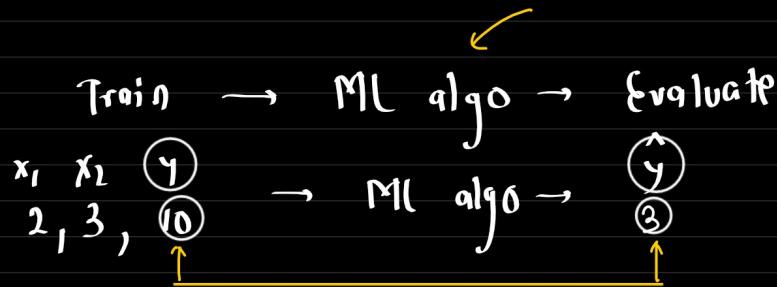
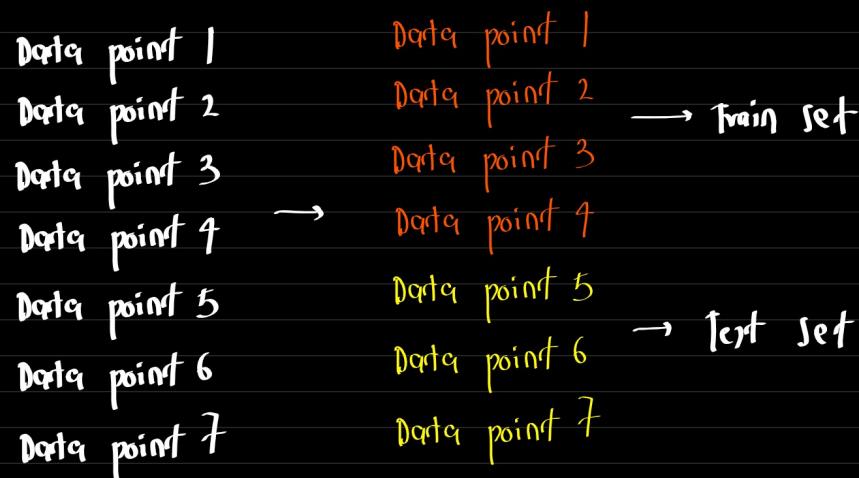


27-09-2025

Agenda: ML - II

ML intro , AI , parts of ML , Types of learning , Data split

Completed



difference → we require a method to tell ML algo about the error, so it can adjust itself.

Evaluation Method

whether our model is learning or not?

whether our model is good or bad?

y (real value)	\hat{y} (ML model prediction)	$(y - \hat{y})$
10	5	5
15	7	8
20	19	1
40	50	-10

error calculation → each data point → loss function
 → set of data (Train, val) → cost function → 4

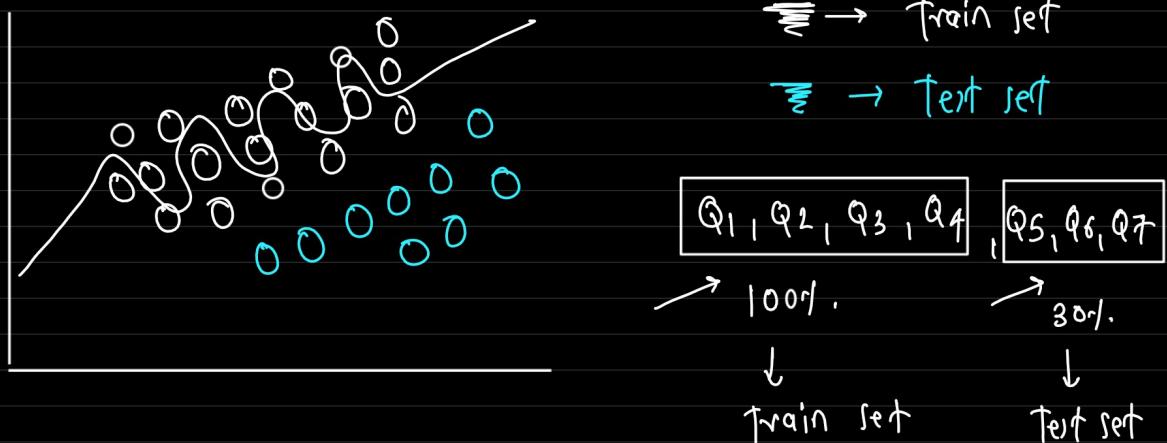
Train data → cost function → 180

Test data → 400

Loss function : Measure the error for a single example

Cost function : Sum of the loss functions over the dataset (Train/Val/Test)

If Training Loss <<<< Test/validation Loss, means this is a overfitting case



Train \rightarrow Error \rightarrow (obj) function
(formula) \rightarrow number

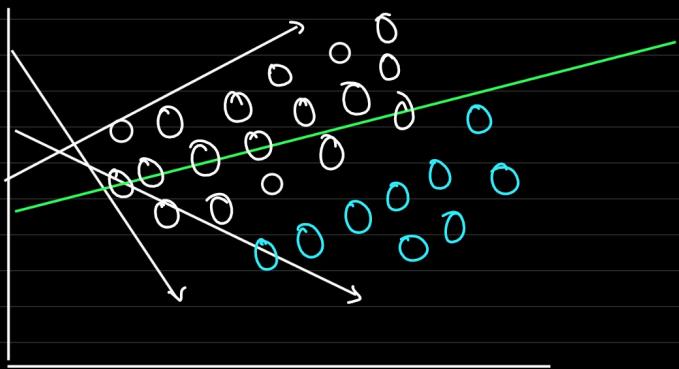
data[i] \rightarrow error \rightarrow loss function \rightarrow number
(formula)

If Training Loss === Test/validation Loss, means this is a "Ideal case"

Train loss \rightarrow 30
Test loss \rightarrow 32

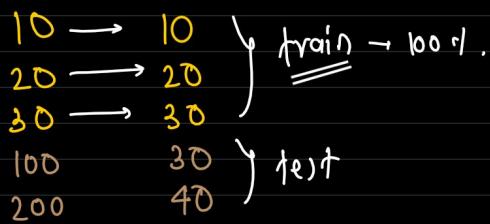
Q_1, Q_2, Q_3, Q_4, Q_5
 $\downarrow \qquad \qquad \qquad \downarrow$
92%. \longleftrightarrow 87%.
 $\uparrow \qquad \qquad \qquad \uparrow$
train test

If Training Loss and Test/validation Loss is bad, means this is a "underfitting case"

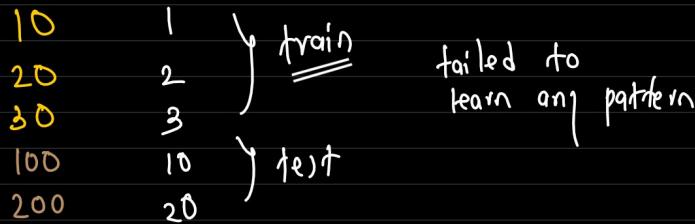


Train \rightarrow 100 \rightarrow 30% \leftarrow
sample is small
Test \rightarrow 5 \rightarrow 50% \nearrow
 \rightarrow

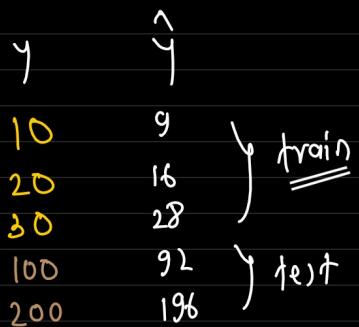
Overfitting $\rightarrow y \hat{y}$



Underfitting $\rightarrow y \hat{y}$



Good fitting $\rightarrow y \hat{y}$

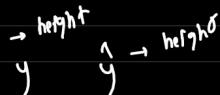


Evaluation Metrics: (Regression)

$$(1) \sum_{i=1}^n (y - \hat{y}) \quad (2) \perp \frac{1}{n} \sum_{i=1}^n (y - \hat{y})$$

$$(3) \text{MSE} \rightarrow \perp \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Mean Square Error



$$\begin{array}{l} \rightarrow \text{height} \\ \rightarrow 10 \end{array} \quad \begin{array}{l} \rightarrow \text{height}' \\ \rightarrow 8 \end{array} \quad \begin{array}{l} \rightarrow 2^2 \\ \rightarrow 4 \end{array}$$

$$\begin{array}{l} \rightarrow \text{height} \\ \rightarrow 20 \end{array} \quad \begin{array}{l} \rightarrow \text{height}' \\ \rightarrow 14 \end{array} \quad \begin{array}{l} \rightarrow 6^2 \\ \rightarrow 36 \end{array}$$

$$\begin{array}{l} \rightarrow \text{height} \\ \rightarrow 30 \end{array} \quad \begin{array}{l} \rightarrow \text{height}' \\ \rightarrow 20 \end{array} \quad \begin{array}{l} \rightarrow 10^2 \\ \rightarrow 100 \end{array}$$

English Math Science

more
attention

$$\begin{array}{l} (-1)^2 \rightarrow 1 \\ (1)^2 \rightarrow 1 \end{array} \quad \begin{array}{l} \text{making error positive} \\ \text{making error positive} \end{array}$$

$$\begin{array}{l} (1) \rightarrow 2 \rightarrow 1 \\ (2) 3 \rightarrow 1 \\ (3) 5 \rightarrow 1 \end{array} \quad \begin{array}{l} 1+1+1 = 3 \\ 1+1+1 = 3 \\ \sum_{i=1}^n (y_i - \hat{y}_i)^2 \end{array}$$

$$(4) \text{ MAE} \rightarrow \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Absolute Error

$$(5) \text{ RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

(6) Modified MSE

$$\frac{1}{n} \sum_{i=1}^n \underline{(y_i - \hat{y}_i)}^2 \xrightarrow[\text{derivative}]{} 2 \cdot \boxed{(y_i - \hat{y}_i) \cdot \frac{\partial}{\partial \theta} (y_i - \hat{y}_i)}$$

$$\rightarrow \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \cancel{\frac{1}{n}} \cancel{\sum_{i=1}^n} \cdot \frac{\partial}{\partial \theta} (y_i - \hat{y}_i)$$

	y	\hat{y}
Dog	0	0
Cat	1	1
Elephant	2	2

Evaluation Metric: (classification)

TP, TN, FP, FN (confusion matrix)

Accuracy: $\frac{(TP + TN)}{(TP + TN + FP + FN)} / \frac{\text{All true cases}}{\text{All cases}}$

: Overall correctness. $= 0.96 \rightarrow 96\%$.

Precision: $\frac{(TP)}{(TP + FP)} \rightarrow 0.857 \rightarrow 85.7\%$.

: How often we are correct

Email \rightarrow ML \rightarrow spam
not spam

Predicted		
Actual	spam	Non-spam
spam	600 TP	300 FN
Non-spam	100 FP	900 TN

$$\text{Recall} : \frac{\text{TP}}{\text{TP} + \text{FN}}$$

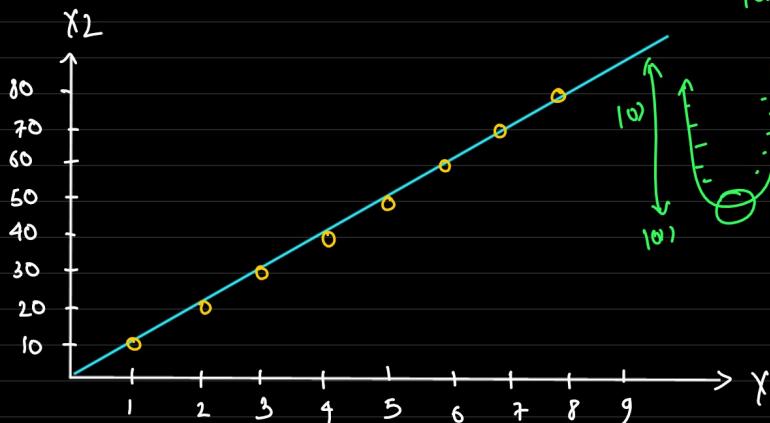
: $0.6667 \rightarrow 66.67\%$.

: How many actual spam caught.

F1 Score: Harmonic mean \rightarrow balance precision & Recall

$$: \frac{2(\text{Precision} * \text{Recall})}{\text{Precision} + \text{Recall}}$$

Linear Regression



X_1	X_2
1	10
2	20
3	30
4	40
5	50
6	60
7	70
8	80

$$y \leftarrow x$$

linear regression :

$$y = mx + c$$

m reduce complexity

add or subtract parameters \rightarrow model \rightarrow linear model

$$\rightarrow y = mx + c \rightarrow \text{element not multiplied}$$

\rightarrow not raised to a power

$$\rightarrow y = m_1x_1 + m_2x_2 + m_3x_3 + c$$

$$\rightarrow y = m_0 + m_1x_1 + m_2x_1^2 + x_2$$

Non-linear:

$$y = e^{m_0 + m_1x_1}, \quad y = \frac{m_0 + m_1x_1}{1 + m_2x_2}, \quad m_1(x_1)m_2(x_2) + c$$

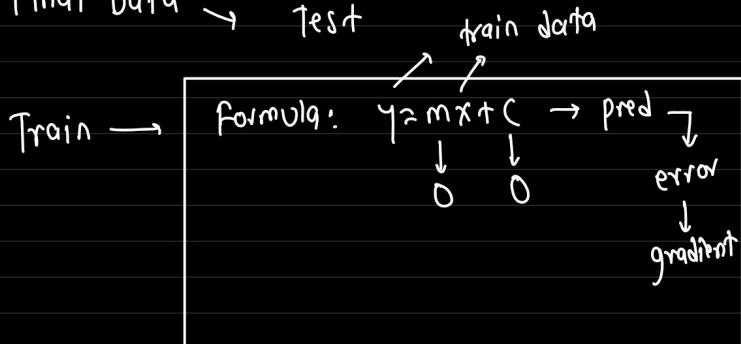
- linear regression
- polynomial regression ← linear features

x_1	x_2	y
1	4	120
2	8	240

$$\begin{aligned}
 120 \cdot x_1 &\rightarrow y \quad | \text{ linear} \\
 60x_1 + 2x_2 &\rightarrow y \\
 \downarrow & \downarrow \\
 m_1 = 60 & \quad m_2 = 2 \\
 \underline{\underline{=}} & \quad \underline{\underline{=}}
 \end{aligned}$$

Data → EDA → FE → Final Data

Final Data $\xrightarrow{\text{train}}$ $\xrightarrow{\text{test}}$



x_1	x_2
1	10
2	20
3	30
4	40
5	50
6	60
7	70
8	80

x_1	y
1	10
2	20
3	30
4	40
5	50
6	60
7	70
8	80

x_1	x_2	y
1	2	
2	3	
4	5	
6	7	
8	9	

$$(m_1, m_2)$$

$$x_1 m_1 + x_2 m_2 = \hat{y}$$

\downarrow \downarrow

∞ ∞

$$\text{error} = \hat{y} - y$$

$$\text{final loss} \rightarrow m_1, m_2$$

$$y \rightarrow 10, x = 1$$

$$m, c \rightarrow 0$$

$$\hat{y} = mx + c \rightarrow 0$$

$$\text{error} \rightarrow y - \hat{y} = 10$$

Dataset

Year-of-exp	Completed-project	Salary-k
1	0	39.764
2	1	48.400
3	1	56.978
4	2	68.240
5	3	77.867
6	4	85.022