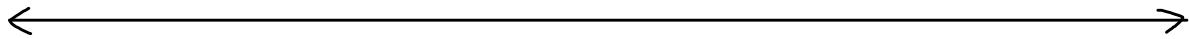
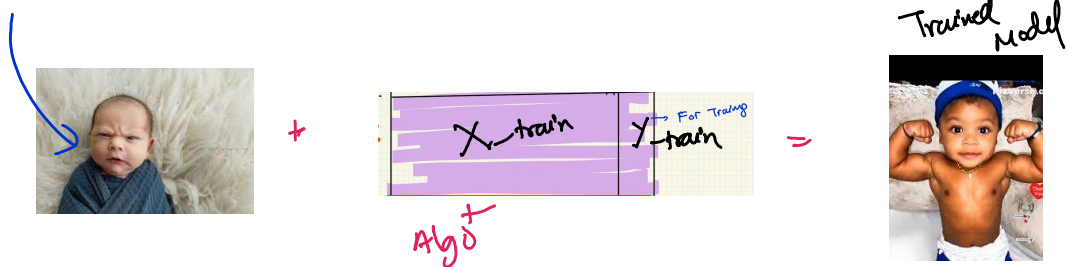


Linear Regression - 1

23 July 2025 20:52

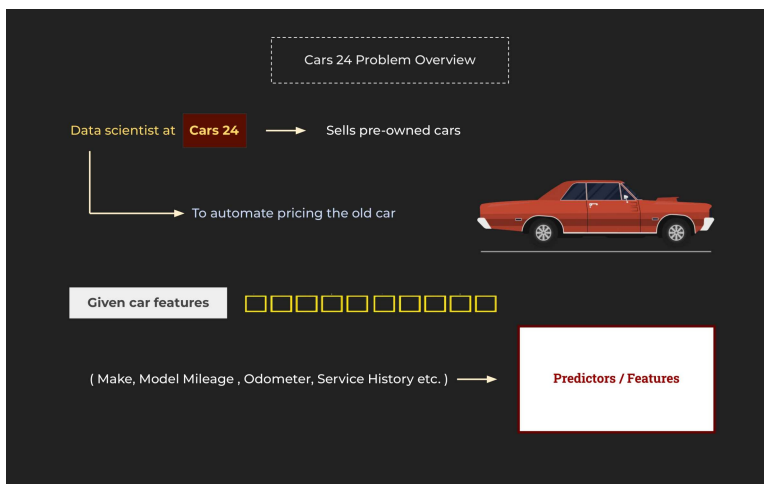
Feature	Regression	Classification
Output Type	Continuous values (e.g., price, temperature)	Discrete labels (e.g., cat/dog, spam/not spam)
Objective	Predict a numerical value	Assign a class/category
Examples	House price prediction, stock price forecasting	Email spam detection, disease diagnosis

Linear Regression	Points (Experience)
Classes	4 [Assumptions, Mathematics]
Interview	<20LPA (Very Imp Algo.)
Characteristics	Simple, Elegant, Explainable
Fields	Finance, Sales, Business (CEOs)
Age	Old, (1920s)

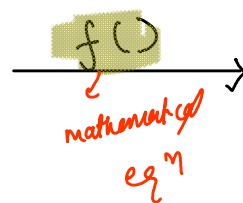
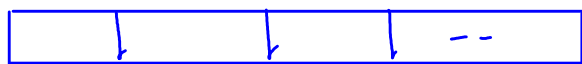


- Agenda →
- I ~~cars~~ 24 Dataset
 - clean [prepare]
 - prepare for further processing
 - II ML- concepts (cont.) [General]
 - III Linear Regression (Basic Theory)
 - cont. in next classes
 - IV [Apply L-R on our data to predict car prices]

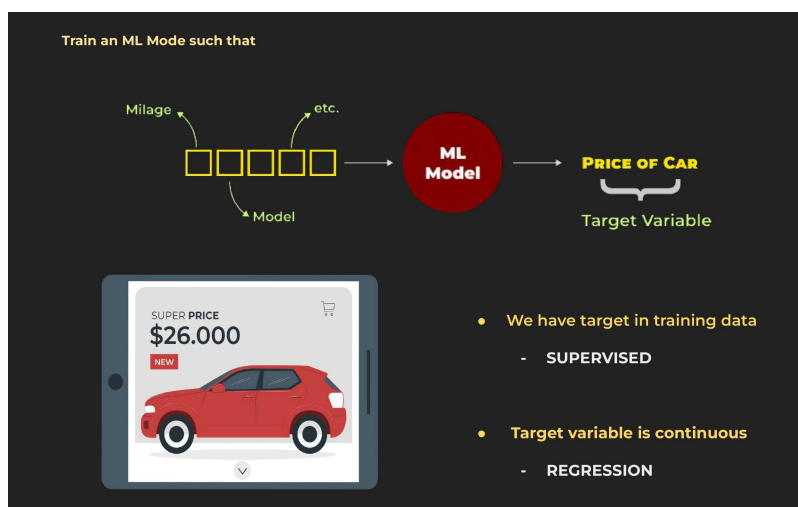
Feature	Regression	Classification
Output Type	Continuous values (e.g., price, temperature)	Discrete labels (e.g., cat/dog, spam/not spam)
Objective	Predict a numerical value	Assign a class/category
Examples	House price prediction, stock price forecasting	Email spam detection, disease diagnosis



Features



Price
(Predicted)



(Training Data)

Features	Price

Target Encoding

Revision - Target Variable Encoding

Using One Hot Encoding will result in a lot of dimensions -> **curse of dimensionality**

Hence, Target Encoding \otimes Target(Y)

32336

Target encoding replaces

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

MEAN

Maruti - 1.2

Hyundai - 3.3

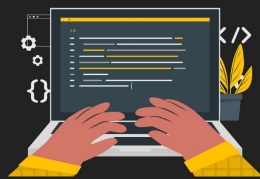
Ford - 5.76

Replace

Ford $\rightarrow \frac{5 \cdot 2 + 1 \cdot 2}{2}$

Line of Code

```
df['make'] = df.groupby('make')['selling_price'].transform('mean')
df['model'] = df.groupby('model')['selling_price'].transform('mean')
```



M →
M →
M →

Alb
Alto
Alb

Revision - Scaling the Data

$$X_{Scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Mix - Max Scaler



Km_Driven
10,000
5000
2000

Km_Driven
1
0.375
0

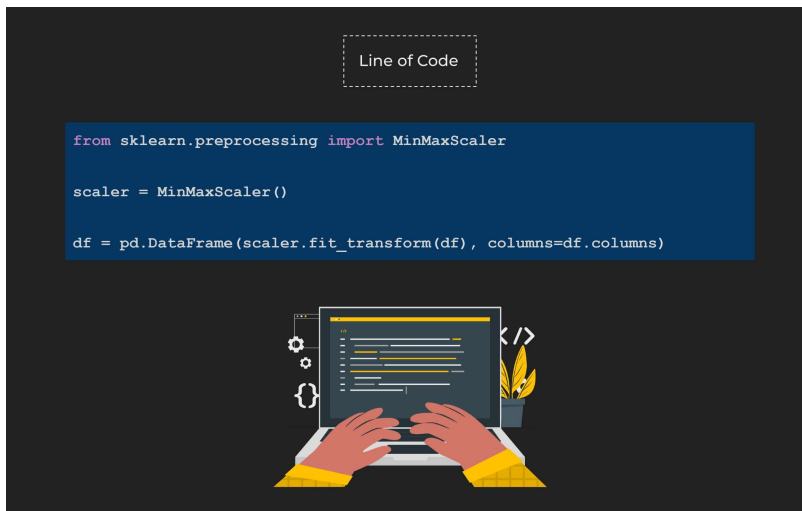
$$\begin{bmatrix} X_{min} = 2000 \\ X_{max} = 1000 \end{bmatrix}$$

Note:

Since all the features are in different ranges, we need to scale the data.

Later we will see why is scaling important

$$x' = \left(\frac{x - \mu}{\sigma} \right)$$



Before

	selling_price	km_driven	mileage	engine	max_power	age	make	model	Individual	Trustmark Dealer	Diesel	Electric	LPG	Petrol	Manual	5	>5
0	1.20	120000	19.70	796.0	46.30	11.0	4.684721	1.180000	1	0	0	0	0	1	1	1	0
1	5.50	20000	18.90	1197.0	82.00	7.0	4.458819	4.818750	1	0	0	0	0	1	1	1	0
2	2.15	60000	17.00	1197.0	80.00	13.0	4.458819	3.394000	1	0	0	0	0	1	1	1	0
3	2.28	37000	20.92	998.0	67.10	11.0	4.684721	2.242676	1	0	0	0	0	1	1	1	0
4	5.70	30000	22.77	1498.0	98.58	8.0	5.858258	6.777576	0	0	1	0	0	0	1	1	0
...
19815	6.50	69480	23.59	1364.0	87.05	6.0	10.532763	7.075000	0	0	1	0	0	0	1	1	0
19816	9.25	18000	17.50	1373.0	91.10	4.0	4.684721	7.128571	0	0	0	0	0	1	1	0	1
19817	4.25	87000	21.14	1498.0	103.52	8.0	7.182067	4.454000	0	0	1	0	0	0	0	1	1
19818	12.25	380000	16.00	2178.0	140.00	7.0	7.315421	8.096522	0	0	1	0	0	0	0	1	0
19819	12.00	13000	18.00	1497.0	117.60	4.0	5.879802	11.014286	0	0	0	0	0	1	0	1	0
...

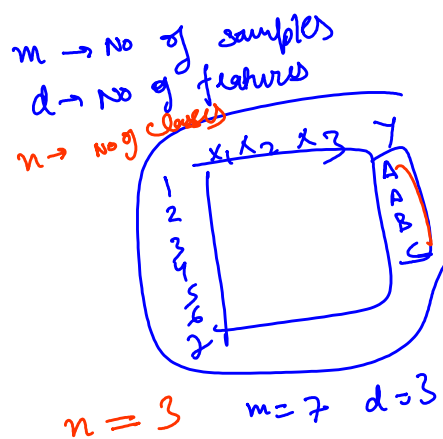
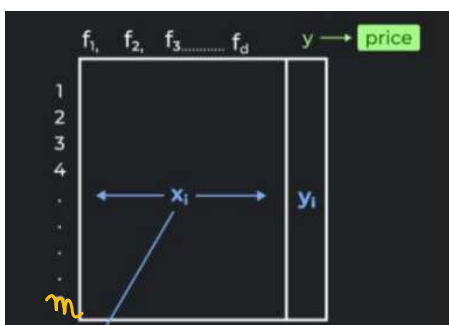
After

	selling_price	km_driven	mileage	engine	max_power	age	make	model	Individual	Trustmark Dealer	Diesel	Electric	LPG	Petrol	Manual	5	>5
	0.043684	0.031553	0.135345	0.117891	0.066506	0.310345	0.194048	0.041550	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	0.0
	0.252397	0.005237	0.128448	0.177281	0.123994	0.172414	0.232517	0.218382	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	0.0
	0.089795	0.015764	0.112069	0.177281	0.120773	0.379310	0.232517	0.149143	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	0.0
	0.095134	0.009711	0.145862	0.147808	0.100000	0.310345	0.194048	0.093193	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	0.0
	0.262104	0.007869	0.161810	0.221860	0.150709	0.206897	0.252367	0.313574	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0

No norm in scalts

Data Notation

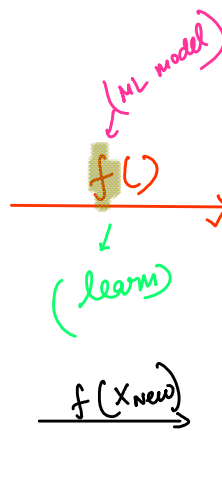
Tabular Data



Goal of Machine Learning

km_driven	mileage	engine	max_power	age	make	model	Individual	Trustmark Dealer	diesel	Electric	LPG	Petrol	Manual	>5
0.031553	0.135345	0.117891	0.066506	0.310345	0.194048	0.041550	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0
0.005237	0.128448	0.177281	0.123994	0.172414	0.232517	0.218382	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0
0.015764	0.112069	0.177281	0.120773	0.379310	0.232517	0.149143	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0
0.009711	0.145862	0.147808	0.100000	0.310345	0.194048	0.093193	1.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0
0.007869	0.161810	0.221860	0.150709	0.206897	0.252367	0.313574	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0
...
0.018258	0.168879	0.202014	0.099919	0.137931	0.484670	0.328028	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0
0.004711	0.116379	0.203347	0.138647	0.068966	0.194048	0.330632	0.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0
0.017606	0.147759	0.221860	0.158647	0.206897	0.318156	0.206656	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0
1.000000	0.103448	0.322719	0.217391	0.172414	0.324782	0.377671	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0
0.003395	0.120690

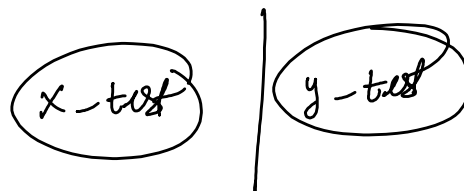
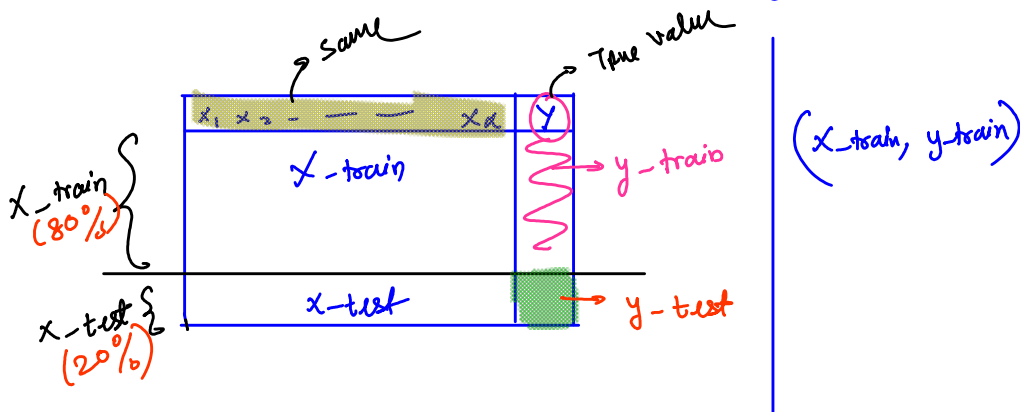
selling_price
0.043684
0.252397
0.089795
0.095134
0.262104
...
0.300934
0.434413
0.191724
0.580027
0.567892



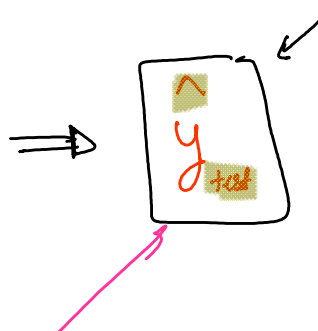
Goal: To generalise and so that we can make prediction on new cars

Train-Test Split

→ we don't use all data for Training



X_{test} ⇒  This has been trained on X_{train}, y_{train} ...





X_{train}

Y_{test}

\hat{Y}_{test} given by models

X_{test}	Y_{test}	\hat{Y}_{test}
① Alto	2 lakh	2.1 lakh
② Maruti	1.5 lakh	1.6 lakh
③ MG Hector	22 lakh	21 lakh

How well my model is performing

Solution:

Because we don't know the ground truth (y) for New Sample

DO NOT use whole data for TRAINING

Typically,

80%

20%



Use this for training the data

Use this as a proxy for unseen data

'y' known

Check if: $\hat{y} \approx y$

For evaluating model performance



Mean

(substitute)

Train test

X_{train}

Model	Model	S.P
S/3	A → 2	2
S/3	A	2
3	B	3
5	C	5
0	D	0
0	D	0

we have done

Target encoding

A^*

$$\Rightarrow \frac{2+2+1}{3} \Rightarrow (5/3)$$

X-test

5	C	0
0	D	0
2	E	2
5/3	A	1

✓
✗
⇒ (5/3)

Correct A → $\frac{2+2}{2} = 2$

	male	price
X-train	A A B C B	
X-test	B C A E	y-test

→ Formulas

	Age	Age scaled
Train	2 3 5 10	
Test	8 20	

$\frac{8 - x_{min}}{x_{max} - x_{min}}$
 $x_{min} = 2$
 $x_{max} = 10$

$$\frac{20 - 2}{10 - 2} = \text{fit}$$

- transform

① fit

Age → x_{min}
 Km-dr → x_{max}
 → x_{min}

② transform

test - data

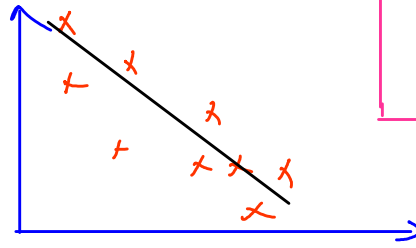
← (Linear Regression) →

↑ x

x_1	y
Age	Price
2	10
3	6

① Age

Price
(y)



2	6
3	6
5	3

(x) Age

$$y = mx + c$$

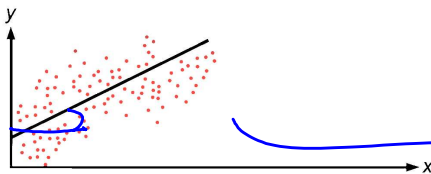
$$\text{Price} = m (\text{Age}) + c$$

$$\text{Price} = w (\text{Age}) + c$$

Geometrically → we are trying to fit a line through our data

① → Uni Variate LR

Linear Regression



$$\text{Price} = f(\text{Age})$$

$$\text{Price} = w (\text{Age}) + c$$

Linear → Linear fn

$$y = f(x)$$

(Linear)

$$y = mx + c$$

$y \neq \max + cx$

(linear)

$$y \neq mx + c$$

$$y = \sin(x)$$

y is a function of x

$$x^i \xrightarrow{f()} y^i$$

$$\hat{y}^i = f(x^i)$$

f() can be $x+2$, x^2+3 , $\sin(x)$, etc

1 predictor LR = Straight line

$$y^i = mx + c$$

$$y^i = w_1 x + w_0$$



Geometrically,

Linear Regression is trying to find a line such that most of the points are close to the line

Learning??

Estimating best values of <w> vector

Sk learns
(Uni-varial LR)

Linear Regression using Sklearn

Univariate

Fitting the model on 1 predictor

```
X1=X[['model']]
X1_train = X_train[['model']]
X1_test = X_test[['model']]

from sklearn.linear_model import LinearRegression
model = LinearRegression()
model.fit(X1_train, y_train)
```

LinearRegression
LinearRegression()

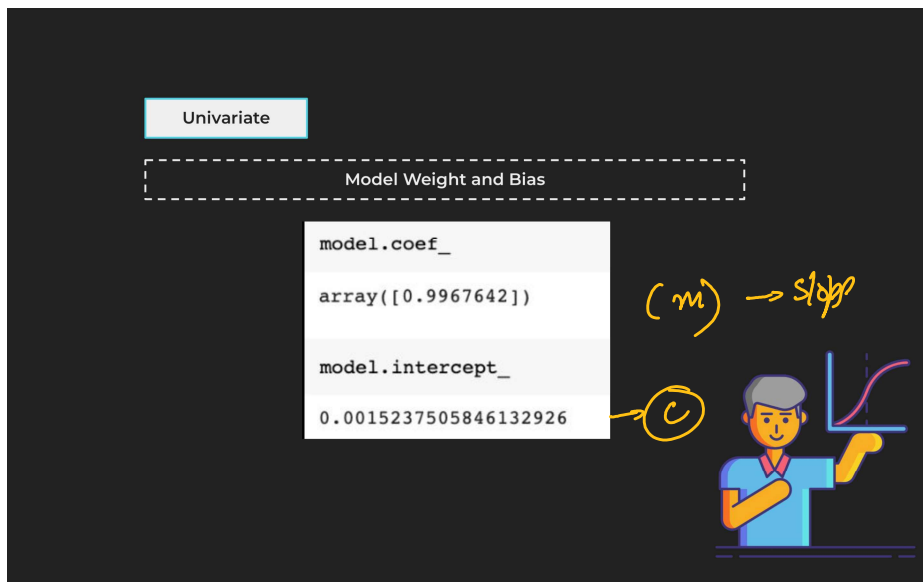
scale of fit param()

$$y = mx + c$$

$$y = mx + c$$

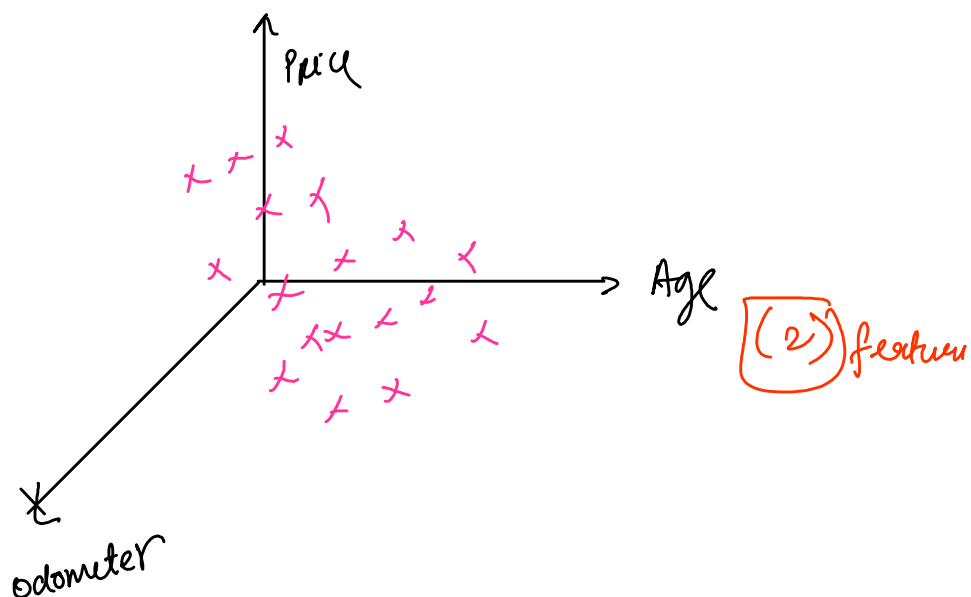
Intercept (Bias)

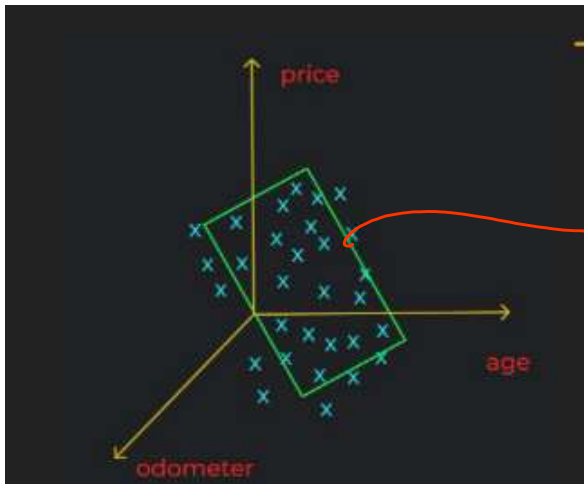
coeff



Multiple variables → MLRM

(2)





2 variables

$$\text{Price} = w_1 (\text{Age}) + w_2 (\text{odometer}) + w_0 \quad (\text{Bias})$$

Handwritten notes: 'feature' points to the terms in parentheses; 'wt' (weight) points to w_1 and w_2 ; 'Bias' points to w_0 .

3 features \rightarrow 4D $\{ \text{Hyperplane} \}$

Univariate linear regression - 1 feature

Bivariate linear regression - 2 features

Multivariate linear regression - d features

Straight line 2D

Plane 3D (Target)

d+1 hyperplane to visualize

Learning in Linear Regression

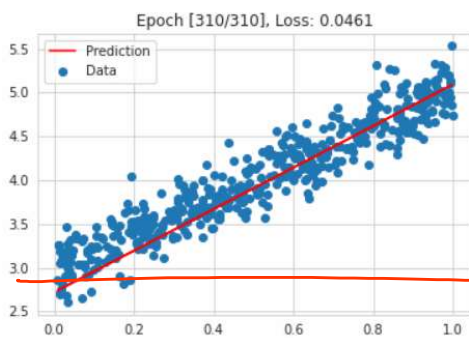
Estimating best values of $\langle w_0, w_1, w_2, \dots, w_d \rangle$

Which means to find the weights of best fitting line

3D

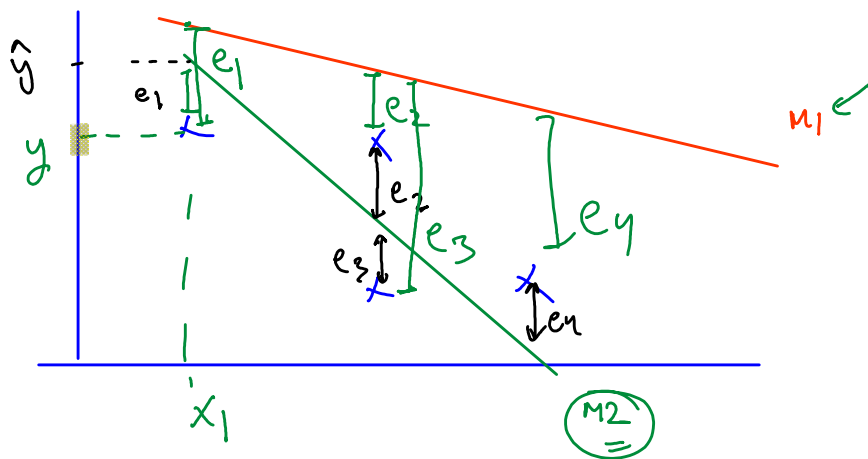
d+1

Machine Learning

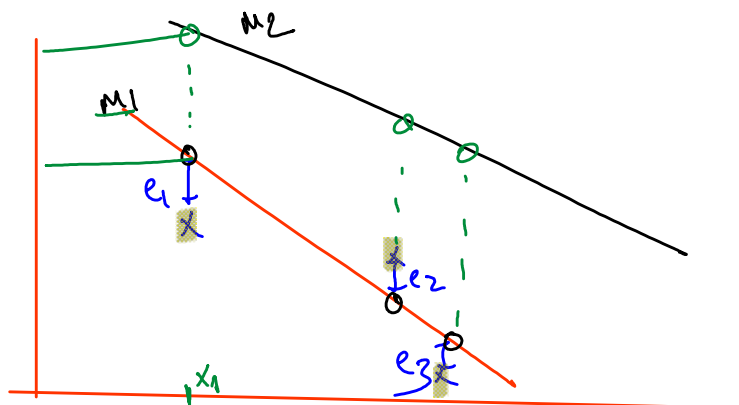
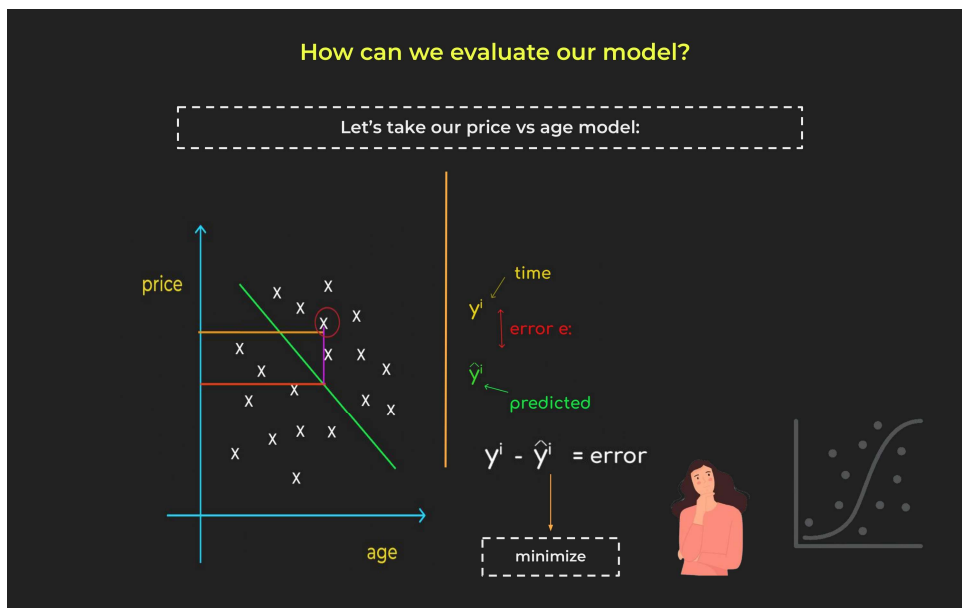


good or Bad?

Evaluation



$$e = y - \hat{y}$$



Red vs Black

$$\sum e_i = e_1 + e_2 + e_3$$

$$M_1 \rightarrow \sum e_i = 10$$

$\left[\sum e_i = e_1 + e_2 + e_3 \right]$
 $M_1 \rightarrow \sum e_i = 10$
 $M_2 \rightarrow \sum e_i = 15$

good or Bad? \rightarrow $e_i = y - \hat{y}$

So, for multiple points..

It will be the sum of the error divided by m : $\frac{1}{m} \sum_{i=0}^m (\hat{y}^{(i)} - y^{(i)})$

Can we use sum of errors as a metric?? \rightarrow NO

	Error	Value
1	e1	3
2	e2	4
3	e3	-7

$\sum_{i=1}^m e^i = 0$
 \downarrow
 Perfect model?? **WRONG!!**

$\sum (y - \hat{y})$

$\sum |y - \hat{y}|$
 m
 Mean Absolute Error

$y = |x|$

$y = x^2$ (Square)
 $\sum |y - \hat{y}|^2$
 $\sum e_i^2$
 m

Can you compare two models?

Best among two models ?

	MSE
Model 1	21
Model 2	31

M1 IS BETTER!!

Issue with MSE or MAE,

- > it gives a value between 0 to ∞
- > value depends on the range of values predicted.
Can be in 100s, 1000s, or 10000s

How small of an MSE is good enough?

Is 1 a good mse score, or 10, or 100?

This makes a problem since there is no common scale/range to measure the performance.

Don't Recog R^2

SOLUTION ??


Mod $\rightarrow \frac{1}{m} \sum_{i=0}^m |\hat{y}^{(i)} - y^{(i)}|$

Square $\rightarrow \frac{1}{m} \sum_{i=0}^m (\hat{y}^{(i)} - y^{(i)})^2$

MAE ABSOLUTE

MSE SQUARED

LOW ERROR = BETTER MODEL



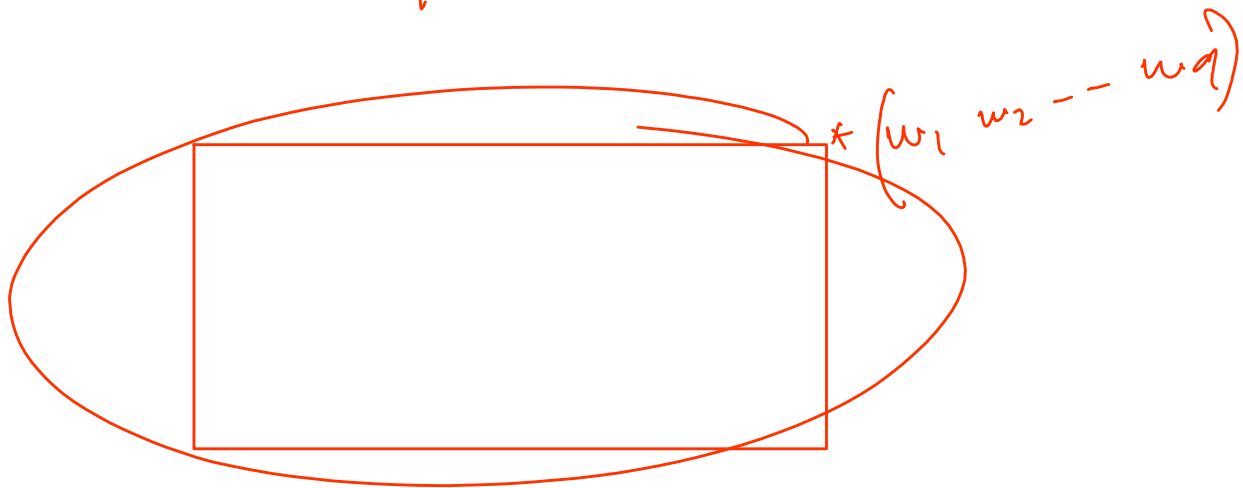
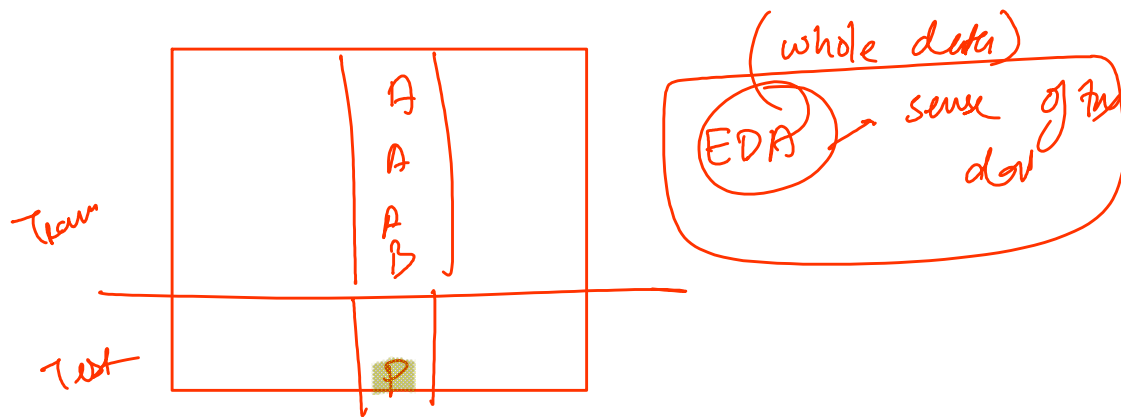
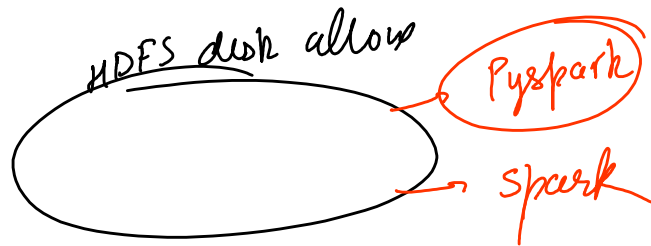
(Lim Reg 2)

Doubts \rightarrow

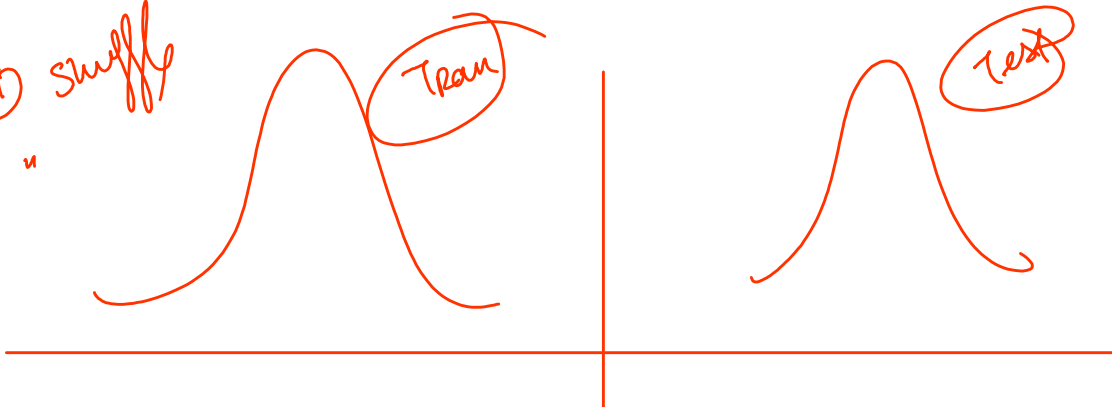
let \rightarrow find the parameter

$fit \rightarrow$ find the parameter

$(fit - transform) \rightarrow$ find param + transform

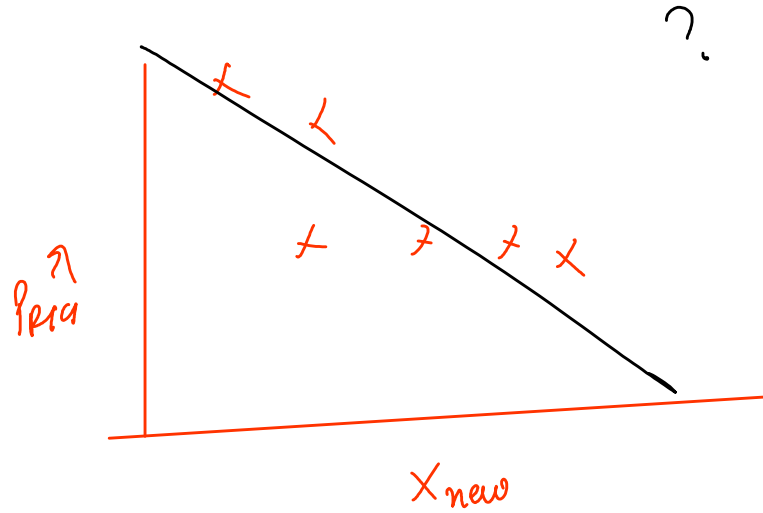


① shuffle

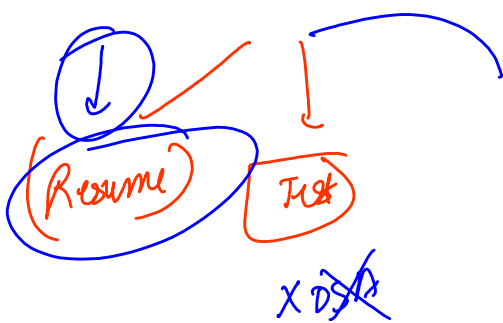


Label encoding → Business Units → Stores Test

$\begin{pmatrix} S \\ M \\ H \end{pmatrix} \rightarrow \begin{matrix} 1 \\ 3 \\ 5 \end{matrix}$



Scaler - fit - transform (of train)



Retrain *

Data Augment