

1 Linear

Regression

Agenda

- ① Linear Regression
- ② Cars24 Case
- ③ Intuition lin. reg
- ④ Maths \rightarrow Algebraic
- ⑤ Sklearn \rightarrow Code

Cars 24 Problem Overview

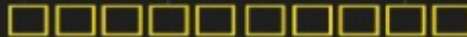
Data scientist at **Cars 24** → Sells pre-owned cars

resale price

→ To automate pricing the old car



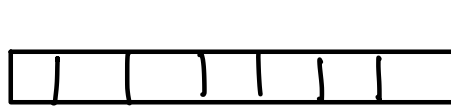
Given car features



*↓
info.*

(Make, Model Mileage , Odometer, Service History etc.) →

Predictors / Features



input var.
independent



Price of Car

output /
dependent
target

$y \in \mathbb{R}$

What do you think about the nature of Car Resale price prediction?

"Supervised"



110 users have participated

A	Regression	85%
B	Classification	9%
C	Clustering	5%

ONE → 3300 New cols/features
"Curse of Dimⁿ"

⇒ Target Encoding

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

Target
Encoding

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

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

MEAN

Maruti - 1.2

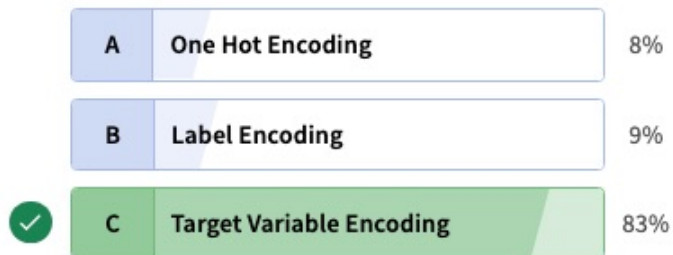
Hyundai - 3.3

Ford - 5.70

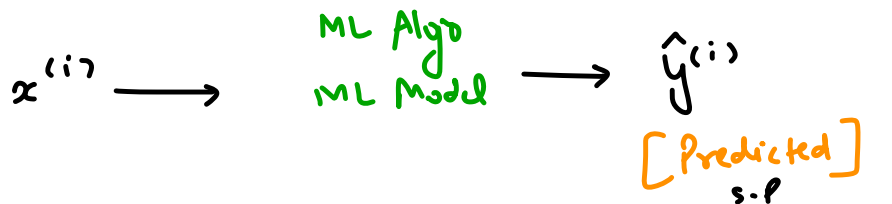
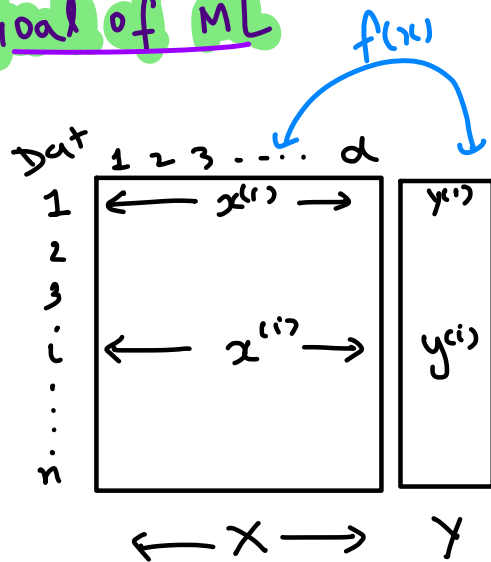
Replace

How do you think we should handle the large number of categories in make and model column?

117 users have participated



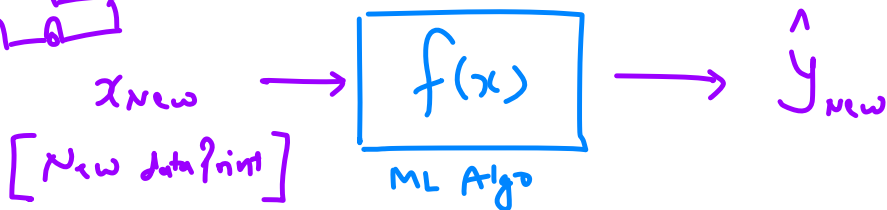
Goal of ML



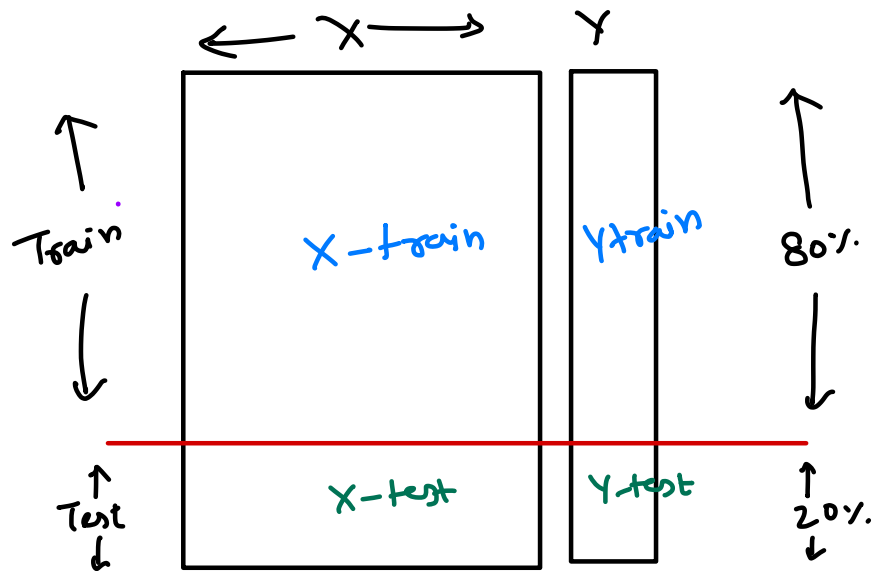
ideally,

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

original predicted.



Train Test Split



ML Train
↓
(X-train, Y-train)
 $f(x) / h(x)$

"Evaluation"
 $X_{\text{-test}} \rightarrow h(x) \rightarrow \hat{Y}_{\text{pred}}$
↓ Compare
 $Y_{\text{-test}}$

Intuition L.R

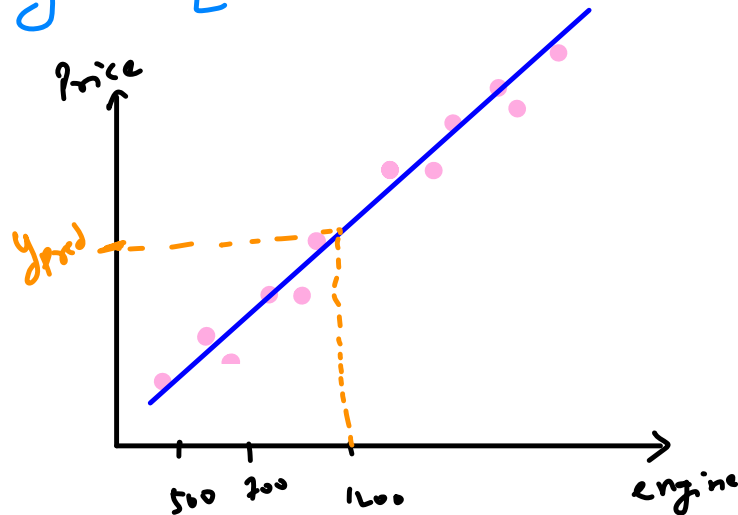
- Univariate lin. Reg
- Multivariate lin. Reg

[1 input var]
[> 1 input var]

input → output
engine → Price

New car (1200 cc)

↓
input this in line eq.
↓
get Price



line eq
 $h(x)$
 $x_{new}^{(i)}$ → $\hat{y}_{new}^{(i)}$ (predicted)

St. line

$$y = mx + c$$

$$y = w_1 \cdot x + w_0$$

weight

engine

intercept/bias

For. eg (Assume)

$$\hat{y} = \overset{w_1}{200} \cdot x + \overset{w_0}{50000}$$



(1200 cc)

$$\begin{aligned} \text{Price} &= 200 \times 1200 + 50K \\ &= 2.9L \end{aligned}$$

2-input features

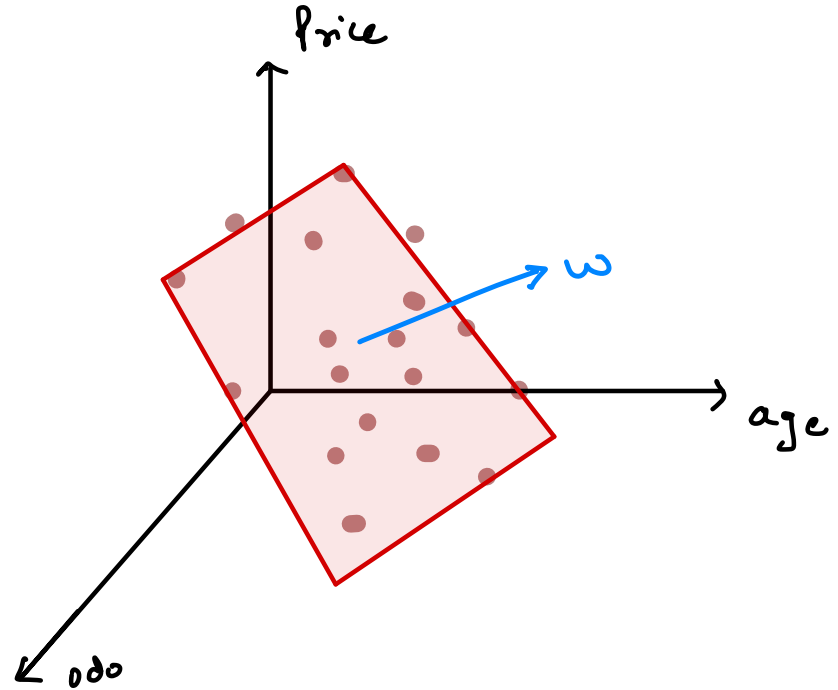
age, odometer

x_1	x_2	
age	odo	y

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

Annotations for the equation:

- x_1 is labeled "age" (red arrow)
- x_2 is labeled "odo" (red arrow)
- w_1 is labeled "weight" (green arrow)
- w_2 is labeled "weight" (green arrow)
- w_0 is labeled "intercept" (green arrow)



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

d-features

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

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

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

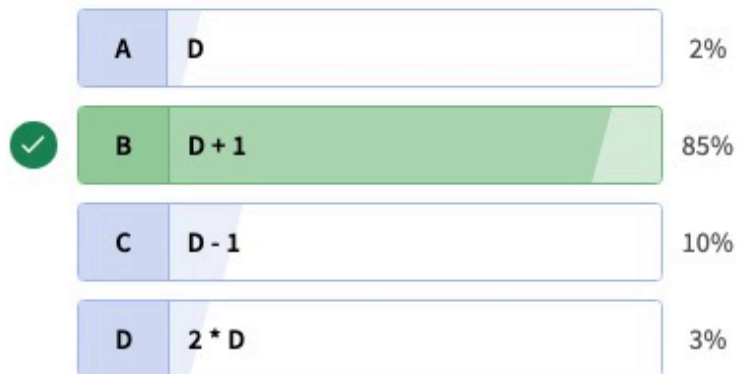
1 feature \rightarrow 2D line

2 features \rightarrow 3D plane

\vdots
 \vdots
d feature \rightarrow hyperplane (d+1)

If your data contains d features, how many dimensions will be required to fit the hyperplane through that data?

98 users have participated



Evaluation Metric

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

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

⋮

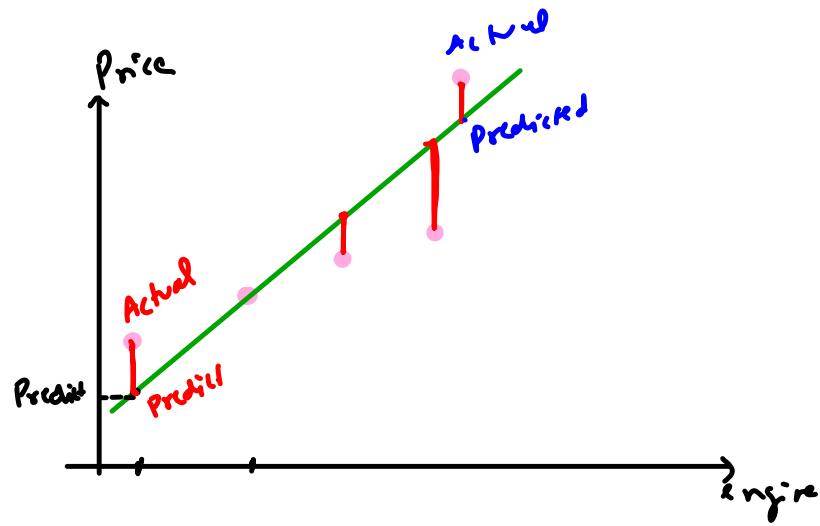
$$x^{(5)} \rightarrow y^{(5)} - \hat{y}^{(5)} = e^{(5)}$$

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

$$\text{Total error} = e^{(1)} + e^{(2)} + \dots + e^{(5)}$$

$$= 3 + 0 + (-2) + (-4) + (3)$$

$$= 0$$



$$|3| = 3$$

$$|0| = 0$$

$$|-2| = 2$$

$$|-4| = 4$$

$$|+3| = 3$$

$$\text{Error} = 12$$

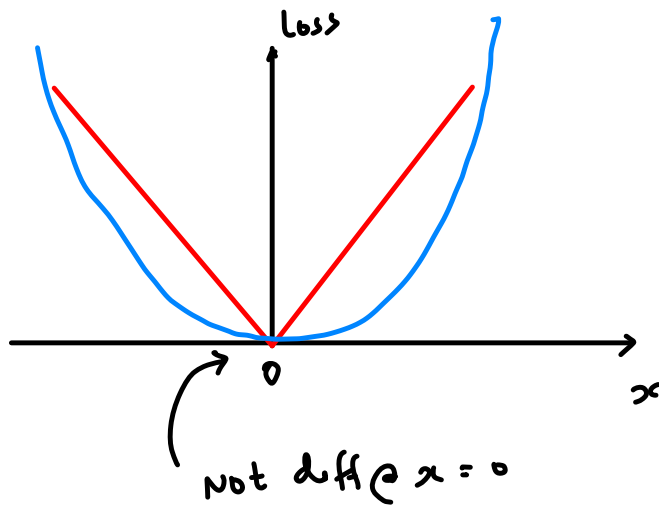
\Rightarrow

$$\text{Error} = \sum_{i=1}^3 |e^{(i)}| = \sum_{i=1}^3 |y^{(i)} - \hat{y}^{(i)}|$$

$$\text{MAE} = \frac{1}{3} \sum_{i=1}^3 |y^{(i)} - \hat{y}^{(i)}|$$

Mean
Sq. Error

$$\text{MSE} = \frac{1}{3} \sum_{i=1}^3 (y^{(i)} - \hat{y}^{(i)})^2$$



MSE

M_1

9.62

≡ Better

M_2

15.41

R2 Score → Next

In linear regression, if the MSE value is 0, it indicates:

85 users have participated



A

The predicted values perfectly match the actual values.

69%

B

The model has no predictive power and fails to explain the dependent variable.

15%



The model has high bias and underfits the data.

8%



D

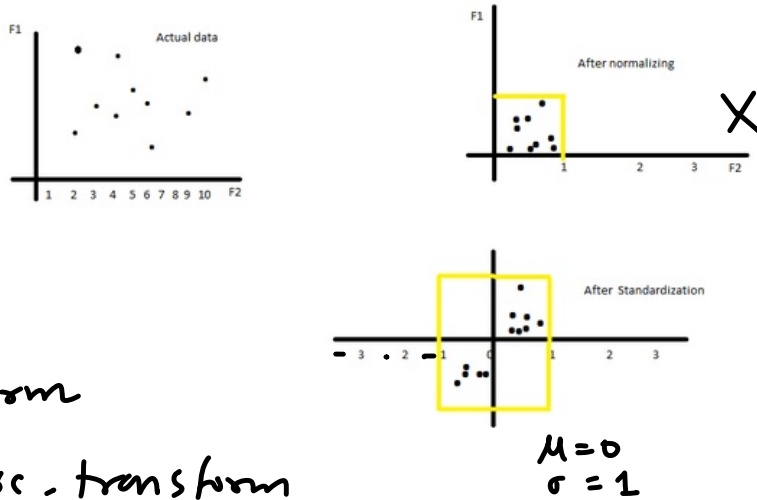
The model has high variance and overfits the data.

7%

[End Quiz Now](#)

MinMaxScaling \rightarrow images $[0-255]$ $[0-1]$

StandardScaling $\rightarrow \frac{x - \mu}{\sigma}$



. transform

. inverse - transform

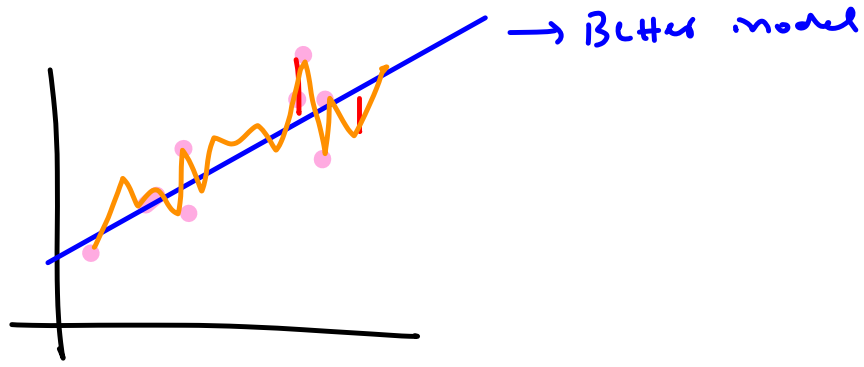
$$x' = \frac{x - \mu}{\sigma}$$

$$[x = x' \cdot \sigma + \mu]$$

$\rightarrow 1.4$
 \downarrow
8 yrs

age $\rightarrow -0.8$
 \downarrow
4 yrs

$\rightarrow -2.1$
 \downarrow
1 yr.



$$MSE \in (0, \infty)$$

Rmse

$$y = \omega x + \omega_0$$

$$(0.1) \cdot x + (0.1)$$



100

d

2^{100}

2^d

$$0 + 1 + 1 + 1$$

2

4

8

