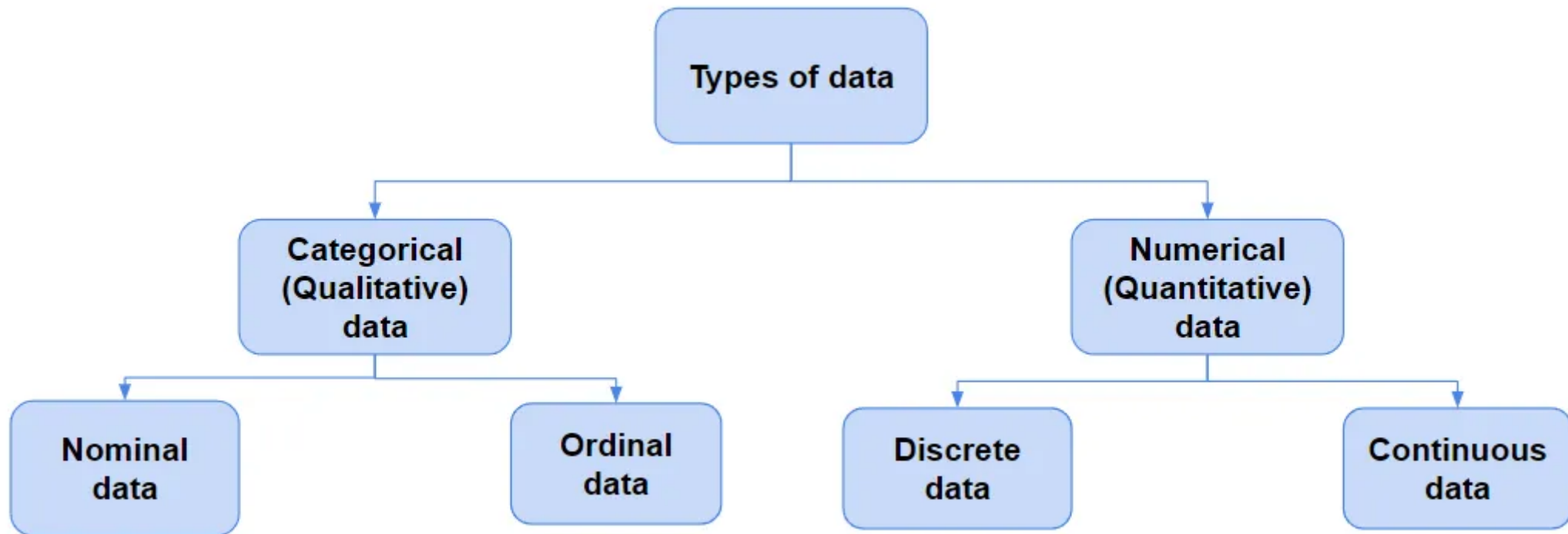
Balachandar K

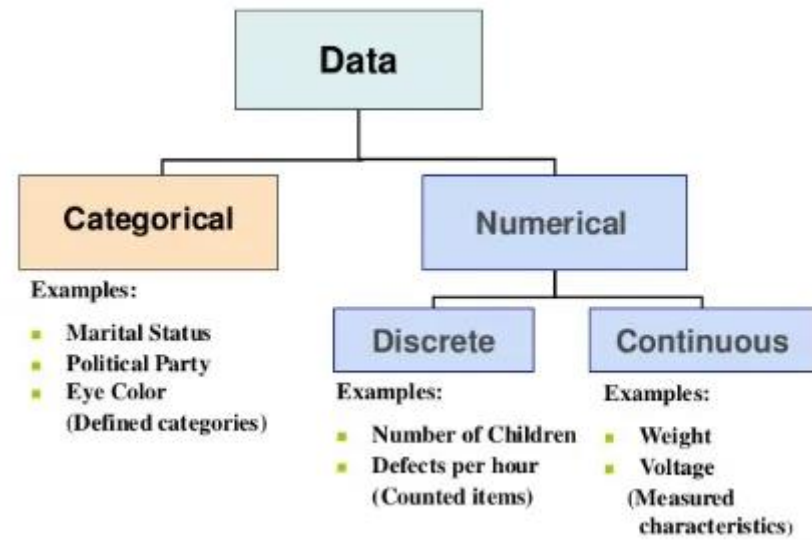
# Important one....

- *Data Types*
- *Probability & Bayes' Theorem*
- *Measures of Central Tendency*
- *Skewness*
- *Kurtosis*
- *Measures of Dispersion*
- *Covariance*
- *Correlation*
- *Probability Distributions*
- *Hypothesis Testing*
- *Regression*

# What is data ?







# Data Types

## Qualitative data

**Nominal**

Hair Color  
Gender  
Name  
Nationalities

**Ordinal**

Groups  
Economic St.  
Ratings

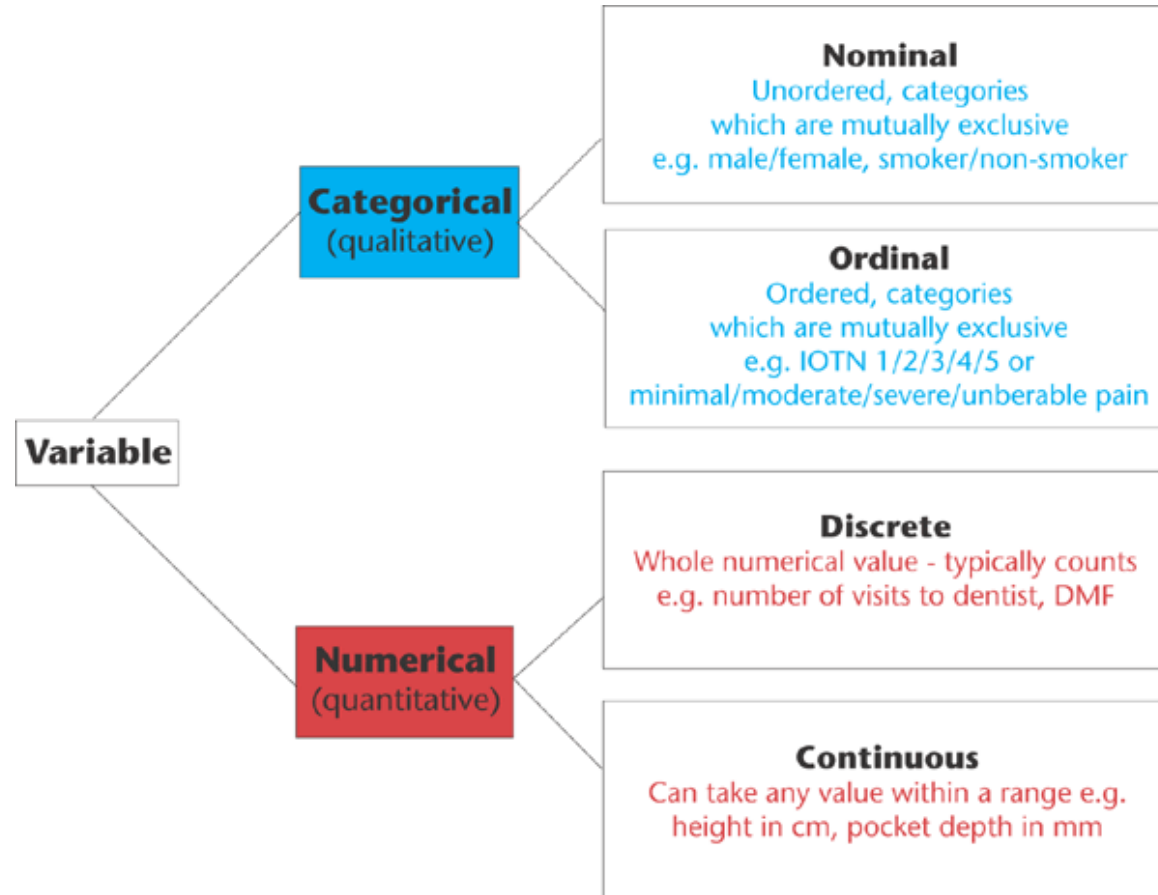
## Quantitative data

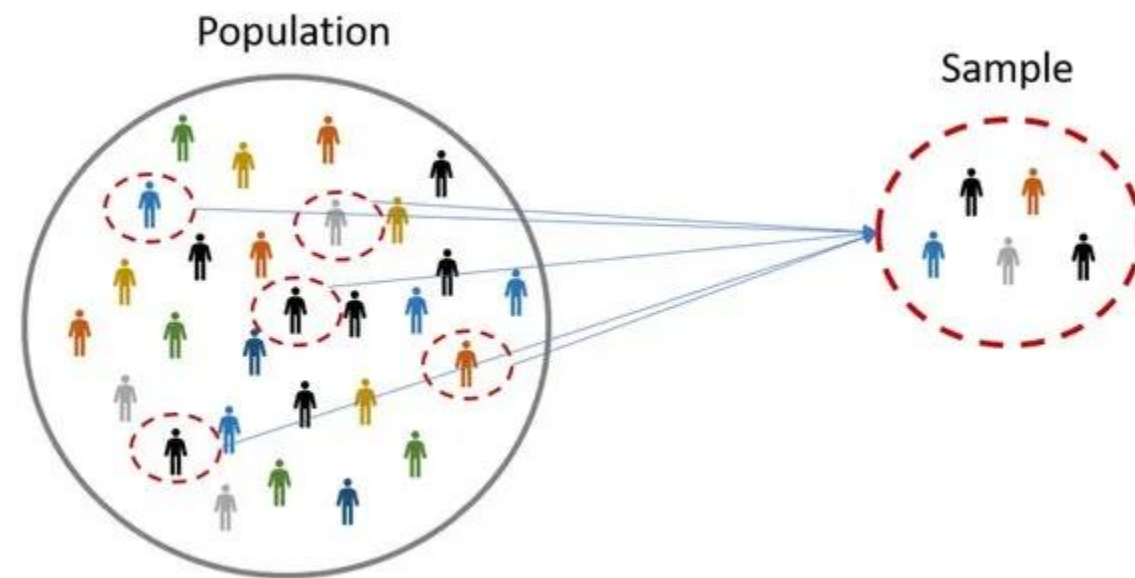
**Discreet**

Shoe Size  
No. of emp  
No. of children  
(Counted Items)

**Continuous**

Weight  
Height  
Temperature







# Measure of central tendency

## Mean

Where you add up all the numbers and then divide by the amount of numbers you added.

Example :

13, 18, 13, 14, 13, 16, 14, 21, 13.

$$(13 + 18 + 13 + 14 + 13 + 16 + 14 + 21 + 13) \div 9 = 15$$

## Median

It is the "middle" value in the list of numbers.

Example :

13, 18, 13, 14, 13, 16, 14, 21, 13

Arrange the number

13, 13, 13, 13, 14, 14, 16, 18, 21

There are nine numbers in the list.

$$(9 + 1) \div 2 = 10 \div 2 = 5\text{th number.}$$

13, 13, 13, 13, 14, 14, 16, 18, 21

The median is 14.

## Mode

The number or numbers that occur most often in a set of numbers.

Example :

13, 18, 13, 14, 13, 16, 14, 21, 13

The number that is repeated more often is 13.

## Range

The difference between the highest and the lowest numbers in a set of numbers.

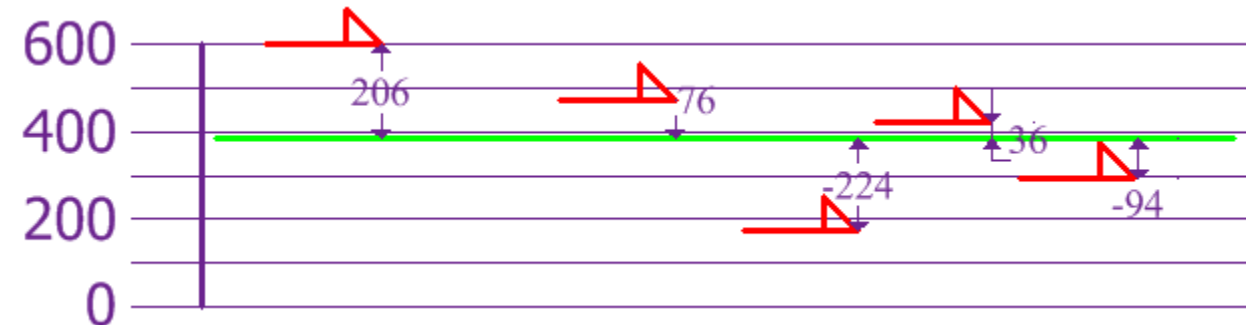
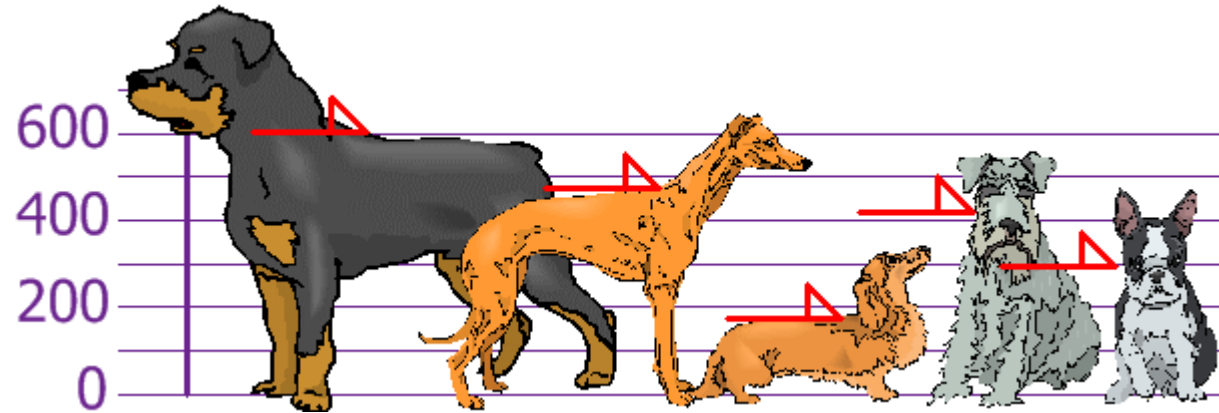
Example :

13, 18, 13, 14, 13, 16, 14, 21, 13

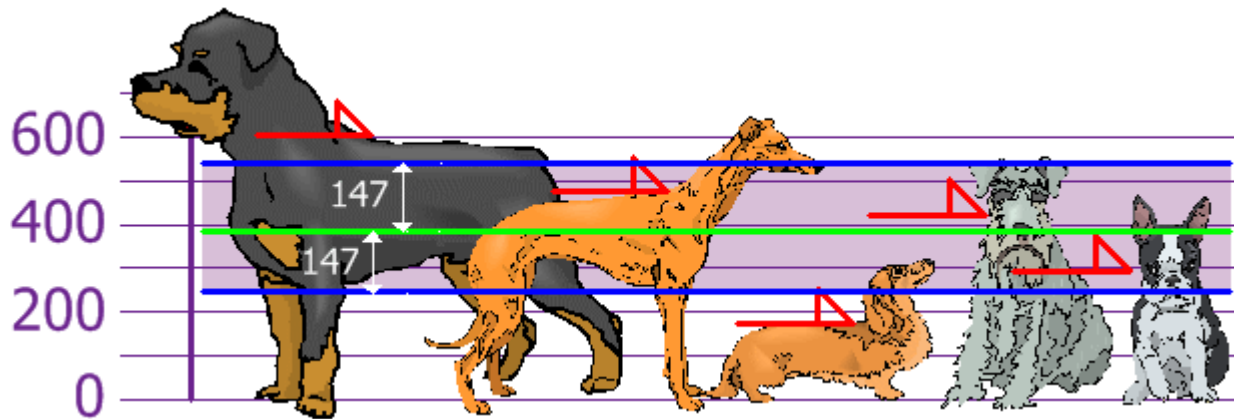
The largest value in the list is 21, and the smallest is 13,

$$\text{so the range is } 21 - 13 = 8.$$

# Height & Average Height (mean)



# Variance & SD



$$\text{Variance } \sigma^2 = (206^2 + 76^2 + (-224)^2 + 36^2 + (-94)^2) / 5$$

$$\Rightarrow (42436 + 5776 + 50176 + 1296 + 8836) / 5 \Rightarrow 108520 / 5 \Rightarrow 21704$$

So the Variance is 21,704

$$\text{Standard Deviation } \sigma = \sqrt{21704} = 147.32... = 147$$

# Statistics methods....

```
statistics.mean([1,2,3,4,4])
```

2.8

```
statistics.median([1,3,5])
```

3

```
statistics.median([1,3,5,7])
```

4.0

```
statistics.mode([1,1,2,3,3,3,3,4])
```

3

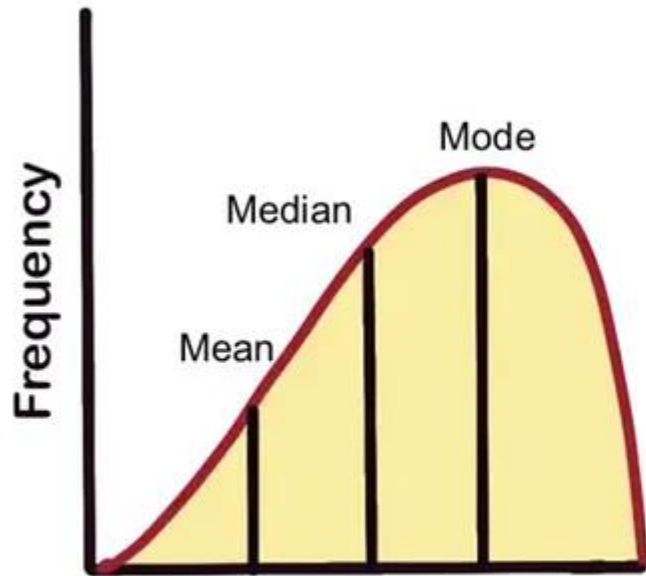
```
statistics.mode(["red","blue","red","red"])
```

'red'

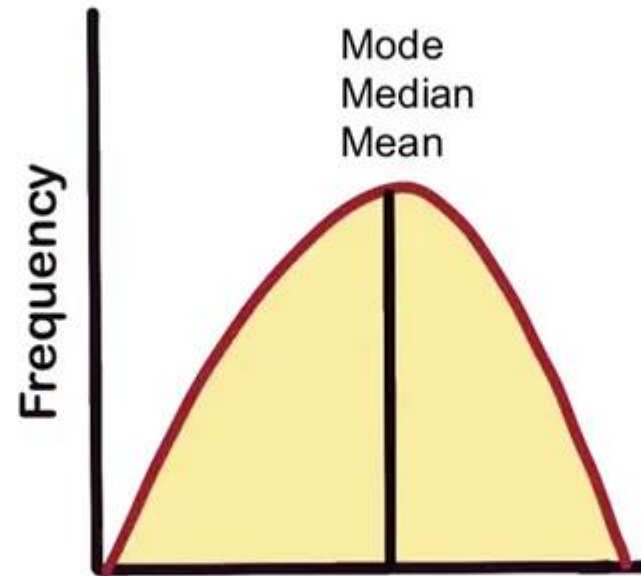
# Measures of central tendency

Type of Variable	Best measure of central tendency
Nominal	Mode
Ordinal	Median
Interval/Ratio (not skewed)	Mean
Interval/Ratio (skewed)	Median

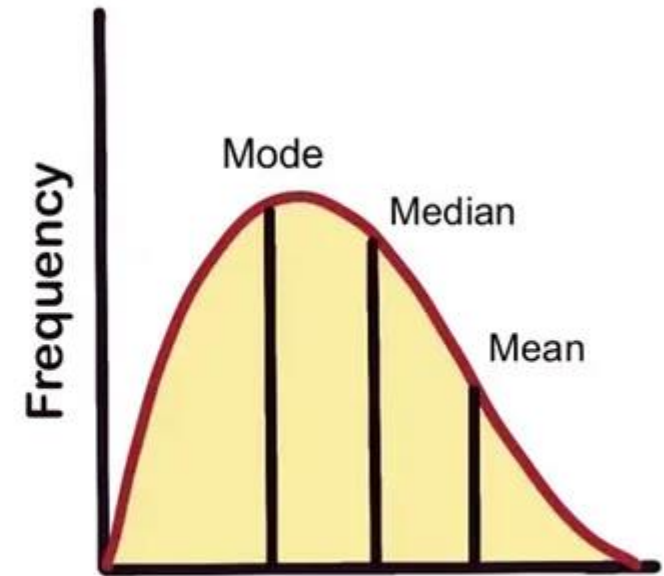
# Frequency Distribution



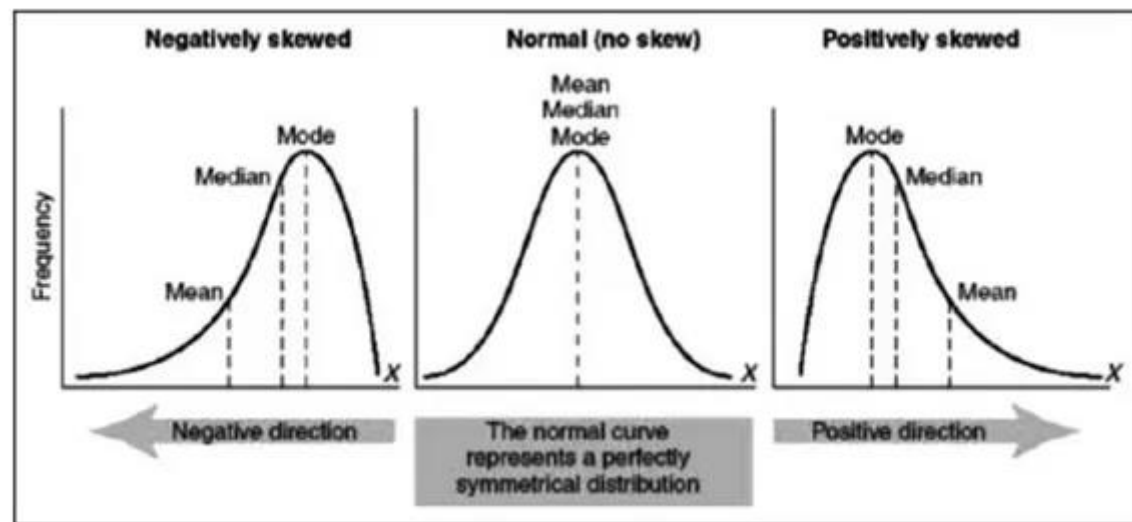
Left Skew  
 $\text{Mean} < \text{Median}$



Normal Distribution  
 $\text{Median} = \text{Mean}$



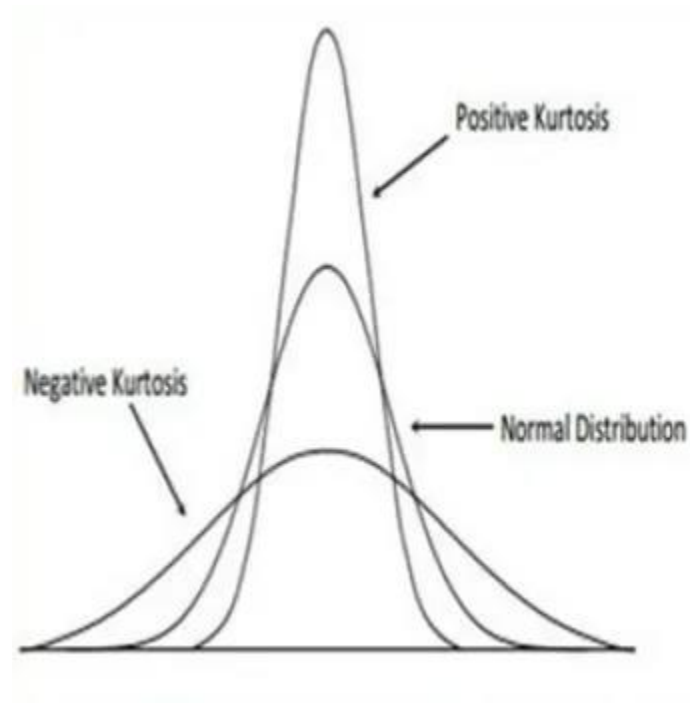
Right Skew  
 $\text{Mean} > \text{Median}$



- Positive Skewness: It occurs when the  $\text{Mean} > \text{Median} < \text{Mode}$ . The tail is skewed to the right in this case, i.e outliers are skewed to the right.
- Negative Skewness: It occurs when the  $\text{Mean} < \text{Median} < \text{Mode}$ . The tail is skewed to the left, i.e the outliers are skewed to left

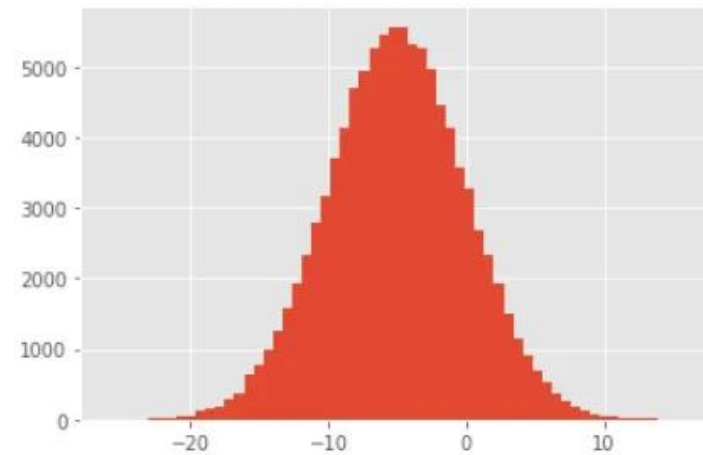


# Kurtosis



```
import numpy as np
from scipy.stats import kurtosis
from scipy.stats import skew
import matplotlib.pyplot as plt
plt.style.use('ggplot')
data = np.random.normal(-5, 5, 100000)
plt.hist(data, bins=60)
print("skew : ",skew(data))
print("kurt : ",kurtosis(data))
```

skew : -0.005168246772588942  
kurt : 0.010248125871068048



# Measures of Dispersion

```
data=[1,2,3,4,5,7,9]  
#sample variance  
statistics.variance(data)
```

7.9523809523809526

```
#population variance  
statistics.pvariance(data)
```

6.816326530612245

```
#square root of sample variance  
statistics.stdev(data)
```

2.819996622760558

```
#square root of population variance  
statistics.pstdev(data)
```

2.610809554642438

---

### quartiles, deciles, percentiles

These are three common measures used in statistics to divide an ordered data set into equal parts.

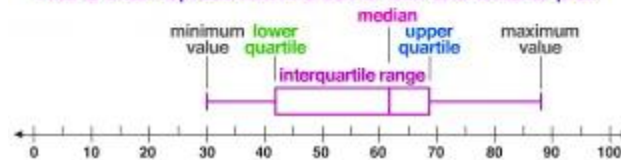
Quartiles = 4 equal parts, Deciles = 10 equal parts,  
Percentiles = 100 equal parts.

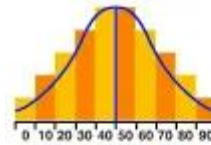


### five number summary

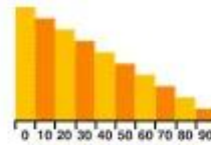
The five number summary gives the minimum value, lower quartile, median, upper quartile, and the maximum value.

This is often represented in a box or box-and-whisker plot.

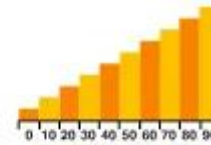




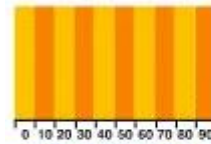
normal distribution  
unimodal, symmetric,  
aka 'bell curve'



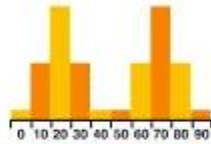
skewed distribution  
positively skewed,  
skewed right



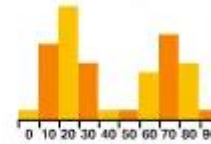
skewed distribution  
negatively skewed,  
skewed left



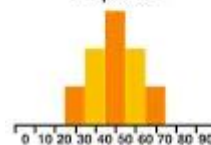
uniform distribution  
equally spread,  
no peaks



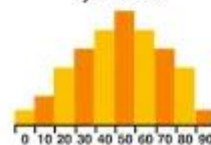
bimodal distribution  
two modes,  
symmetric



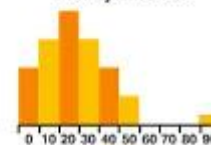
bimodal distribution  
two modes,  
non-symmetric



spread  
narrow range



spread  
wide range

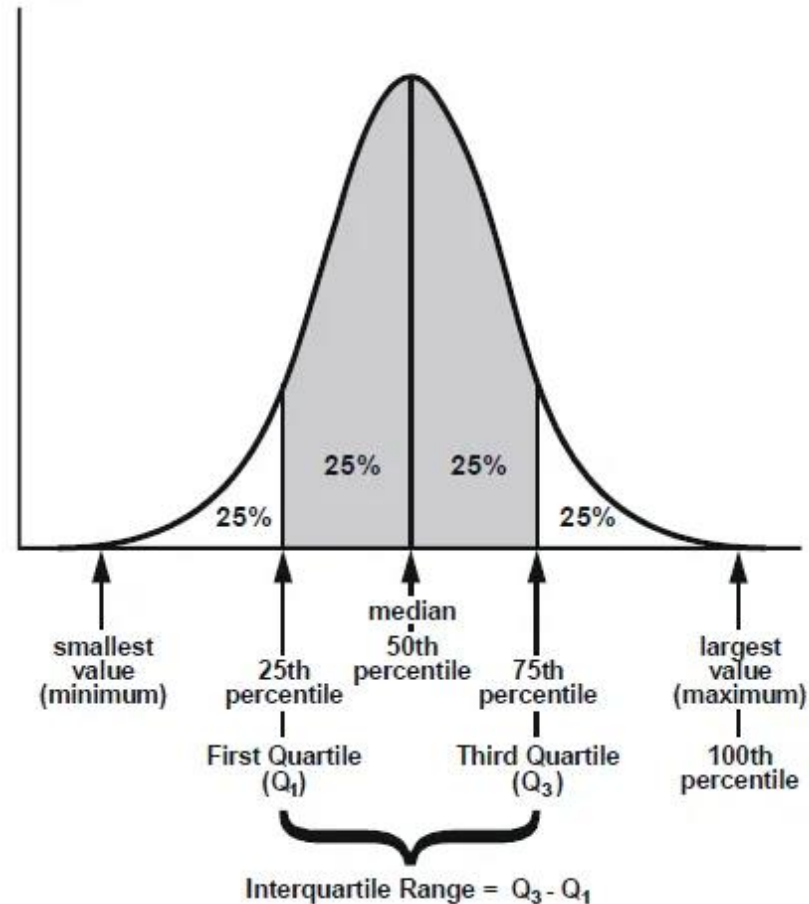


spread  
outlier

- `import numpy as np`
- `arr1= np.array([[12,43,56],[78,88,95],[79,89,43], [101,34,67]])`
- `arr2 = np.array([5,6,7,12,34,67,89])`
- `#Mean function`
- `print("Mean:", np.mean(arr2))`
- `#Median function`
- `print("Median:",np.median(arr2))`
- `#Standard Deviation Function`
- `print("Standard Deviation:", np.std(arr2))`
- `#Variance Function`
- `print("Variance:",np.var(arr2))`
- `#Average Function`
- `print("Average:",np.average(arr2))`
- `#Percentile Function`
- `print("Percentile:",np.percentile(arr2,5,0))`
- `#Minimum Function`
- `print("Minimum element:",np.amin(arr))`
- `#Maximum Function`
- `print("Maximum element:",np.amax(arr))`

- Mean: 31.428571428571427
- Median: 12.0
- Standard Deviation: 31.409084867994768
- Variance: 986.530612244898
- Average: 31.428571428571427
- Percentile: 5.3
- Minimum element: 12
- Maximum element: 101

# Measure of Spread





```
arr = [31, 35, 45, 49, 59, 69, 74, 79, 80, 81, 89, 94, 96, 99, 101, 104, 112,  
117,119,127,134]
```

```
# First quartile (Q1)  
Q1 = np.median(arr[:12])
```

```
# Third quartile (Q3)  
Q3 = np.median(arr[12:])
```

```
# Interquartile range (IQR)  
IQR = Q3 - Q1
```

```
print(IQR)
```

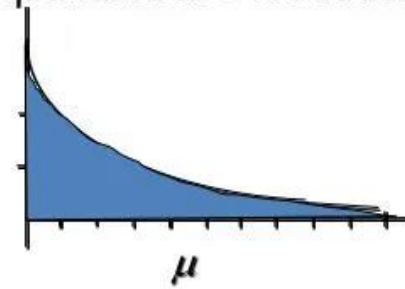
## Continuous Probability Distributions

---

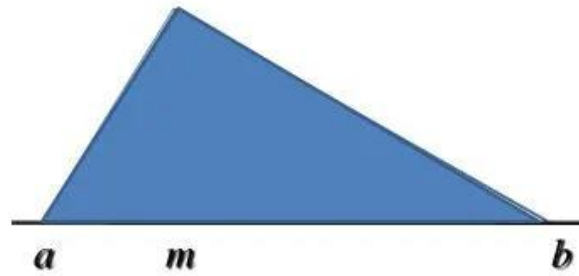
Uniform Distribution



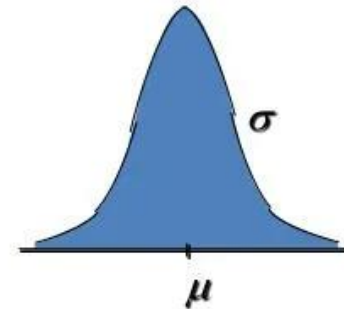
Exponential Distribution



Triangular Distribution



Normal Distribution



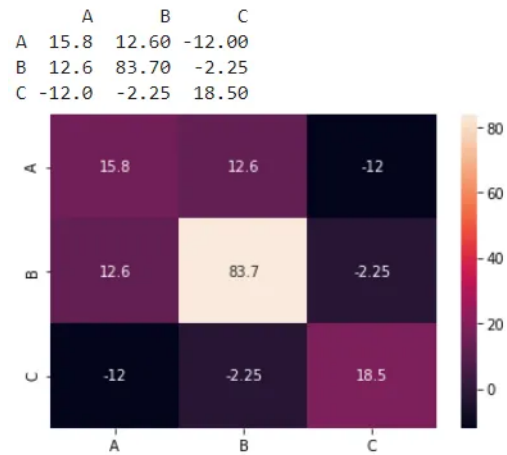
# Covariance

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

data = {'A': [55,47,52,45,49],
        'B': [48,41,36,38,23],
        'C': [20,25,27,31,22]
        }

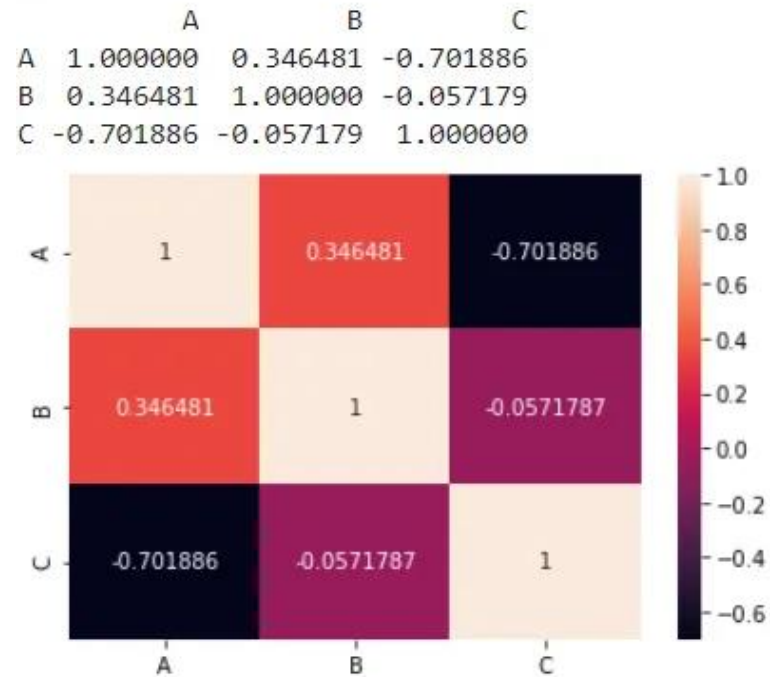
df = pd.DataFrame(data,columns=['A','B','C'])

covmtx = pd.DataFrame.cov(df)
print (covmtx)
sns.heatmap(covmtx, annot=True, fmt='g')
plt.show()
```



# Correlation

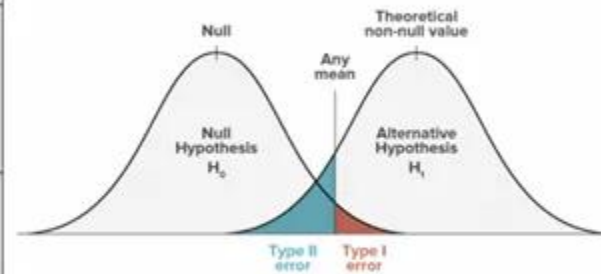
```
corrmtx = pd.DataFrame.corr(df)
print (corrmtx)
sns.heatmap(corrmtx, annot=True, fmt='g')
plt.show()
```



Covariance	Correlation
Measure of how much two random variables vary together	Indicates how strongly two variables are related
Involves the relationship between two variables or datasets	Involves relationship between multiple variables as well
Lies between $-\infty$ and $+\infty$	Lies between -1 and +1
Measure of correlation	Scaled version of covariance
Provides direction of relationship	Provides direction as well as strength of relationship
Dependent on scale of variable	Independent of scale of variable
Has dimensions	Dimensionless

# Type 1 & 2

		Conclusion	
		Accept the Null	Reject the Null
The True State of the Nature	$H_0$ is True	Correct	False Positive Type I Error $\alpha$
	$H_0$ is False	False Negative Type II Error $\beta$	Correct



# Linear Regression

The diagram illustrates the linear regression equation  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ . It includes labels for the dependent variable, independent variable, and the components of the equation: Y intercept, Slope Coefficient, and Error Term.

Dependent Variable  
(Response Variable)

Independent Variable  
(Predictor)

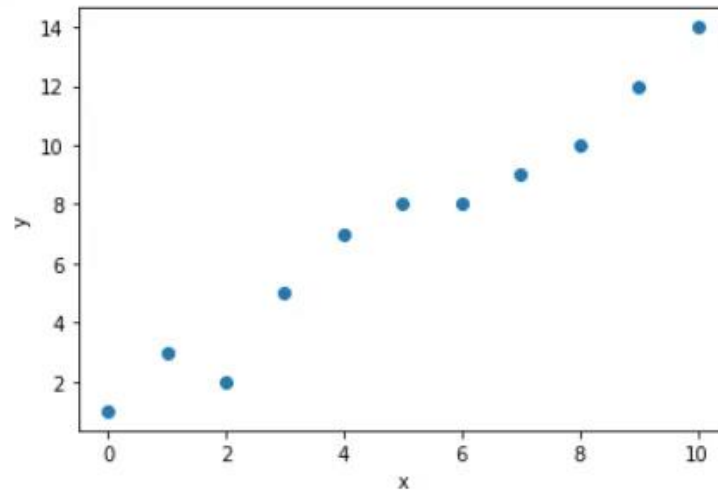
$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Y intercept

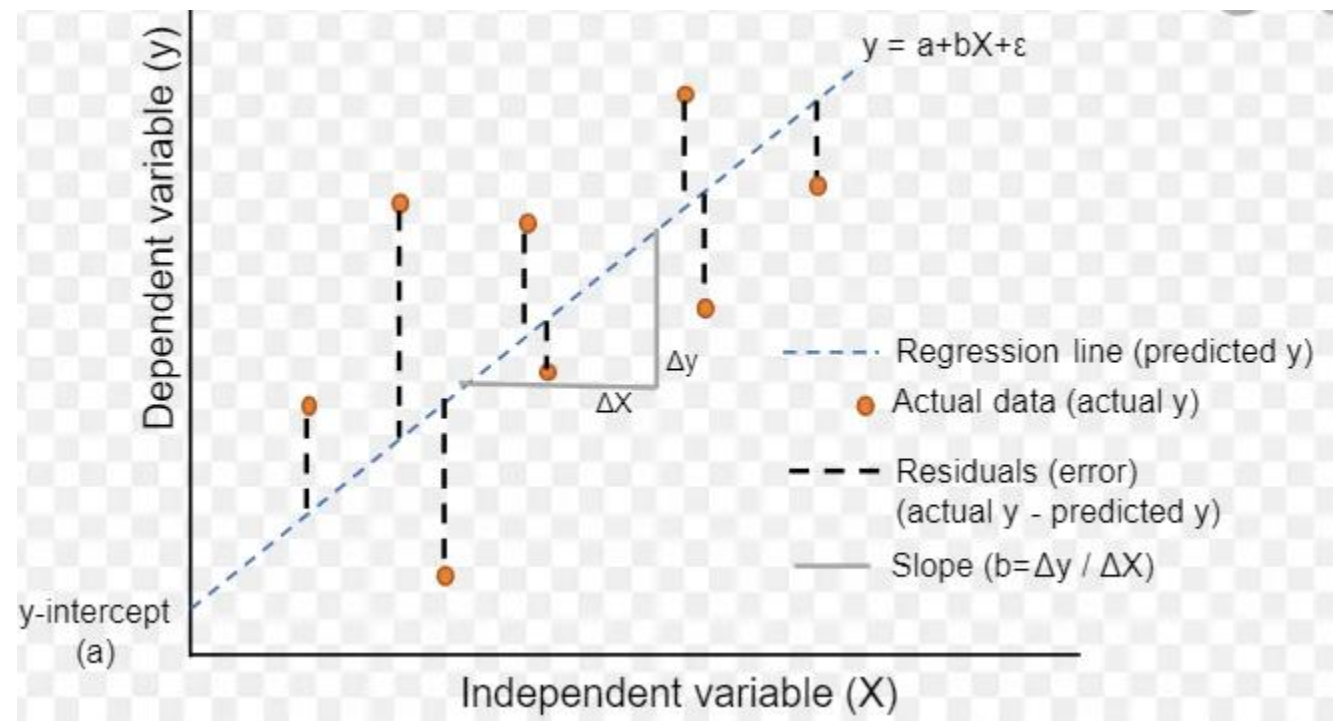
Slope  
Coefficient

Error Term

```
import numpy as np
import matplotlib.pyplot as plt
#sample dataset
x = np.array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
y = np.array([1, 3, 2, 5, 7, 8, 8, 9, 10, 12, 14])
#scatter plot of dataset
plt.scatter(x,y)
plt.xlabel('x')
plt.ylabel('y')
plt.show()
```





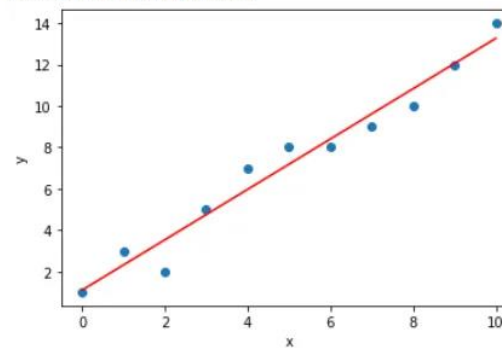


```

n = np.size(x)
mx = np.mean(x)
my = np.mean(y)
# cross-deviation and deviation about x
SSxy = np.sum(y*x) - n*my*mx
SSxx = np.sum(x*x) - n*mx*mx
# estimating regression coefficients
b1 = SSxy/SSxx
b0 = my - b1*mx
print("Estimated coefficients:\nb_0 = {} \
      \nb_1 = {}".format(b0, b1))
plt.scatter(x, y)
y_pred = b0 + b1*x
# plotting the regression line
plt.plot(x, y_pred, color = "r")
plt.xlabel('x')
plt.ylabel('y')
plt.show()

```

Estimated coefficients:  
 $b_0 = 1.09090909090909$   
 $b_1 = 1.21818181818183$



# Multiple Regression

The diagram illustrates the Multiple Regression equation: 
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \varepsilon$$
 Each component of the equation is labeled with an arrow pointing to it: 

- Dependent Variable (Response Variable)** points to  $Y$ .
- Independent Variables (Predictors)** points to  $X_1$  and  $X_2$ .
- Y intercept** points to  $\beta_0$ .
- Slope Coefficient** points to  $\beta_1$  and  $\beta_2$ .
- Error Term** points to  $\varepsilon$ .