

* SLR

- ↳ gradient Descent
- ↳ Evaluation metrics

Scikit Linear Regression

$\left\{ \begin{array}{l} R\text{ square} := \cdot f \cdot Y \text{ by } X \\ \text{adjusted } R\text{-square} := \end{array} \right.$

MSE

R MSE

MAE.

\rightarrow SLR \rightarrow two datasets.

Agenda

- ↳ Assumptions of Linear Regression
- ↳ Multiple linear Regression
- ↳ Polynomial regression.
- ↳ Lasso, ridge, Elastic Net.
- ↳ Cross validation & Hyperparameter tuning
- ↳ Logistic Regression

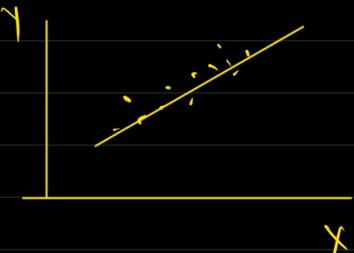
Multiple linear Regression

| Perch # of rooms | Locality | Area of house | Price of house | % of syllables completed | No. of houses Studied | Marks obtained |
|------------------|----------|---------------|----------------|--------------------------|-----------------------|----------------|
| . | . | - | - | ; | 5 | 58 |
| . | . | - | - | ; | ; | , |
| . | . | - | - | ; | ; | , |
| . | . | - | - | ; | ; | , |
| . | . | - | - | ; | ; | , |

* All the information γ is not captured by X , in order to capture more & more information we require more X .

MLR \Rightarrow more the IIV is used.

$$\text{IIV} \Rightarrow \underline{y_{\text{pred}}} = \theta_0 + \theta_1 x_1 \rightarrow$$



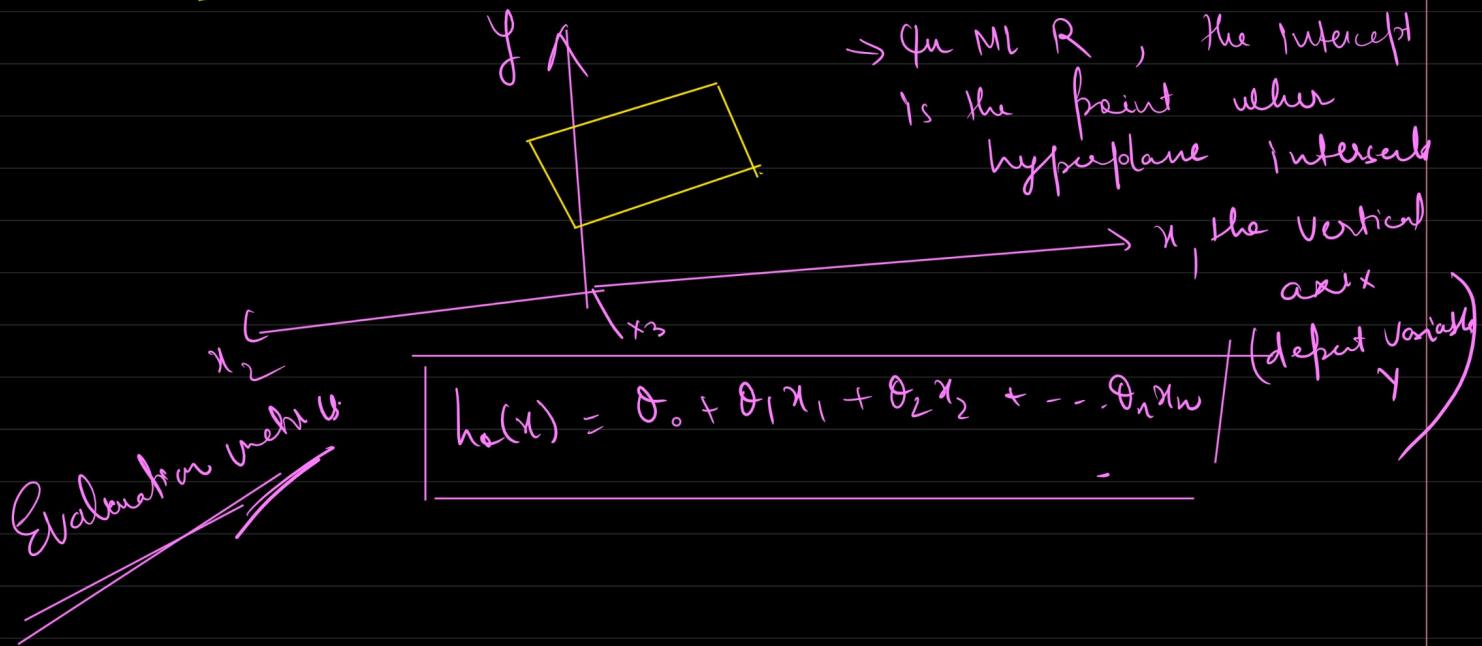
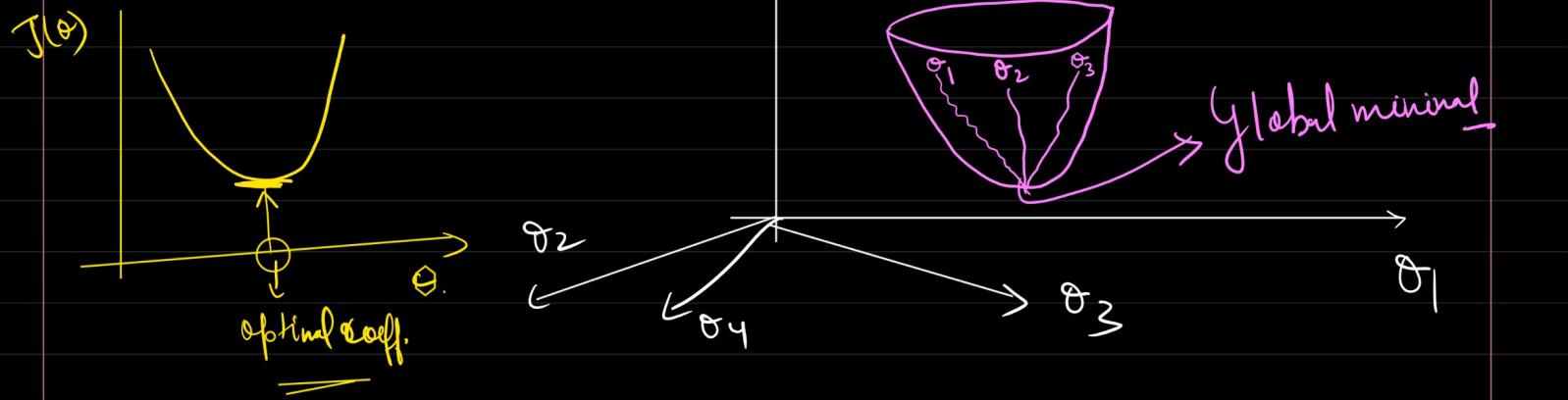
$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 - \dots - \theta_n x_n$$

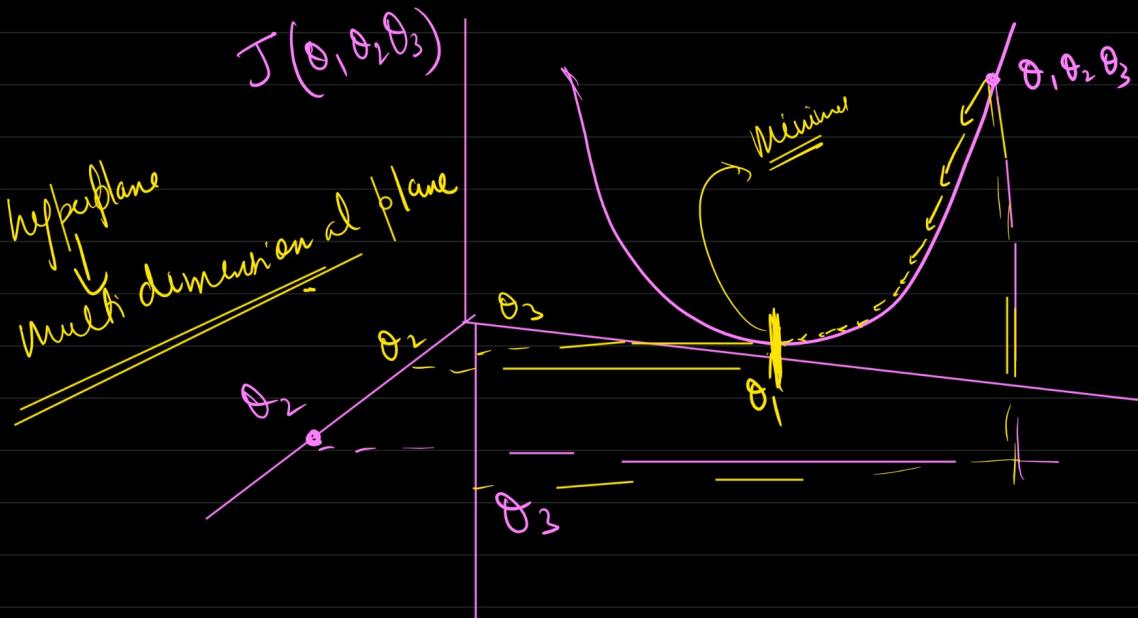
$$CF = \frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - \underline{y_{\text{pred}}})^2$$

$$(y_{\text{act}} - (\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n))^2$$

$\theta_1, \theta_2, \dots, \theta_n$ are optimal coeff.

Earlier in SLR.





$$\theta_j : \theta_j - \eta \frac{\partial C_F}{\partial \theta_j}$$

$$X_{fit} = x_1, x_2, x_3, x_4, x_5$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_0$$

\oplus

of hours studied

-
-
-
-

Marks obtained by a student even if he didn't study

$$y = 50 + 2.5x$$

$$q_{33} \rightarrow q_{50}$$

$\rightarrow \theta_1 \approx$ With 1 unit increase in # of hours studied, the marks increase by 2.5 units.

$$y = 120 + 3.2x$$

\rightarrow with 1 unit increase in x , y increases by 3.2 units

\ominus Sp of car

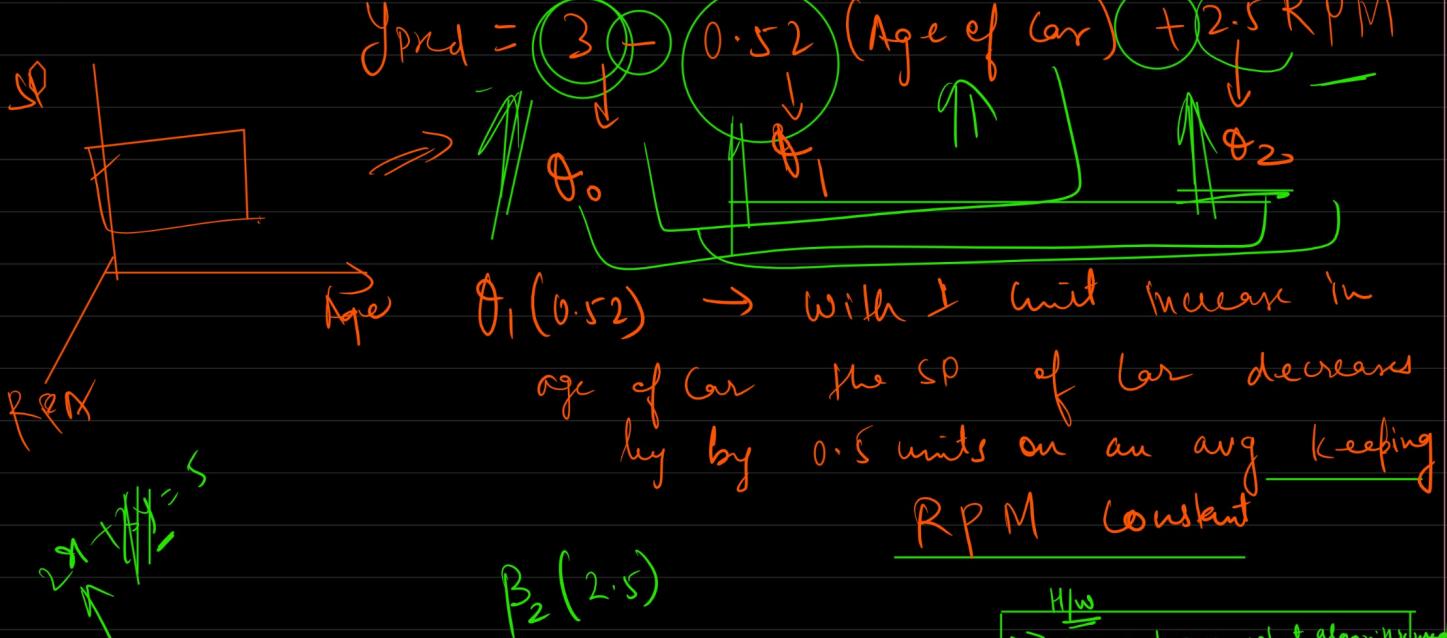
$$y_{pred} = 120 - 3.2(Age)$$

Unit increase in x , y decreases by 3.2 units

interpretation of 1.2 doesn't make sense

$Age = 0$
Sp of car $\rightarrow 120$

MLR



* Polynomial Regression

H/w

- Suchil learn what algorithm used
- Explain | Unexplain Varian in Regress
- Statsmodel
- Feature importance in linear regression

→ Simple Linear Regression, $h_0(x) = \underline{\theta_0 + \theta_1 x}$
for MLR : $h_0(x) = \underline{\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n}$



* Poly nomial regression

Simple LR (1 DV - 1 IV)

① Simple Polynomial regression (1 DV, 1 IV)

Polynomial degree 0, $h_0(x) = \theta_0 x^0$

$$= \theta_0 \quad \Rightarrow \quad (1)$$

Polynomial degree 1: $h_0(x) = \theta_0 x^0 + \theta_1 x^1$

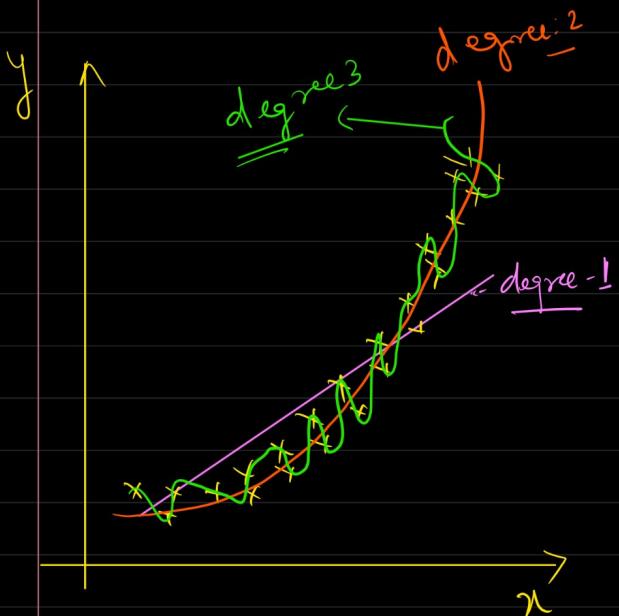
$$\Rightarrow \theta_0 + \theta_1 x^1$$

↓
Simplify LR.

Polynomial degree 2

$$h_0(x) = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2$$

$$h_0(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2$$



* As you increase the degree, you might get an overfitting model

→ Non-linear relationship:

Decision Trees

Polynomial degree 'n' = $h_0(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3 + \dots + \theta_n x^n$

for multiple X (for 3 X)

Polynomial degree 2:

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 +$$

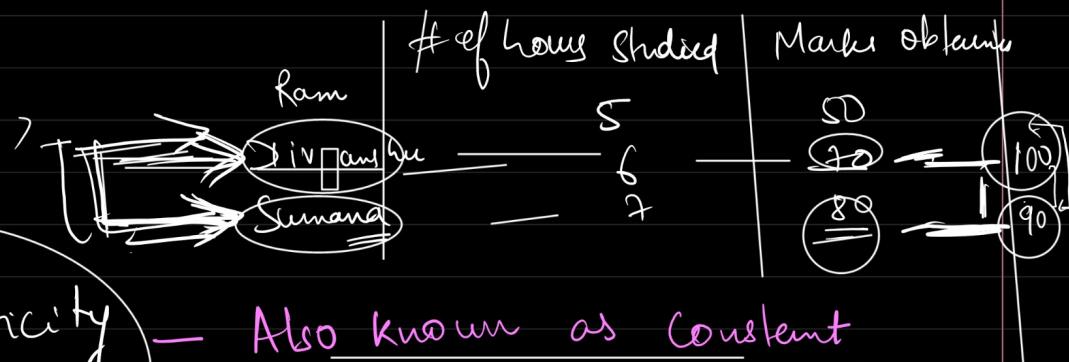
$$\theta_4 x_1^2 + \theta_5 x_2^2 + \theta_6 x_3^2 + \theta_7 x_1 x_2 + \theta_8 x_2 x_3 + \theta_9 x_3 x_1$$

$x_1 x_2$
↓
Cross product
keeping power in mind

$x_1 x_2$ $x_2 x_3$ $x_1 x_3$ x_1 x_2 x_3

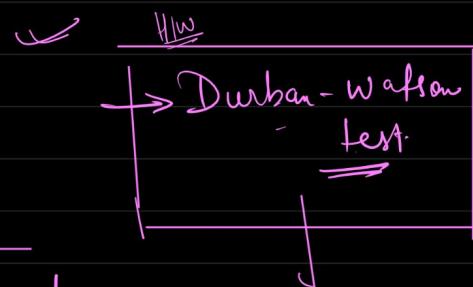
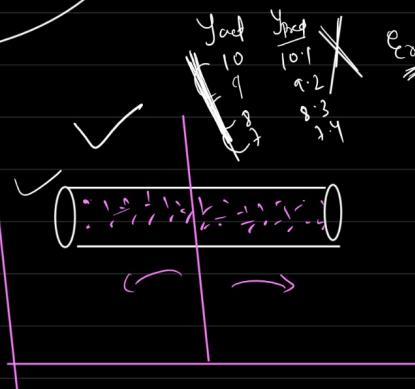
* Assumption of Linear Regression

- ① Linearity → X and Y should have linear relationship.
- ② Independence → Observations (rows) are independent of each others.
→ Errors should be independent

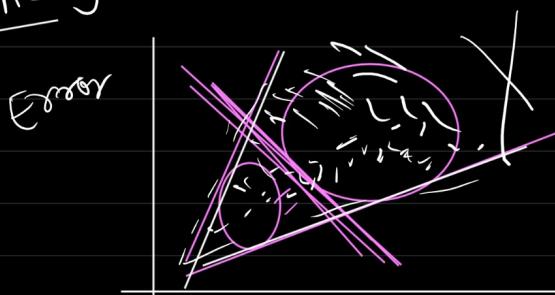


③ Homoscedasticity

- Also known as Constant Variance. The variance of errors are constant.



Heteroscedasticity



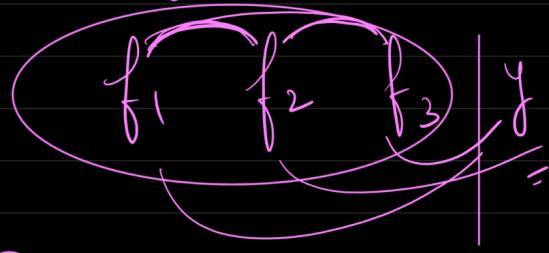
$$\text{Var}A = \text{Var}B$$

Homoscedastic

④ Normality of Error

→ Errors should be normally distributed.

⑤ The feature should not be related or
should have relation



$$f_1 \quad f_2 \quad f_3 \quad y \quad f_1 \rightarrow y$$

$$\left\{ \begin{array}{l} x_1 = x_2 \\ 8x_1 + 2x_2 = y \\ 10x_1 \\ 10x_2 \\ 2x_1 + 8x_2 \end{array} \right.$$

Multi - Col - linearity
many together linear relation

$$\left\{ \begin{array}{l} x_1 - x_2 \\ x_1 - y \end{array} \right\} \rightarrow \text{Correlation}$$

$$\left\{ \begin{array}{l} \cancel{x_1} \approx (x_2, x_3) \\ \cancel{x_1} \approx (x_2, x_3, x_4) \end{array} \right\}$$

Multicollinearity → where a feature exhibits a linear relationship with more than one variables.

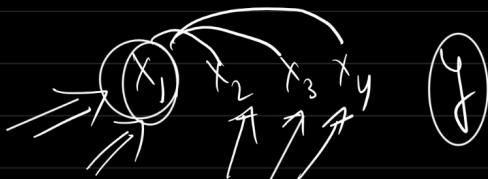
$$\left\{ \begin{array}{l} x_1 \sim x_2 \sim x_3 \\ x_1 \sim x_2 \sim x_4 \\ x_5 \sim x_1 + x_2 + x_3 \end{array} \right.$$

Multicollinearity

What → One feature explained by other features

Why → You can not interpret correctly what is the contribution of each individual feature wrt Y
Severe? → No effect on prediction

~~Problem~~ → Interpretability
 → Computationally expensive - model training

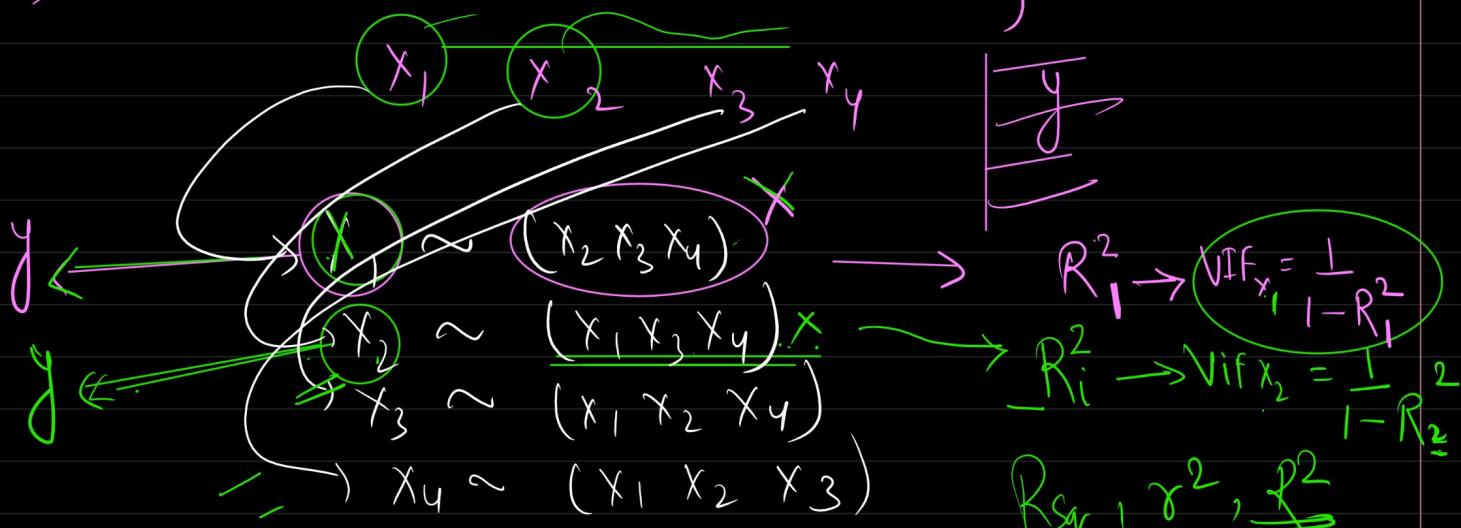


How? → VIF (Variance Inflation factor)
 ↳ a measure of amount of multicollinearity in regression

$$VIF_i = \frac{1}{1 - R_i^2}$$

→ R^2 → % age variation in y explained by x

$$x_1 \sim (x_2, x_3, x_4, \dots, x_n)$$



* Since VIF doesn't impact prediction, you should always ask business team before dropping any feature based on VIF

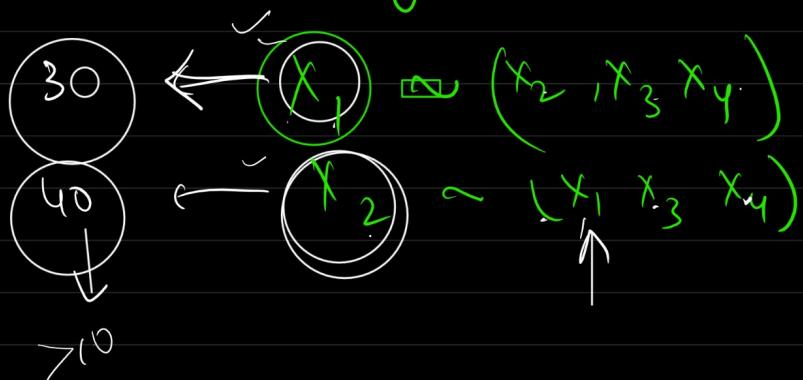
~~x_1~~ $\sim (x_2, x_3, x_4)$ — 200
 $\checkmark \rightarrow x_2 \sim (x_1, x_3, x_4)$ — 5
 $x_3 \sim (x_1, x_2, x_4)$ — 8
 $x_4 \sim (x_1, x_2, x_3)$ — 2

VIF $\sim 0 \text{ to } \infty$

$\text{VIF} > 10$, then drop features one by one.

- 1 Why $\text{VIF} > 10$?
- 2 Drop feature one by one?

(Starting with highest value)



dropping one feature can change VIF of other features

$\text{VIF} > 10$

$\checkmark (x) \sim (x_2, x_3, x_4)$

$$10 = \frac{1}{1 - R^2}$$

$$1 - R^2 = \frac{1}{10}$$

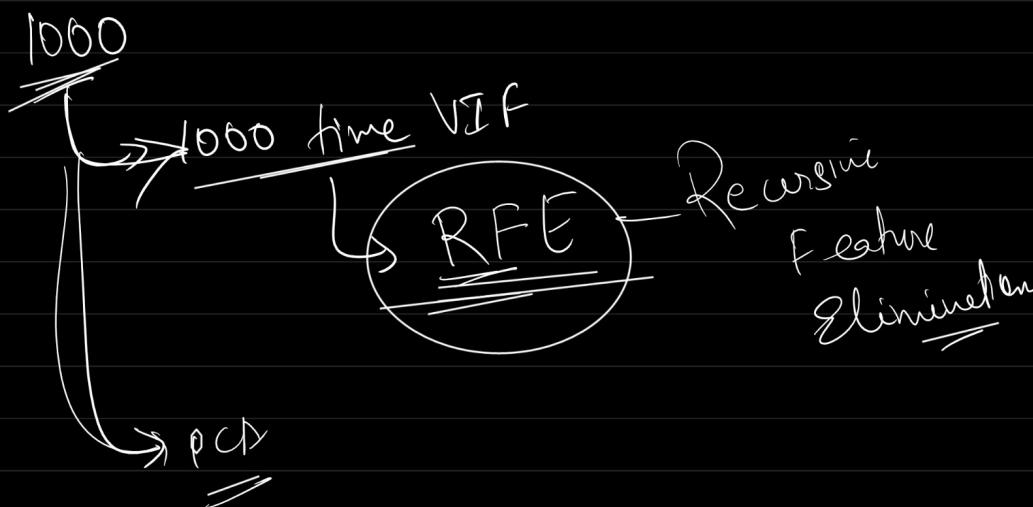
- Only variance in x is explained by x_2, x_3, x_4

$$R^2 = 1 - \frac{1}{10} = \frac{10 - 1}{10} = \frac{9}{10} = 0.9$$

| Feature | VIF |
|---------|-----|
| X_1 | 13 |
| X_2 | 12 |
| X_3 | 8 |
| X_4 | 7 |

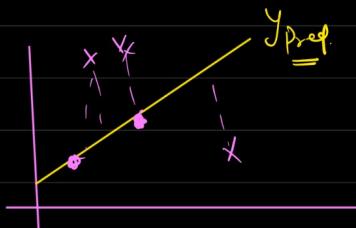
> 10

first drop $X_1 \rightarrow$ Again calculate
VIF \Rightarrow



* Regularisation \Rightarrow To add something | to regularisation

To regularize
 \downarrow
To penalize

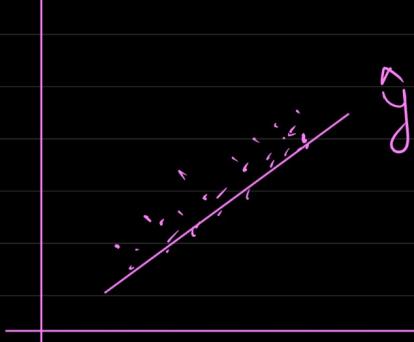


data
train
 \Rightarrow test
 \Rightarrow overfitting

Acc \uparrow
(low bias)

Acc \downarrow
(High Variance)

CF

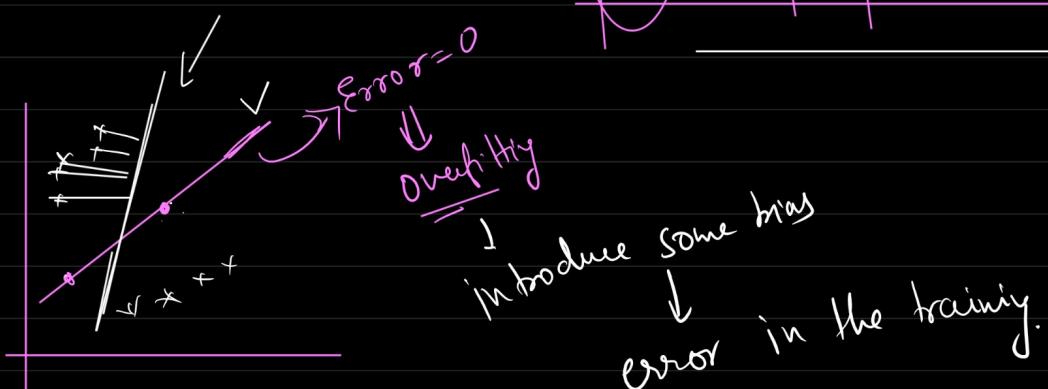
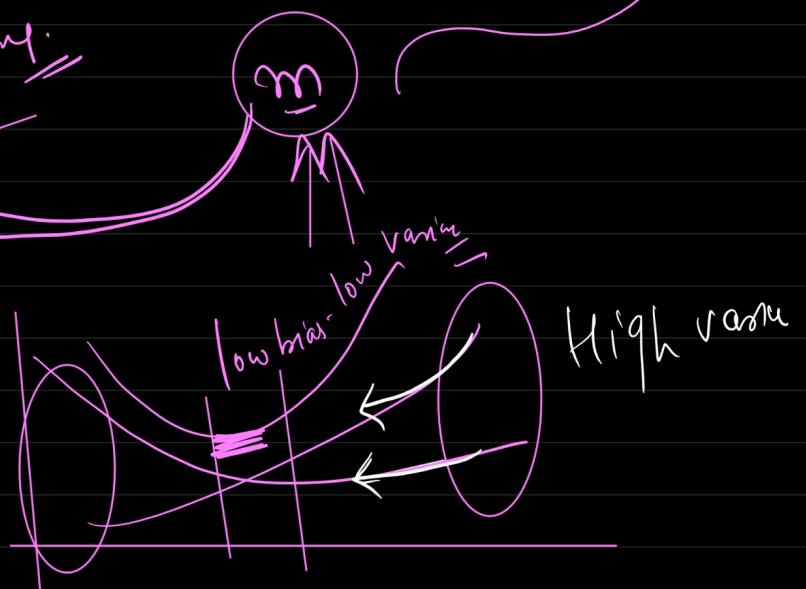
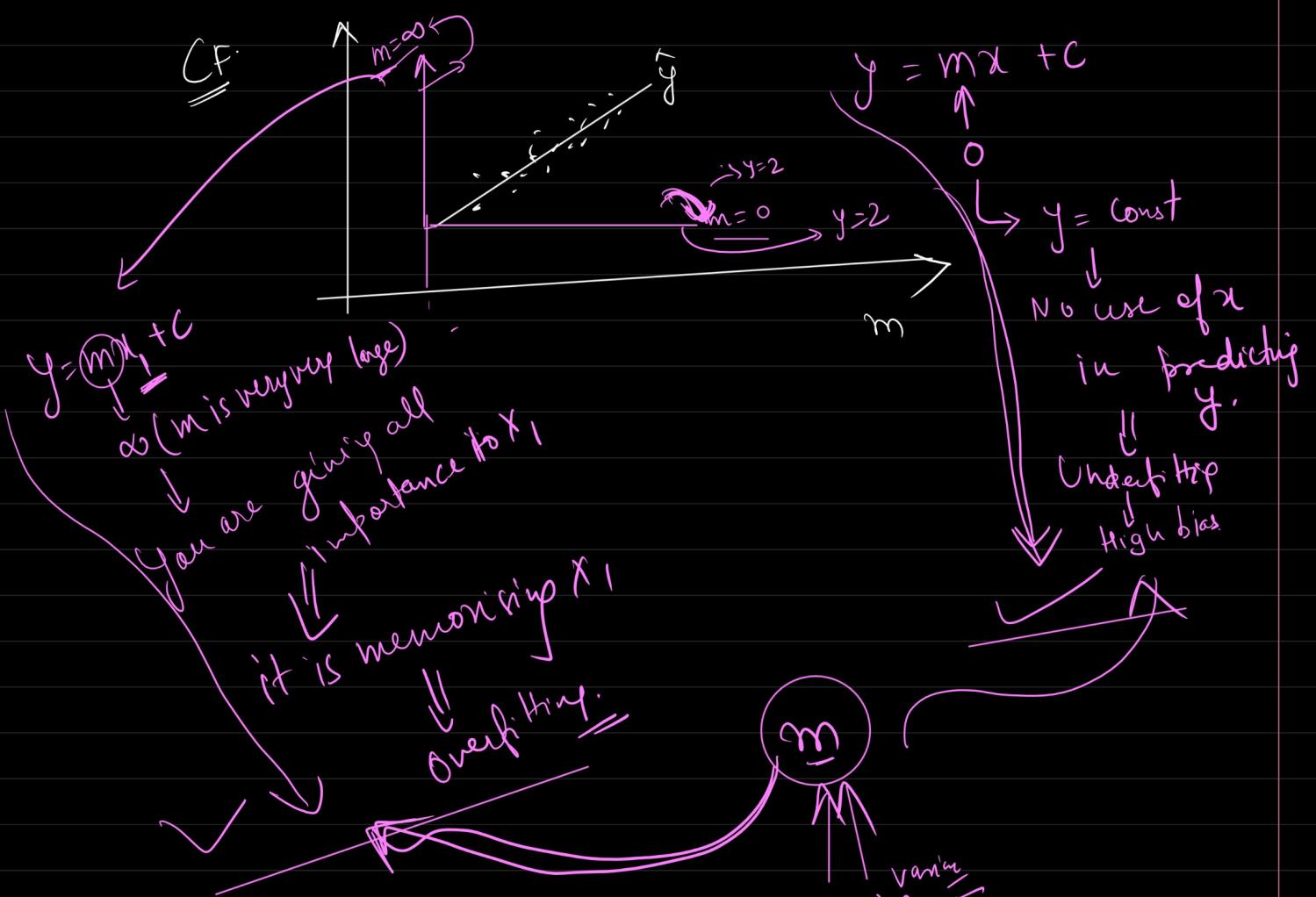


m is very very high \rightarrow
model has memorised
the data \rightarrow it means
overfitting

$$y = 3.2 + 2.5x$$

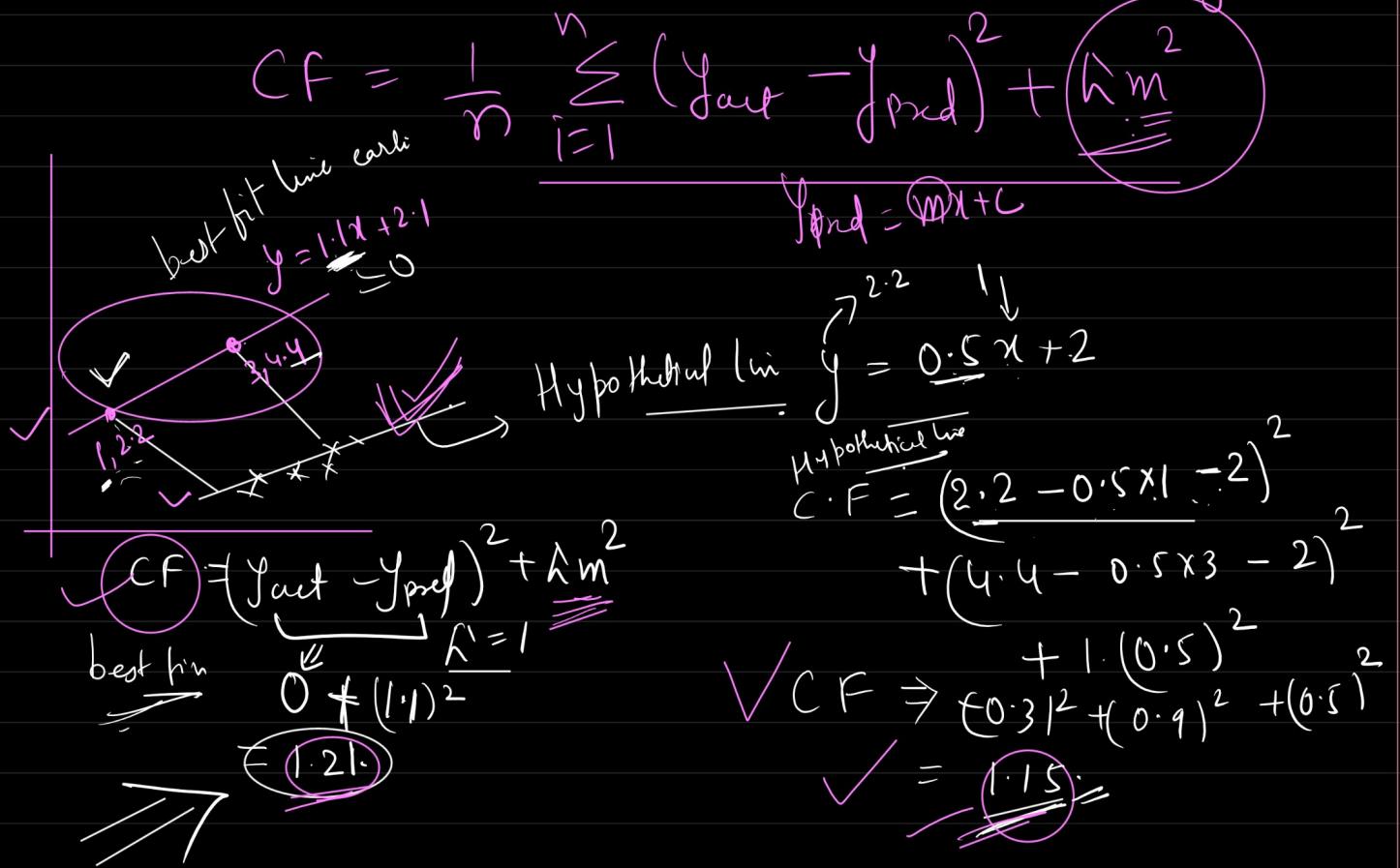
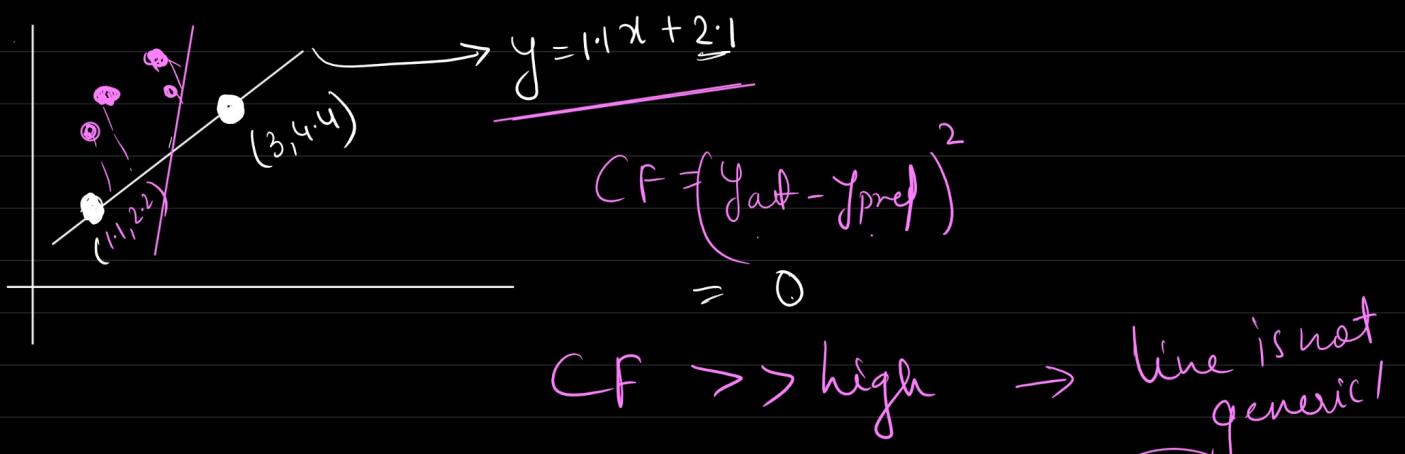
$$y = 4.1 + 2.50x$$

m .



Cost fn = $\frac{1}{n} \sum_{i=1}^n (\text{fact } f(x_i) - c)^2 + \lambda^2 m^2$

Error = 0



$$CF =$$

* Hypothetical line will be selected even if it is performing bad with train data because CF is lower as compared to other line.

With them Ridgeless Estm. \rightarrow Logistic

$$\left\{ \begin{array}{ll} \textcircled{1} \text{ Ridge} & - CF + \lambda m^2 \\ \textcircled{2} \text{ Lasso} & - CF + \lambda |m| \\ \textcircled{3} \text{ Elastic net} & - CF + \lambda m^2 + \lambda |m| \end{array} \right.$$
