# Linear Regression

# Agenda

1. Linear Regression

2. Cars24 Case

3. Intuition lin. reg

4. Maths → Algebraic

5. Sklearn → Code

```
[ | | | | | | | ]  ⟶   ┌─────────┐   ⟶   Price of car
                        │ ML      │
                        │ Model   │
                        └─────────┘
```

input var.
independent

output
dependent
target

OHE ⟶ 3300 New Cols.

"Curse of Dimn"

"Target Encoding"



| Make | Selling Price |
|------|--------------|
| Maruti | 1.2 |
| Hyundai | 5.5 |
| Hyundai | 2.15 |
| Hyundai | 2.26 |
| Ford | 5.70 |

→ avg

Target Encoding →

↓

Target encoding replaces the categories with a number representing the average target value associated with each category.

MEAN

Maruti - 1.2

Hyundai - 3.3

Ford - 5.70

| Make | Selling Price |
|------|--------------|
| 1.2 | 1.2 |
| 3.3 | 5.5 |
| 3.3 | 2.15 |
| 3.3 | 2.26 |
| 5.7 | 5.70 |

Replace

# Min Max Scaling (Normalization) / Standardisation

$$X_{Scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$
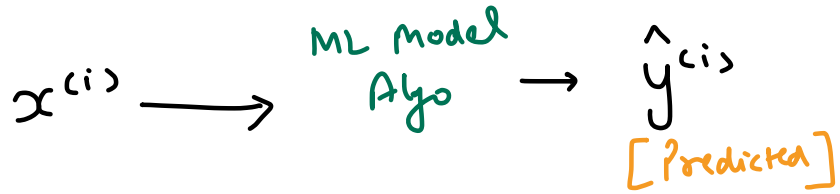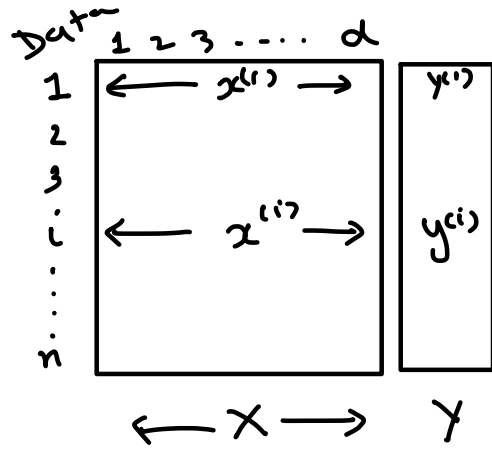
**Mix - Max Scaler**

| Km_Driven |
|:---:|
| 10,000 |
| 5000 |
| 2000 |

**Scaling** →

| Km_Driven |
|:---:|
| 1 |
| 0.375 |
| 0 |

→ [0, 1]

$$X_{min} = 2000$$
$$X_{max} = 10000$$

# Goal of ML

Data

$$
\begin{array}{c}
1 \ 2 \ 3 \ldots \ d
\end{array}
$$

$\xleftarrow{x^{(1)}}$   $y^{(1)}$

1
2
3
$i$   $\xleftarrow{x^{(i)}}$   $y^{(i)}$

$n$

$\xleftarrow{\quad X \quad}$   $Y$

$x^{(i)} \longrightarrow$ ML Model Algo $\longrightarrow \hat{y}^{(i)}$

[Predicted]

ideally,

$y^{(i)} \approx \hat{y}^{(i)}$

original   Predicted

$x_{NEW} \longrightarrow$ ML Algo $\longrightarrow \hat{y}_{NEW}$

# Train | Test



ML Train
↓
( X-train , Y-train )

Evaluate
↓
( X-test , Y-test )

# Intuition L.R

$\rightarrow$ Univariate Lin. Reg. [ 1 input var.]

$\rightarrow$ Multivariate Lin. Reg. [ > 1 input var ]

input
engine $\longrightarrow$ output
Price

New Car (1200 cc)
$\downarrow$
input this in line eq.
$\downarrow$
get Price

$$x^{(i)} \xrightarrow{\substack{\text{line eq} \\ f()}} \hat{y}^{(i)}$$

St. **line** →

$$y = mx + c$$

$$y = w_1 x + w_0$$

$w_1 x$ → weight
$w_1$ → Engine
$w_0$ → intercept

For eg. →

$w_1$   $w_0$

$$\hat{y} = \boxed{200} x + \boxed{50000}$$

(1000cc)
New Car

Price = $200 \times 1000 + 50000$

= 2.5 L

# 2-input features.

age , odometer

$x_1$    $x_2$

| age | odo | y |
|-----|-----|---|
|     |     |   |

$$y = w_1 x_1 + w_2 x_2 + w_0$$

age → $w_1 x_1$
odo → $w_2 x_2$

weight    weight    intercept



Price

$w$

age

odo

$$y = -10000 x_1 - 10 \cdot x_2 + 5000$$

# d-features

$$y = w_1 x_1 + w_2 x_2 + \ldots + w_d x_d + w_0$$

$$\boxed{y = w^T x + w_0}$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

1 feature $\longrightarrow$ line 2D

2 features $\longrightarrow$ plane 3D

$\vdots$

d features $\longrightarrow$ hyperplane "d+1"
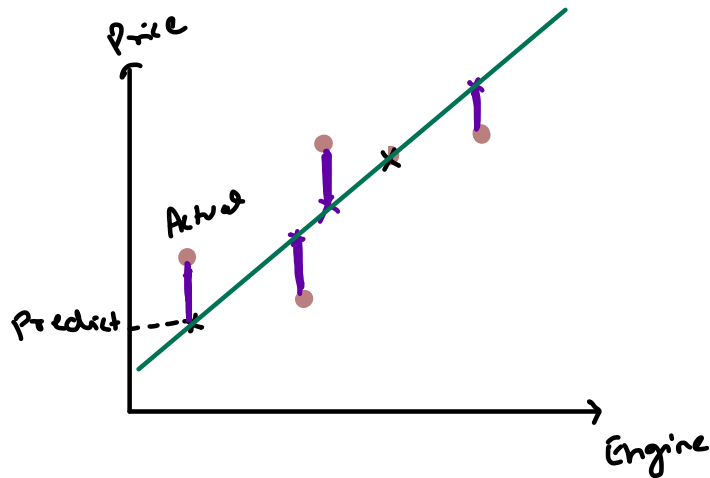
# Evaluation Metrics

$$x^{(1)} \rightarrow y^{(1)} - \hat{y}^{(1)} \rightarrow e^{(1)}$$

$$x^{(2)} \rightarrow y^{(2)} - \hat{y}^{(2)} \rightarrow e^{(2)}$$

$$x^{(3)} \rightarrow y^{(3)} - \hat{y}^{(3)} \rightarrow e^{(3)}$$

$$x^{(4)} \rightarrow y^{(4)} - \hat{y}^{(4)} \rightarrow e^{(4)}$$

$$x^{(5)} \rightarrow y^{(5)} - \hat{y}^{(5)} \rightarrow e^{(5)}$$



$$\text{Total error} = \sum_{i=1}^{n} e^{(i)}$$

$$\text{Total error} = e_1 + e_2 + e_3 + e_4 + e_5$$

$$\text{Total error} \rightarrow 5 + (-4) + (+3) + (0) + (-1)$$
$$\Rightarrow 3$$

$$|5| \rightarrow 5$$

$$|-4| \rightarrow +4$$

$$|3| \rightarrow 3$$
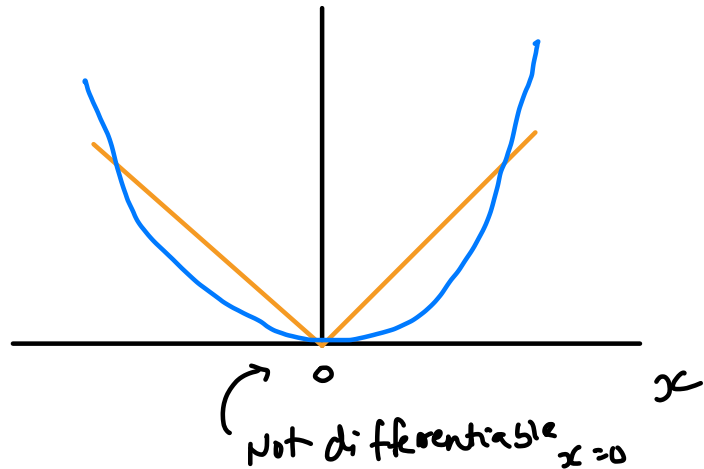
$$|-1| \rightarrow +1$$

$$\text{Error} = \sum_{i=1}^{n} |e_i| \Rightarrow \sum_{i=1}^{n} |y^{(i)} - \hat{y}^{(i)}|$$

Mean Absolute Error

$$\boxed{MAE = \frac{1}{n} \sum_{i=1}^{n} |y^{(i)} - \hat{y}^{(i)}|}$$

Mean Sq Error

$$\boxed{MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \hat{y}^{(i)} \right)^2}$$
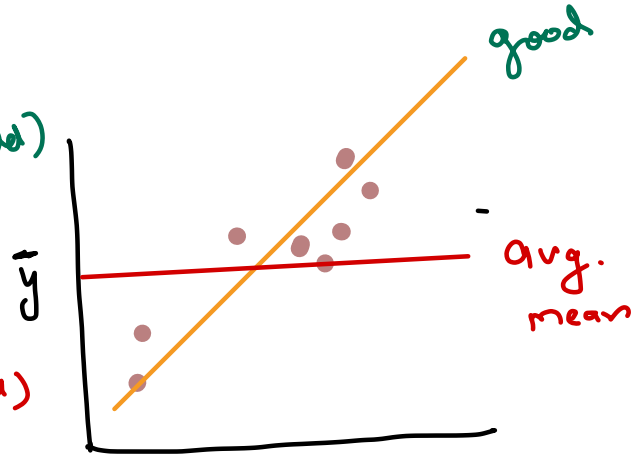


Not differentiable $x = 0$

MSE

$M_1$      9.62 ⌣ Better Model

$M_2$      15.41
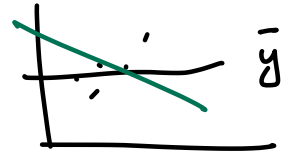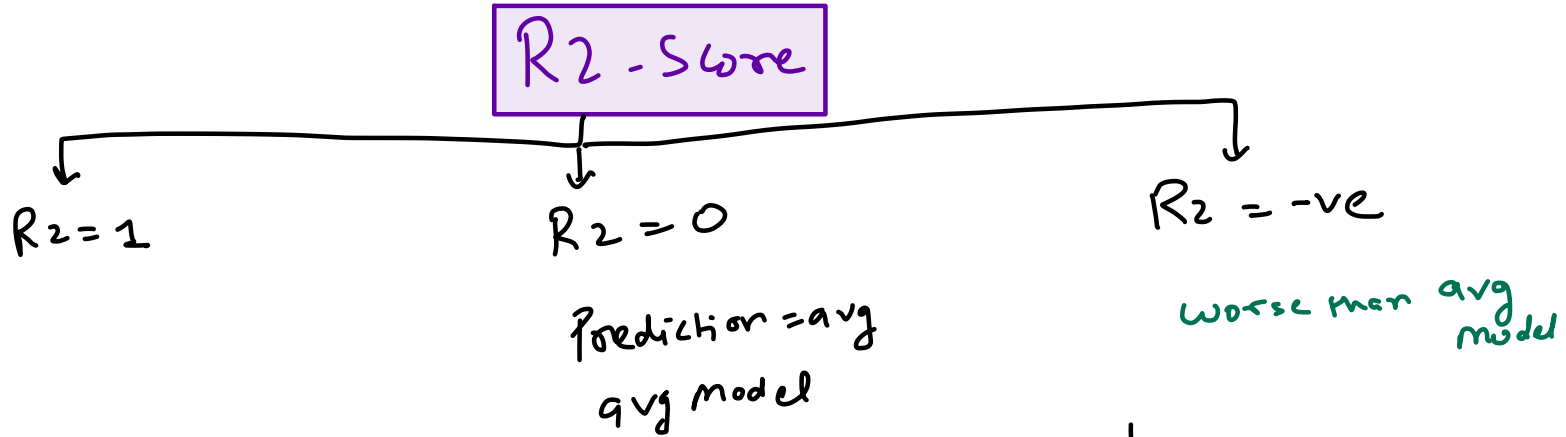
# R2_Score / R-squared / Coeff. of Determination

~~accuracy~~ → Performance

$$R2\text{-}Score = 1 - \frac{SS_{res} \text{(MSE of good model)}}{SS_{total} \text{(MSE of avg. model)} \text{ mean}}$$

$$R_2\text{-Score} \Rightarrow \quad 1 - \frac{\frac{1}{n}\sum_{i=1}^{n}\left(y^{(i)} - \hat{y}^{(i)}\right)^2}{\frac{1}{n}\sum_{i=1}^{n}\left(y^{(i)} - \bar{y}\right)^2}$$

$$\boxed{R2\text{-Score}}$$

$R_2 = 1$

$R_2 = 0$

Prediction = avg
avg model

$R_2 = -ve$

worse then avg
model

$\bar{y}$

$$R2\text{-}Score = 1 - \frac{SS_{res}}{SS_{total}}$$

Generally $\rightarrow$ $[\ 0,\ 1\ ]$

Bad

Best

# Quiz time!

🕐 Quiz Ended!

**In a multiple linear regression with five features, the coefficient of determination R2 is found to be 0.85. What does this value indicate about the model's performance?**

33 users have participated

| ✓ | A | The model explains 85% of the variation in the target variable | 27% |
| ✗ | B | The model's predictions are 85% accurate | 24% |
| ✗ | C | The model has an 85% probability of making correct predictions | 24% |
| ✗ | D | The model is 85% confident in its predictions | 24% |

https://colab.research.google.com/drive/1ajDuNR_-K_9mozthz5BJuYBZYgcXcsiZ?usp=sharing