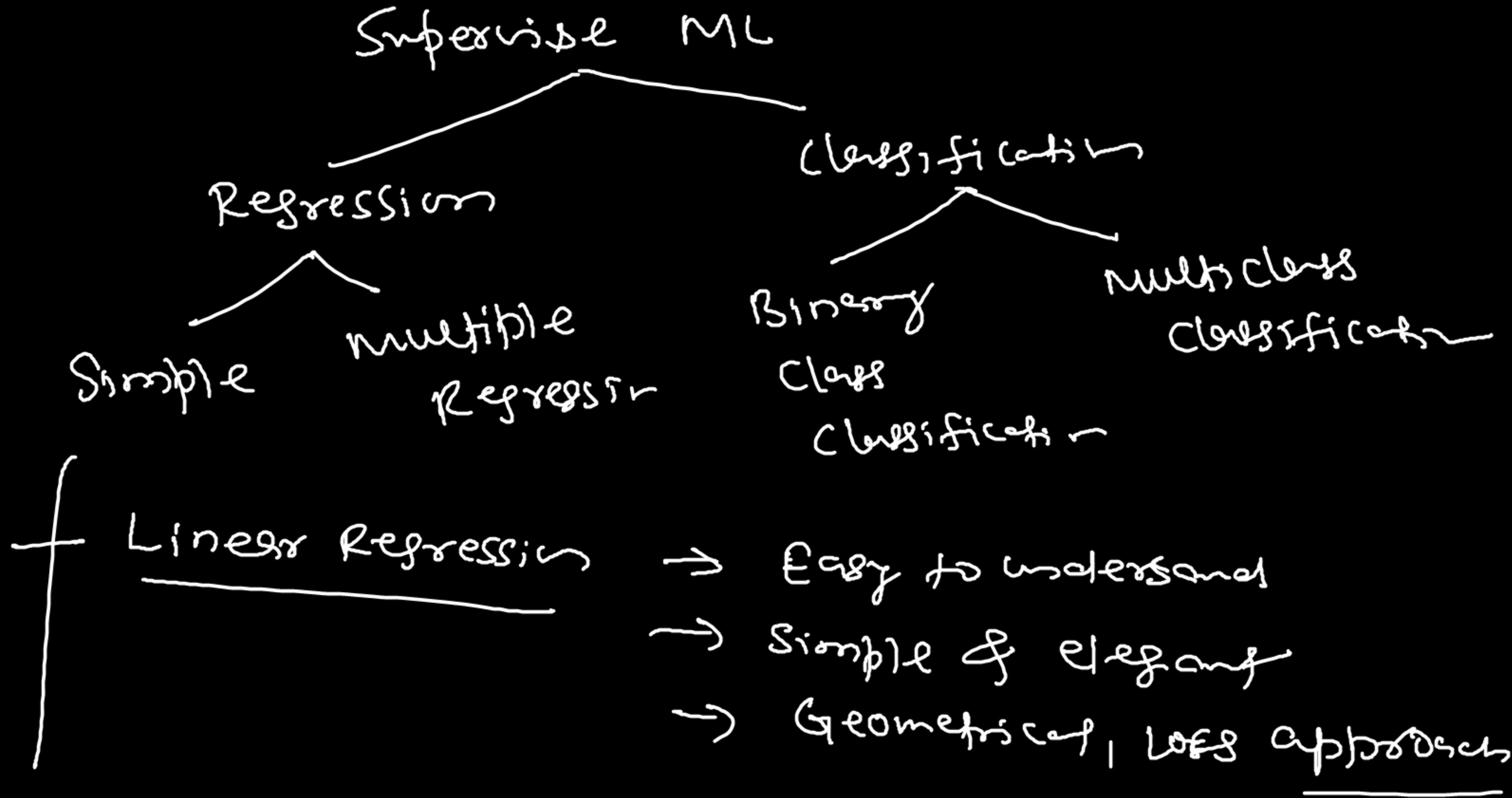


You are screen sharing

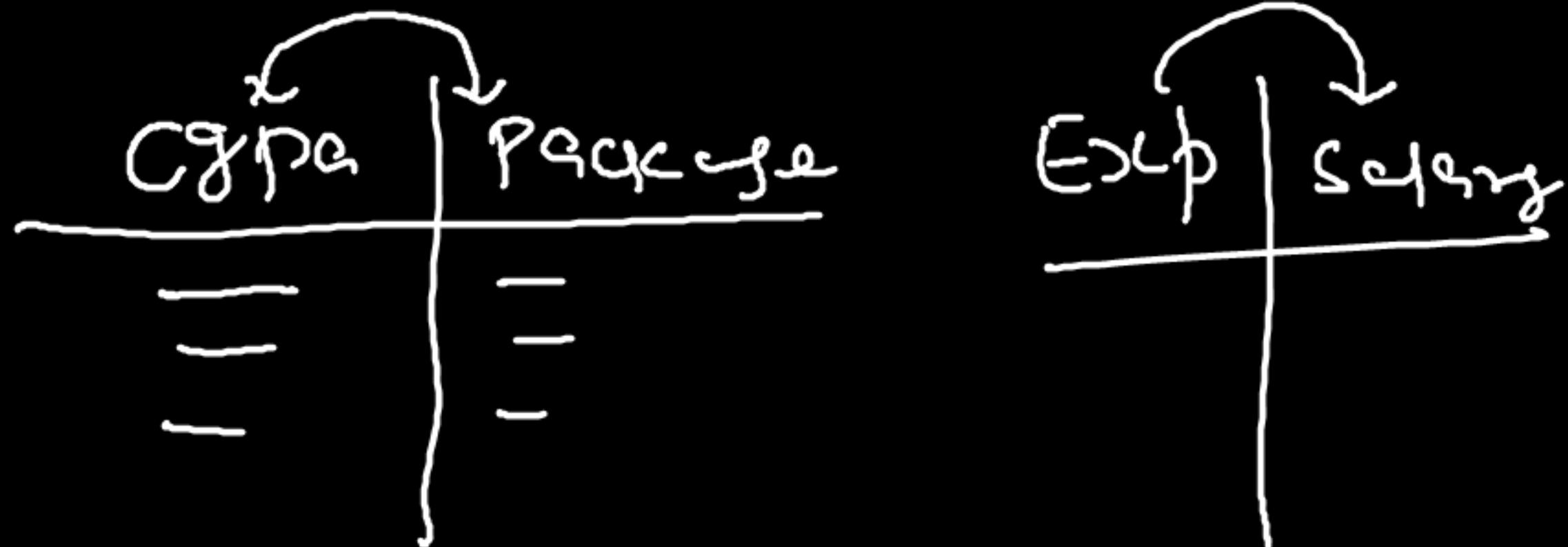
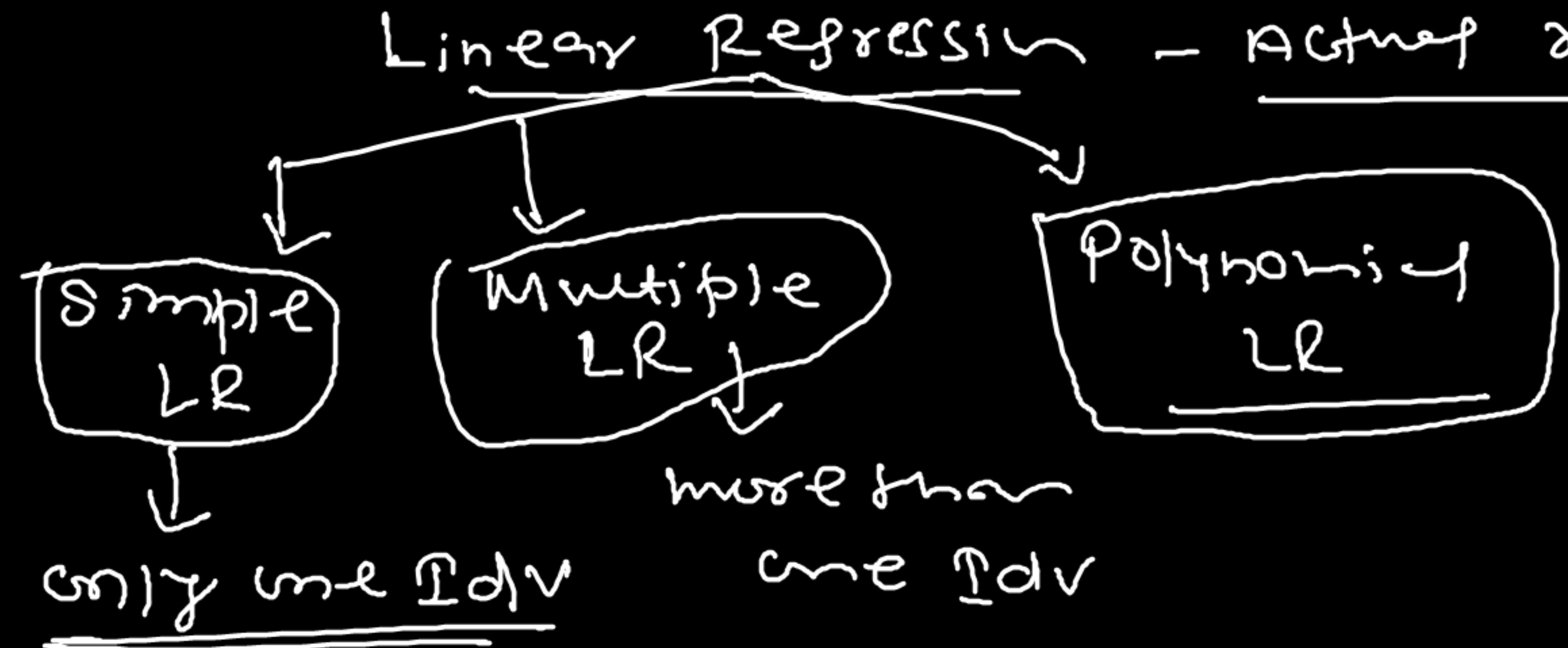
Stop Share



You are screen sharing

Stop Share

Linear Regression - Actual regression problem

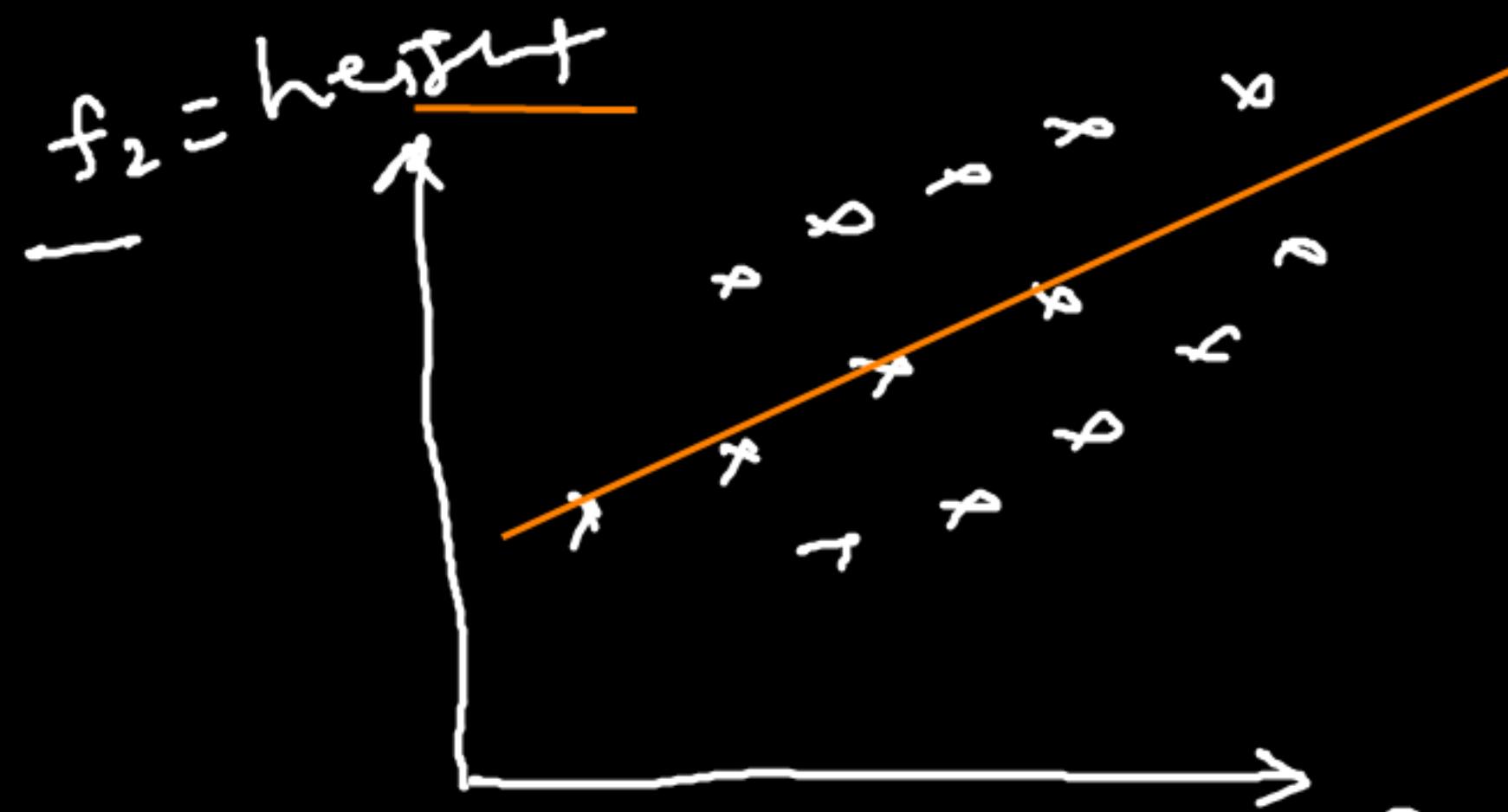


You are screen sharing

Stop Share

Geometric Intuition of Linear Regression

e.g.: Predict height $\in \mathbb{R}$ given weight, Gender, ethnicity, haircolor etc.



Task:— find a line/plane that fits the given data

$$\text{height} = w_0 + w_1 * \text{weight}$$

$y = mx + c$

$Ax_1 + By_1 + C = 0$ — General eqn.

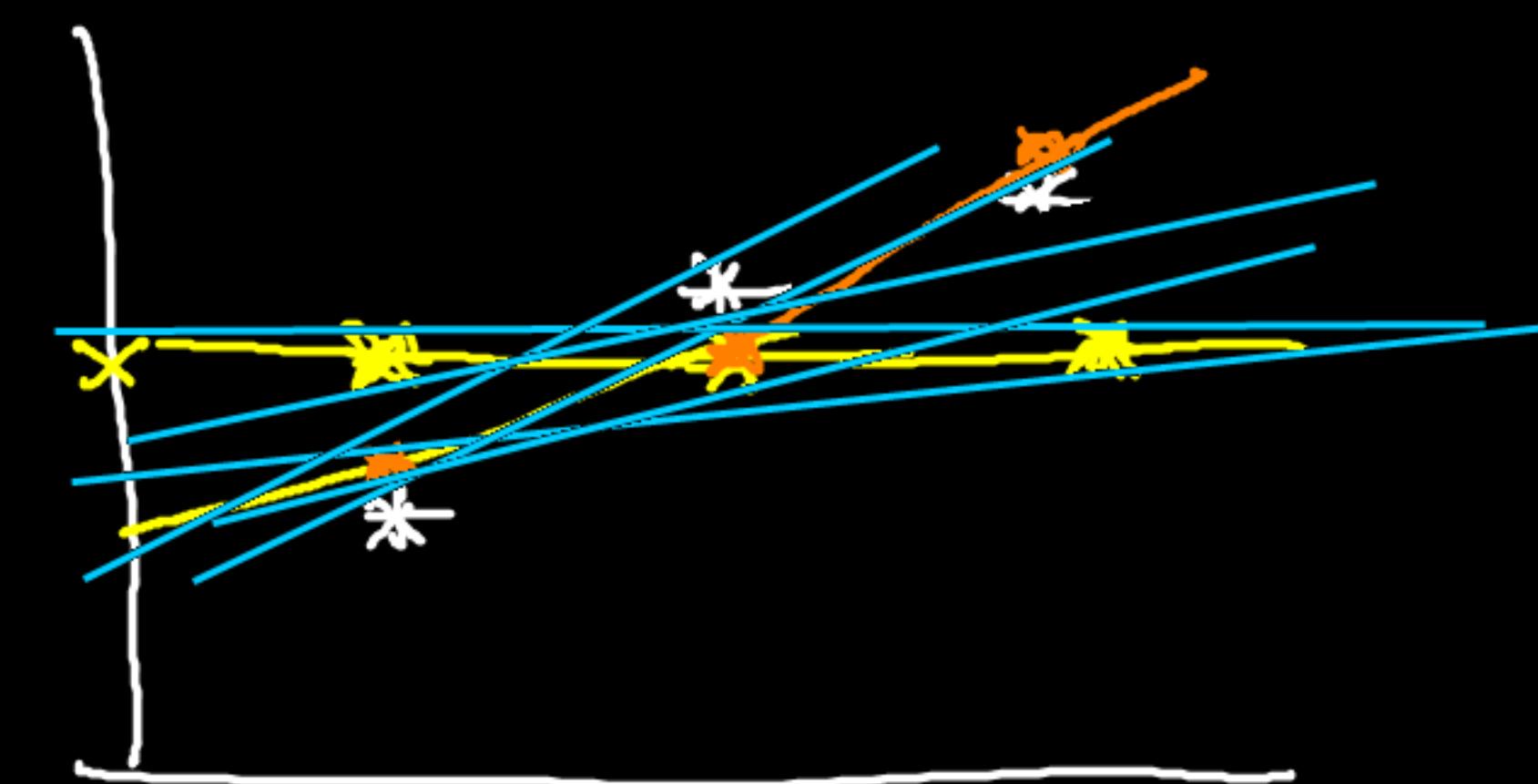
$$Ax_i + By_i + \epsilon$$

$$\begin{cases} x_i, y_i = \text{var} \\ A, B = \text{co-eff} \end{cases}$$

$$C = \text{Int} = \text{Const}$$

<u>(x)</u>	<u>(y) - predict</u>	
Exp	Salary	Avg
2	10800	12500
3	13000	12500
4	14500	12500
✓ 0	?	
✓ 10	?	
✓ 5	?	

$$\text{Salary} = \underline{\text{Intercept}}(w_0) + \text{slope}(w_1) * \text{Exp}$$



y slope = ?

Intercept = ?

$$\text{Avg Salary} = 10 + 13 + 1f.$$

- eqn.

$$\text{Salary} = ?$$

$$Y = (w_0) + w_1 * X$$

$$\underline{w_0} + \underline{w_1} = ?$$

$$\checkmark \text{Exp} = \text{Given}$$

$$\text{Slope/Coefficient} = \frac{\sum_{i=1}^n \underline{\text{cov}(x)} * \underline{\text{cov}(y)}}{\sum_{i=1}^n \underline{\text{cov}(x)^2}} = \frac{4500}{2} = 2250$$

$$\underline{\text{cov}(x)} = (x_{\text{actual}} - x_{\text{avg}})$$

$$\text{cov}(y) = (y_{\text{act}} - y_{\text{avg}})$$

$$C = 12500 - 2250 * 3$$

$$y = mx + c$$

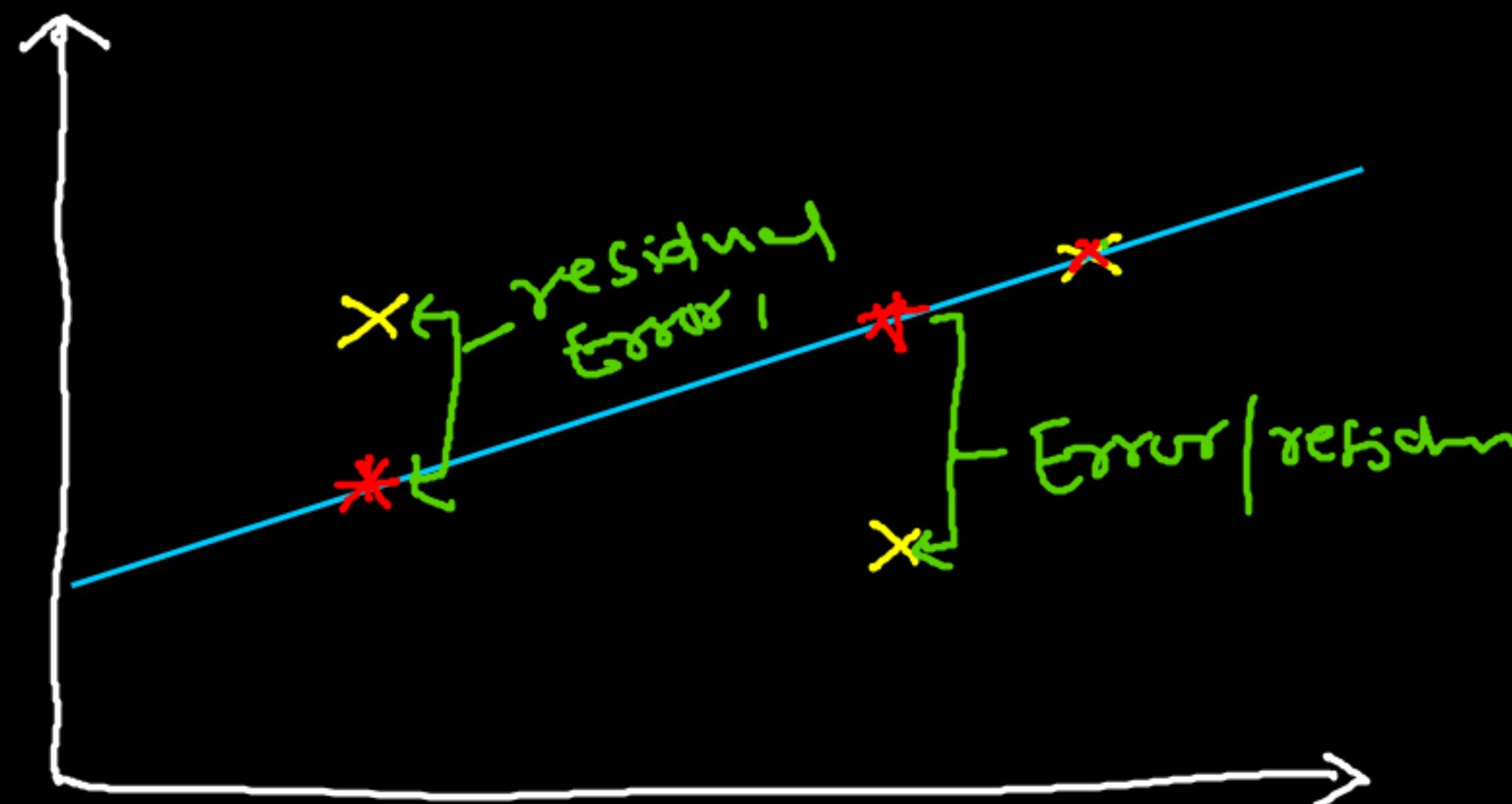
$$c = y - mx$$

$$y_{\text{avg}} = 12500$$

$$x_{\text{avg}} = 3$$

$$m = 2250$$

Find a line/plane that best fits the data points



Yellow color = Actual data

Red color = Prediction

Error 1 = Actual - Prediction

$(\text{Error 1})^2 = \underline{\text{the}} - \underline{\text{the}}$

Error 2 = Actual - Prediction

$(\text{Error 2})^2 = \underline{-\text{the}} - \underline{-\text{the}}$

Error 3 = Actual - Prediction

$(\text{Error})^2 = \underline{\text{zero}} - \underline{\text{zero}}$

To remove errors, we have to use square

Gradient

$$\begin{cases} |x_1 - x_2| = \text{abs} & \frac{\partial |x|}{\partial x} = \text{Can't find it} \\ (x_1 - x_2)^2 = 59 & \frac{\partial (x^2)}{\partial x} = 2x \end{cases}$$

Gradient descent

$$\text{new weight} = \text{old weight} - \gamma \times \frac{\partial L}{\partial \text{weight}}$$

hyperparameter
learning rate
(0 to 10)

head slope

$$w_0, w_1$$

Int

ΣSlope

\rightarrow epochs : iteration

OIS

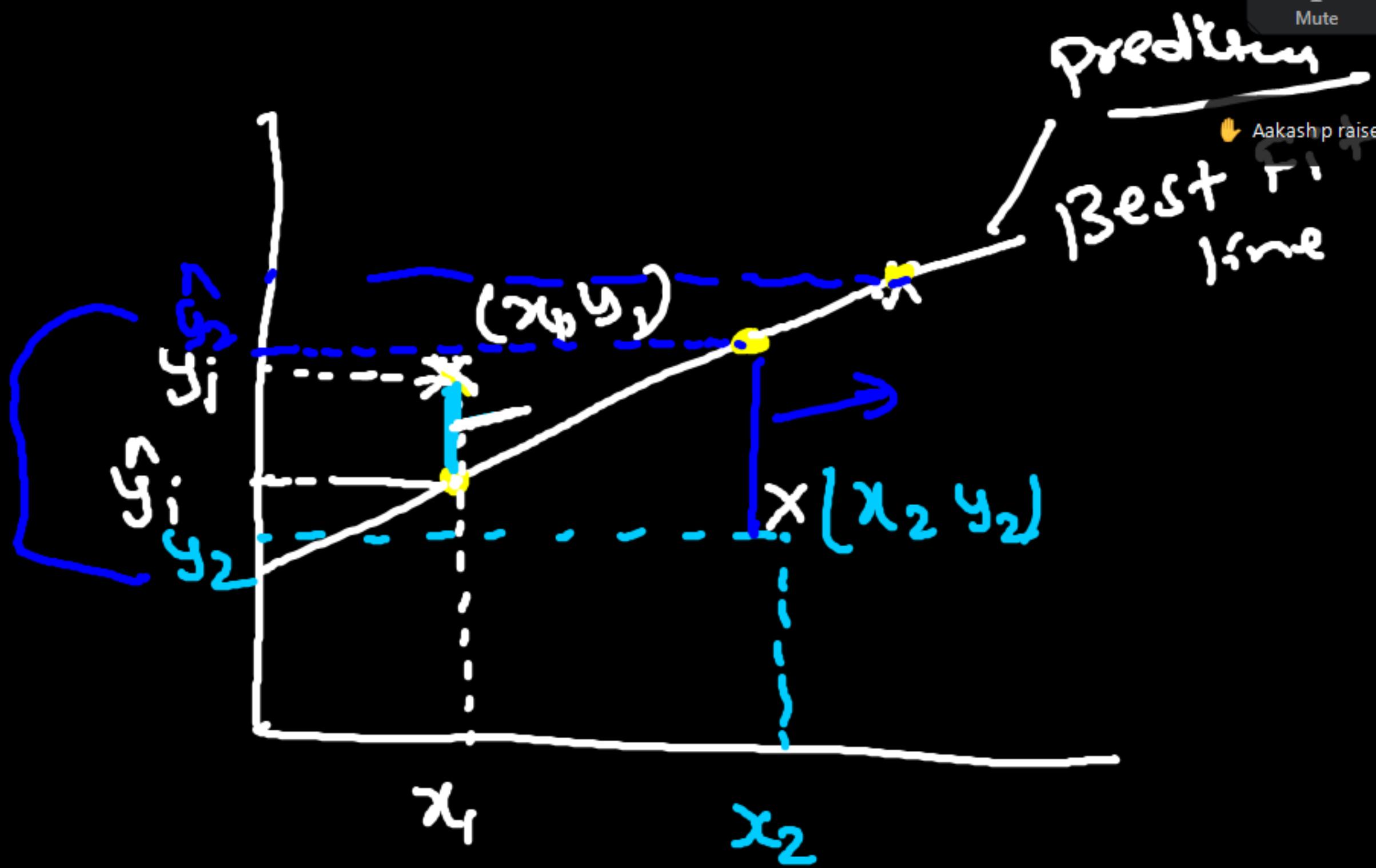
error = 500
 $\sqrt{\text{error}} = 300$ - best fit
min Sum of error

across our training data

2250 / 2000 / 1800 - slope

You are screen sharing

Stop Share



error (residual 1) = $(y_1 - \hat{y}_1)^2$

error (residual 2) = $(y_2 - \hat{y}_2)^2$

error 3 = $(y_3 - \hat{y}_3)^2$

= 0

(7)

Concept = Min Sum of Error = best-fit line

\hat{y}_i

$$Y = w_1 x_1 + w_0$$

Int = w_0

Indv = x_1 = SLR

Int = ~~w_0~~

Slope = w_1, w_2, w_3

$Y = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$) Indv = x_1, x_2, x_3 ~ MLR

You are screen sharing

Linear Regression - Algorithm or OLS = Ordinary Least Square LLS = Linear Least Square

$$w_0, w_1 = \underset{w_0, w_1}{\text{arg min}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

error is min,
best fit line

scaler
(Intercept)
vector
(slope)

Inv = given
find = $w_0 \& w_1$

$$\hat{y}_i = w_1^T x_i + w_0$$

$$y = \underline{w} x + c$$

$$\underline{w} = w_1, c = w_0$$

Dot Product

$$[w_1, w_2, w_3, w_4, \dots, w_n] [x_1, x_2, x_3, \dots, x_m]$$

$$w^T x_i$$

$$\hat{y}_i = w^T x_i + w_0$$

Notation

$$[w_1, w_2, w_3, w_4, \dots, w_n] [x_1, x_2, x_3, \dots, x_m] \\ w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n$$

Optimization method

$$(w_1^*, w_0) = \underset{w^0, w_1^*}{\arg \min} \sum_{i=1}^n \left\{ y_i - (w^T x_i + w_0) \right\}^2$$

SQ-LOSS function

Part 1 - Regularization

LASSO ✓
Ridge ✓

ElasticNet ✓

Picked scope
+ case study
Multiple
Linear Reg

Part 2 - Gradient descent

Part 3 :- Assumption

- Linearity
- Normality of Residual
- Homoscedasticity
- No Auto Correlation
- No or Little Multicollinearity
- No endogeneity Problem

Linear Reg

OLS / SE-loss -

R-Square, Adj-R-Square

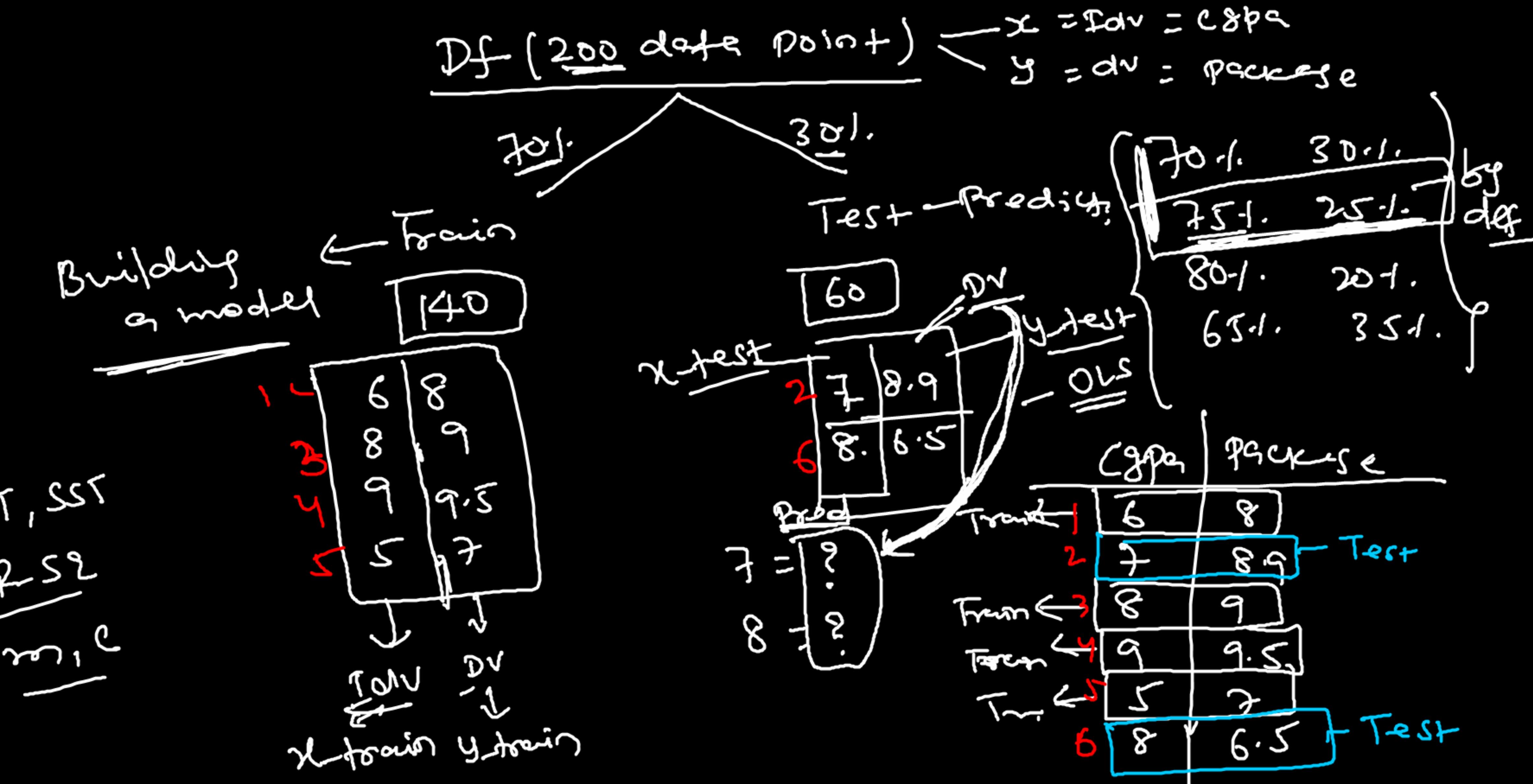
P-value

Life Cycle of model

- ① collect / extract data
- ② Preprocessing of the data
- ③ Divide the dataset into train & test
- ④ EDA - Perform Descriptive analysis
- ⑤ Build Regression model
- ⑥ Estimate regression parameters (Int , slope , R^2)
- ⑦ Perform regression model diagnosis
- ⑧ Validate the model by test data
- ⑨ Deployment on the cloud -

You are screen sharing

Stop Share



ML - Problems

① underfitting Problem | high bias

② overfitting Problem | high variance

training	test	variance	
~90%	~70%	20%	overfitting Prob high variance more than 10% - high var high variance
~75%	~90%	15%	
~50%	~50%	less than 70%	underfitting Problem high bias
~90%	~90	Perfect fit	Perfect fit ✓

Data leakage