11- OCT- 2025

Agenda:

    Logistic Regression

    Implementation of LR

    ROC , AUC

    Cross - validation

Logistic regression:    $y = mx + c$

  → classification

  ↳ Activation $f^n$ / logistic function → sigmoid → input range

                    ↳ $-\infty$ to $\infty$

               → output range → 0 to 1

    $z = mx + c$ → sigmoid ⤳ $\begin{matrix}0\\1\end{matrix}$ } threshold > 0·5 → $\begin{matrix}0\\1\end{matrix}$

$$Z = mx + c$$

height $\longrightarrow$ short $\longrightarrow$ 0 $\Big\}$ binary
$\searrow$ tall $\longrightarrow$ 1

$$0.2 \longrightarrow \text{prob of tall is } 20\%.$$

sigmoid $(z) \to \hat{y} \longrightarrow$ prob $\nearrow$

$\searrow$ $0.8 \longrightarrow$ prob of tall is $80\%.$

## linear regression.

$$\begin{matrix} y = 80 \\ \hat{y} = 91 \end{matrix} \Big\} \text{mse} \longrightarrow (y - \hat{y})^2$$

## logistic regression.

$$0 - 1$$

$$Z = 80$$
sigmoid $(z) \to 0.8 \leftarrow \hat{y}$ $\nearrow$
$\Big\}$ mse will not work
$1 \leftarrow y$
$\downarrow$
$0 \ 1$

To solve this prob, we need a different cost function.

## Modified MSE (Linear Regression)

$$\frac{1}{2m} \sum_{j=1}^{m} (y_i' - \hat{y_i})^2$$

Logistic Regression:

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^{m} \left( y_i \cdot \log(\hat{y_i}) + (1 - y_i) \cdot \log(1 - \hat{y_i}) \right)$$

Cost function =

$m = 1$, $y = 0$, $y\_pred = 0.9$ (case 1)

$$= -\left( y_i \cdot \log(\hat{y_i}) + (1 - y_i) \cdot \left( \log(1 - \hat{y_i}) \right) \right)$$

$$= -\left[ 0 \cdot \log(0.9) + (1 - 0) \cdot \left( \underline{\log(1 - 0.9)} \right) \right)$$

$$= -\left[ 0 + 1 \cdot \log(0.1) \right]$$

$$= -\log(0.1)$$

$$= -(-2.3025)$$

$$= 2.3025$$

$$\boxed{MSE = 0.81}$$

$m = 1$, $y = 0$, $y\_pred = 0.1$ (case 1)

$$= -\left[ 0 \cdot \log(0.1) + (1 - 0) \cdot \log(1 - 0.1) \right]$$

$$= -\left[ 0 + 1 \cdot \log(0.9) \right]$$

$$= -\log(0.9)$$

$$= -(-0.10536)$$

$$= 0.10536$$

Case 1 $(y=0, y\_pred = 0.9)$ : Cost $\approx$ 2.3025 (High Cost)
Case 2 $(y=0, y\_pred = 0.1)$ : Cost $\approx$ 0.10536 (Low Cost)

Flow of logististic regression.

1) $x_1, x_2, x_3, x_4$ $\longrightarrow$ linear equation for intermediate o/p $\rightarrow z$

2) $z \rightarrow$ sigmoid $\rightarrow \hat{y}$

3) loss-function gradient $(y, y\_pred) \rightarrow$ update weight/parameters.
$$\downarrow$$
$$m, \beta$$
$$\downarrow$$
learable parameters

4) calculate final loss using loss/cost function.

5) Iterate until satisfice

6) Done

Out of content topic.

Data Card:
    $\rightarrow$ info about data
    $\rightarrow$ sources
    $\rightarrow$ rows & columns
    $\rightarrow$ features $\rightarrow$ info $\rightarrow$ range of each feature, type
    $\rightarrow$ Missing values

Model card:
    $\rightarrow$ performance
    $\rightarrow$ link of test set on which we did evaluation
    $\rightarrow$ sklearn
    $\rightarrow$ confusion matrix
    $\rightarrow$ date

Imputation = [ 1, 2, 3, 4, 6 ]

numerical = [ 1, 2, 3, 4, (5), 6 ]

transfer = [ ] → object of sklearn pipeline

if imputation :

    ①    pipeline → imputation on Imputation list
             → scaling on (imputation list)
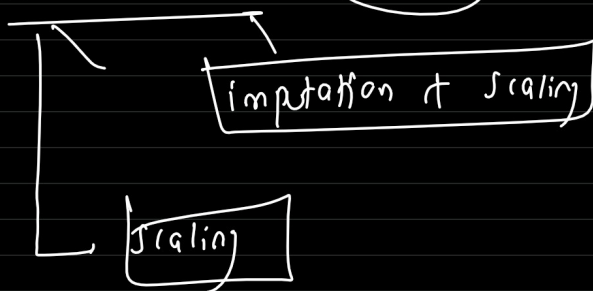
    append ① to transformr list

    ②    [ if numerical is not present ] → [ 5 ]
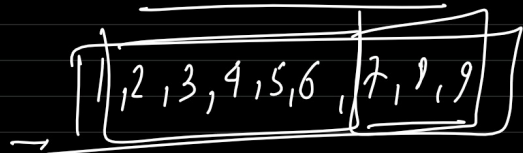
        pipeline → scaling → [ 5 ]

else,

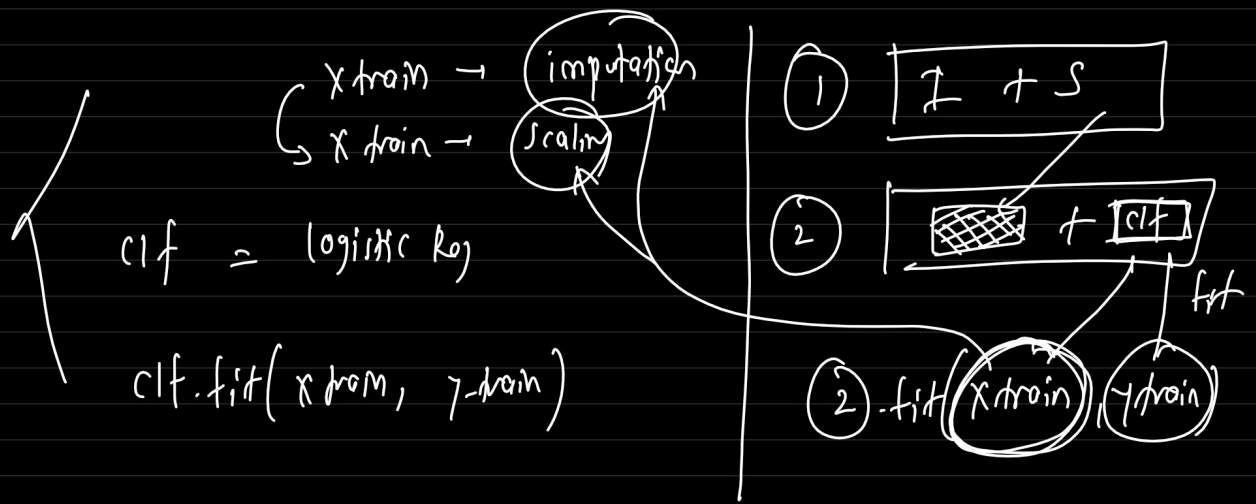    pipeline → scaling → [ numerical (all) ]
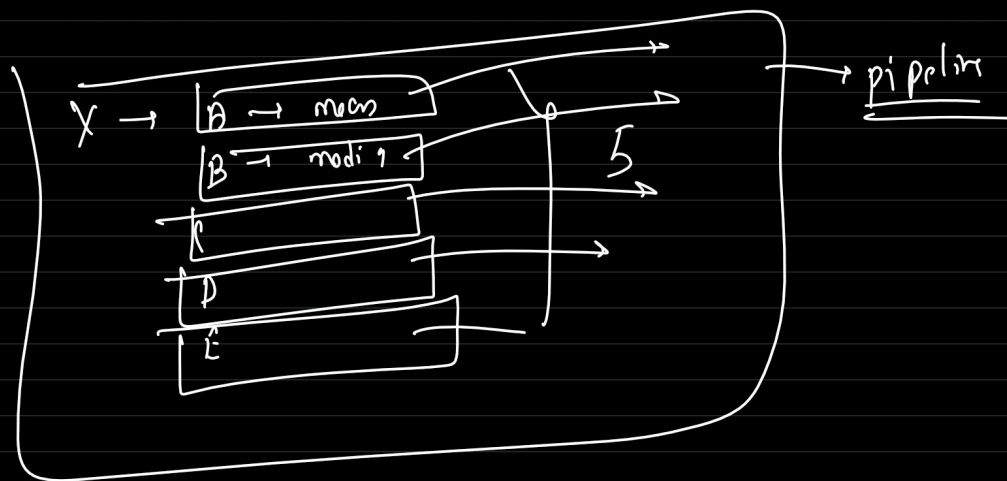
1, 2, 3, 4, 5, 6, (7, 8, 9)   remainder = "drop")

          imputation + scaling

        scaling

      transform ()

remainder = "passthrough"

1, 2, 3, 4, 5, 6, 7, 8, 9

X train → imputation

X train → scaling

① $\boxed{I + S}$

② ▨ + [clf]

fit

②.fit (X train) (y train)

clf = logistic Reg

clf.fit( x train, y train)

fit.

transform

X → [A → mean]

[B → nodi 1]

[C]

[D]

[E]

→ pipeline

predict

$$Z = X_{train} \cdot dot(weights) + bias$$

$$P = sigmoid(Z)$$

$$\hat{y} = (P >= 0.5) \; 1 \; else \; 0$$