

Data Science questions

- ① Give me a name of all the Regression Algo that you know.
- ② Linear Regression
- ③ Polynomial Regression
- ④ Ridge Regression
- ⑤ Lasso Regression
- ⑥ Elastic Net Regression
- ⑦ Support vector Regression
- ⑧ Decision tree Regression.
- ⑨ Random Forest Regression
- ⑩ Gradient Boosting Regression
- ⑪ xgboost Regression
- ⑫ LightGBM Regressor
- ⑬ CatBoost Regressor
- ⑭ K-Nearest Neighbour Regressor.
- ⑮ Bayesian Ridge Regressor
- ⑯ Huber Regressor
- ⑰ Theil-Sen Regressor
- ⑱ Quantile Regressor
- ⑲ Tweedie Regressor
- ⑳ NN Regressor

③ Give me a list of all the classification algorithm that you know

① Logistics Regression.

② K-Nearest Neighbour

③ Support Vector machine

④ Decision tree

⑤ Random Forest

⑥ Gradient Boosting.

⑦ XG Boost

⑧ Light GBM

⑨ Cat Boost

⑩ Naive Bayes

⑪ NN Classification

⑫ CNN

⑬ RNN

⑭ LSTM

⑮ GRU

⑯ Ada Boost

⑰ Bagging classifier

⑱ Extra tree classifier

⑲ LDA

⑳ QDA

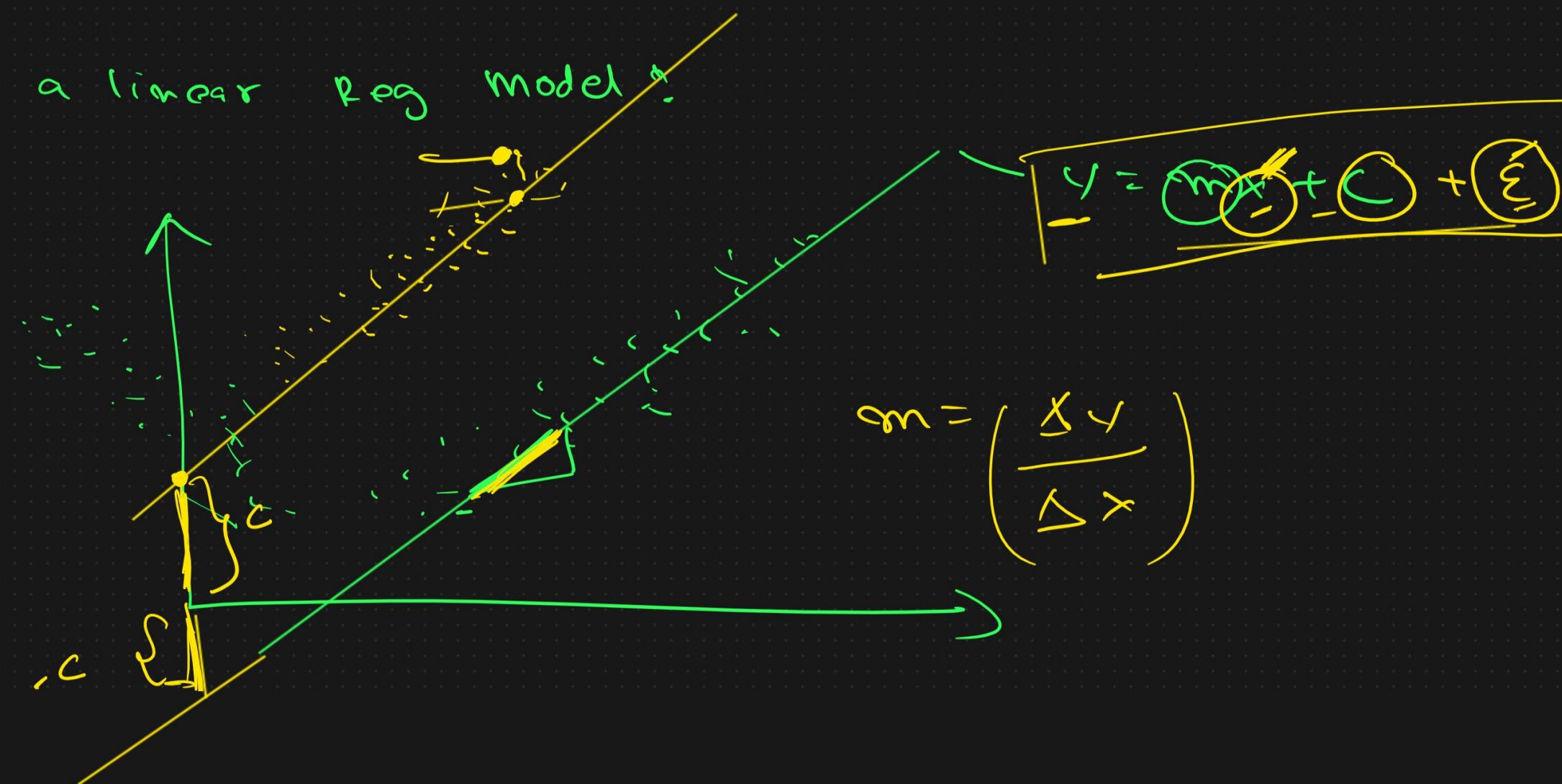
㉑ SGD

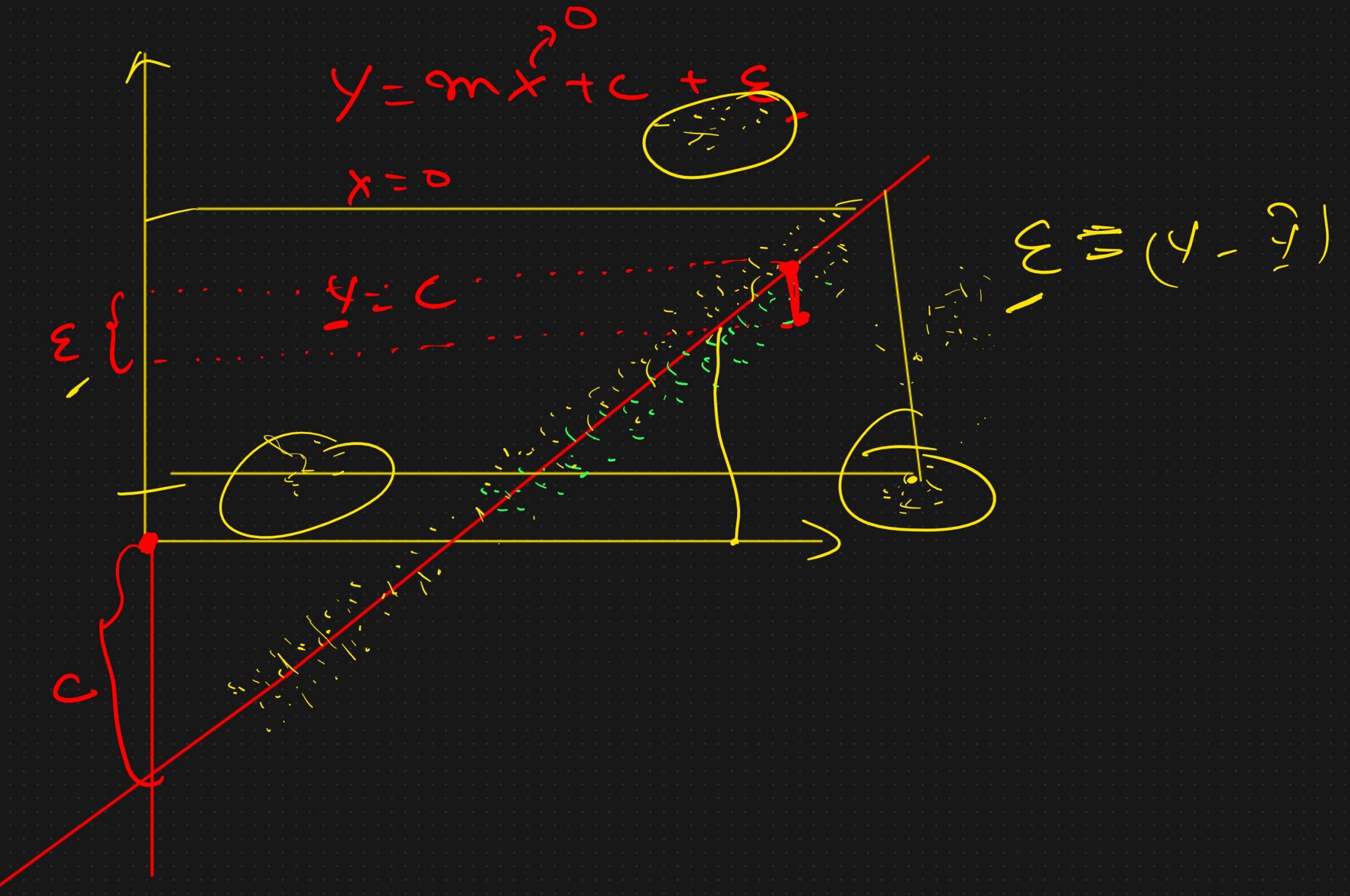
㉒ (OVR) & (OVO)

③ what is your understanding about Regression ?



④ what is a linear Reg model ?





⑤ When you're tuning linear Regression what are the factors inside data you have considered.

① Linearity.

② No Autocorrelation (Independence of Error) →

③ Homoscedasticity. (Constant Variance of Error)

④ Normality in Error →

⑤ No multicollinearity.

⑥ Exogeneity

⑥ what do you understand by Ridge Regression? (L₂)

⑦ What is overfitting

$$J(\beta) = \frac{\sum (y_i - \hat{y}_i)^2 + \lambda \sum \beta_j^2}{\text{SSE}}$$

↓ ↓
SSE L₂ Penalty term.
λ → hyper parameter

Whenever we have multiple co-related features

⑧ L1 Lasso

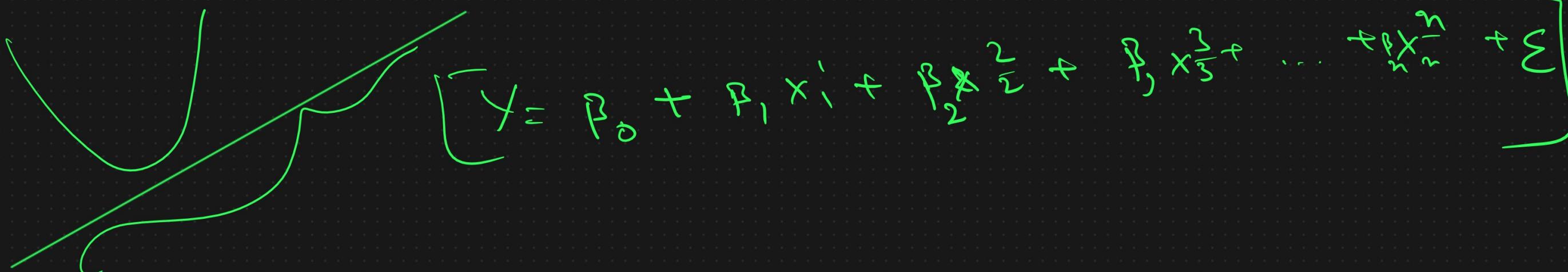
$$J(\beta) = \sum (y_i - \hat{y}_i)^2 + \lambda \sum |\beta_j|$$

L1 penalty term

Feature Selection.

$$y = \alpha + \beta_1 x_1 + \epsilon$$

⑨ What do you understand by Poly Regression?



⑩ Lets talk about logistic Regression . . .

$$P(Y = k | X) = \frac{1}{1 + e^{-(m_x + c)}} \quad \rightarrow \text{(Sigmoid)}$$

$$\geq 0.5 \rightarrow 1 \checkmark$$

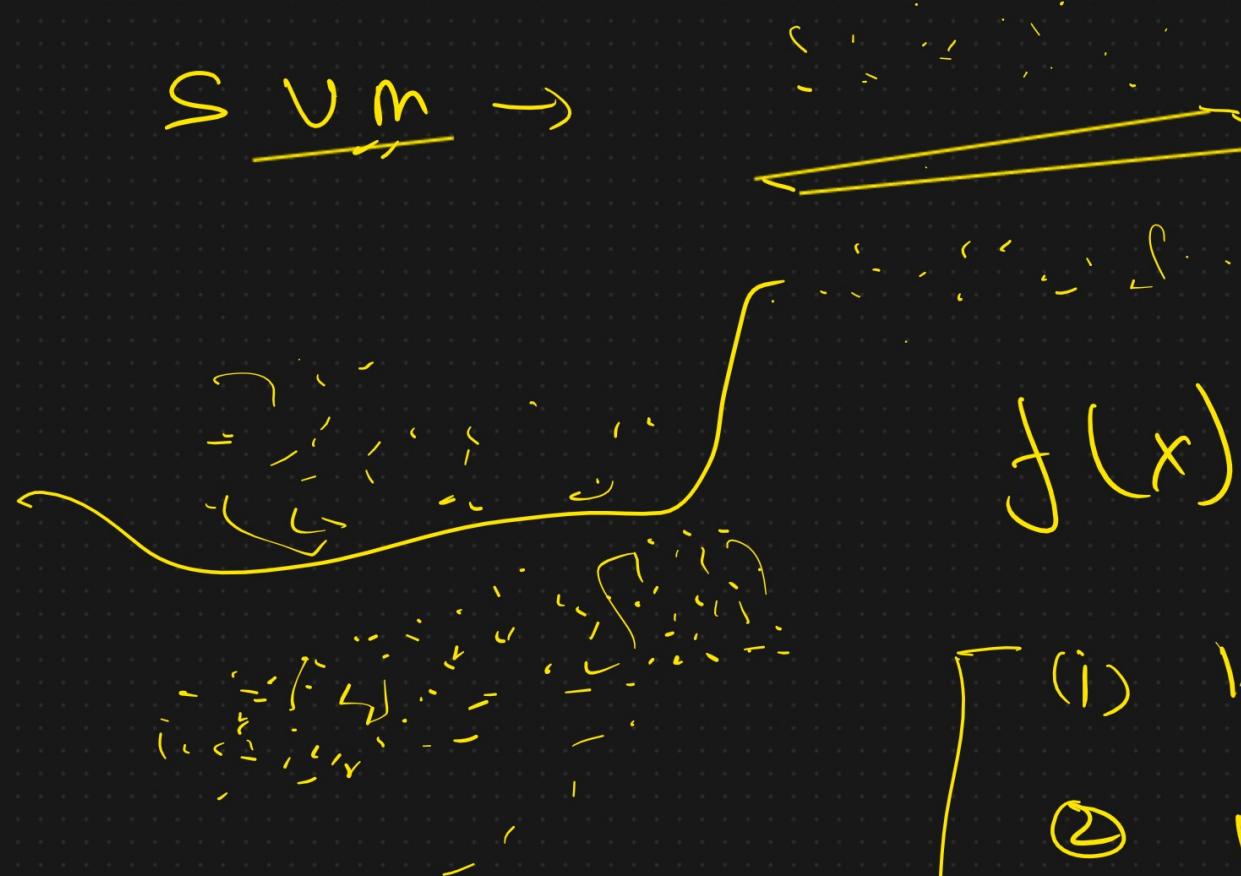
$$< 0.5 \rightarrow 0 \checkmark$$

(One vs Rest) / Done or can \rightarrow softmax

$$P(Y = k | X) = \frac{e^{kx}}{\sum e^{jx}}$$

①

$$\sum \underline{w} \cdot \underline{x} \rightarrow$$



$$f(x) = (\underline{w} \cdot \underline{x} + b)$$

(i) linear SUM

② non linear SUM \rightarrow

② SUR

Kernel trick

(RBF, Polynomial)

① multiple Features

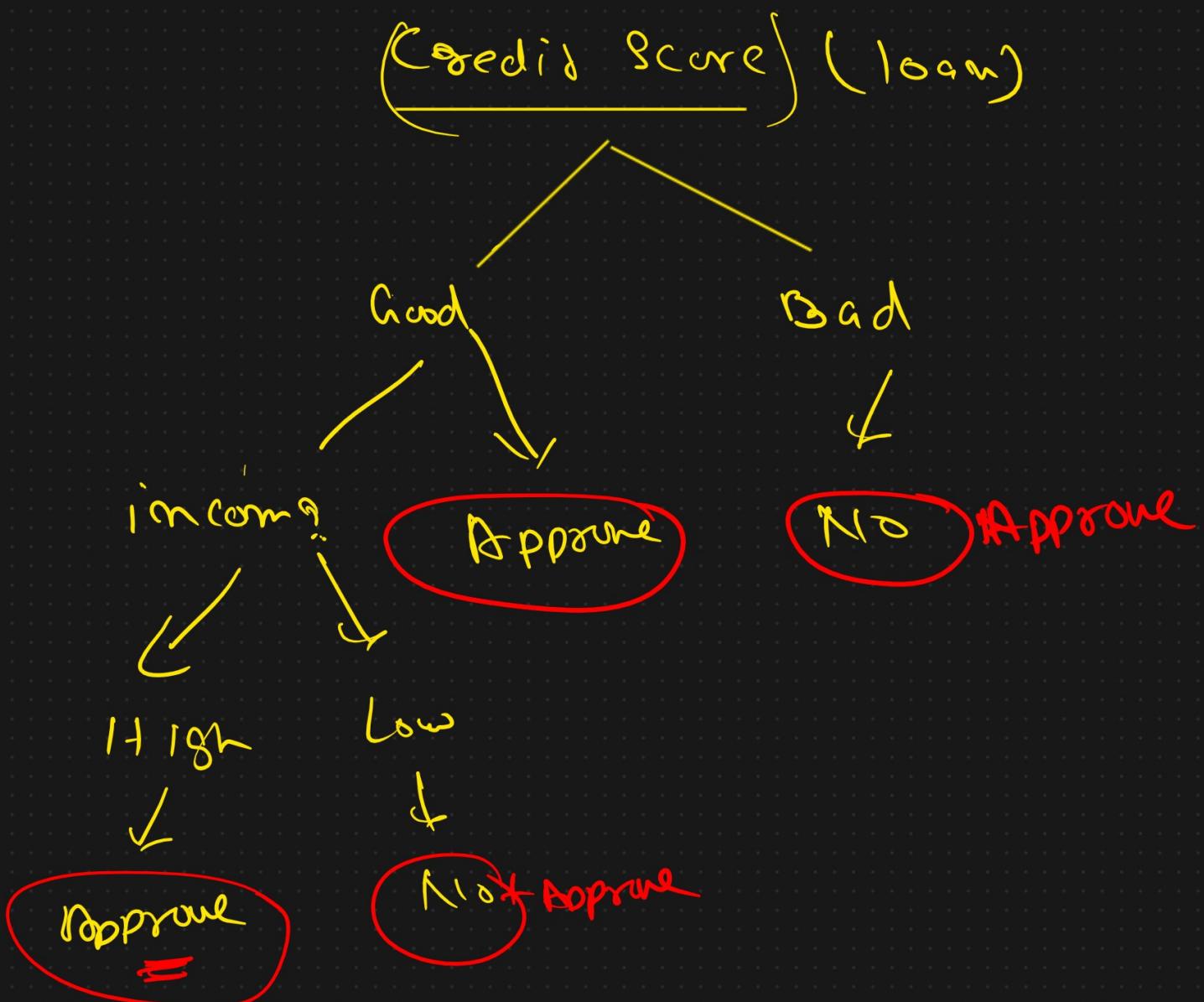
② Data is not linearly separable.

③ Small and huge data

④ high - dimension data.

12

Decision tree ?



$$\text{Gini} = 1 - \sum p_i^2$$

low Gini = Better Split.

$$\text{Entropy} = - \sum p_i \log_2(p_i)$$

lower Entropy = more pure Node

(13) Random Forest ?
=

(14) Boosting algo \rightarrow XGBoost

(15) XGBoost (Extreme Gradient Boosting)

① For Boosting it uses Gradient Descent.

② Prevent overfitting

③ Parallel Computation.

④ Missing Value Imputation

(5)

Stop automatically..

(6)

MAE \rightarrow

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{y}_i|$$

↳ Absolute error

(7)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

(8)

$$\text{RMSE} = \sqrt{\text{MSE}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

⑯ R^2 (R-Squared)

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{Total}}}$$

$$\underline{SS_{\text{res}}} \rightarrow \text{Residual Sum of Square} \rightarrow \sum (y - \hat{y})^2$$

$$SS_{\text{Total}} \rightarrow \text{Total Sum of Square} = \sum (y - \bar{y})^2$$

R^2 is -ve \Rightarrow

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2) (n-1)}{n-p-1}$$

↑ No of data points
↓ Number of independent Variable

Confusion matrix

Type I

Type II

$A P \rightarrow$	$I(I)$	$O(E)$	
$I(T)$	TP	FN	T_{v-2}
$O(E)$	FP	TN	T_{v-1}

