## Agenda

- Cars24 Data
- Data Notation
- Goal of ML
- Linear Regression Intuition
- SKLearn Implementation
- ==Loss Function and Evaluation Metric==

---

### Cars24 Dataset

- Goal : To build a ML Model that can predict price of used cars accurately

Target → Sell_Price

Features : 17

Ordinal Encoding

① HE ⊃ X

Target Encoding

↓.

replace each of the category
with. mean of that category

↓

(Target)

---

**Feature Scaling**

---

① Min-max Scaler ⊃ [1, 1]

② Standar Scaler ⊃  mean = 0
std = 1

① Feature Scaling make model
training Faster

② Weights are more Interpretable

$$[ \; 10 \quad 20 \quad -5 \; ]$$

$$[ 1, 1 ]$$

max_val = 20
min_val = -5 $\Rightarrow$ $x_i = \dfrac{x_j - min}{max - min}$

- Target_Var : Selling_price

$$\downarrow$$

$$[a, b]$$

$$\downarrow$$

$$\mathbb{R}$$

- Regression

|   | $F_1$ | $F_2$ | $F_3$ | $F_4$ | - . . - | $Y$ |
|---|---|---|---|---|---|---|
| 1 | | | | | | |
| 2 | | | | | | |
| 3 | | | | | | |
| | | | | | $x_i$ | |
| 5 | | | | | | |

$$x_{3\ominus} \begin{bmatrix} & & & & \end{bmatrix}$$

$$y_i = \text{Scalar}, \quad y_3 = 1000 \text{ dollars}$$

$$\boxed{X_{n \times d}, \quad Y / y}$$

Training $(X, Y)$

$X_3$ ——— Prediction ———→ $\boxed{\text{ML Model}}$ ———→ $\hat{y}_3$

Predicted Price

$$\hat{Y} \quad SS \quad Y$$

Predicted SS Actual

Train → 80 %

Test_set: 20 %

$$1 \begin{cases} 0.8 \rightarrow 80\% \\ 0.2 \rightarrow 20\% \end{cases}$$

Train_set (80%)

Training

X   Y

(20%)

Test_set
X → [ M2 ] → $\hat{Y}_{test-act}$ (20%)

$$Y_{test\_set} - \hat{Y}_{test\_set}$$

## Linear Regression Intuition

⊙ Linear Model

Linear Hyperplane



$$y = w_1 x + w_0$$

$$\hat{Y}_{2019} \Rightarrow w_1 \times 2019 + w_0$$

Price (y-axis), Year (x-axis), 2012, 2019, 2020

$$y = mx + c$$

$$\downarrow$$

$$y = w_1 x + w_0$$

$3_i = 2$
$3_0 \Rightarrow 5$   $\Sigma$ Training

$y_i = 2x + 5$

$x$ (2012) $\rightarrow$   $\hat{y} \Rightarrow 2 \times 2012 + 5 \Rightarrow 4029$ dollars

$y = 3500$   ②

Diff $\Rightarrow 4029 - 3500$
$\Rightarrow 529$



$y = w_1 x_1 + w_2 x_2 + w_0$

$$\left.\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array}\right\} \text{Training} \quad \begin{array}{c} 2 \\ 3 \\ -2 \end{array}$$
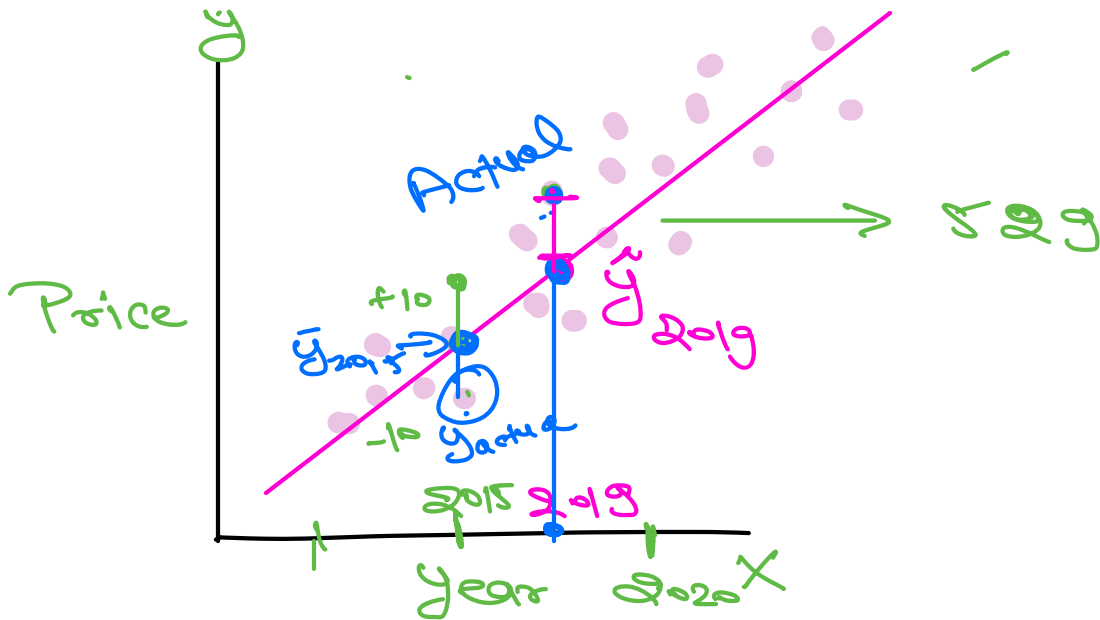
$$x_1 \Rightarrow 2012 \qquad y = 2 \times 2012 + 3 \times 1000$$
$$x_2 \Rightarrow 1000 \qquad\qquad\qquad -2$$

$$\boxed{2 \text{ features}} \qquad 1 \text{ target}$$

How many Parameters $\Rightarrow 3$

Dimension y Hyperplane $\Rightarrow 2d$

Distance of Each of the prediction
to their actual Values

③

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ ??? \end{bmatrix}$$

20% test_set

$e_i \ni \hat{y_i} - y_i$

③ Sum ③    $e_1 + e_2 + e_3 - - - -$

✗



$e_1 + e_2 \ni 0$ ✗

Ⓓ ✓ Sum of absolute

$|e_1| + |e_2| - - - -$

$$|e_1| + |e_2| \to 20$$

For a Good Model

min ( Sum of absolute Error )

$\downarrow$

MAE → Mean Absolut Error

min $\boxed{\dfrac{1}{n} \sum_{i=0}^{n} |e_i|}$

③ Sum of errors Square

Total Error = $e_1^2 + e_2^2 - - -$

$\downarrow$

$$\boxed{MSE \to \dfrac{1}{n} \sum_{i=0}^{n} e_i^2}$$

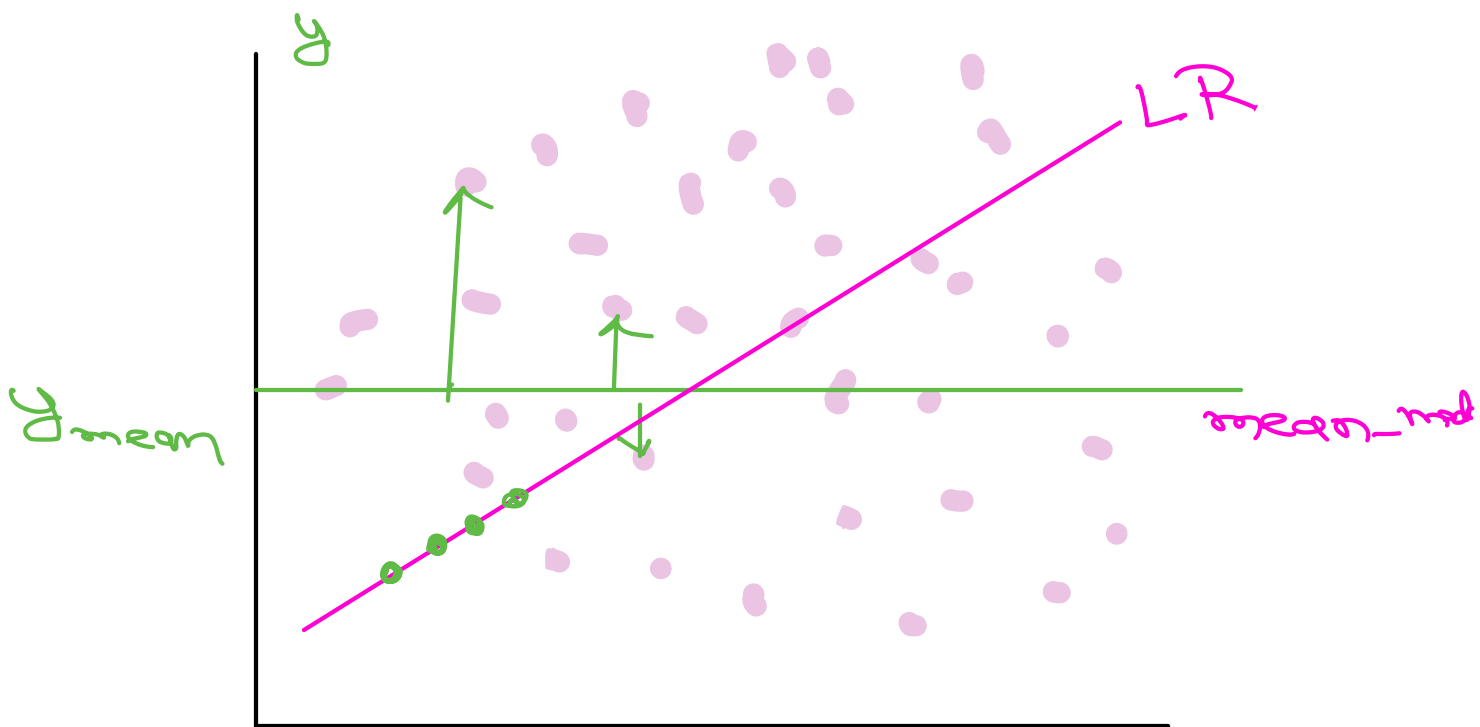| MAE | MSE |
| --- | --- |
|  |  |

Fill this H.W

1) In Regression, Loss function itself can be used as Evaluation metric

MSE and MAE ↓↓ 0 minimum

$m_1$ → | 2000 |
MAE
$m_2$ → | 1500 |    $m_2$ is better than $m_1$

To compare any model with mean-model as Baseline

Error

$$SS_{Total} \text{(Mean)} = \sum_{i=0}^{n} \left(y_i - y_{mean}\right)^2$$

$$SS_{res} \text{(LR)} = \sum_{i=0}^{n} \left(y_i - \hat{y}_i\right)^2$$

$R^2$
(R2_score)

$$1 - \frac{SS_{Res}}{SS_{mean}}$$

How well Linear Regression
Model is doing Compared
to mean Model

Todo: Calculate R2_score

Case 1) LR has same Error
as Mean

Case 2) LR has 0 Error

Case 3) LR has ∞ Error

# ① Ordinal Encoding

### Education_status

B. Tech $\Rightarrow$ 3

12th $\Rightarrow$ 2

10th $\Rightarrow$ 1

# ② OHE

| Make | Maruti | BMW | Ford |
|---|---|---|---|
| Maruti | 1 | 0 | 0 |
| BMW | 0 | 1 | 0 |
| Ford | 0 | 0 | 1 |

# ③ Target Encoding

| Make | |
|---|---|
| Maruti | Ymean of Maruti |
| BMW | Ymean of BMW |
| Ford | Ymean of Fords |