

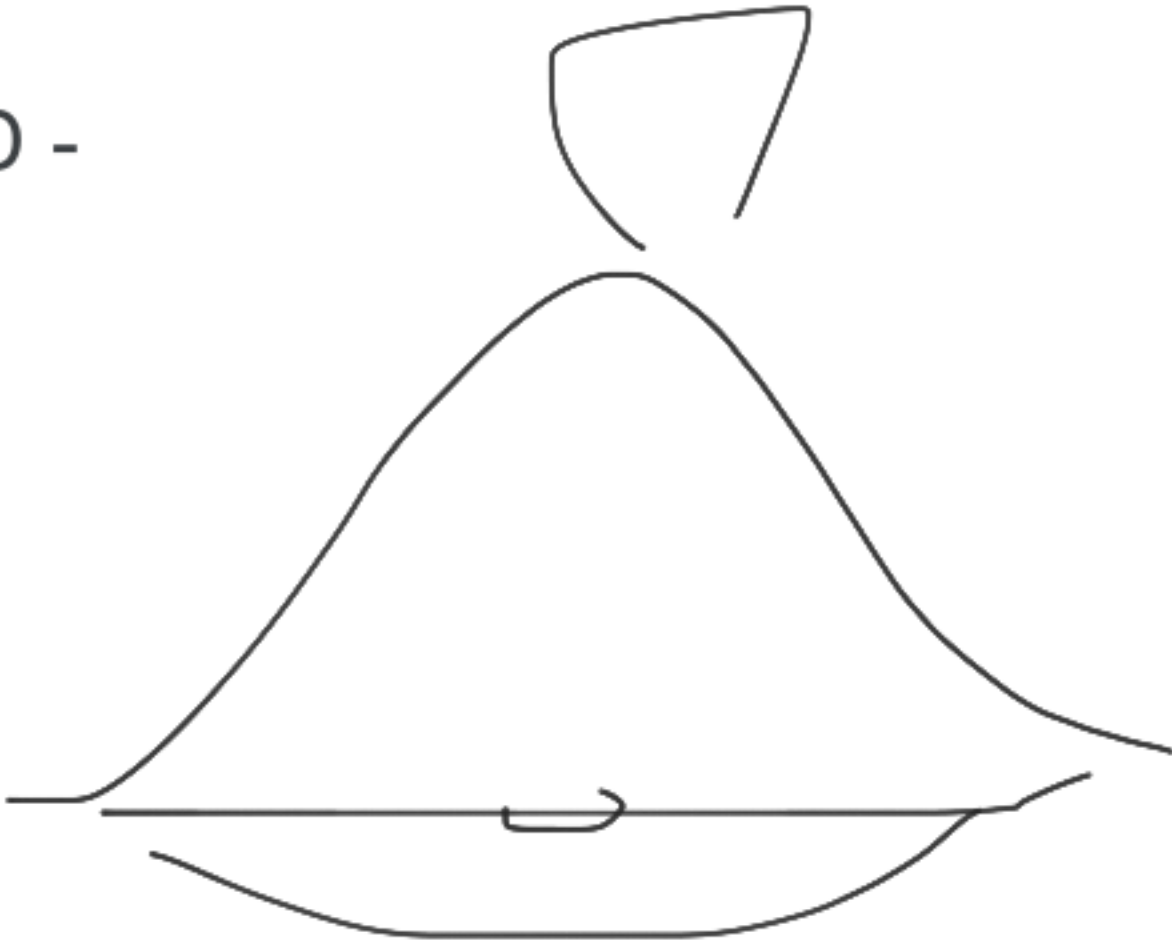
Stats

Type of Analysis

Types of Stats

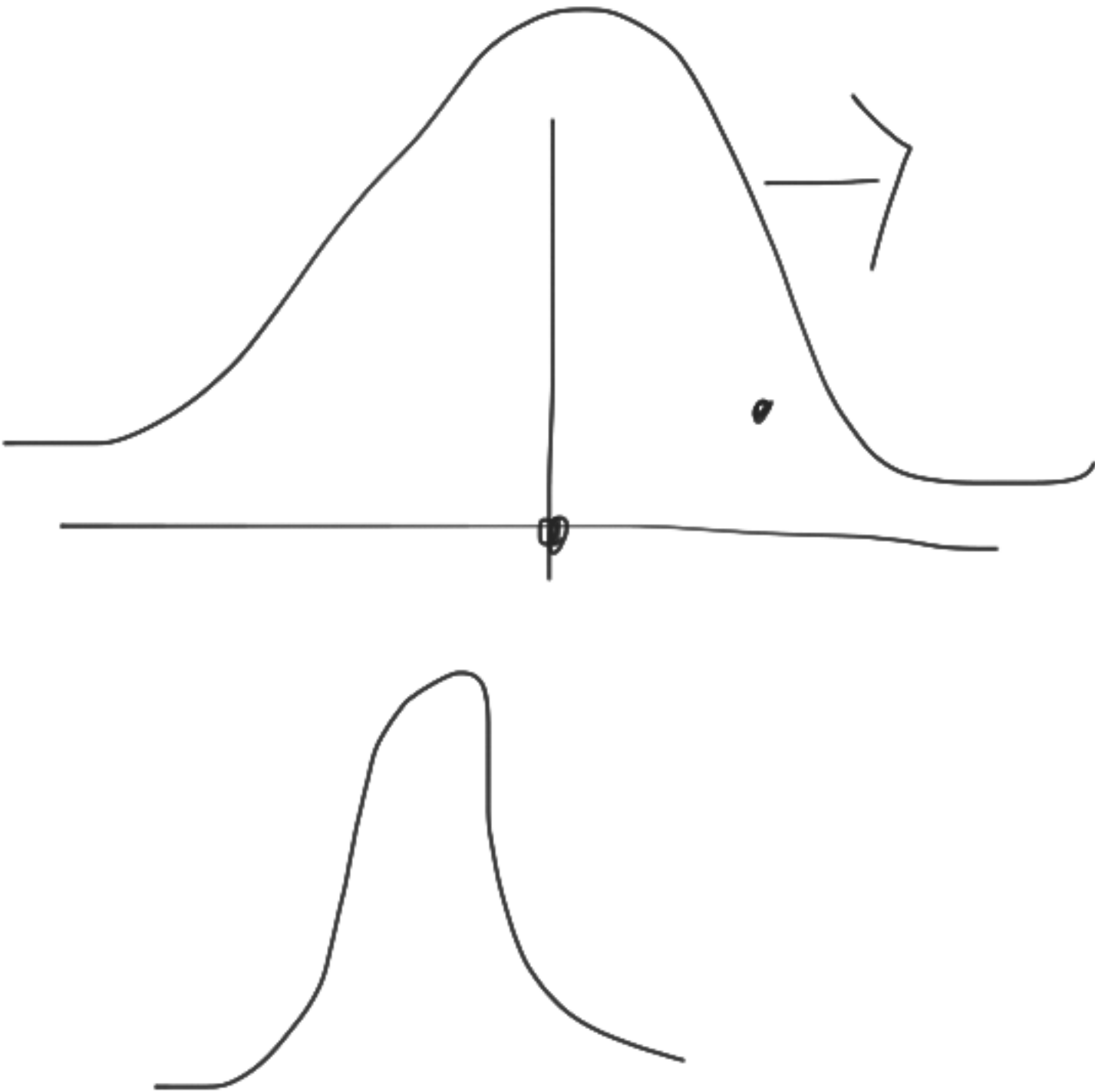
Measure of central Tendency - Mean, Median, Mode

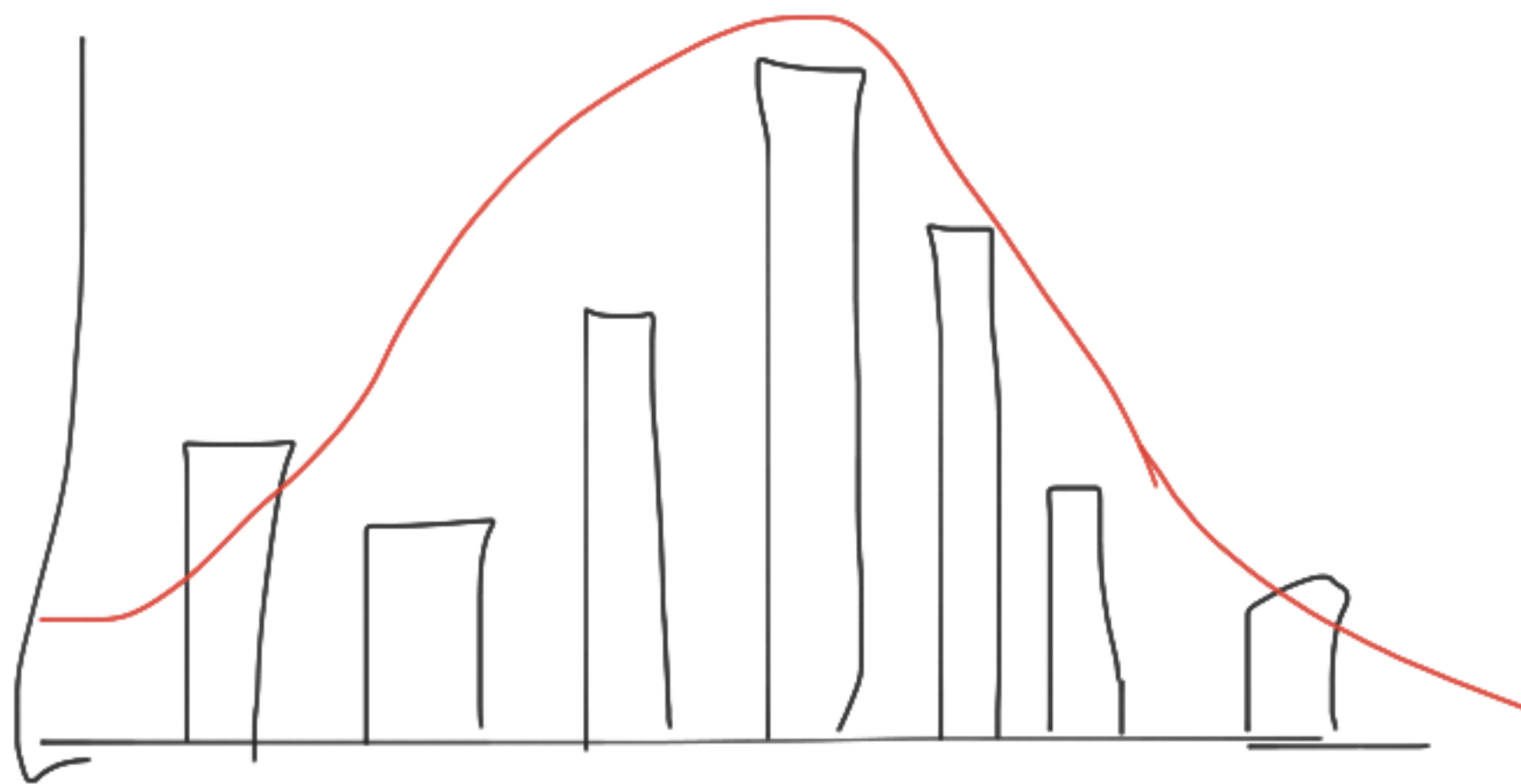
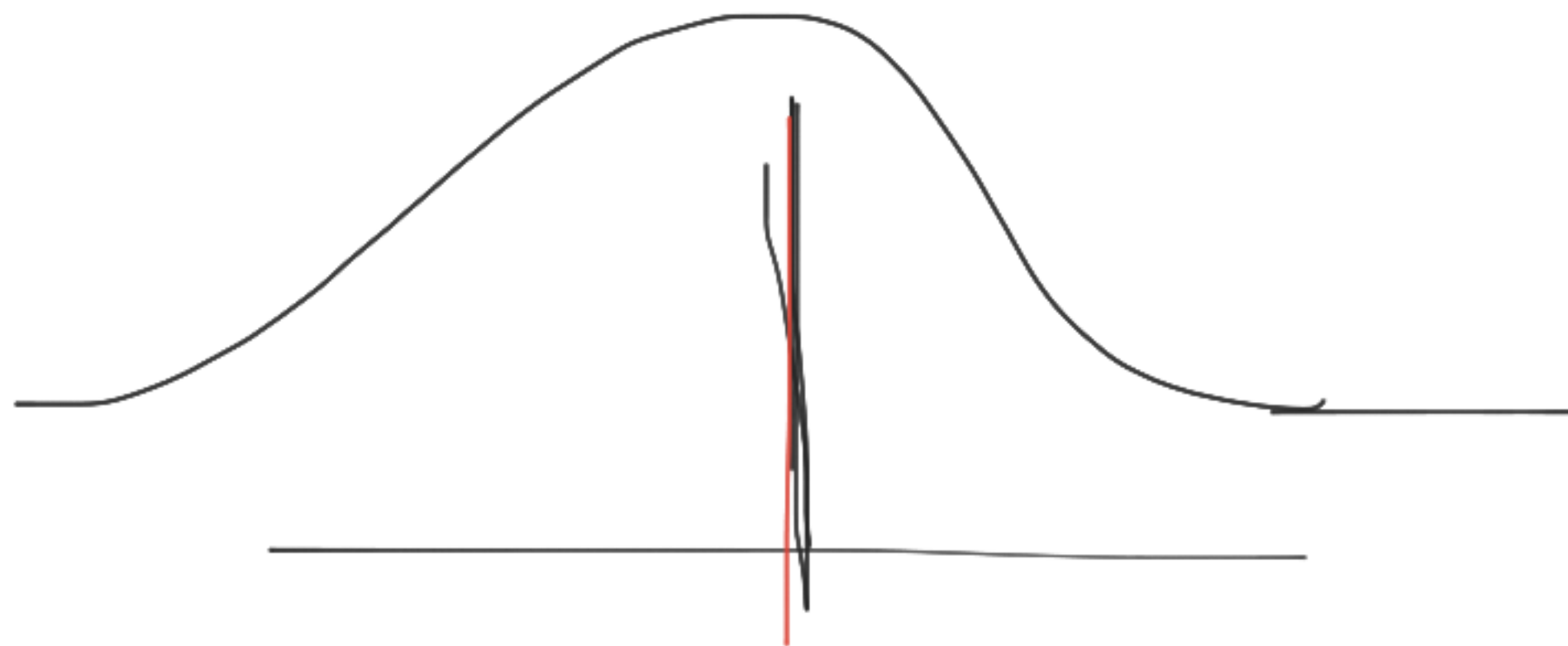
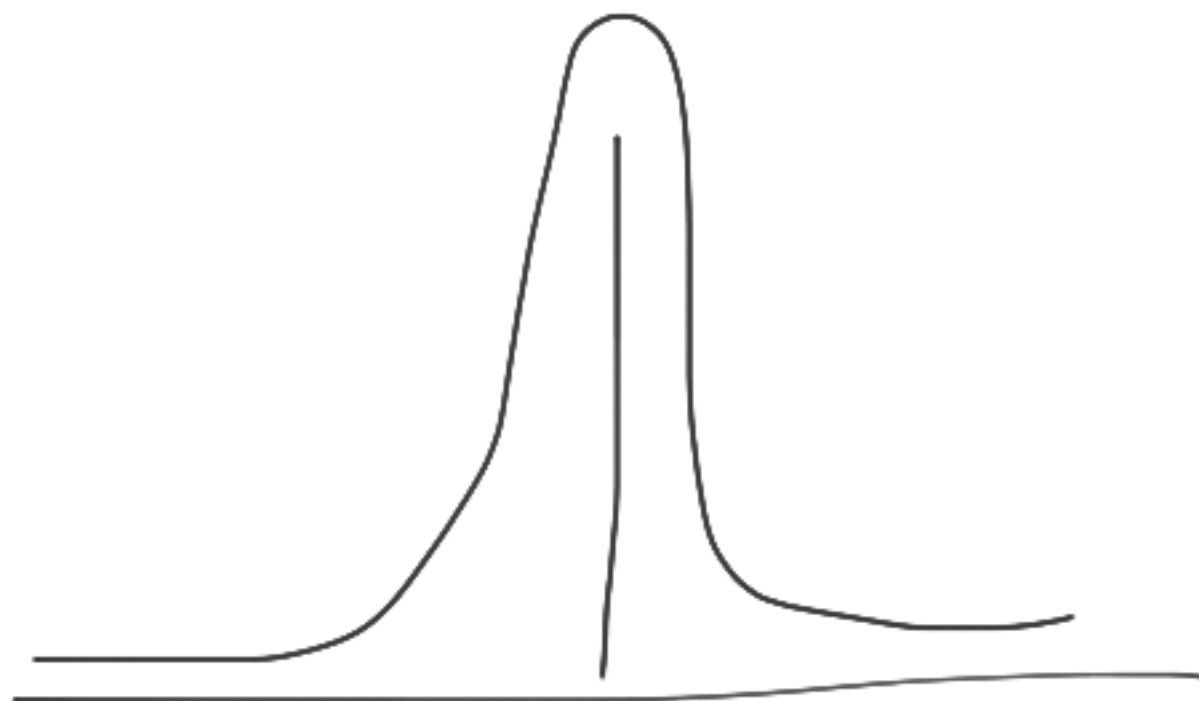
SD -



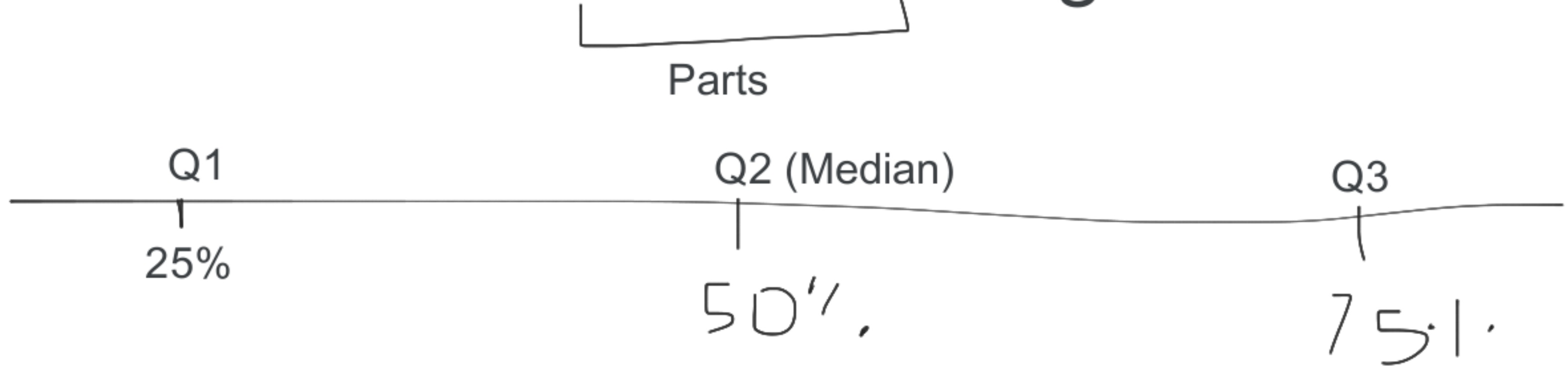
Measure of Variance

Vary





Inter Quartile Range IQR



IQR - It is a measure based on the dividing the dataset into quartiles (Parts). It also gives an idea where the bulk of the values lies. Plot -> Box Plot, Violin Plot

*It helps us to detect the outliers

$$\begin{aligned} \text{IQR} &= Q3 - Q1 \\ Q1 &= 25\% \\ Q2 &= 50\% \\ Q3 &= 75\% \end{aligned}$$

Steps to find out the IQR

1. Arrange the data in the form of ascending order
2. Find Q2 (Median)
3. Find Q1 & Q3
4. Find $IQR = (Q3 - Q1)$
5. Upper Bound: $Q3 - (1.5 * IQR)$
6. Lower Bound: $Q1 - (1.5 * IQR)$

iPhone 13
iPhone 12 Mini
iPhone 14 Pro
iPhone SE
iPhone 13 Pro Max
Samsung S22 Ultra

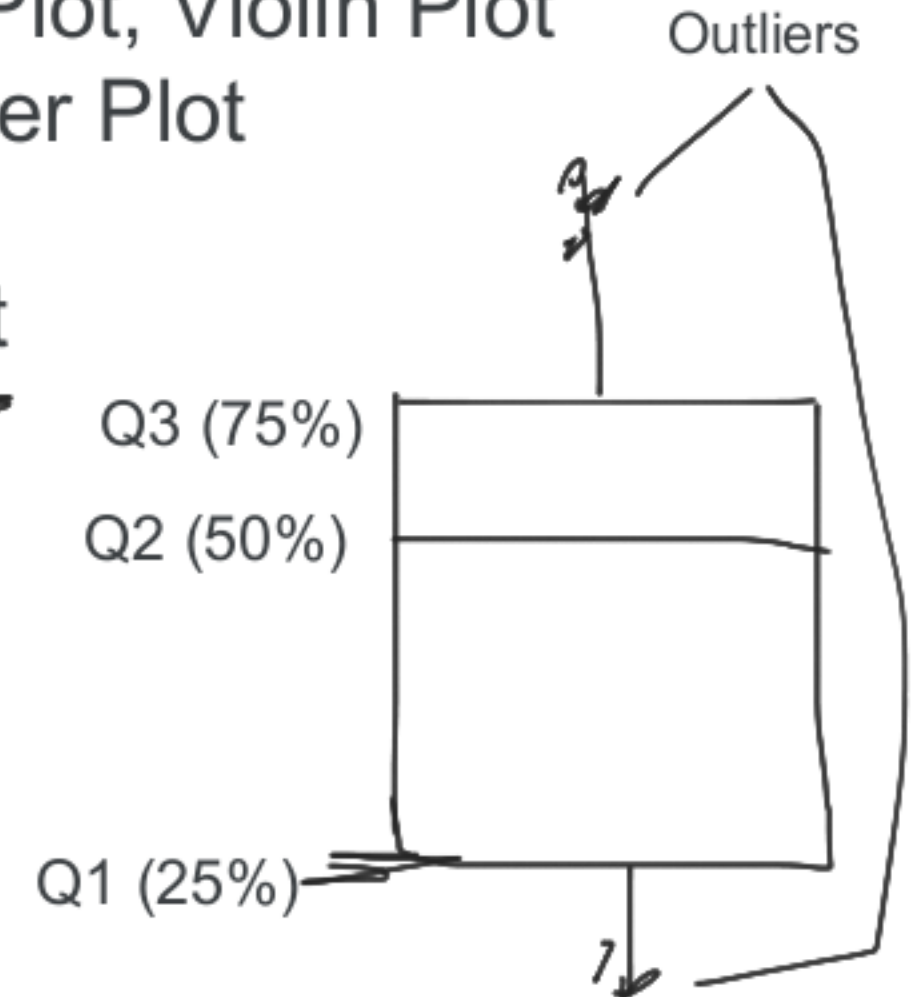
1
2
3
4
5
6
7
8
9
250

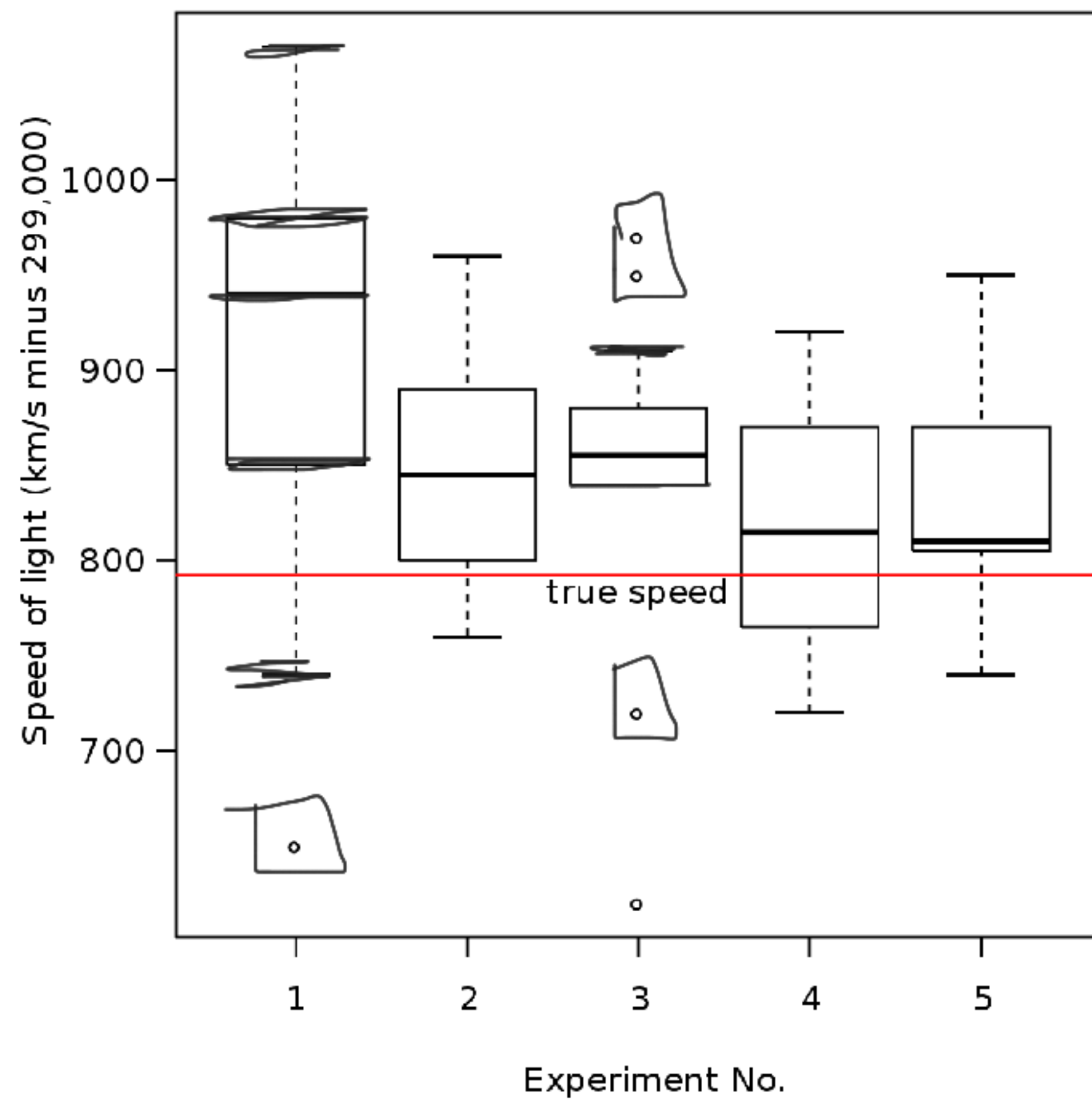
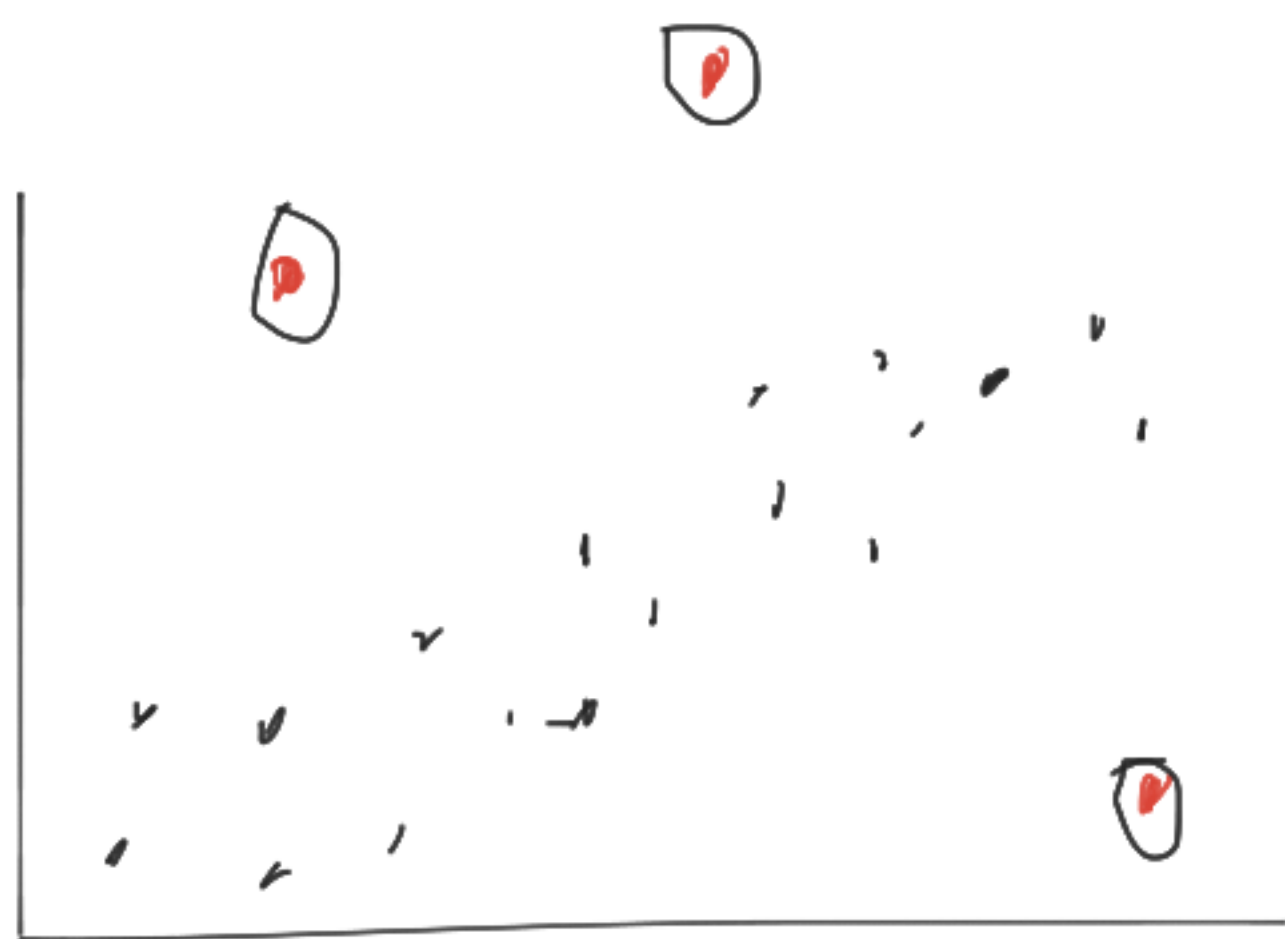
11, 20, 23, 24, 34, 25, 46, 90

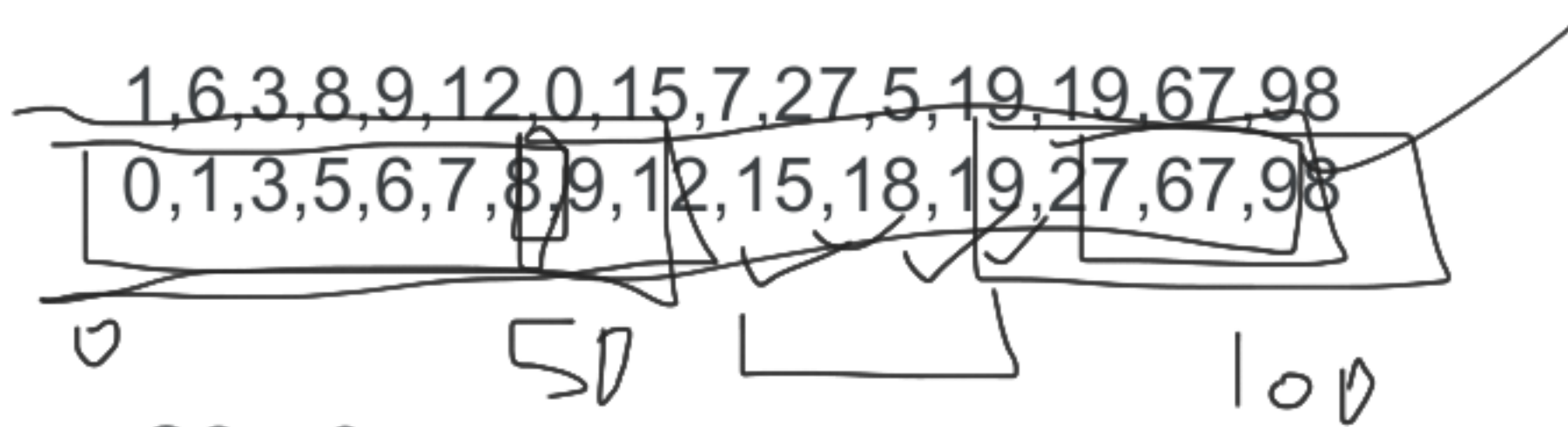
Note: Any datapoint that lies outside the Lower & Upper Bound is an OUTLIER.

Ways to find out the outlier

1. Box Plot, Violin Plot
2. Scatter Plot
3. IQR
4. z-test







$$Q2 = 9$$

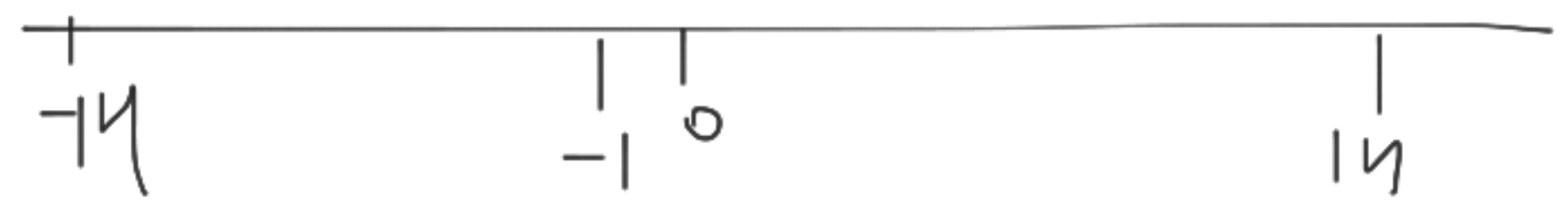
$$Q1 = 5.5$$

$$Q3 = 18.5$$

$$IQR = Q3 - Q1 = 18.5 - 5.5 = 13$$

$$UB = 18.5 - (1.5 \cdot 13) = -1$$

$$LB = 5.5 - (1.5 \cdot 13) = -14 = +14$$



Random Variables

P

It is a variable that stores value that depends upon the outcomes of a random phenomeon.

Types of random Variable

Numerical
RV

✓ $x = 25$
✓ $x = 0.5$
✓ $x = 92.9$

Categorical
RV

$x = \text{'Male'}$
 $x = \text{'Electronics'}$
 $x = \text{'Domestic Animal'}$

Discrete RV

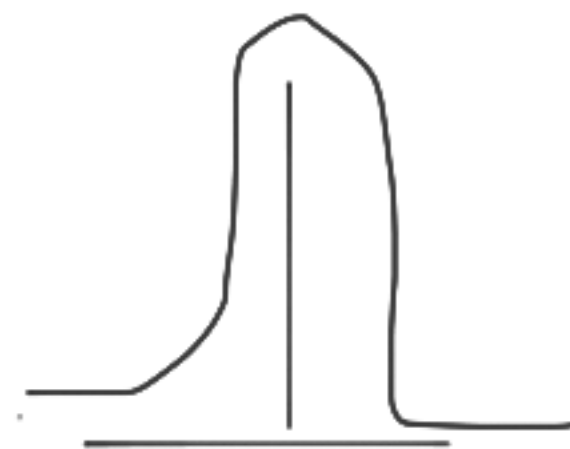
Should be a Whole
number and cannot
be -ve

Continuous
RV

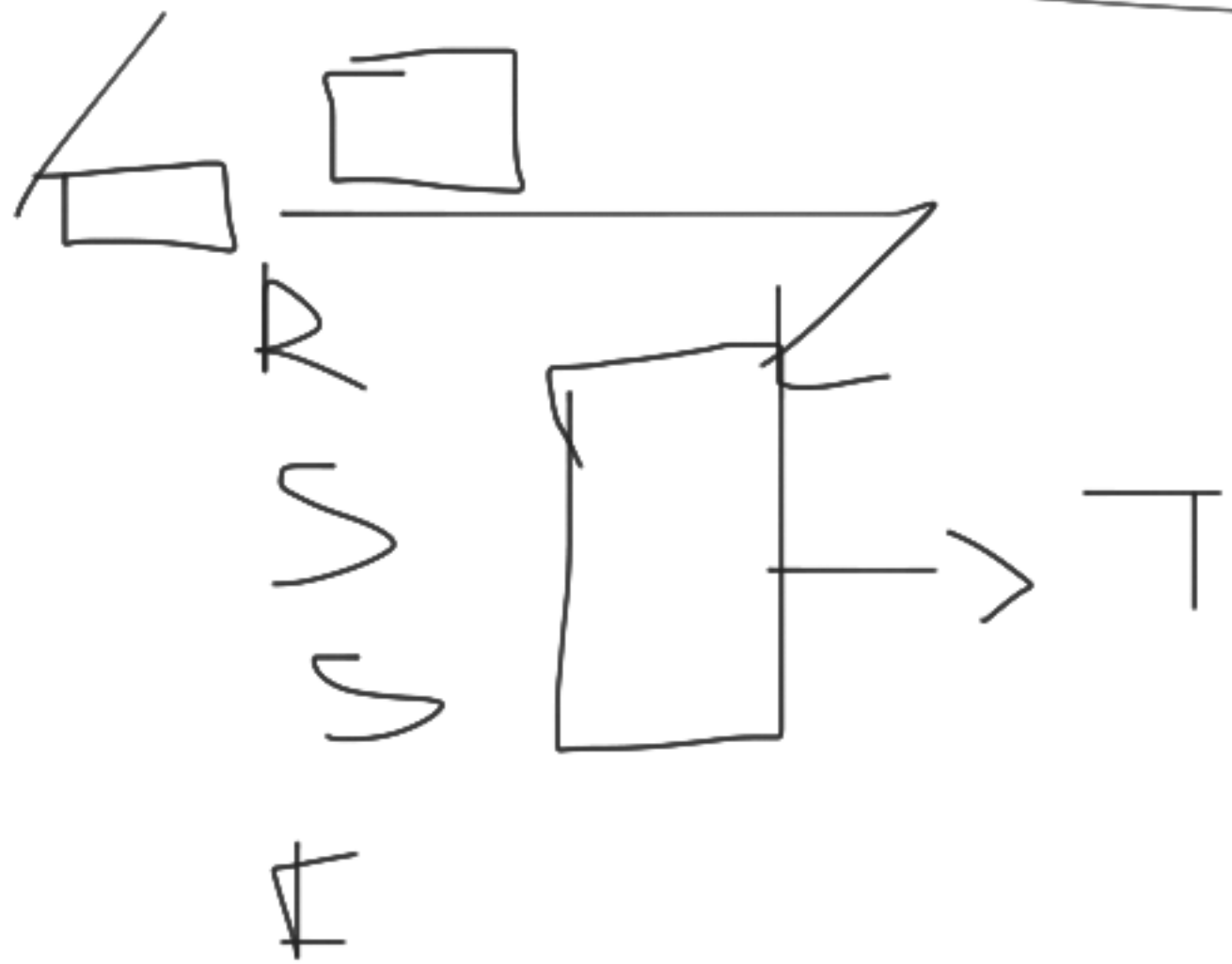
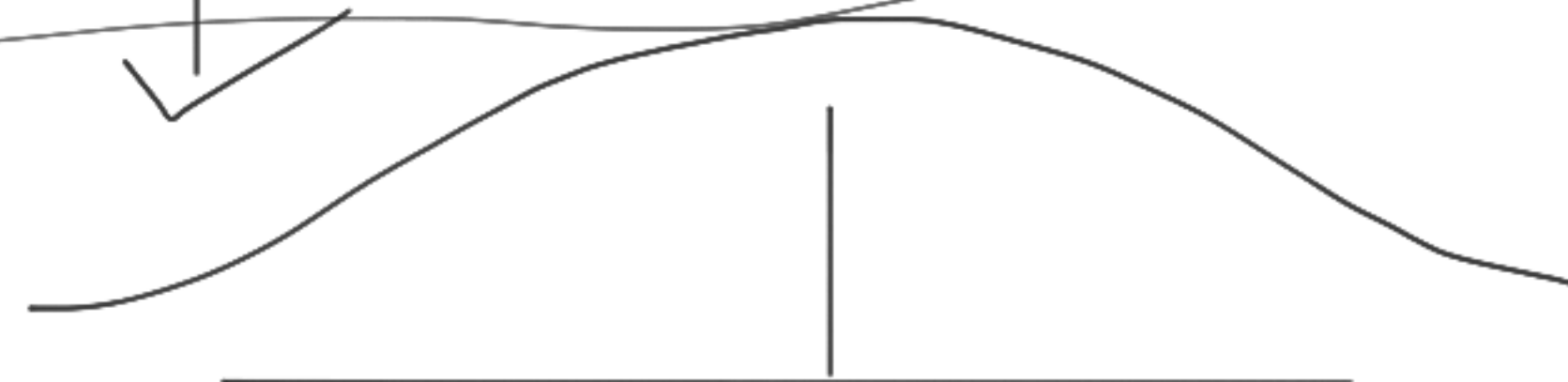
Can have any number
that includes floating,
decimal

Feature Engineering (Feature Selection)

Feature
Target



A perfect model is a model which is having low bias & low variance



Covariance
Karl Pearson Correlation coefficient
Spearman Rank correlation

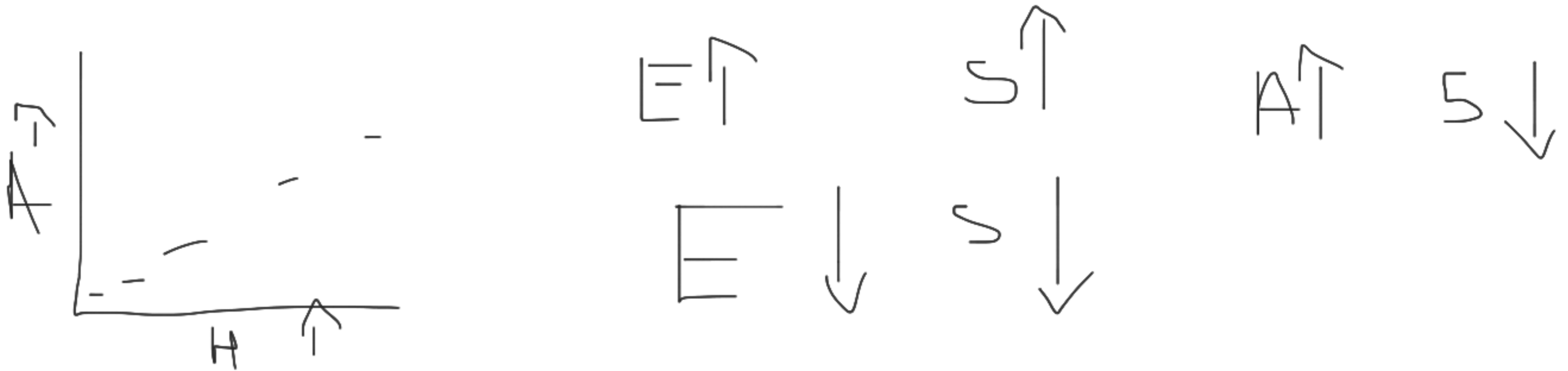
Covariance

It is a measure to show the relationship between two random variables.
It will only help to find out the direction of the relationship (Disadvantage)

It can take any +ve & -ve values

Positive Covariance - When the outcome of the two RV tend to move in the same direction

Negative Covariance - When the outcome of the two RV tend to move in the different direction



$$\text{Cov}(x,y) = \sum \left(\frac{[x(i) - \text{mean}(x)][y(i) - \text{mean}(y)]}{n} \right)$$

```
def cov(x,y):  
    mean_x = np.mean(x)  
    mean_y = np.mean(y)  
    count = 0  
    for i,j in zip(x,y):  
        result = (i - mean_x)*(j - mean_y)  
        count += result  
    return count/len(x)
```

np.cov(df.X, df.Y)

Pearson Correlation Coefficient

-1 to 1

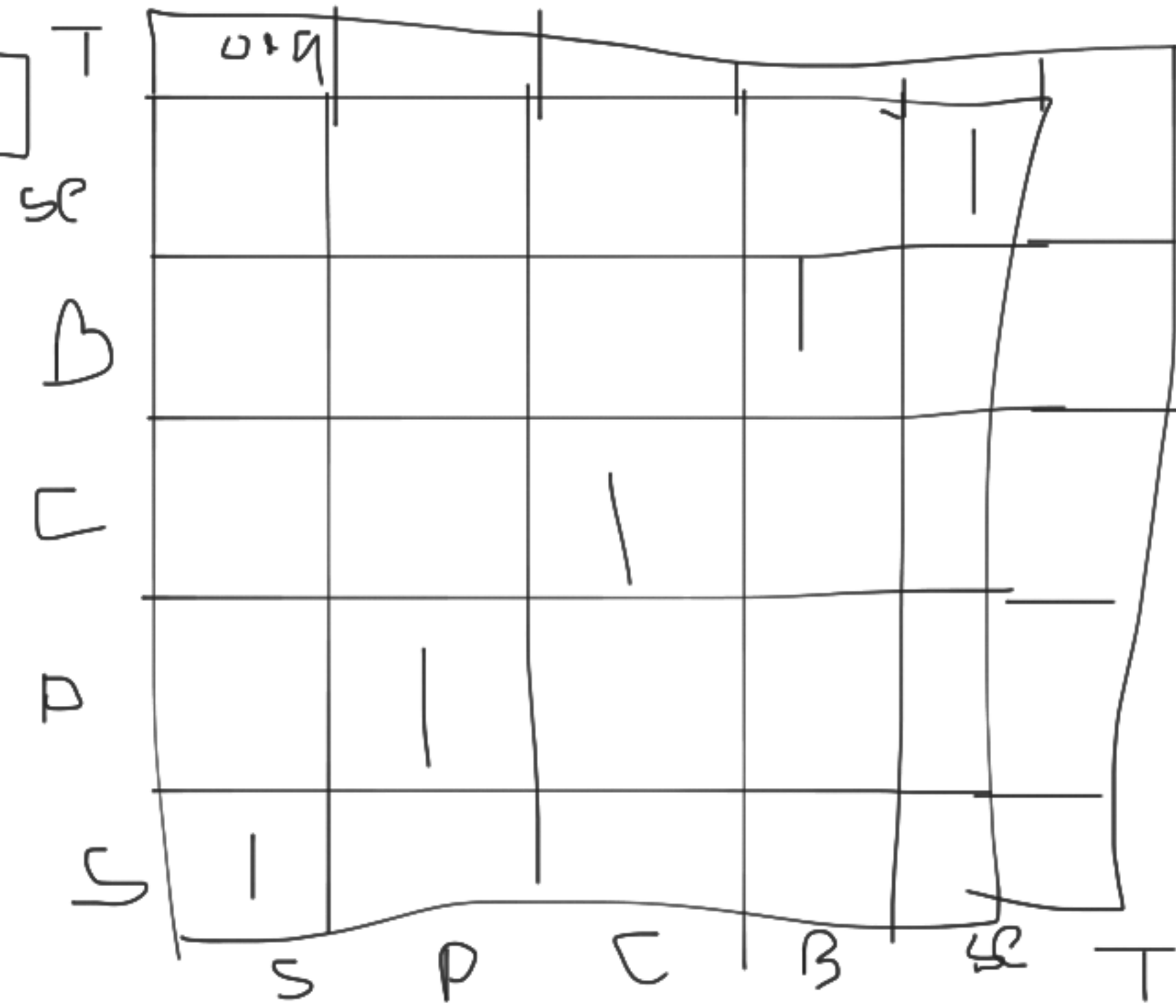
It is a measure to show the relationship between two random variables. It not only helps to find out the direction of the relationship, but also the **strength**.

Apple -> New budget phone

Top 5

Screen
Processors
Camera
Battery
Security

- ✓ 1. Screen
- ✓ 2. Processors
- ✓ 3. Camera
- ✗ 4. Design
- ✓ 5. Battery
- ✗ 6. Storage
- ✗ 7. Heating reduction ✗
- ✓ 8. Security
- ✗ 9. Color
- ✗ 10. Screen refresh rate



$$\text{Corr}(x,y) = \frac{\text{Cov}(x,y)}{\text{SD}(x) * \text{SD}(y)}$$

`np.corrcoef(df.Age, df.Height)`

`df.corr()`

1 -> Perfect Positive Correlation

0.99 - 0.79 -> High +ve Corr

0.78 - 0.5 -> Moderate +ve Corr

> 0.5 - 0 -> Low +ve corr

```
def corrcoef(x,y):  
    mean_x = np.mean(x)  
    mean_y = np.mean(y)  
    count = 0  
    for i,j in zip(x,y):  
        result = (i - mean_x)*(j - mean_y)  
        count += result  
    cov = count/len(x)  
    SD_x = np.std(x)  
    SD_y = np.std(y)  
    denominator = SD_x * SD_y  
    return cov/denominator
```

Variable Measurement Scales

It is important because different types of data allow for different types of Data Analysis

1. Nominal - Qualitative Data. It is the data that split into categories-> Gender, Color, Car Type
2. Ordinal - Qualitative Data. It is the data where the order matters but the distance b/w the values doesn't matter.
3. Interval - Quantative Data. It is the data where the order matters and the distance b/w the values matters, and a natural zero (0) is not present.
4. Ratio - Quantative Data. It is the data where the order matters and the distance b/w the values matters, and a natural zero (0) is present



Spearman Rank Correlation

Data: Ordinal, Interval, Ratio

It is a measure to show the monotonic relationship between two random variables. It not only helps to find out the direction of the relationship, but also the strength.

A monotonic relationship is where:

1. One variable increases and the other also increases, or
2. One variable decreases and the other also decreases

	A	B	C	D	E	F	G
		Physical activity	Blood presure	Physical activity	Blood presure	d	d ₂
1	Name	(min)	(mm Hg)	(rank)	(rank)		
2	Alan	60	118	1	9	-8	64
3	Carl	55	117	2	10	-8	64
4	David	25	120	8	7	1	1
5	Don	50	121	3	6	-3	9
6	John	40	119	5	8	-3	9
7	Matt	45	122	4	5	-1	1
8	Mike	35	123	6	4	2	4
9	Neal	10	124	10	3	7	49
10	Rick	30	125	7	2	5	25
11	Rob	20	126	9	1	8	64
12							290
13	Spearman correlation		-0.757575758	=CORREL(D2:D11, E2:E11)			
14			-0.757575758	=1-(6*G12/(10*(10^2-1)))			

$$\rho = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

$$r = \frac{\text{cov}(r(x), r(y))}{\text{SD}(x) * \text{SD}(y)}$$

import scipy

scipy.stats.spearman(a,b)

r(x) -> rank of x

r(y) -> rank of y

```
def spearman(x,y):  
    rank_x = x.rank(ascending=False)  
    rank_y = y.rank(ascending=False)  
    covariance = cov(rank_x, rank_y)  
    SD_x = np.std(x)  
    SD_y = np.std(y)  
    deno = SD_x * SD_y  
    return covariance/deno
```