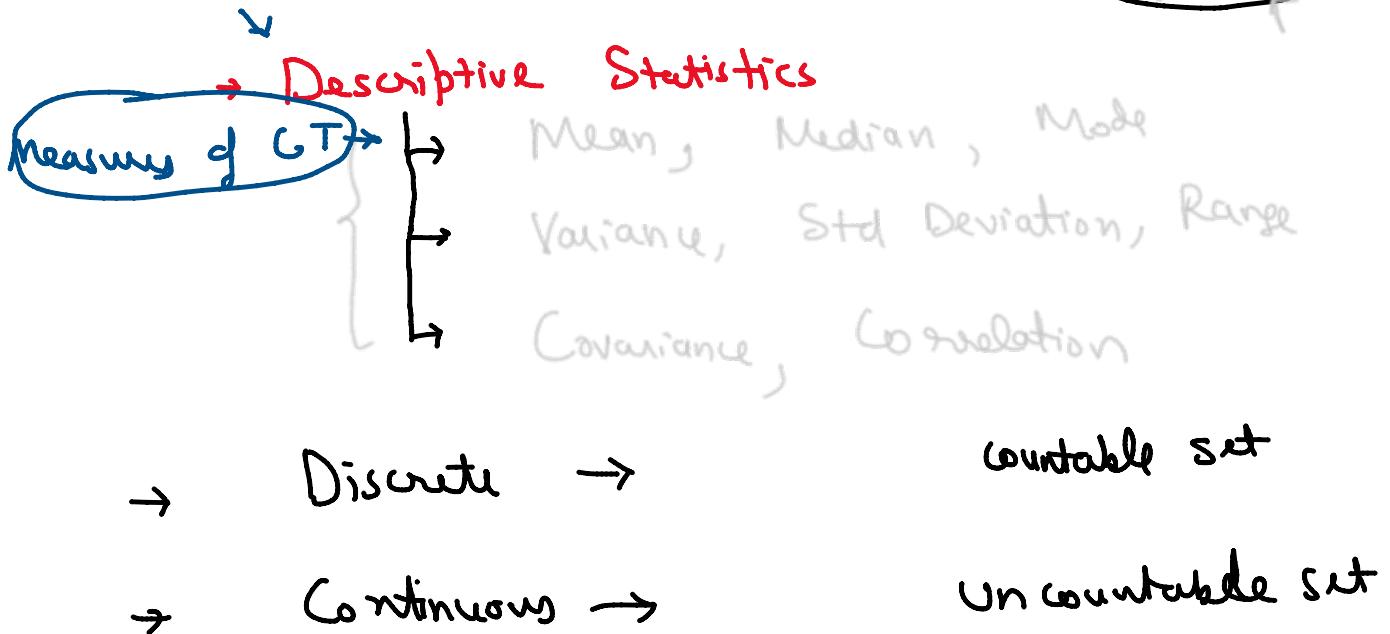


Starting @ 4:07 PM

Topics ↴

→ Statistics



1. ✓
Measure's of Central Tendency

1. Measures of Central Tendency

1. Mean

(Average)

H	W	G
Height (cm)	Weight (kg)	Gender
160	50	Female
180	78	Male
146	45	Female
162	51	Female
184	80	Male
180	60	Male

NOTE: Assume that all the students in the class come from same geographical region and fall under same age group

$$\underline{\mu_H} = 168.67$$

$$\underline{\mu_W} = 60.66$$

$$\underline{\mu_H} =$$

$$\frac{h_1 + h_2 + h_3 + \dots + h_6}{6} \rightarrow \text{Total no. of obj}$$

$$\underline{\mu_H} = \frac{\sum_{i=1}^n h_i}{n}$$

160, 180, 146, 162, 184, 180

$$\underline{168.67} \checkmark$$

{Central}
[Tendency]

Google form

$$168.67 \quad (\text{without error})$$

$$438.67 \quad (\text{with error})$$

1 error

Extremely small
or
a very large value

$$\begin{matrix} H \\ \downarrow \\ 5.6 \end{matrix} \quad \begin{matrix} W \\ \downarrow \\ 18 \end{matrix}$$

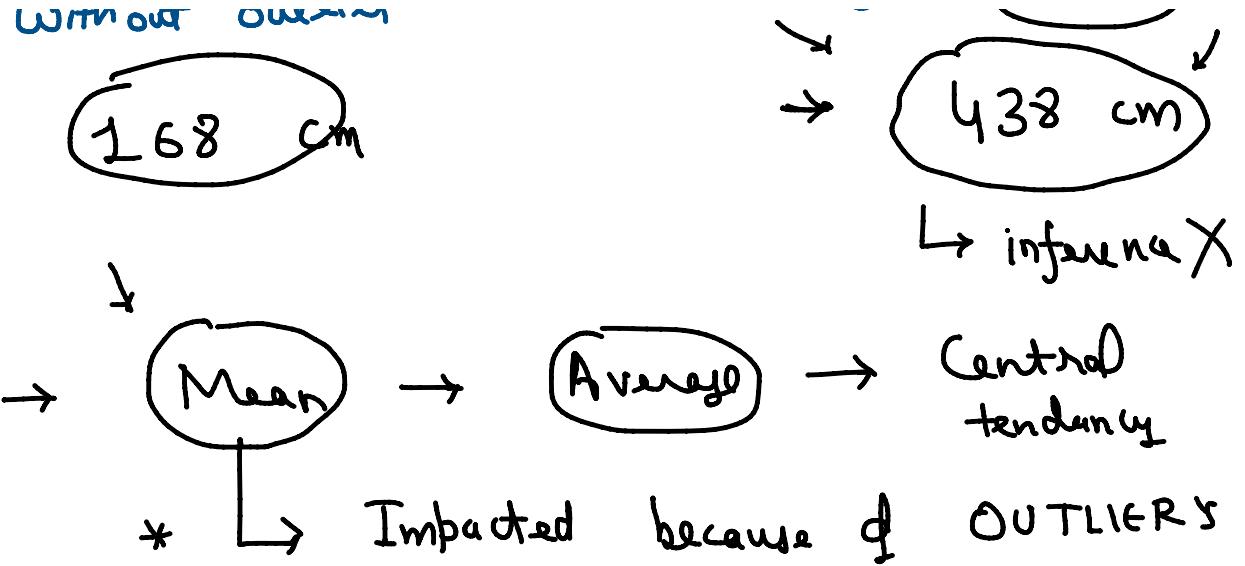
→ outliers

without outlier

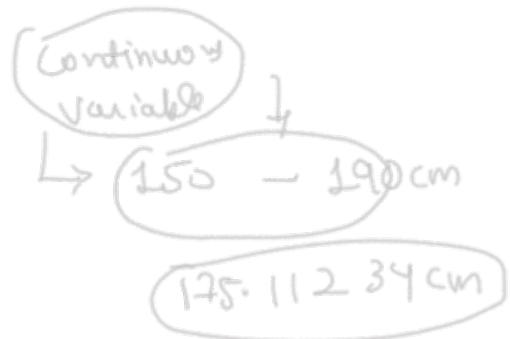
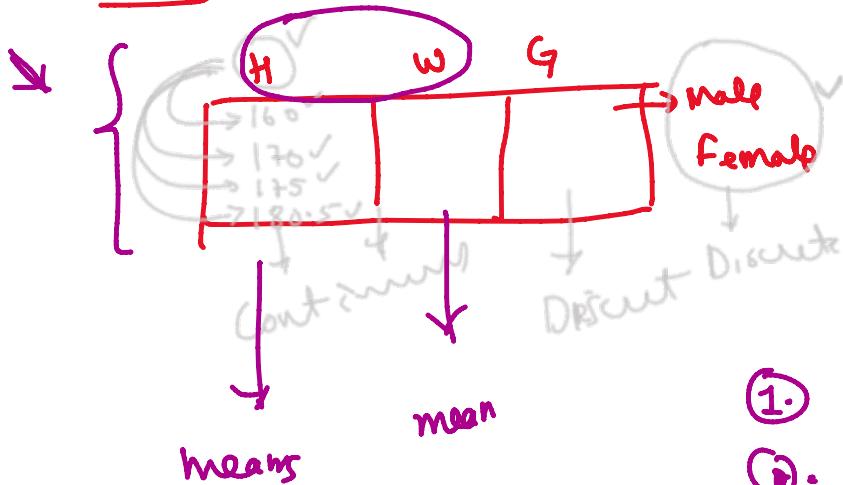
with outlier

1.22 ...

Without outliers



Columns → Variables



- ① 1 variable at time ✓
- ② collection

Univariate Analysis → Mean

Height (cm)	Weight (kg)	Gender
160	50	Female
180	78	Male
146	45	Female
162	51	Female
184	80	Male
180	60	Male

NOTE: Assume that all the students in the class come from same geographical region and fall under same age group

Mean

Univariate Measure of

$$\mu_H = \frac{\sum_{i=1}^n h_i}{n}$$

→ 3 variables

→ H, W → Continuous Variable

→ G → Discrete/ Categorical

analysis central tendency

$$\bar{X}_H = \frac{\sum_{i=1}^n x_i}{n}$$

Impacted in the presence of Outliers

↳ Solution → Drop the Outliers ✓

Americans →

50,000 \$

Avg.
~~35 LPA~~

IIT B

Admission →

17 LPA

1. Mean

2. Median ↴

160, 180, 146, 162, 184, 180

1. Sort the numbers ✓

2. Pick the middle value

merge sort
→ nlogn



$$\frac{162 + 180}{2} \Rightarrow 171$$



① sort

② → Middle

$$\frac{5+3-2}{2}$$

1. Sort the values in increasing order
2. If no. of observation is
Even \rightarrow median = $\frac{(\frac{n}{2})^{th} + (\frac{n+1}{2})^{th}}{2}$
Odd \rightarrow median = $(\frac{n+1}{2})^{th}$ value

→ Measures of Central Tend ↴

① Mean → 3.5 → O(n)

② Median →
 1. Sorting $\rightarrow O(n \log n)$
 2. Middle $\rightarrow O(1)$
 $\underline{O(n \log n)}$

Measure Central Tendency ↴

160, 180, 146, 162, 184, 1800 → [146 160 162 180 184 1800]



Mean

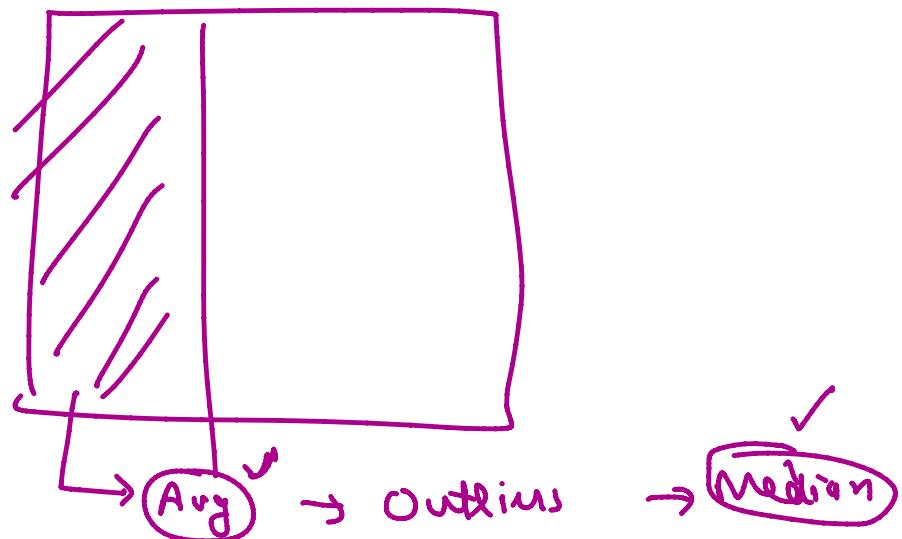
168 cm
(if no outliers)
438 cm
(outliers)

Median

171 cm
(no outliers)

171 cm
(outliers)

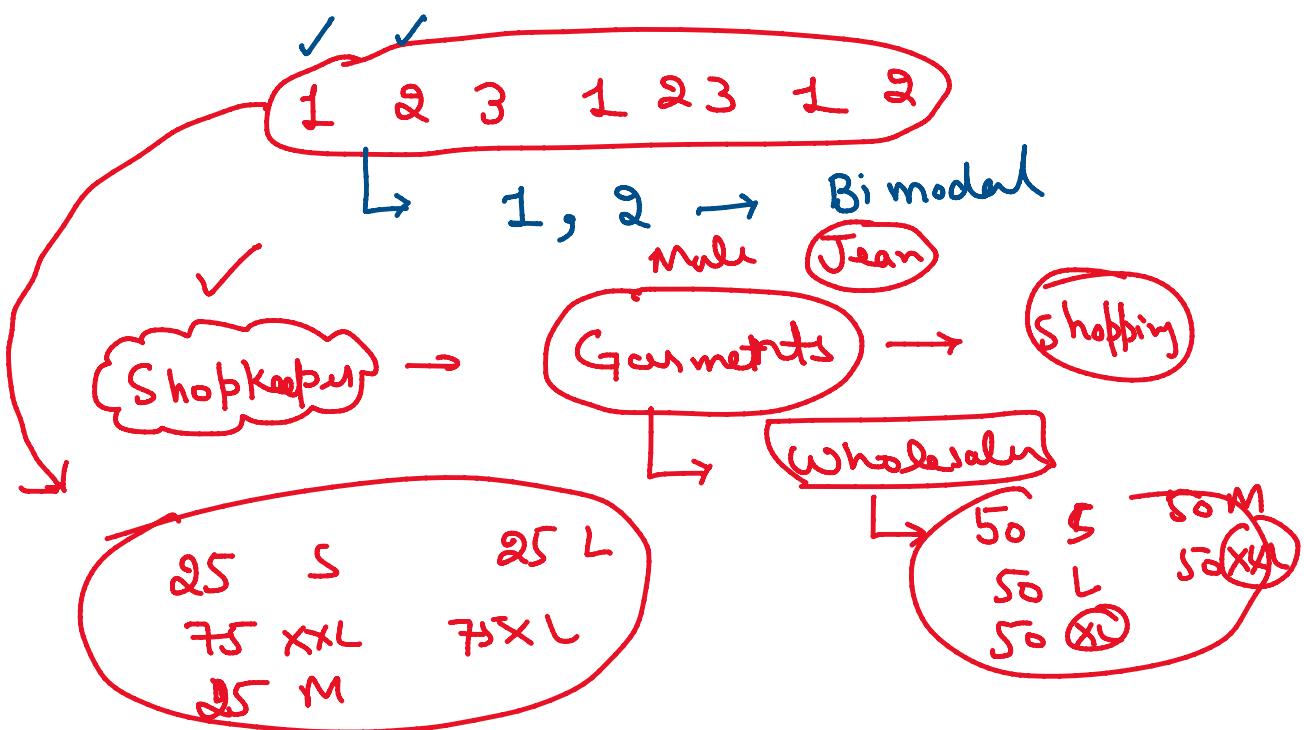
Time ↗



→ Median Package

C T ↴

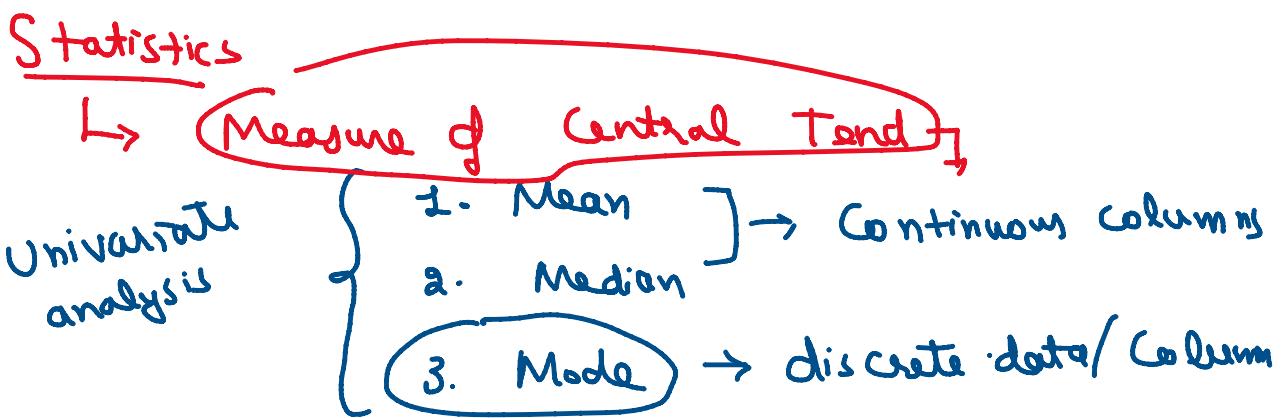
1. Mean
2. Median
3. Mode → Maximum freq. value.



5:18 to 5:23

BREAK

Statistics



Mean → Pick if no outliers exists. Fast to compute

Median → Pick if outliers exists
More time complexity

→ Measure of Spread ↴

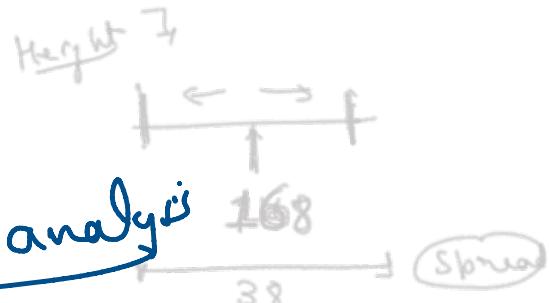
1. Range
2. Variance
3. Std dev

Height ↴

160, 180, 146, 162, 184, 171

1880

Univariate analysis



↳ Range → max - min

$$\rightarrow 184 - 146 \Rightarrow 38$$

with outlier → 1800 - 146 →

171 & 1654

Univariate

2. Variance ↴

σ^2

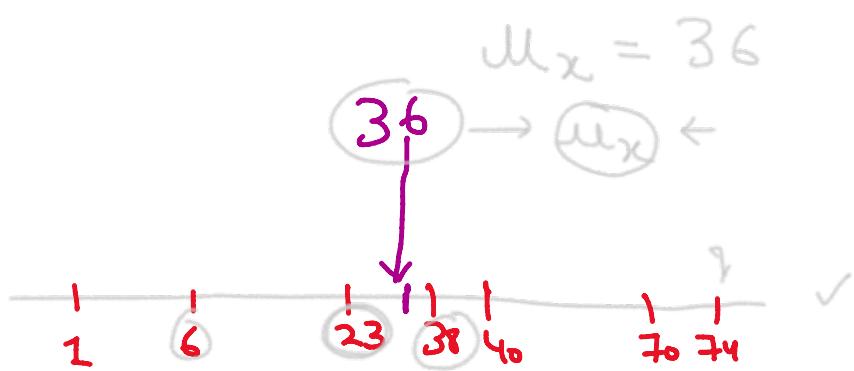
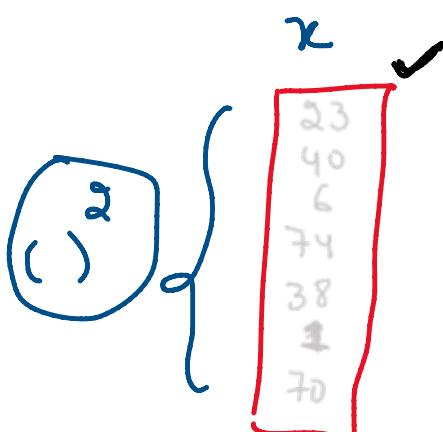
→ Variance measures how far a dataset is spread out around the central

is spreaded out around the central
tendency

$$H \rightarrow \mu_H$$

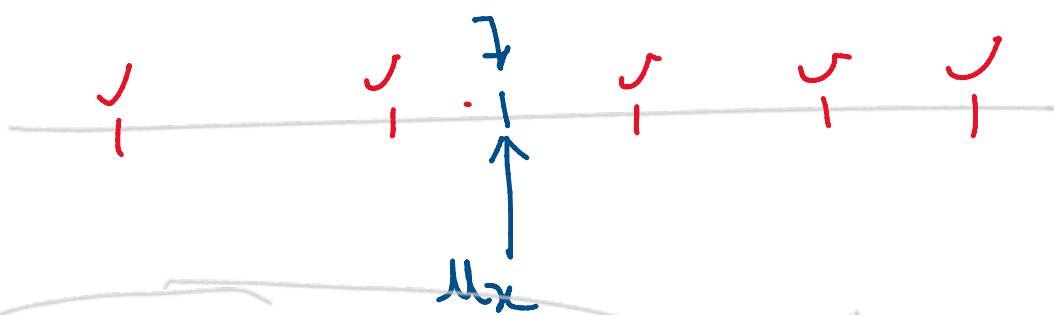
$$\left\{ \begin{array}{|c|c|} \hline H & \mu_H = 160 \\ \hline 160 & \\ 160 & \\ 160 & \\ 160 & \\ 160 & \\ 160 & \\ \hline \end{array} \right. \rightarrow \sigma^2 = 0$$

$\rightarrow \left\{ \begin{array}{l} \text{No variability or zero variance if all the} \\ \text{values are same} \end{array} \right\}$



$$\left\{ (36-0)^2 + (36-0)^2 + (36-23)^2 + (36-38)^2 + (36-40)^2 + (36-70)^2 + (36-74)^2 \right\}$$

\rightarrow Measure how far dataset is spreaded out around mean.



$$\sigma^2 = \sum_{i=1}^n (\mu_x - x_i)^2$$

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

$\mu_x = 36$, $\sigma^2 = 702$

$$\sqrt{\sigma^2} = \sqrt{702}$$

Std. dev $\rightarrow \sigma = 26.49$

- Measure of Central Tendency
 - Mean (impacted by outliers)
 - Median
 - Mode

 - Measure of Spread
 - Range
 - Variance σ^2
 - Std. deviation σ
- Univariate Analysis Statistical tareeke

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

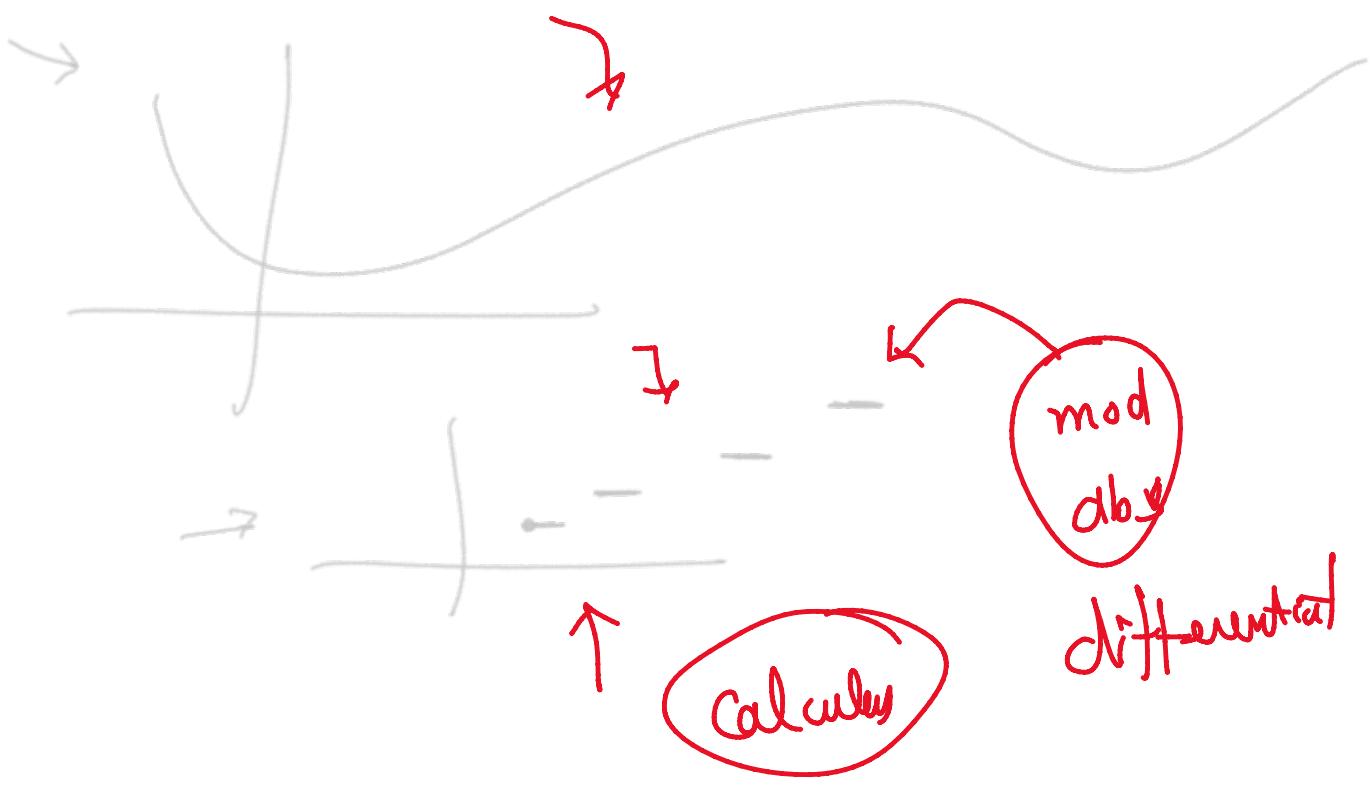
tve OR -ve

Variance \rightarrow Average Squared Distance

Variance → Average Squared Distance
of each point from mean

~~Std dev~~

Continuity & Discontinuous func



Starting @ 4:05 PM

TOPICS ↴

→ Statistics (Part-2)

1. Measure of Central Tendency

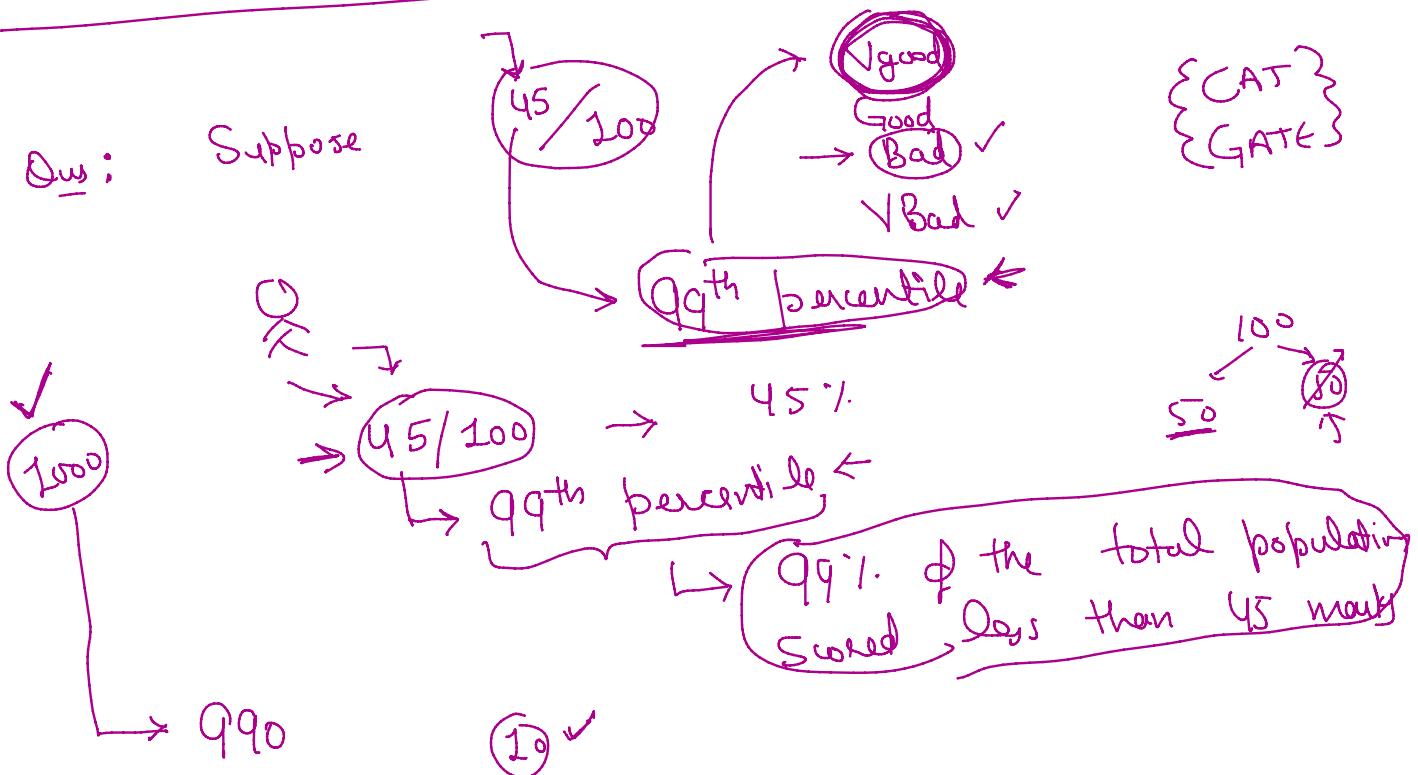
- Mean (Outliers impact)
- Median ✓ (Time complexity)
- Mode ✓

Univariate Analysis

2. Measure of Spread

- Range ✓
- Variance } (outliers impact)
- Std. dev. ✓
- IQR ✓ (Time complexity)

Ques: Suppose

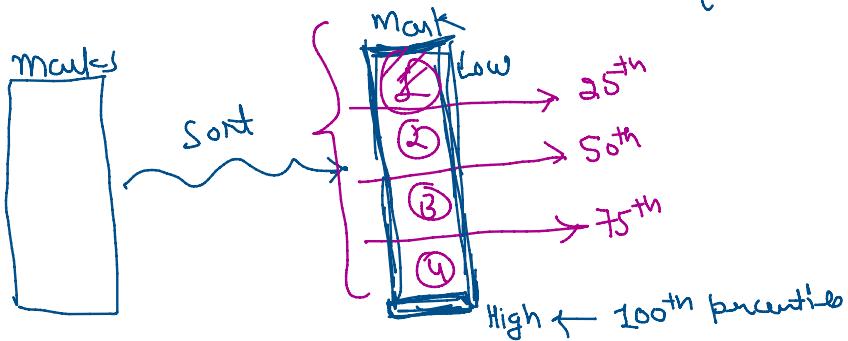


Percentile tells us the rank

90th percentile
90% of the total population scored less than 90

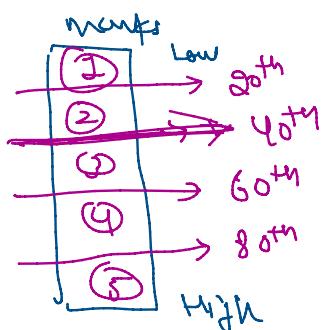
Frequently used Percentiles
1... Divides the data into 4

Frequently used Percentiles → Divides the data into 4 parts

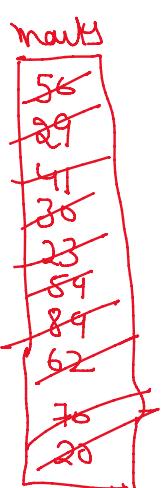


Quartile → 25th, 50th, 75th

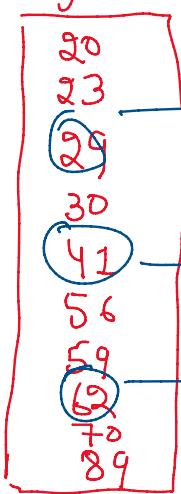
→ Quintiles → Divides the data into 5 parts



Quintiles → 20th, 40th, 60th, 80th



Sort



Quartiles

25th → Q₁

50th percentile → Q₂

75th → Q₃

→ IQR = (Inter Quartile Range)

(Inter

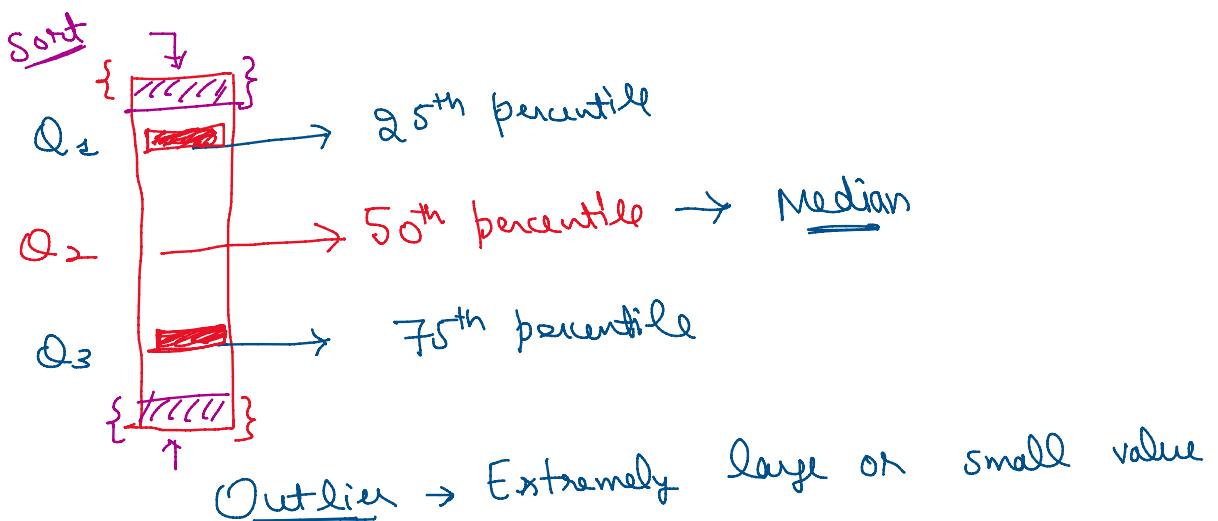
Quartile Range)

Max - min

$$IQR = 75^{\text{th}} \text{ value} - 25^{\text{th}} \text{ value}$$

$$IQR = Q_3 - Q_1$$

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



Univariate Analysis

One Variable at a time

Bivariate Analysis

Two " " " "

Height	Weight	Gender

1. —

2. —

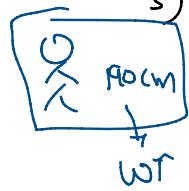
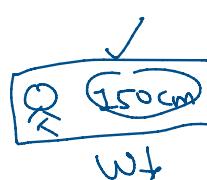
3. Measures of Relationship

(Bivariate Analysis)



Covariance

Correlation



(Covariance)

H	W	G

H↑
H↓

⇒ w↑ OR wt
⇒ w↑ OR wt

$H \uparrow$

$\Rightarrow w \uparrow$ or $w \downarrow$

$H \uparrow \Rightarrow w \uparrow$

$H \uparrow \Rightarrow w \uparrow$

$$\left\{ \begin{array}{l} \text{COV}(H, W) = \\ \text{OR} \end{array} \right. \begin{array}{l} \xrightarrow{\text{fve}} H \propto W \\ \xrightarrow{-\text{ve}} H \propto \frac{1}{W} \end{array} \right\}$$

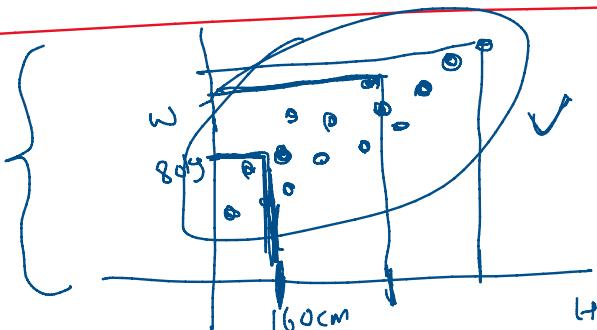
Co Variance

$$\text{Variance} \rightarrow \frac{1}{n} \sum_{i=1}^n \{(u_H - h_i)(u_W - w_i)\}$$

$$\text{COV}(H, W) \rightarrow \frac{1}{n} \sum_{i=1}^n (u_H - h_i) * (u_W - w_i)$$

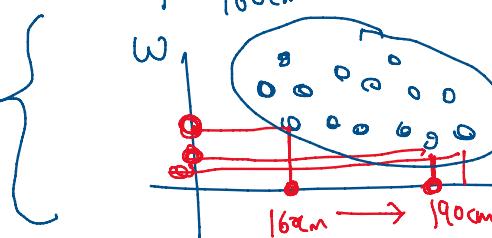
$\text{COV}(H, W) = +\text{ve value} \Rightarrow \overbrace{H \uparrow \Rightarrow w \uparrow}$

$\text{COV}(H, W) = -\text{ve value} \Rightarrow \overbrace{H \uparrow \Rightarrow w \downarrow}$



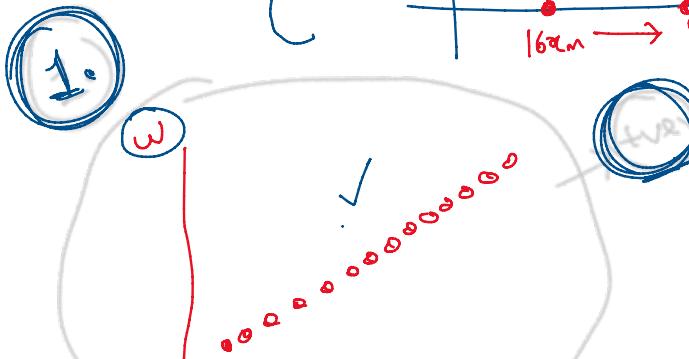
$H \uparrow \Rightarrow W \uparrow$

$\text{COV}(H, W) = +\text{ve value}$



$H \uparrow \Rightarrow W \downarrow$

$\text{COV}(H, W) = -\text{ve value}$

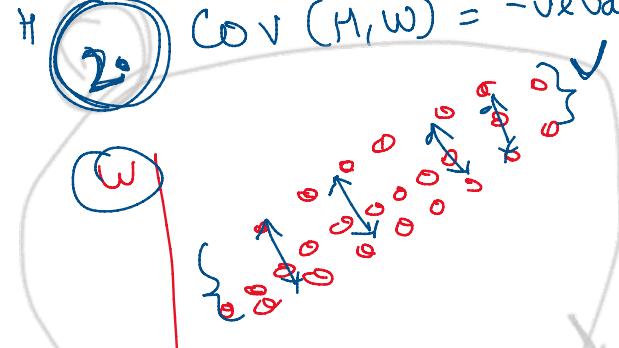


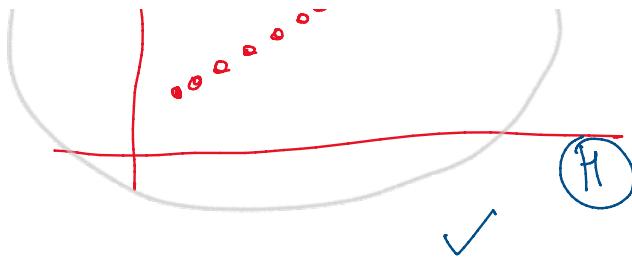
$1.$

W

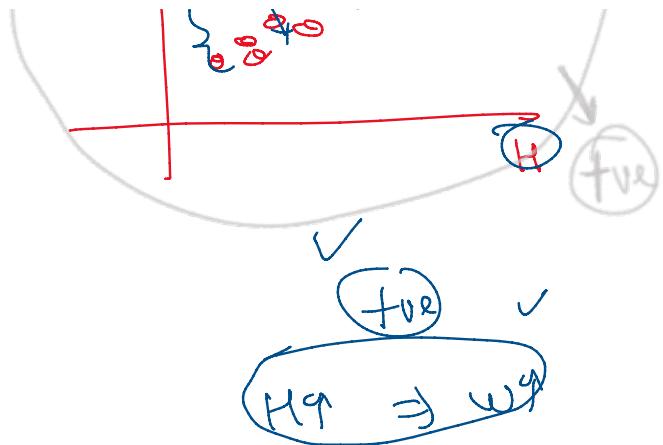
$2.$

W





$$\checkmark \quad H \uparrow \Rightarrow W \uparrow \quad \checkmark$$



$$\checkmark \quad H \uparrow \Rightarrow W \downarrow \quad \checkmark$$

\rightarrow **CORRELATION**

Pearson

Correlation ✓
Spearman Rank Correlation

Pearson Coeff

$$f_{H,W}$$

$$\left\{ -1 \leq f_{H,W} \leq +1 \right\}$$

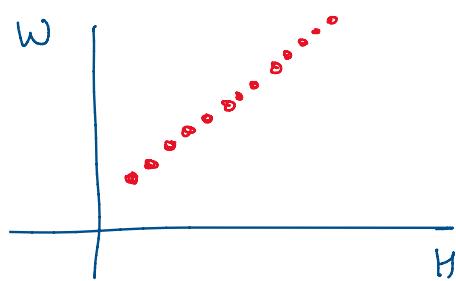
$$f_{H,W} = \frac{\text{COV}(H, W)}{\sigma_H * \sigma_W}$$

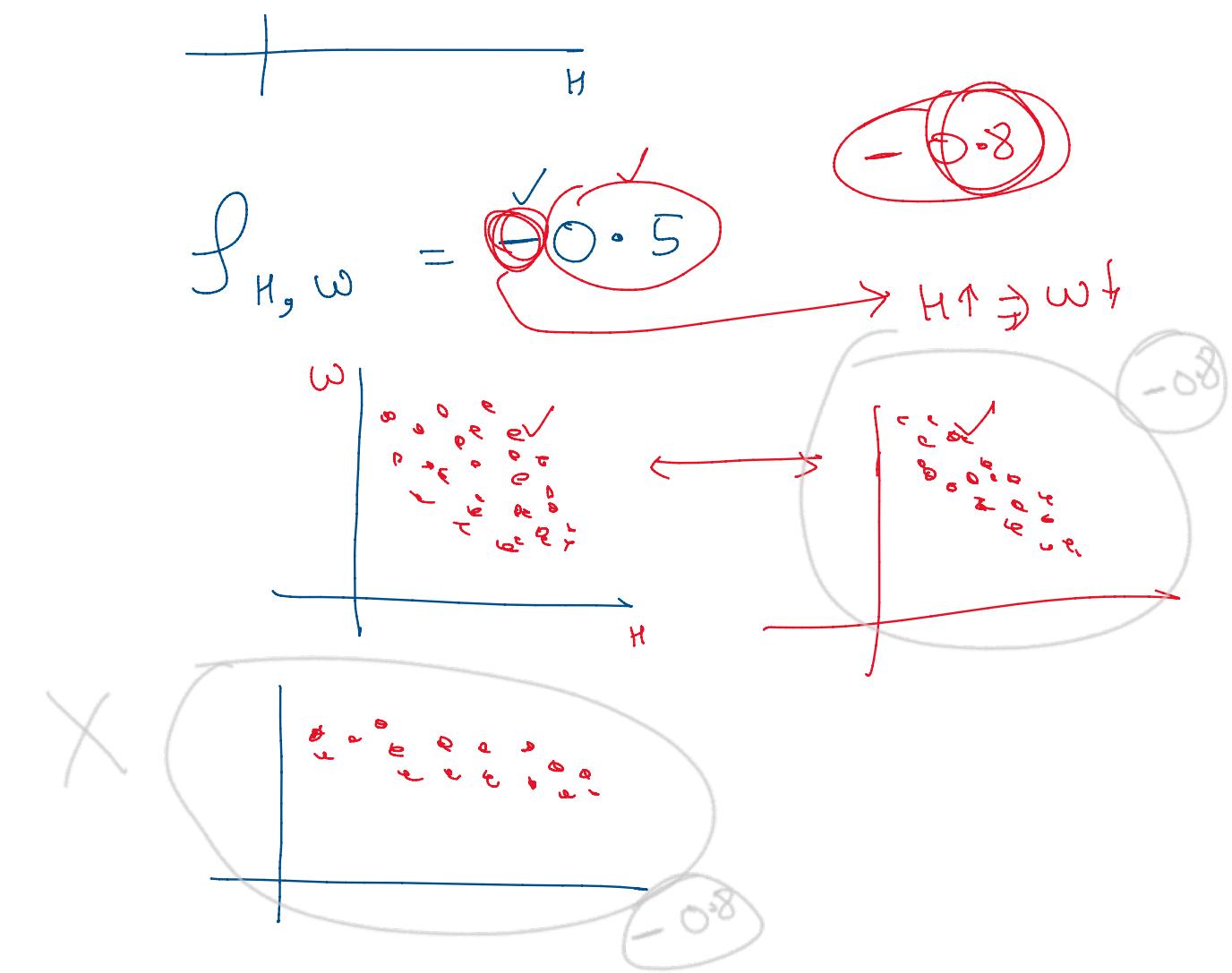
Pearson
 \rightarrow $H \uparrow \rightarrow W \uparrow$ or $W \uparrow$ ← Covariance
 \rightarrow Dependency ✓

$$f_{H,W} = \oplus \boxed{1}$$

$$\oplus \rightarrow H \propto W$$

$$f_H$$

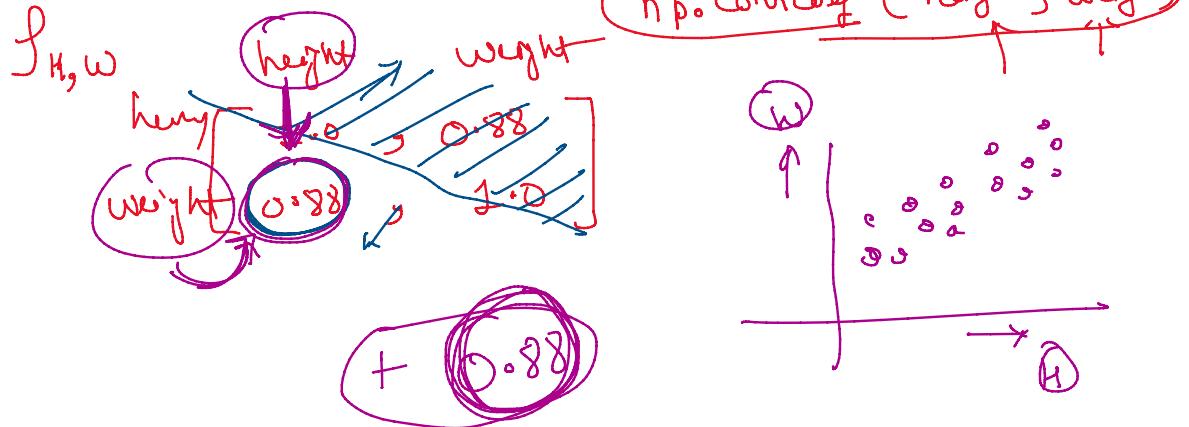




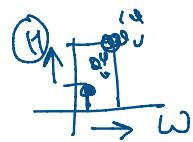
05:03 → 05:08 → 5 min

1. Measures of Central Tendency
 - ↳ Mean → (Outlier Impact)
 - ↳ Median
 - ↳ Mode
 2. Measures of Spread
 - ↳ Range
 - ↳ Var & Std → (Outlier Impact)
 - ↳ IQR
 3. Measures of Relationships
 - ↳ Covariance
 - ↳ Correlation (Pearson) f
- Univariate Statistical Analysis
- Bivariate Analysis

→ Covariance
 → Correlation (Pearson) (f)
 $\boxed{x \uparrow \Rightarrow y \uparrow \Rightarrow T \downarrow}$
 Dependent



$$-1 \leq f_{H,W} \leq +1$$



$f_{H,W} = 0 \Rightarrow$ No relationship b/w H, W
 $f_{H,W} = \pm 1 \Rightarrow H \uparrow \Rightarrow W \uparrow$ & H, W are totally dependent

$f_{H,W} = \frac{1}{\sqrt{2}} \Rightarrow H \uparrow \Rightarrow W \uparrow$ & a "u" shape

$$f_{H,W} = \frac{\text{cov}(H,W)}{\sigma_H * \sigma_W}$$

$\rightarrow 1$

\rightarrow Numpy

\rightarrow Open source free