

MACHINE LEARNING

Machine learning → The subfield of computer science that gives the computer the ability to learn without explicitly programmed.

↓
Patterns in the data

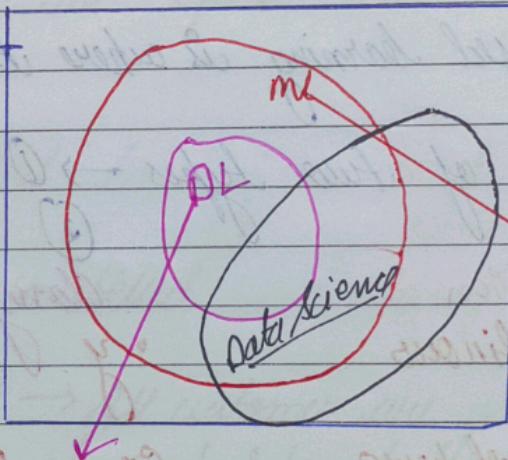
Machines learns pattern from the data and replicate the same in future.

* Why ML is famous these days?

- Advance Processor
- Data is now oil

* Difference b/w AI, ML, DL

AI ←
↓
Describes how computer and technology mimics human intelligence.



DL is specialised ML algo.
that mimics human brain
→ Multilayer Neural Network.

• lot of data
↓
Process this data
 (1) Advance processing

Valuable insights

↓
improve business performance

→ ML focuses on creating algo and statistical model to let computers learn and make prediction without explicitly programmed.

Types of machine learning

- 1) Supervised
- 2) Unsupervised
- 3) Semi-supervised
- 4) Reinforcement learning

1) Supervised ML →

Let's say we have a data set of houses contains
 ① area of house ② no of rooms ③ price of house

Independent Variable

Dependent/Target

Now for new data we predict the target variable
based on independent variable.

Here past data act as supervisor.

Supervised learning is where the data we fed is labeled.

it is of two types → ① Regression
② Classification

Regression

y is continuous

e.g. prices of house

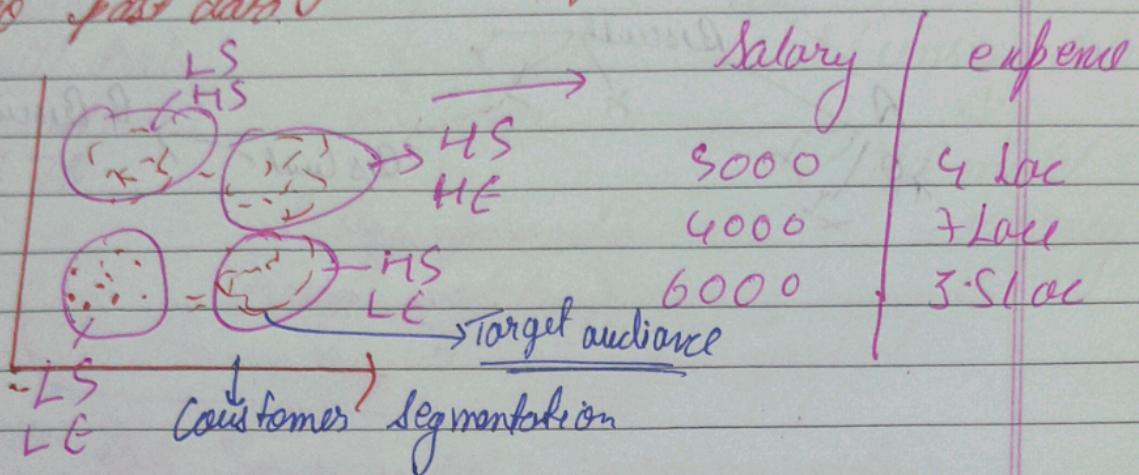
Classification

y is discrete

e.g. Pass/Fail

2) Unsupervised learning.

No supervision \rightarrow No 'y' \rightarrow no target variable \rightarrow that means no 'fair data'.



* We try to find similar groups in data and classify the

ex \rightarrow in above ex created 4 groups (HS, LC) (LS, LC).

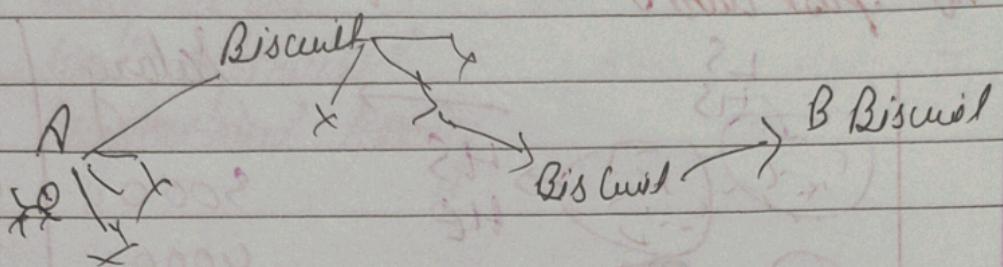
0. we give 20% discount to HS, LC persons so they spend more on expenditure.

4-3) SemiSupervised Learning \rightarrow Combination of both supervised & unsupervised

Netflix/prime \rightarrow All customers are grouped/divided into cluster (USL)
 The they suggest movies (SL)

Vote \rightarrow We convert semi-supervised to either SL or USL

↳ Reinforcement Learning. \rightarrow concerned with how intelligent agent take action in an environment to maximize the award.



Even though I am trying to teach my dog to sit, he just wants to play.

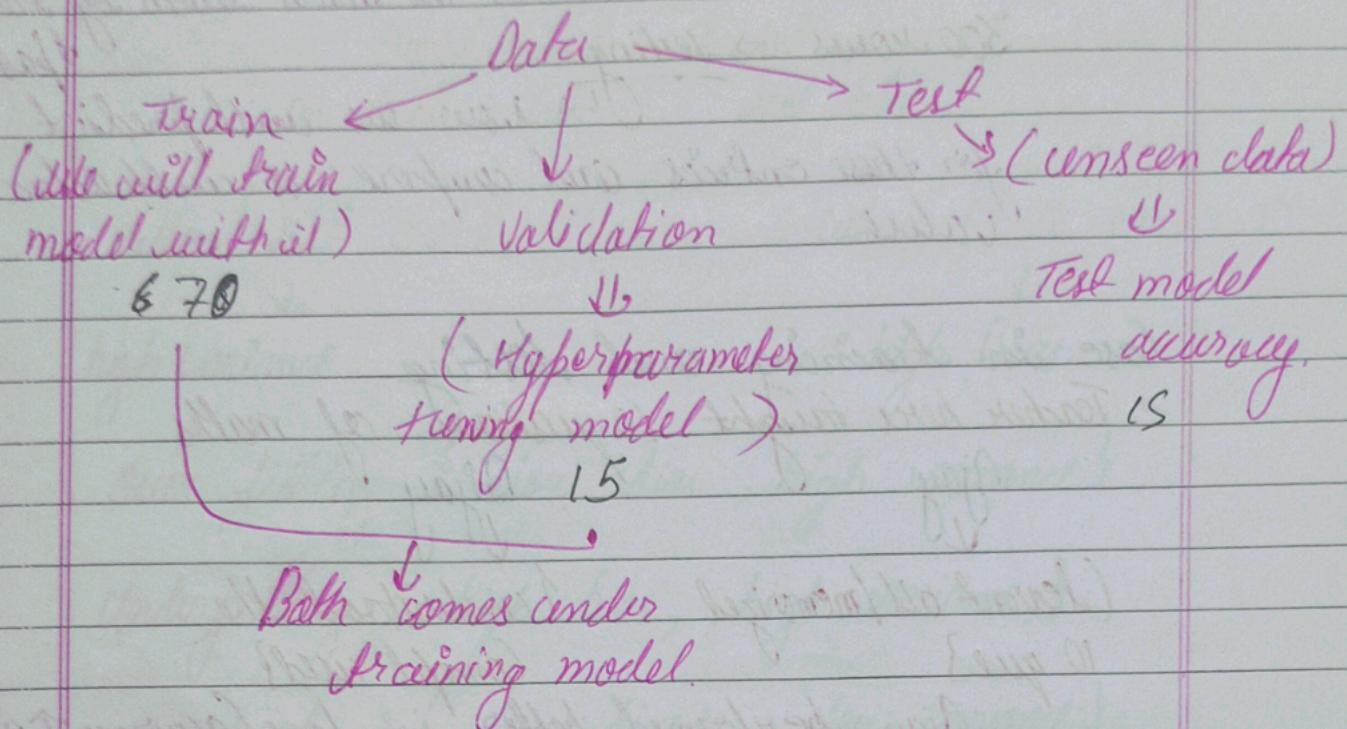
(2) All (1,2,3) above is what we made in training
all 3H or two with one step up. 0
so that dogs do want help with or worked.

→ Reinforcement learning is different from supervised learning.
In supervised learning, we have labels for each data point.

→ Reinforcement learning is similar to supervised learning but without labels.

↳ Reinforcement learning with ML

Train, Validation and Test data



Model training → given the data estimate the prediction function by minimizing error.

area of house	Rooms	Price
2	4	3
3	6	4
4	8	5
5	10	6

} This will generate a pattern,
 $y = f(u)$

with this our model will learn pattern, and predict the price based on different entries.

$$y = mx + c \quad \{ m \text{ is slope} \\ c \text{ is intercept} \}$$

now lets say we have 1000 hours, so we will use 700 hours \rightarrow training

$\downarrow \rightarrow$ The model learn by trend pattern
300 hours \rightarrow testing

$\downarrow \rightarrow$ Now we will predict for these entries and compare it with actual value

Case in Training and Testing

• Teacher here taught 10 question of maths

Ajay

Bijay

(learnt all / memorized 10 que)

{understood the concept used}

Ajay performed better in class (accuracy $\uparrow 98\%$)
10 new que came in board exam

\rightarrow Here ajay have memorized only these 10 que, but when new problem came he failed, scored (70%) while bijay scored (85%).
 \rightarrow Bijay performed better (85%).

\Rightarrow This is case of overfitting

in machine learning, when our model score/gives high accuracy. (model perform well on train data, but failed in test data).

\Rightarrow opposite of this is underfitting

Test accuracy $\rightarrow 40\%$, test accuracy - 35%.
 \rightarrow our model did not learnt anything

⇒ Bias / Variance

Data
Model trained
on train data
Tested on
test data

When test error ↑

high variance

when train error ↑

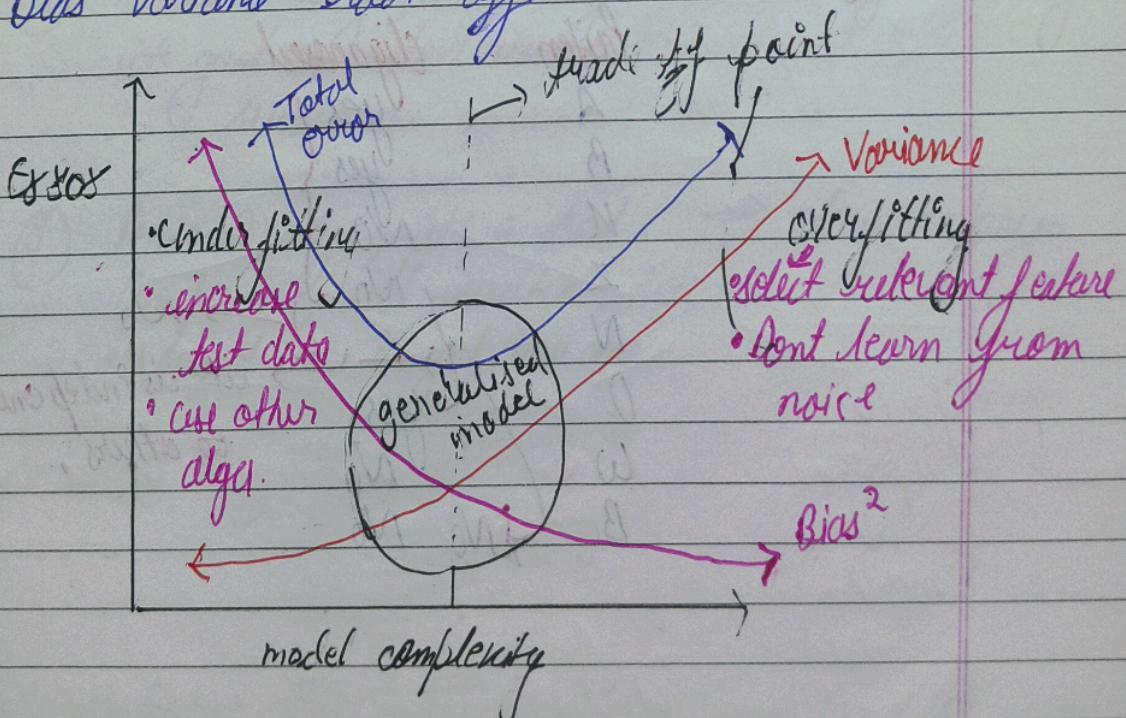
high Bias

over fitting { low bias, high variance }

underfitting { high bias, high variance }

generalised → { low bias, low variance }

⇒ Bias Variance Trade off



- Missing Values And its treatment

missing data → some data points are not filled/missed

Reason → Human error

data is corrupted.

missing data occurs when some info is not saved/shared in our data set.

- Missing Completely at Random (MCAR)

The data missing is independent of observed and unobserved data/missing data.

Survey → if person is diagnosed

patient diagnosed

A yes |

B yes |

C No |

D No |

E / → No

F Yes |

G No |

H No |

it is independent of others.

② Missing at Random (MAR) → missing data depends on observed data not on missing data survey → To income of people. itself

family → income of all family members

lets say one people in family don't want to declare their income due to x-y-z reason

③ Missing data not at Random (MNAR) → missing value depends on missing data itself.

Eg → Thought of employee / satisfied with their job ✕ Income.

① Person 1 says, i am satisfied ⇒ high salary
its good and don't want to switch

② Person 2 says, its a bad company will switch
as soon as i get new one

Missing Value Treatment

A rule says whenever you want to treat missing data, {always as business team}

→ It will treat it like this, it will affect this, should i go on?

Why missing value treatment → For analysis & ML prediction
it is important to deal with missing values

Note → if missing value is less than 1%, and data is huge in amount, drop the missing data

2 • if its greater than >1%, impute/replace the missing data

Imputation

numerical

categorical

if outlier
is present

if not

mode is used
to impute

mean is
used

mean is
used

3 • if any col has missing value greater than 40%.
drop that column.

4) impute the missing value with a random number, or with a value which is not possible

Mark

100
80
60
→ -1/Na

90

70

50

80

60

20

Note → whenever we fill missing value, try to plot a histogram and if the distribution is normal & no many outliers were present, fill with mean in that case.

Imbalanced data.

When we have **1** target variable and it has 2 outcomes
→ binary classification problem.

\rightarrow	Sugar level	Calistrol level	diabetes (g/l)	$\delta(y)$
20	102	0	0	0
38	100	0	0	0
31	105	0	0	0
77	193	0	0	0
89	201	1	1	1

data says we have 90% of the y is 0
or not (non-diabetic)

and you haven't made any model and your
predicted all values to 0 → got accuracy
90% → why?
because 90% is 0

→ This is called class imbalance

im + balanced class

When one class has very high percentage as
compared to other class,

→ in our case → class 0 (non-diabetic) →
class 2 → 1 (diabetic)

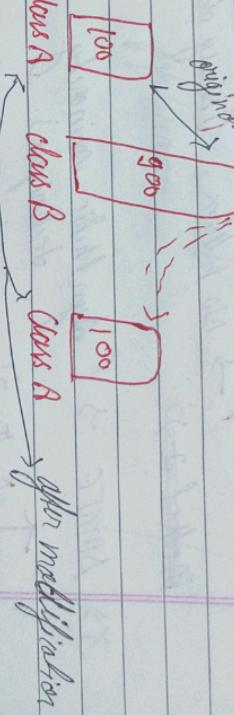
- We don't want our model to learn the only data pattern that is in majority
That's why we need to deal with class imbalance

Class Imbalance → When the difference in the count of both classes is very huge.

Very huge → 80% → Class A → This is not.
20% → Class B → imbalance since

Solution for Class Imbalance

① Under Sampling (DownSampling)



Name: If imbalanced data may be Class B 90%.
We can down sample class 10% may have

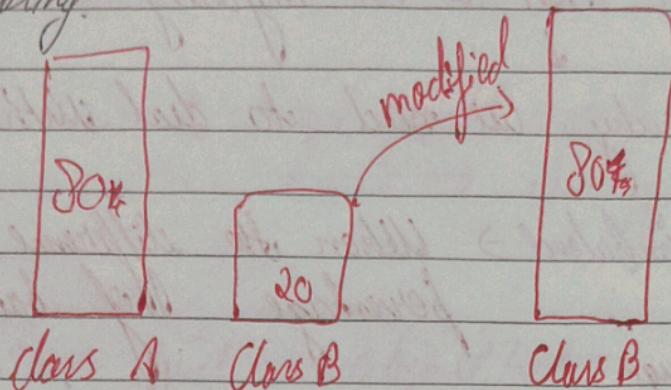
New Class A → 100 values Overall class 50%

Class B → 100 values

Disadvantage → A lot of data is lost.

→ you lose more memory data waste less fine
like hammed about resources use but the

2-5 Ovrsampling.

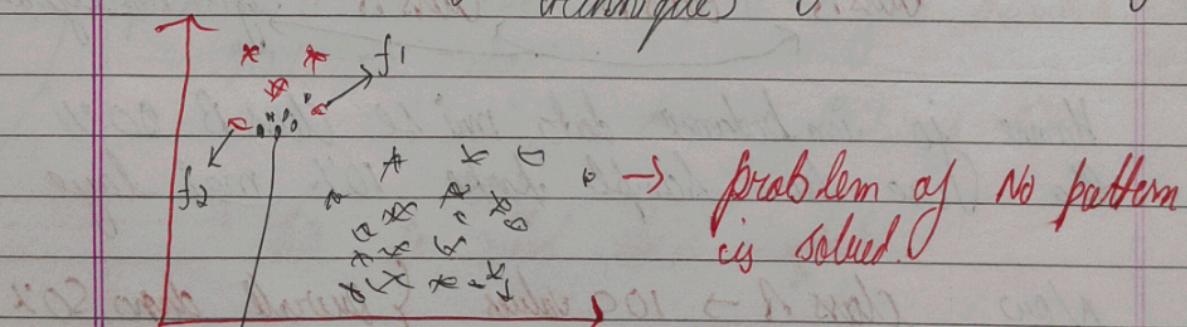


Method 1 → repeat class B to get similar level of class A

Disadvantage → ML is about learning patterns but here data is repeated so no new pattern is learnt from data.
→ No pattern, Noisy data.

Method 2 →

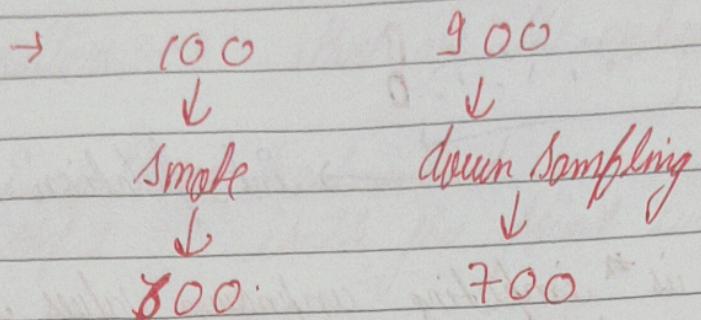
3-5 SMOTE → (Synthetic minority overampling technique)



→ like datapoints will always be together, this is reason
synthetic value is effective.

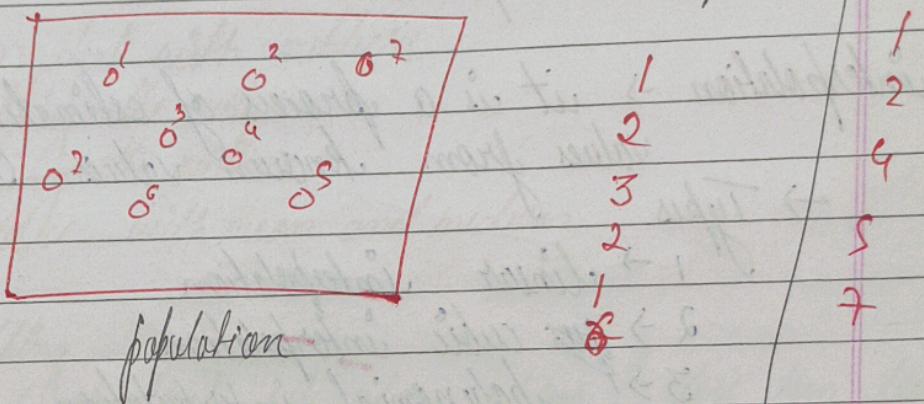
get by aug, ~~max~~ mode, knn etc

Use in industry
How to use Smoke of down sampling in data

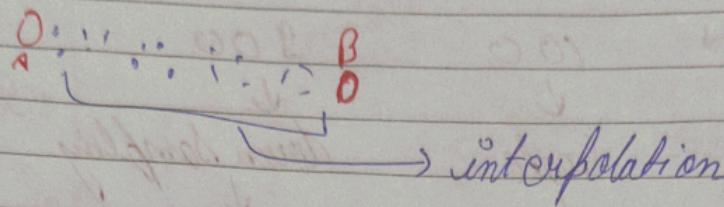


* In practical implementation

With replacement Without



Data Interpolation



- interpolation is ~~a~~ imputing unknown values based on known values.
- it is where we impute new data values b/w two known data points

interpolation → it is a process of estimating unknown values from known values

→ Types

- 1 → linear interpolation
- 2 → ~~cubic~~ cubic interpolation
- 3 → polynomial interpolation

Use Cases → 1 → to fill missing values
2 → to compute new values

Outliers And its treatment

outliers → These are the extreme values

in outlier treatment

Step 1 → check five point summary

To check outlier using plots

① → distplot, boxplot

How to deal with outliers

- ① dropping the outlier
- ② capping the outlier
- ③ replace with mean and median

① dropping

→ drop all values ~~above~~ (less than lower fence & greater than upper fence)

② imputing

→ replacing the outlier values with mean and median

③ capping → capping the outliers with the nearest value which is not an outlier

(lower fence, upper fence)

* Feature extraction and selection

steps we have covered till now

- understanding problem statement
- Data ingestion
- Data preparation / preprocessing

- missing value treatment
- class imbalance
- outlier treatment

next steps ->

- Data encoding
- Feature Selection

Creating new features

Selecting right features

modifying the exist feature

Feature extraction → process of selecting and extracting the relevant features from raw data.

CN → out of 1000 → 600 are relevant.

A question arises out of that → what if it can be a problem

There comes a concept, → curse of dimensionality

Curse of dimensionality

- ↳ with increase in no. of features
- ① Model training becomes computationally expensive
- ② Model interpretation becomes complex

How to solve

- ① Create new features

↳ CN

Distance	Time	Speed = $\frac{D}{T}$
100	1	
110	2	
120	3	
130	4	
...	...	

- ② Modifying existing features

- ① changing data types

Date \rightarrow day / month / year

- ② feature scaling

(optional)

Avg house	rooms	Cost
1500	2	60
2100	5	65
2500	4	100

feature scaling \rightarrow it is the process to scale all the features to same scale

why \rightarrow

- Many of the algorithms are distance based, that's why if high magnitude comes, it becomes computationally expensive.
- interpretation becomes easier.

Types of scaling.

① Standardization

$$Z \text{ score} = \frac{x_i - \bar{x}}{\sigma} \quad \text{for SNO}$$

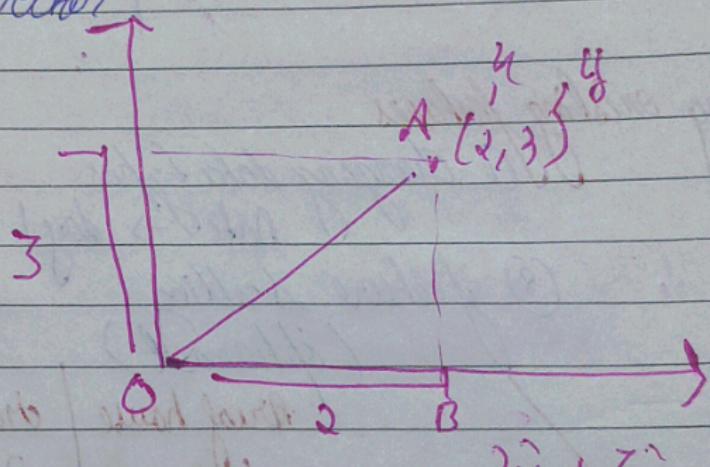
$\bar{x} = 0, \sigma = 1$

② Normalization (min, max Scaler)

\hookrightarrow reduced to [0 to 1]

$$\left\{ \begin{array}{l} \text{Scaled} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \end{array} \right.$$

③ Unit Vector



$$OA^2 = OB^2 + AB^2$$

$$\frac{OA}{||OA||} = \sqrt{3^2 + 2^2}$$

$$\frac{OA}{||OA||} = \sqrt{13}$$

$$\hat{u}_1 = \frac{2}{\sqrt{13}}, \frac{3}{\sqrt{13}}$$

unit vector

(7) Selecting the right features

(1000) \rightarrow (20)

methods

① filter method

calculate correlation with target variable, features with target value will be relevant one

② embedded method.

③ wrapper method

* PCA

df['total bill']
std
~~mean~~ μ, σ $\frac{x_i - \mu}{\sigma}$
{ calculate } $\{ \mu, \sigma \}$ { apply the calculation
on dataset using μ, σ }

This is called as

fit

train data

This is called as

transform

test data

* Scaling should be done after train-test split.

~~Scaling~~ Data encoding.

The machine only understand numerical no. \Downarrow

So we need to convert object data to numerical data.

Types for that we need encoding \Downarrow

① Nominal / One hot encoding

② Label and ordinal

③ Target guided ordinal encoding.

① Nominal / One hot encoding.

Converts categorical to numerical
No order in data.

Status

Single

Married

Engaged

Single

Single

Married

Single \rightarrow 1 0 0 0 1
Married \rightarrow 0 0 1 0 0
Engaged \rightarrow 0 1 0 1 0

disadvantages of One hot encoding

→ A column has many categories

→ The columns we create are called as dummy variables

Note → ~~the~~ n-1 dummy col. can explain n dummy col.

② Label and ordinal encoding

Label encoding →

→ Assign numerical data to each category.

Status	→	Status
M		0
S		1
D		2
E		3
M		0
F		3

disadvantages

→ good for ordinal data, as for nominal data it can learn pattern.

Ordinal

→ High School = 1

College = 2

Master = 3

Job = 4

→ used for ordinal data.

* Target Guided Ordinal Encoding

→ does encoding based on their relationship with the target variable.

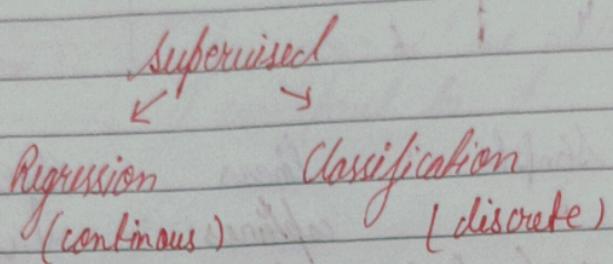
→ useful when we have large no of unique categories in categorical value.

→ categorical groups with mean/median of corresponding target variable.

Regression

Simple Linear Regression

comes under supervised machine learning algorithms

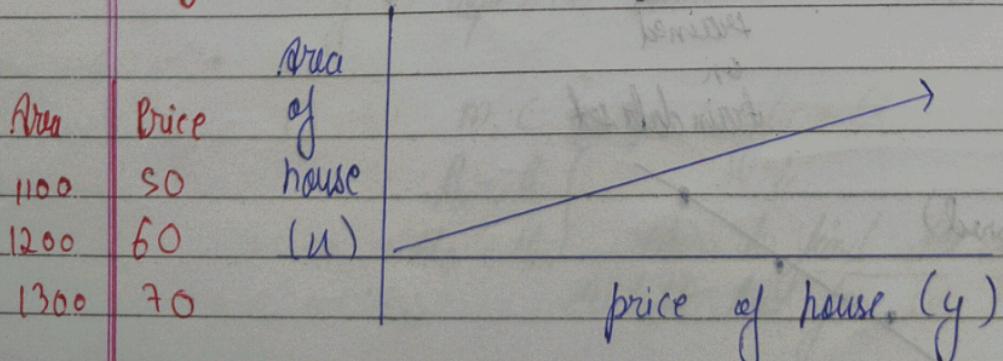


Regression \rightarrow it is to establish relationship b/w two or more than two variables.

Linear \rightarrow The relationship established is linear.

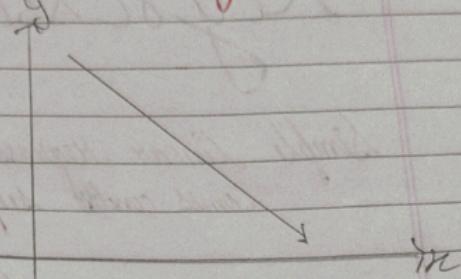
** Simple linear regression attempts to determine the strength and characteristics of relationship b/w one independent variable (x) and another dependent variable (y). **

e.g. \rightarrow predict house of price based on area



e.g. → Selling price of car w.r.t. to years used

x age	y price
1	11
2	10
3	8
4	6



Simple linear Regression
captures understand relationship
Only one (variable) linear relationship

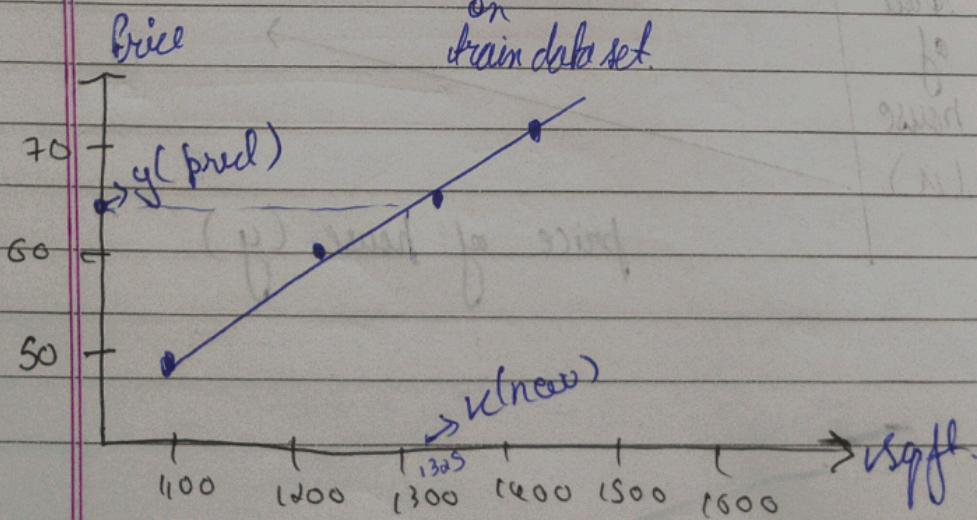
→ There were only 4 data points shall we understand the trend

→ if many data points exist we need to train our model

New no. of rooms → Model → Price of House

trained

on
train data set.



Here the line's eq = $y = mx + c$

$$m \rightarrow \text{slope} = \frac{y_2 - y_1}{x_2 - x_1}$$

$c \rightarrow$ intercept (\rightarrow where it cuts y axis at $y=0$)

using $y = mx + c$ for now 'u' we will get 'y' correspond to it.

$$\hat{y} = mx + c$$

y estimator \rightarrow you have predicted price of house based on m, u, c

\rightarrow This also written as $\hat{y} = \beta_0 + \beta_1 u$

$$\text{or } H_0(u) = \phi_0 + \phi_1 u$$

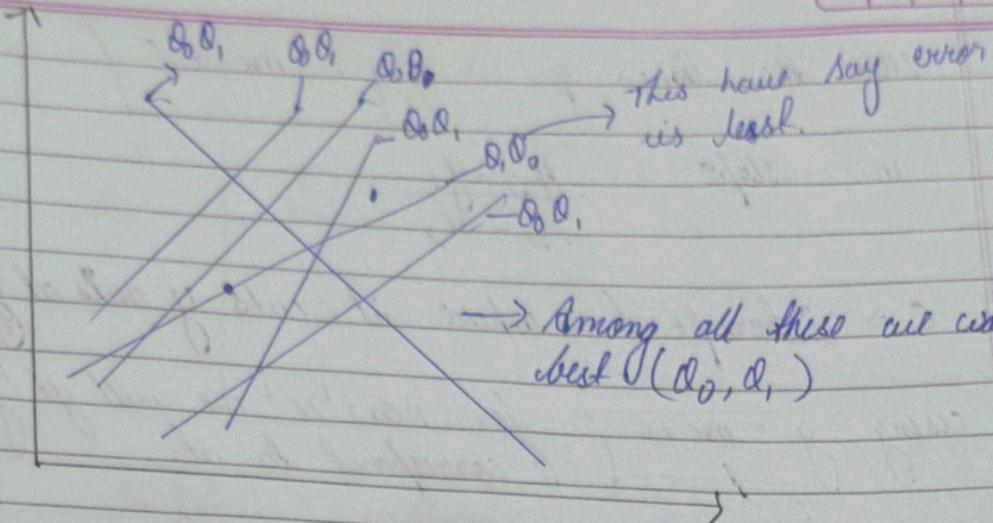
Hypothesis representation
Here

$u \rightarrow$ independent variable

$y \rightarrow$ dependent variable

$m, c \rightarrow$ coefficients
 $\beta_0, \beta_1 \rightarrow$ coefficients
 $H_0, H_1 \rightarrow$ functions

How to find the coefficient suitable best for us



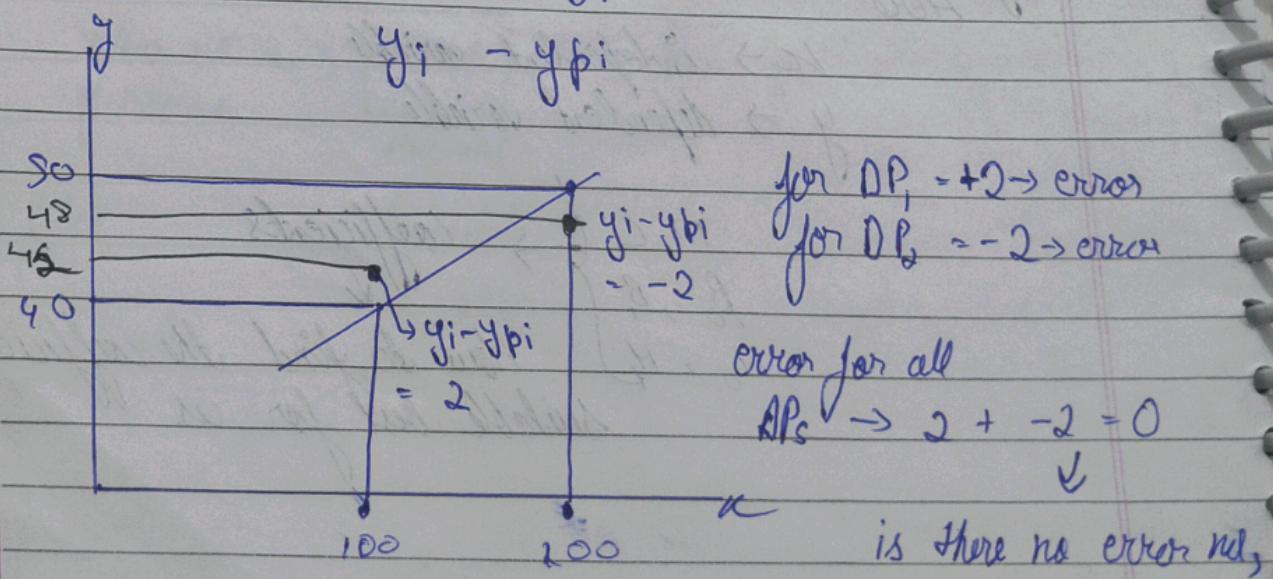
We want that coefficient ($\alpha_0 \alpha_1$, or $\alpha_0 \beta_1$) where error is least.

that coefficient also called as $\alpha_0 \alpha_1$.

for all lines we want $\rightarrow y_{\text{actual}} - y_{\text{pred.}} = \text{last error}$

for $y = mx + c$, (m, c) to be those best optimal (α_0, α_1)

$$\text{error} = y_{\text{act}} - y_{\text{pred}}$$

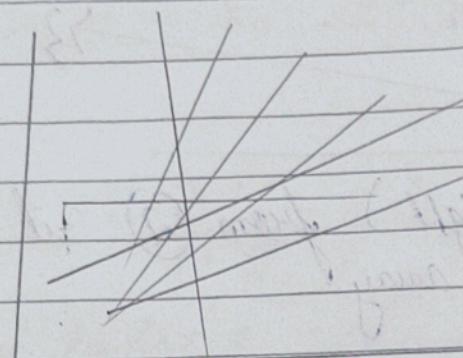


To solve this problem instead of absolute error we measure squared error. $\rightarrow (y_{act} - y_{pred})^2$

$$\text{error} = (df_{\perp 1})^2 + (df_{\perp 2})^2$$

$$= 2^2 + 2^2 = 8$$

$$= 4 + 4 = 8$$



\rightarrow now we will calculate (error) for all the lines and choose the line with least error
[error \rightarrow squared error]

$$\min (y_{act} - y_{pred})^2$$

Best line $\rightarrow \min \sum_{i=1}^n (y_i - \hat{y}_i)^2$

\downarrow

$\hat{y} = mx + c$

Least squares method
or
 $n = \text{no. of lines}$

$$H_0(u) = H_0 u + H_0$$

Ordinary least sq
method

$$\min \sum (y_i - H_0(u))^2$$

$$\min \sum_{i=1}^n (y_i - Q_0 - Q_1 u_i)^2$$

* for best fit line

use least MSE to
be least.

Sum of Sq. error $\rightarrow \min \sum_{i=1}^n (y_i - \hat{y}_i)^2$

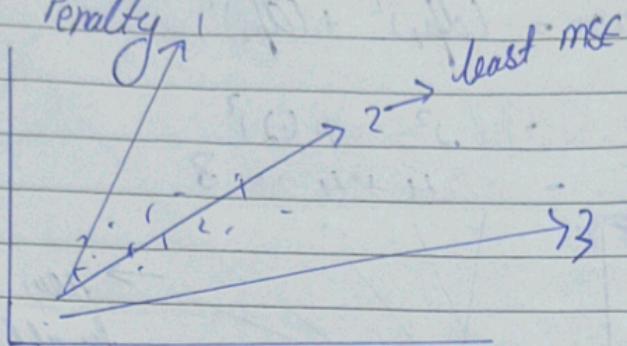
when we divide it by n

it becomes (MSE)

$$\text{mean sq. error} = \frac{\min \sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

→ MSE is used as CF to get optimal (Ch 1,)

why CF is Penalty



deviating (left or right) from ② it is penalty as we are moving away

How to choose the best line

1 → chose right combination of line

2 → find MSE line

3 → find Best Line

The predicted line should be made in such a way that over all sum of error should be minimum or least

Sum of squared error

$$SSE = \sum_{i=1}^n (y_i - \hat{y})^2$$

\downarrow \downarrow
 y_{actual} $\hat{y}_{\text{predicted}}$

$\frac{n}{\downarrow}$ mean squared error

also called as cost function

$$S.S.E = \sum_{i=1}^n (y_i - \hat{y})^2$$

$$\text{Cost function} = J(\theta_0, \theta_1) = \frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n} \quad \text{for } \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

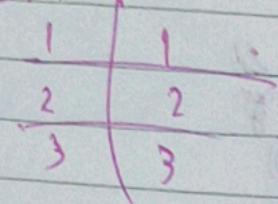
} n = no. of data point
 $y_i \rightarrow \text{actual value}$
 $\hat{y} \rightarrow y_{\text{predicted value}}$
 $\hat{y} = \theta_0 + \theta_1 x /$

$$\min(J(\theta_0, \theta_1)) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 \quad \theta_0 + \theta_1 x / \theta_0 + \theta_1 x$$

Aims \rightarrow minimize cost function to get optimal coefficient.

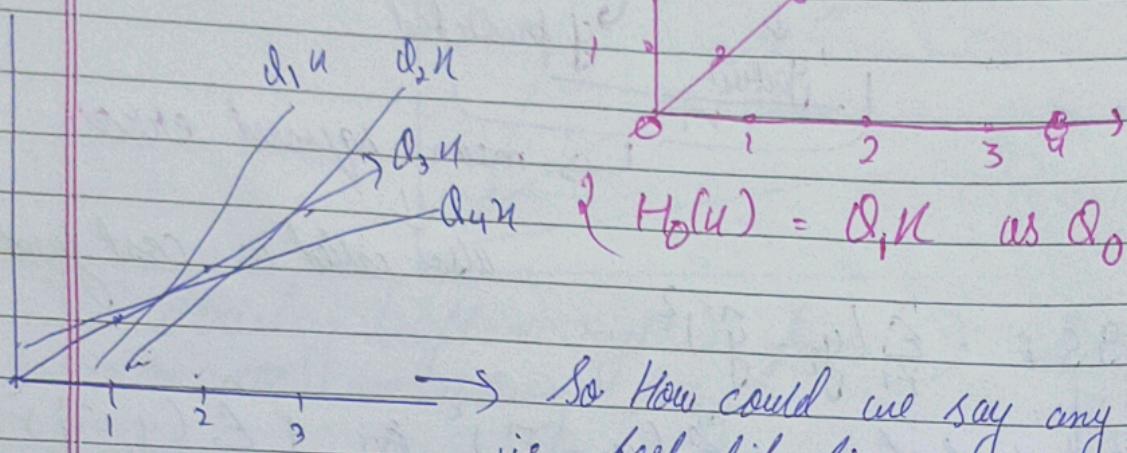
The data point were

like



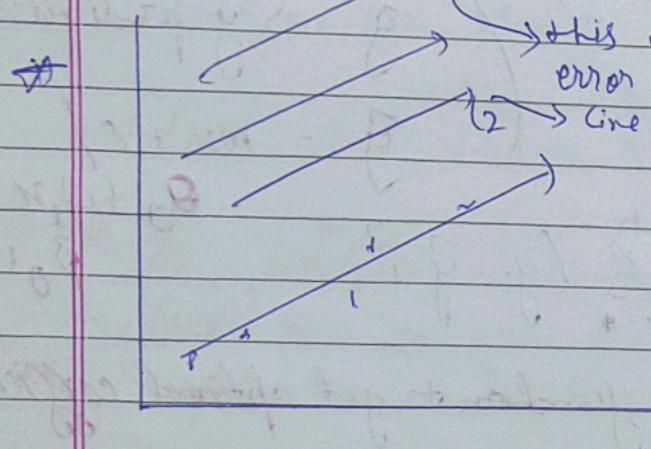
$$\hat{y} = H_0(u) = \theta_0 + \theta_1 u$$

Assume line intersect
at $u=0$



So How could we say any line
is best fit line

\rightarrow we will choose that one line with
the least cost function.



(1) \rightarrow this line has more
error than this

(2) \rightarrow line
in short $m_1(l_1) > m_2(l_2)$
(if at any point ($m_1(l_1)$)
 $> m_2(l_2)$)

then we going again for
best fit line

* Shape of cost function of n line will always
be a parabola.

- for the best fit line we have got an algorithm called as convergence algorithm.

Convergence Algorithm

- keep moving new best fit line in direction of all data point until error (mse / cf) is reduced as compared to previous line

- Repeat until convergence

$$m_{\text{new}} = m_{\text{old}} - \eta \frac{\partial J(m)}{\partial m_{\text{old}}}$$

$$c_{\text{new}} = c_{\text{old}} - \eta \frac{\partial J(c)}{\partial c_{\text{old}}}$$

→ learning rate (beta)

or

Repeat until Convergence (Simultaneously)

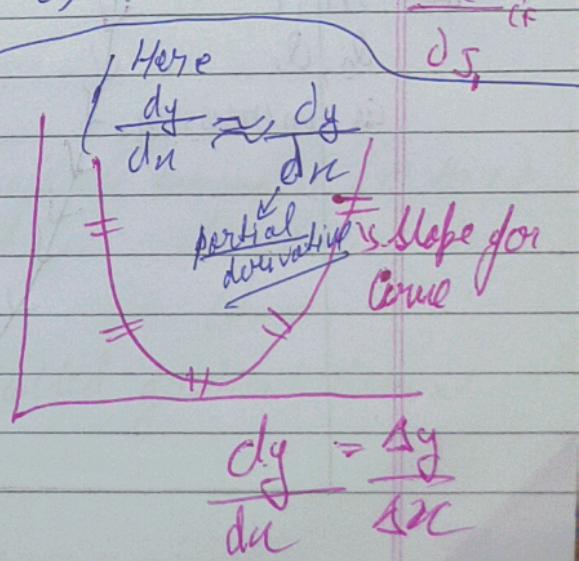
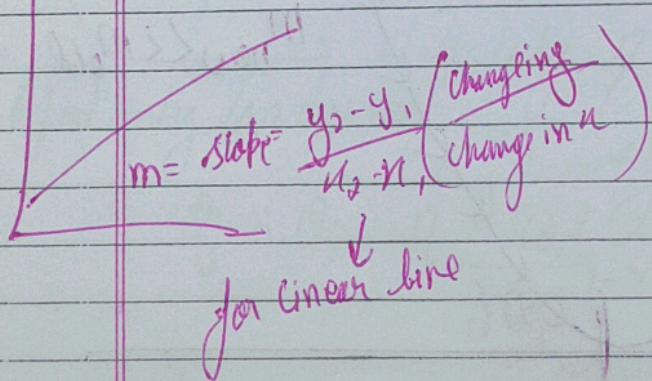
$$\theta_j : \theta_j - \eta \frac{\partial J(\theta)}{\partial \theta_j} \quad \theta_i = \theta_i - \eta \frac{\partial J(\theta)}{\partial \theta_i}$$

$J = \theta_0, \theta_1$

$$\theta_0 = \theta_0 - \eta \frac{\partial J(\theta)}{\partial \theta_0}$$

Here $\frac{dy}{dx} \approx \frac{\Delta y}{\Delta x}$

partial derivative $\frac{\partial y}{\partial x}$ slope for curve



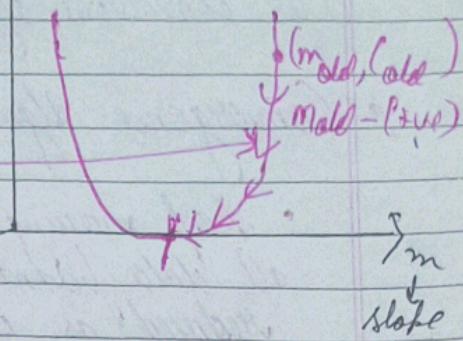
$$\frac{dC^n}{dn} = n \alpha^{(n-1)}$$

$$m_{\text{new}} = m_{\text{old}} - \frac{n \Delta C_p}{2 d_{\text{mold}}}$$

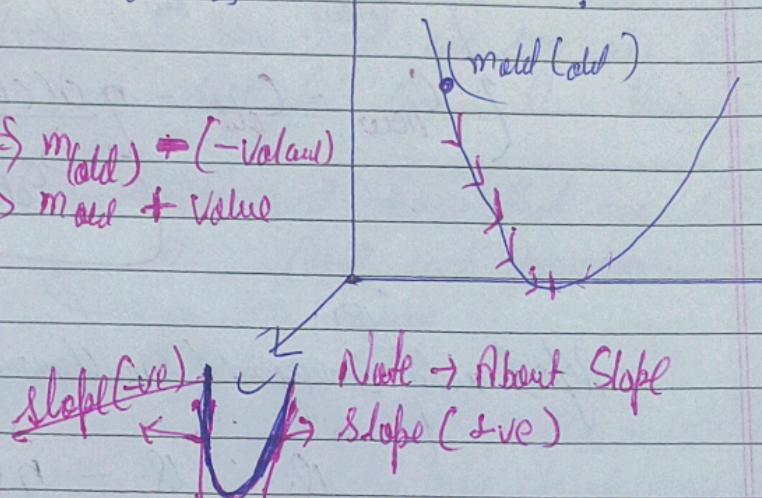
$$C_{\text{new}} = C_{\text{old}} - \frac{n \Delta C_p}{d_{\text{old}}} \quad \text{coefficient}$$

This conversion algo also called as gradient descent

lets say it is $\Rightarrow m_{\text{old}} \leftarrow (-\text{Value})$
 $\Rightarrow m_{\text{old}} + \text{Value}$

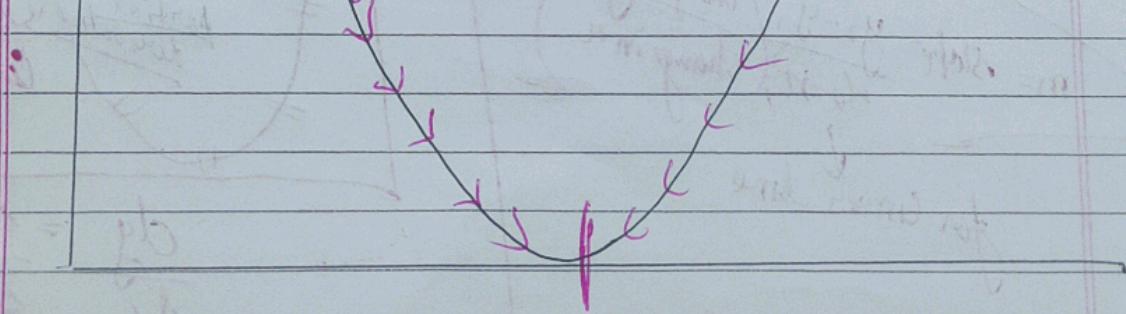


If we initialize at this point?



when slope is -ve
~~m, C~~
~~d_o, d_i~~
 is increasing

Slope is +ve
 when (m, C)
 in getting reduces
~~d_o, d_i~~
 $m_{\text{new}} < m_{\text{old}}$



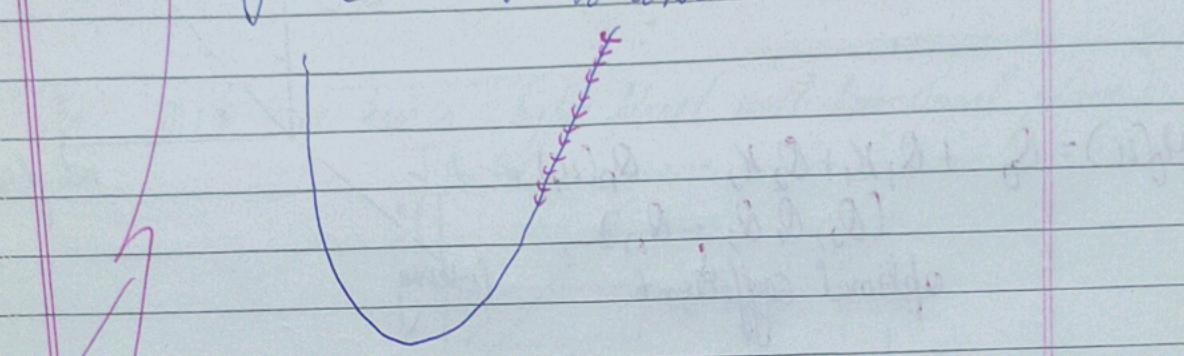
$$m_{\text{new}} = m_{\text{old}} - \eta \frac{\partial C(F)}{\partial m_{\text{old}}}$$

$$C_{\text{new}} = C_{\text{old}} - \eta \frac{\partial C(F)}{\partial C_{\text{old}}}$$

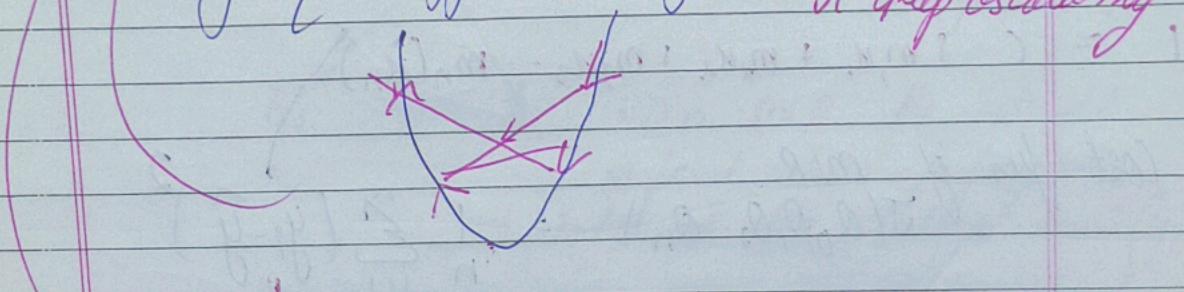
η is learning rate
↓

decreasing conversion speed

if η is very small \rightarrow takes very very high time to convert



if η is bigger (very, very) \rightarrow overshoot the minima or keep oscillating



That's why η should be ~~minimum~~ (0 to 1)
not so large not too big

$$\theta_j : \theta_j - \eta \frac{\partial J(\theta)}{\partial \theta_j} \rightarrow \text{also called gradient descent.}$$

Multiple Linear Regression

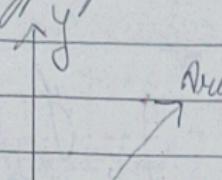
Only one feature can not be enough to predict target variable.

There are multiple independent variables that contribute in predicting target variable and this is called as **multiple linear regression**.

graph of MLR looks like a Hyperplane

Data set

No of Rooms	Area	Parking	Price
1	100	0	10000
2	150	0	15000
3	200	0	20000
4	250	0	25000
5	300	0	30000
6	350	0	35000
7	400	0	40000
8	450	0	45000
9	500	0	50000
10	550	0	55000
11	600	0	60000
12	650	0	65000
13	700	0	70000
14	750	0	75000
15	800	0	80000
16	850	0	85000
17	900	0	90000
18	950	0	95000
19	1000	0	100000
20	1050	0	105000
21	1100	0	110000
22	1150	0	115000
23	1200	0	120000
24	1250	0	125000
25	1300	0	130000
26	1350	0	135000
27	1400	0	140000
28	1450	0	145000
29	1500	0	150000
30	1550	0	155000
31	1600	0	160000
32	1650	0	165000
33	1700	0	170000
34	1750	0	175000
35	1800	0	180000
36	1850	0	185000
37	1900	0	190000
38	1950	0	195000
39	2000	0	200000
40	2050	0	205000
41	2100	0	210000
42	2150	0	215000
43	2200	0	220000
44	2250	0	225000
45	2300	0	230000
46	2350	0	235000
47	2400	0	240000
48	2450	0	245000
49	2500	0	250000
50	2550	0	255000
51	2600	0	260000
52	2650	0	265000
53	2700	0	270000
54	2750	0	275000
55	2800	0	280000
56	2850	0	285000
57	2900	0	290000
58	2950	0	295000
59	3000	0	300000
60	3050	0	305000
61	3100	0	310000
62	3150	0	315000
63	3200	0	320000
64	3250	0	325000
65	3300	0	330000
66	3350	0	335000
67	3400	0	340000
68	3450	0	345000
69	3500	0	350000
70	3550	0	355000
71	3600	0	360000
72	3650	0	365000
73	3700	0	370000
74	3750	0	375000
75	3800	0	380000
76	3850	0	385000
77	3900	0	390000
78	3950	0	395000
79	4000	0	400000
80	4050	0	405000
81	4100	0	410000
82	4150	0	415000
83	4200	0	420000
84	4250	0	425000
85	4300	0	430000
86	4350	0	435000
87	4400	0	440000
88	4450	0	445000
89	4500	0	450000
90	4550	0	455000
91	4600	0	460000
92	4650	0	465000
93	4700	0	470000
94	4750	0	475000
95	4800	0	480000
96	4850	0	485000
97	4900	0	490000
98	4950	0	495000
99	5000	0	500000
100	5050	0	505000



$$H_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

($\theta_0, \theta_1, \theta_2, \dots, \theta_n$)

optimal coefficient.

$$y_{\text{pred}} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

$$l = C + m_1 x_1 + m_2 x_2 + m_3 x_3 + \dots + m_n x_n$$

Cost fun of MLR

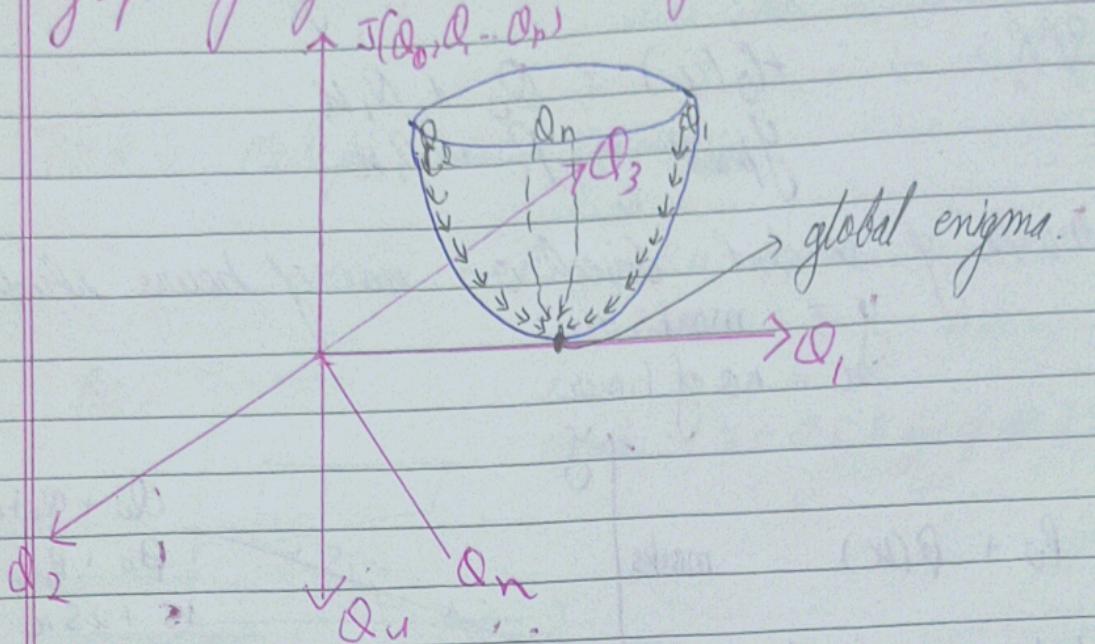
$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

$$(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)$$

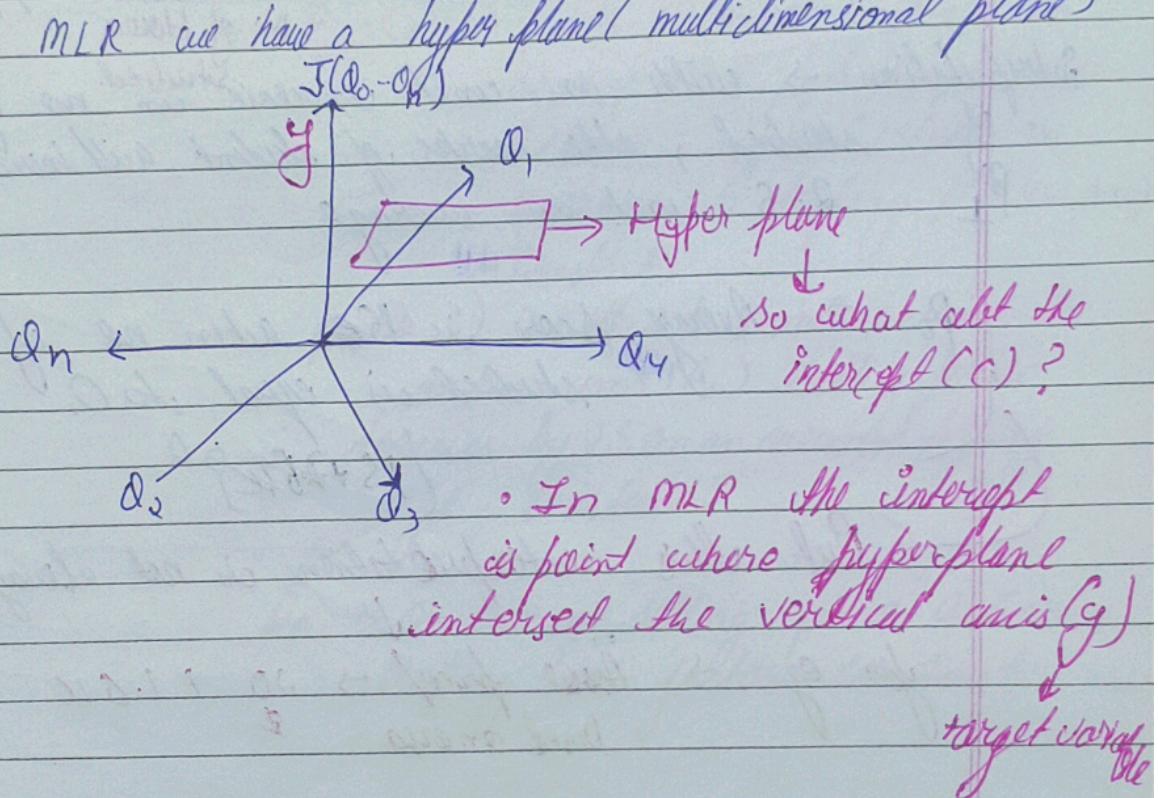
↓

minimize the existing
gradient descent to get optimal
coeff ($\theta_0, \theta_1, \theta_2, \dots, \theta_n$)

graph of gradient decent of MLR



in MLR we have a hyper plane (multidimensional plane)



In MLR the intercept is point where hyperplane intersected the vertical axis (y)

target variable

Interpretation of SLR slope coefficient

$$SIR \rightarrow Y_{\text{SLR}} = \beta_0 + \beta_1 X$$

$$Y_{\text{pred}} = P_0 + P_1 X$$

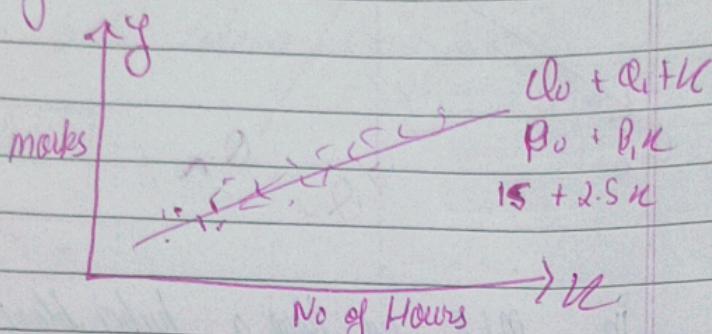
marks of student based on no. of hours studied.

y = marks

x = no. of hours

$$Y_{\text{SLR}} = P_0 + \beta_1 X$$

~~with hours studied~~



interpretation \rightarrow with one unit increase ~~studied~~ in no. of hours studied, the marks of student will increase by β_1 2.5 unit on average

$P_0 \rightarrow$ Average score is 15 when no. of hours studied is equal to 0.

$$(15 + 2.5 \times 0)$$

\rightarrow But this interpretation is not always true.

for ex \rightarrow House price \rightarrow 20 + 1.6 X
based on area \downarrow

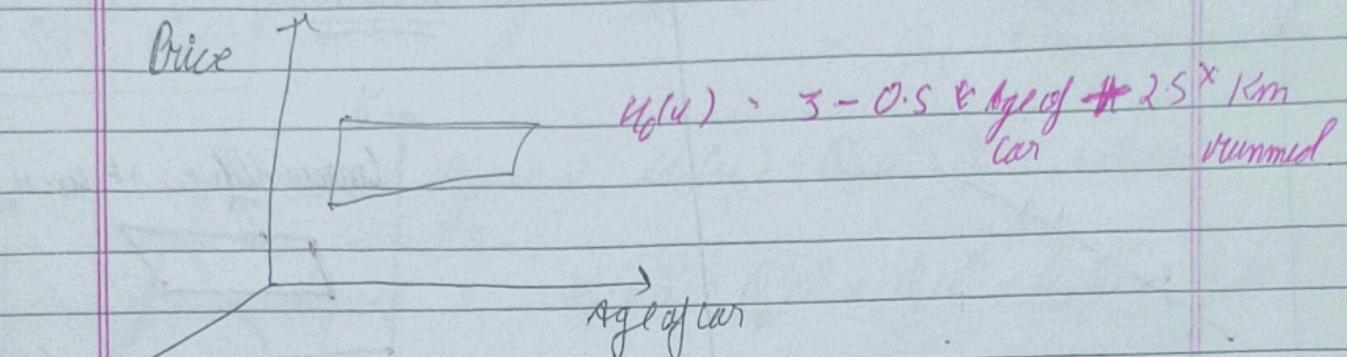
• i.e. Then if $P_0 = 20$ lac on average when area is 0,

But is area is 0, how house price is 20 lac?

The correct interpretation \rightarrow There are multiple interpretations like this is on average price of car.

MLR \rightarrow Multiple linear regression

$$Y_0(u) = \beta_0 + \beta_1 u_1 + \beta_2 u_2 + \dots + \beta_n u_n$$



$\beta_1(0.5)$ \rightarrow with 1 unit increase in age of car, selling price decrease by 0.5. On an average keeping $\beta_2(\text{km runned})$ constant

$\beta_2(2.5)$ \rightarrow with 1 unit increase in km runned (β_2) selling price decreases by 2.5 on an average keeping age constant

β_0 \rightarrow on average selling price is 3 or keeping all factors constant.

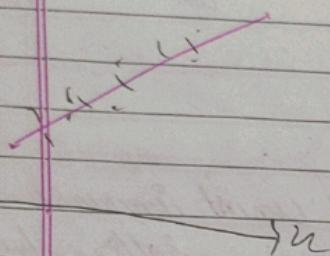
Polynomial ~~Regressions~~ Regression \rightarrow

Till now

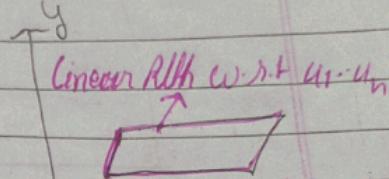
$$SLR = f(u) = P_0 + P_1 u$$

$$MLR = Y_0(u) = P_0 + P_1 u_1 + P_2 u_2 + \dots + P_n u_n$$

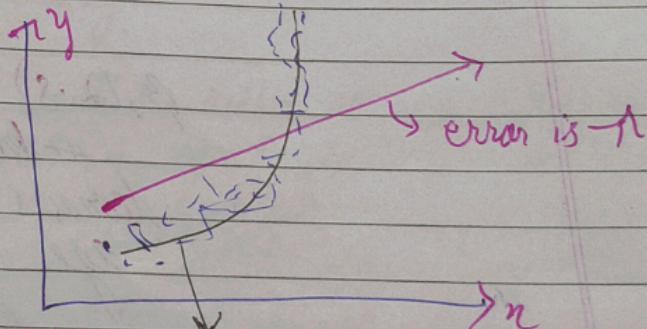
SLR



MLR



if data is like



try to make non linear
best fit line

Polynomial Regression

1 → Simple Polynomial Regression → (1 Dep. Variable
+ Indep. Variable)

$$\text{polynomial degree } 0: y(u) = \alpha_0 = \alpha_0 \times 1 = \alpha_0 u^0$$

$$\text{polynomial degree } 1: y(u) = \alpha_0 u^0 + \alpha_1 u^1 = \alpha_0 + \alpha_1 u$$

+ SLR

$$\text{polynomial degree } 2: y(u) = \alpha_0 u^0 + \alpha_1 u^1 + \alpha_2 u^2$$

$$\text{degree } 3: y(u) = \alpha_0 u^0 + \alpha_1 u^1 + \alpha_2 u^2 + \alpha_3 u^3$$

→ degree-2

→ degree-3

Note → as degree ↑↑ you
might get everything
model says it tries to connect
all points

2 → Multiple polynomial Regression (n → indep. variable)

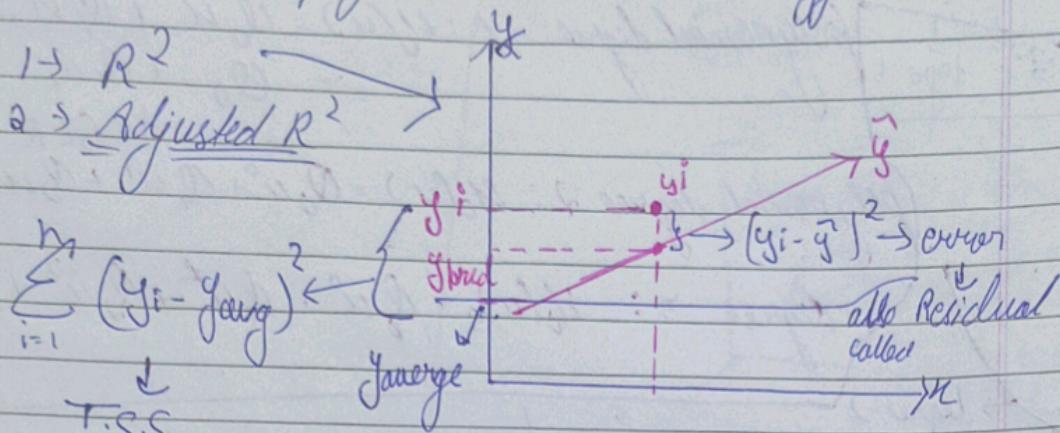
u_1, u_2, u_3, y

$$\text{polynomial degree } 2 \rightarrow y(u) = \alpha_0 + \alpha_1 u_1 + \alpha_2 u_2 + \alpha_3 u_3 + \alpha_4 u_1^2 + \alpha_5 u_2^2 + \alpha_6 u_3^2 + \alpha_7 u_1 u_2 + \alpha_8 u_2 u_3 + \alpha_9 u_1 u_3$$

* Original I.V & power 2 ($u_1^2 + u_2^2 + u_3^2$)
+ Cross product for all ($u_1 u_2, u_2 u_3, u_1 u_3$)

Till now we have learnt ~~to~~ about SLR, MLR, PR
but how we interpret our model performs good?

for that now we have evaluation metric \rightarrow to know performance we have different methods:



Total Sum of Square.

$$x \cdot \frac{\text{error}}{\text{R.S.C}} / \frac{\text{T.S.S.E}}{n} \rightarrow \text{MSE}$$

$\sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow \text{Residual Sum of squares}$

• why (\bar{y}) y average?

$\bar{y} \rightarrow$ if you don't build any model, the output value is the baseline solution.

e.g. → what's the rent of room in your neighbour hood.

• SSR (Sum of square due to regression) \Rightarrow explained variation in y by best fit line

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

SSE (R.S.S) \rightarrow unexplained variation

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

NOTE = R^2 closer to 1 \rightarrow model is close to perfect
 R^2 closer to 0 \rightarrow model is not good.

- TSS (Total Sum of Squared Error) \rightarrow Total variation in y / diff. in y
 w.r.t. + to baseline (\bar{y})

$$TSS = \frac{\text{Explained Variation}}{SSR} + \frac{\text{Unexplained error}}{RSS / S.S.E}$$

~~$R^2 = 1 - \frac{RSS}{TSS}$~~ or ~~$R^2 = \frac{RSS}{TSS}$~~

R^2 \rightarrow coefficient of determination \rightarrow out of total variation,

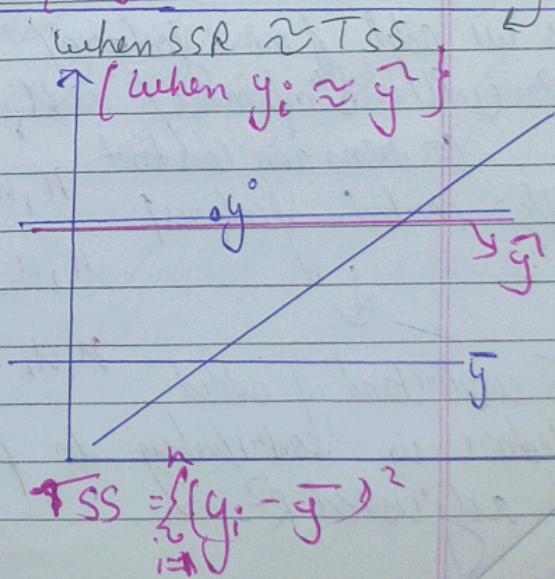
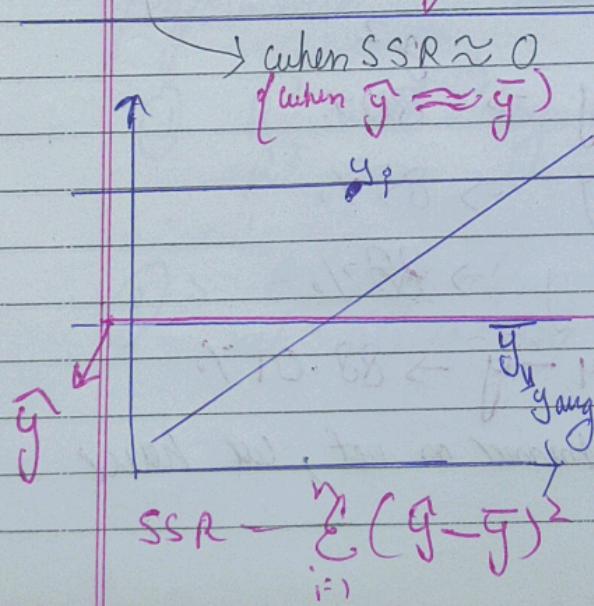
R^2 is the variation which is explained.

or \rightarrow Total percentage of variation explained by model.

$$R^2 = \frac{SSR}{TSS} \times 100$$

min value of $R^2 = 0$ when $\frac{SSR}{TSS} = 0$

max value of $R^2 = 1$ when $\frac{SSR}{TSS} = 1$



Can R^2 can be negative
 ↓

yes → when → if best fit line is very far from \bar{y} (baseline model)
 ↓
 which is not possible

↓

So R^2 can be negative but it is not possible due to nature of algorithm

2) Adjusted R^2

$\hat{y}^2 \rightarrow$ % explained variance in y due to x

u_1 area of house	u_2 No of rooms	u_3 city distance from house	y price of house
---------------------------	-------------------------	--------------------------------------	-----------------------

as we add more features

\hat{y}^2 will improve / $u_i - y \rightarrow 80\%$

or remain constant
after certain features

$$u_1 u_2 - y \rightarrow 85\%$$

$$u_1 u_2 u_3 - y \rightarrow 88\%$$

$$u_1 u_2 u_3 u_4 - y \rightarrow 88.01\%$$

To understand if added feature is contributing to performance or not, we have adjusted R^2 .

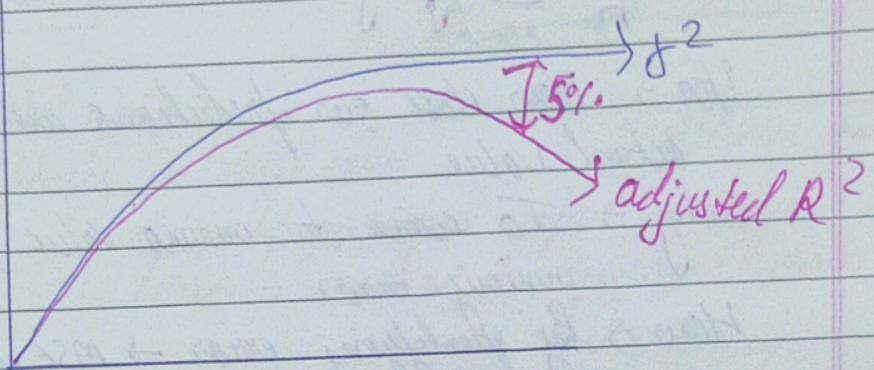
To calculate $\{\text{adjusted } R^2\}$ we have a formula

$$\frac{1 - (1 - \text{R square})(N-1)}{(N-p-1)}$$

Rg.

$n \rightarrow$ no of data points

$p \rightarrow$ no of independent variables



\rightarrow at some point R^2 becomes constant but when the $R^2 \rightarrow$ goes constant \rightarrow adjusted R^2 start decreasing.

if difference b/w (R^2 & adjusted R^2) that means it is not contributing much to performance. \rightarrow that added feature

① R^2 will always be ~~less~~ greater than adjusted R^2 .

② \Rightarrow diff. b/w R^2 & adjusted R^2 should not be more than 5%.

- MSE (Mean Squared Error)

$$\text{Sum of errors} = \sum_{i=1}^n (y_i - \hat{y})^2$$

Mean squared error = Sum of errors

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

Idea \rightarrow How close our predictions are close to actual value.

why \rightarrow To come to minima value, we need to minimize error

How \rightarrow By quantifying error \rightarrow MSE is used here

~~Error~~ Lower MSE, Higher performance

MSE measure average of squared error b/w predicted and actual values

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

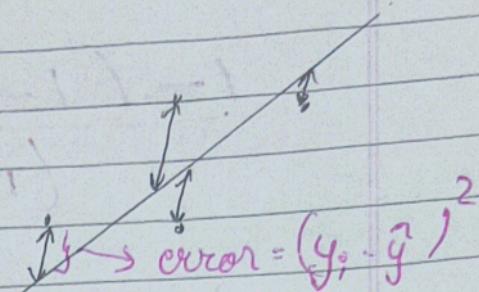
Why Squaring difference?

\rightarrow it ensure the value do not cancel out each other ($(-3)^2 + (73)^2 = 9$)

$$(-3 + 73 = 0)$$

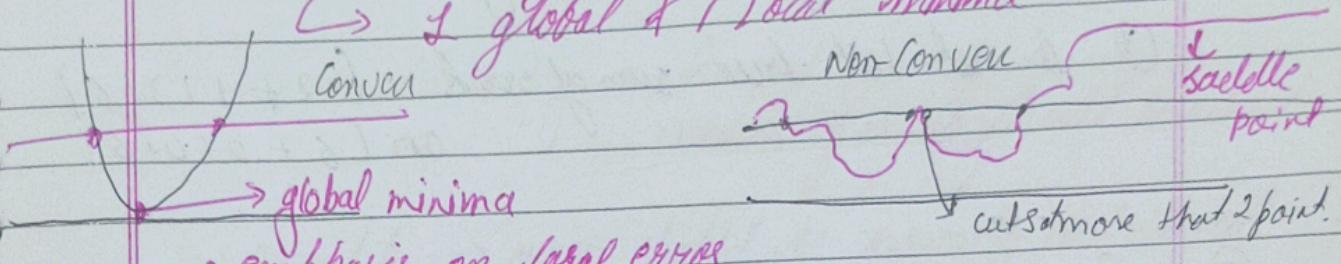
↙

Not error - x ?



Advantages of MSE

- Errors don't cancel out each other
- MSE is used as CF
- MSE is differentiable $\frac{\partial J}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$
- used as gradient descent as CF
- it is convex function
 - ↳ 1 global & 1 local minima



Disadvantage

- Not robust to outliers $(y_{out} - \hat{y}_{pred})^2$
- Not in same unit.
 - ↳ unit's squared.

③ Mean Absolute error → measure absolute difference b/w actual & predicted value

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantage →

- does not square the unit
- less sensitive to outliers
- more interpretable → easy cal

• Disadvantage

① Convergence usually takes more time, optimizing is complex

② Time consuming, Model training

③ Do not tell true sum of error $(2+1+3+6)$
or $[6+0+0+6]$

④ RMSE Root mean sq. error

$$\sqrt{MSE}$$

Advantage:-

- units not changes
- differentiable
- less sensitive to outlier
- magnifies larger error

Disadvantage

- interpretation → complex

Assumptions of Linear Regression

1 → Linearity → linear relationship (y)

2 → Independence → Rows (Observation) should be independent.

3 → Homoscedasticity → Constant variance

• Variance of error should be constant

4 → Normality → error should be normally distributed.

5 → Features should not overlap or should have least correlation. → multicollinearity.

M.L Pipeline

- ① Data ingestion / Reading dataset
- ② EDA
- ③ Data Preparation
- ④ Divide data into dependent & independent
- ⑤ Train test split
- ⑥ Scale the data (optional)
- ⑦ Model Train
- ⑧ Model evaluation

Pickling of model

pickle → allows you to store python object on a disk

python pickle module is used for serializing and de-serializing a python object structure.

Serializing → it is a process of converting data structure or object into the format that can easily stored, transmitted or persisted.
During serialization, objects state converted into a stream of bytes or other that can be stored on the disk, sent over a network or stored in memory

de-serializing → it is a process of reconstructing a serialized object back to the python object or what ever its primary form was.

The serialized object/data is read and converted back into a data structure or object that matches original

Multi-Collinearity

In stat \rightarrow correlation is How two values are correlated,
like when $x \rightarrow y$, y increases/decreases
by some % chance (This is the interpretation
by this [correlation])

$$x \rightarrow y \text{ by } \approx 94\%$$

if we have multiple features: $x_1, x_2, x_3, \dots, x_n$

\rightarrow Here one feature can correlate to several features

Corr. b/w two (Ind. Vs) at a time

$$\text{Corr}(x_1, x_2) = 95$$

$$\text{Corr}(x_2, x_3) = 80$$

Two points

• if x has high correlation with $y \rightarrow$ the model
will be good \rightarrow good thing

1 \rightarrow what will happen if feature among themselves
has good correlation?

2 \rightarrow To measure relationship b/w two variables/features
we have correlation. But how to measure correlation
among multiple features?

① Correlation among features itself →

if $\text{corr}(v_1, v_2) \approx 95\%$.

→ by this in terms of variation we can say
 $v_1 \approx v_2$

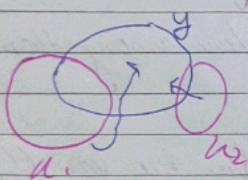
if $v_1 = v_2 \rightarrow$ the interpretation of MLR becomes
tough { which feature is contributing
to y will be tough to decide
as v_1 & v_2 are same }

② Measure of correlation among multiple features

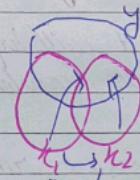
for that we have multiple multicollinearity.

many together linear relation

multicollinearity → either a feature exhibits a linear
relationship with more than two variables.



So both v_1, v_2 are corr with y
no multicollinearity



→ multicollinearity.

Why multicollinearity?
why it is a concern
→ affects interpretation

→ increase computation

→ it increases variance

$VIF \rightarrow$ defining some
outfitting increases

Solution →

To measure multi
inflation pattern
with high

RFE (Recursive)

VIFs → is a measure

VIFs

What's multicollinearity?

→ why it is a concern?

→ affects interpretation $\kappa_1 = \kappa_2$

$$\delta_{\kappa_1} + \delta_{\kappa_2} \xrightarrow{?} y$$

κ_1 or κ_2 ?

→ increase computation if $\kappa_1 = \kappa_2$, then why do we need to use both (κ_1, κ_2) as our feature

→ it increases overfitting if alone κ_1 can explain y then why not when we use κ_2 → defining some then it memorise and chances of overfitting increases

Solution →

To measure multicollinearity we have VIF (Variation Inflation Factor) and drop feature one by one with high VIF.

• RFE (Recursive Feature Elimination)

$\kappa_1 - \kappa_n$ (correlation/heating)

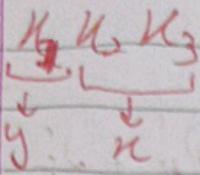
VIF is a measure of amount of multicollinearity in regression.

$\underbrace{\kappa_1, (\kappa_2, \kappa_3, \kappa_4, \kappa_5)}$

$$VIF_p = \frac{1}{1 - R_p^2} \quad \left\{ \begin{array}{l} R^2 \rightarrow \% \text{ Variation in } y \\ \text{explained by } y \end{array} \right.$$

dependent $\rightarrow \kappa_1 \rightarrow y$

Independent $\rightarrow \kappa_2, \kappa_3, \kappa_4, \kappa_5 \rightarrow K$

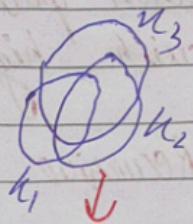
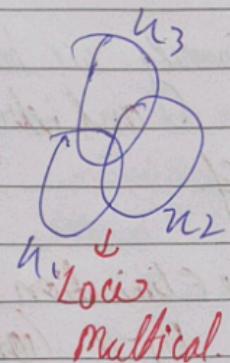
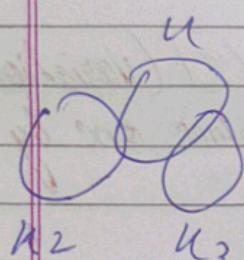


$R^2 = \text{% of variation explained on } y \text{ by } u.$

Case I $\Rightarrow u_1 \approx (u_2, u_3) \rightarrow \text{VIF}_{u_1} \Rightarrow \frac{u_1}{y} \text{ or } \frac{u_2, u_3}{n}$] We don't need to clean this for model by ourself

$$\text{VIF}_{u_1} = \frac{1}{1 - R_{u_1}^2}$$

Case II $\Rightarrow u_2 \approx (u_1, u_3) = \text{VIF}_{u_2} = \frac{1}{1 - R_{u_2}^2}$] Their will be a variation in Inflation Factor
Case III $\Rightarrow u_3 \approx (u_1, u_2) = \text{VIF}_{u_3} = \frac{1}{1 - R_{u_3}^2}$] there are just two all features it will give us all possible cases



Note when $\text{VIF} > 10 \rightarrow$ High Multicollinearity

If $\text{VIF} > 10 \rightarrow$ High $\rightarrow \text{VIF} = \frac{1}{1 - R_i^2} \Rightarrow 10 = \frac{1}{1 - R_i^2}$] High VIF \rightarrow Drop the feature only one

$$10 = 10R_i^2 \rightarrow R_i^2 = \frac{10 - 1}{10}$$

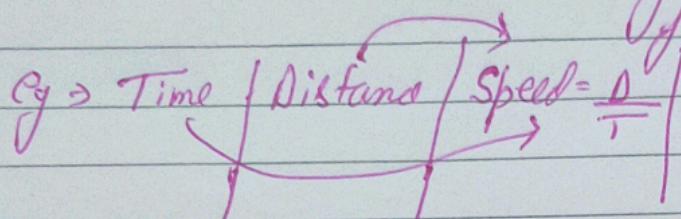
$$R_i^2 = \frac{10-1}{10} = \frac{9}{10} = 0.90$$

Kinds of multicollinearity

There are two kinds of multicollinearity.

- 1) Auto-based collinearity. → present in data itself due to its nature
 - eg → Latitude, Longitude

- 2) Structured based multicollinearity → caused due to new feature from existing feature



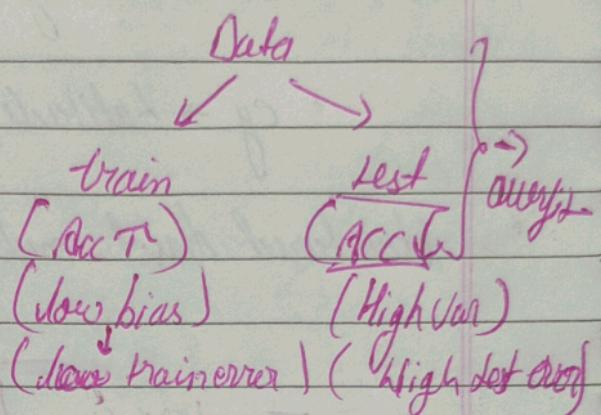
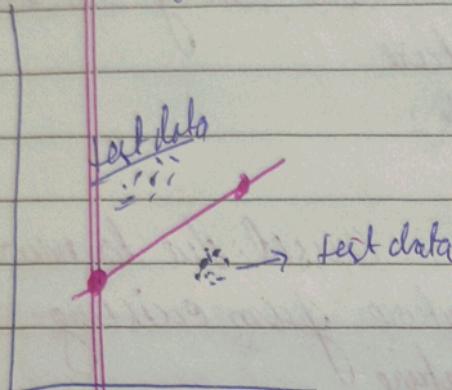
- * We talked that using VIF drop feature one by one what if there are 1000 features

- There comes RFE → Recursive Feature Elimination
 - it will make model with all 1000 features
 - start dropping the feature which has VIF > 10
 - one by one
 - Does it recursively until desired no. of features is arrived.

Regularization

↓ To add something to reduce anything overfitting

To regularise / To penalise



Since best fit line passes through both data point, so the model is overfitted model.

→ overfitted model means model is performing well on train data but on test data

Analogy behind this:

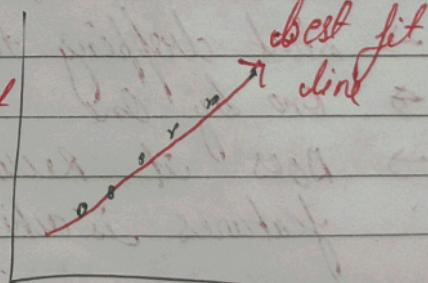
→ lets say this is our best fit line for our train data

→ we got best fit line using gradient descent

→ Now our test accuracy is 100%

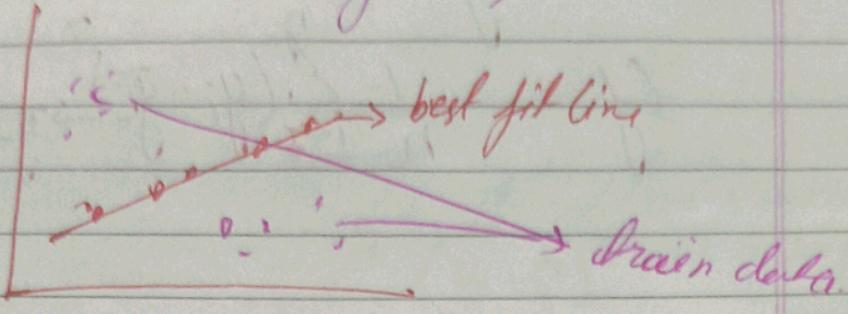
→ as our data points lies on the best fit line

- This is case of overfitting.



because train data fits good, but in test data.

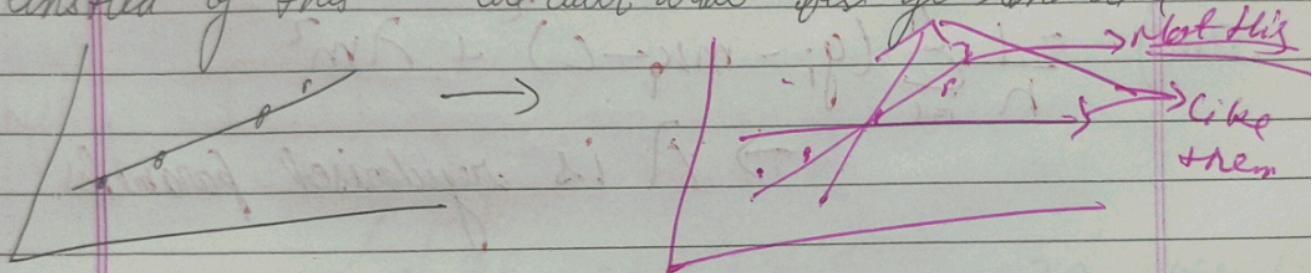
data points were not in test fit line.



So in order to make sure our model performs well on test data (unseen) as well, we will introduce some error / penalize while training model itself.

instead of this

we will draw best fit line like



→ While training we introduced some small error in order to increase the test accuracy. Because of this concept of regularization came into picture.

of Regularization → To add something to reduce overfitting, other something is error.

Standard Cost function

$$J = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 \rightarrow \hat{y} \rightarrow y_{pred}$$

$\hookrightarrow m\mathbf{x} + C$
 $\hookrightarrow Q_1\mathbf{x} + Q_0$

How do introduce error off this?

Regularised Cost function / Regularised linear model

$$J = \frac{1}{n} \sum_{i=1}^n (y_{act} - \hat{y}_{pred})^2 + \lambda m^2$$

$\hookrightarrow m\mathbf{x} + C$

$$J(m) = \frac{1}{n} \sum_{i=1}^n (y_i - mx_i - C)^2 + \lambda m^2$$

$\rightarrow \lambda$ is regularised parameter.

if error = 0 as accuracy = 100% for overfitted model

$$\sum_{i=1}^n (y_i - \hat{y}_{pred})^2 = 0$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{pred})^2 + \lambda \theta^2$$

$\hookrightarrow 0$

$$J(\theta) \rightarrow \theta + \lambda \theta^2$$

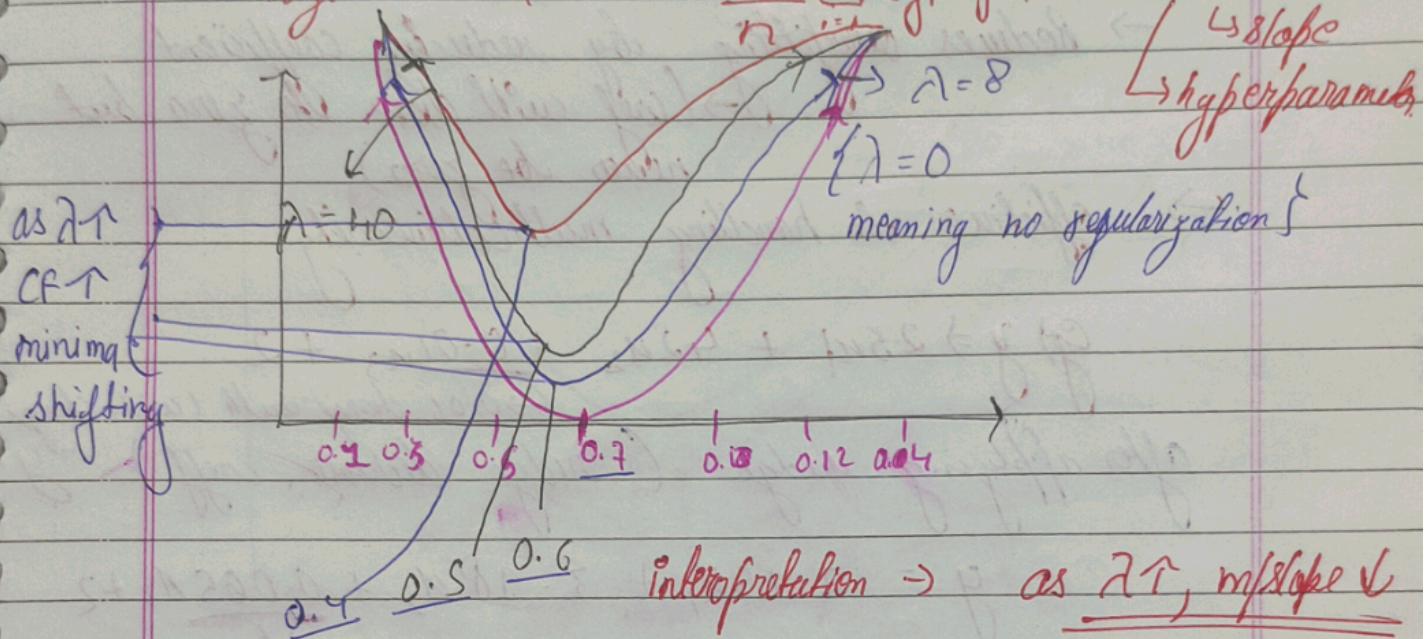
This error is introduced here

Different types of regularization

→ Ridge Regression (also called as L2 penalty, L2 norm or Tikhonov regularization)

$$\text{In SLR, } J(\alpha) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 \xrightarrow{\substack{\text{bow-2 stat'scny} \\ \text{L2 regularization}}} \quad \quad \quad$$

$$\text{Ridge} = C_F(J\alpha) = \frac{1}{n} \sum_{i=1}^n (g_i - \hat{y})^2 + \lambda \alpha^2 \quad \quad \quad \begin{cases} \downarrow \text{slope} \\ \downarrow \text{hyperparameter} \end{cases}$$



Note → Slope (m) decreases but never become zero in Ridge

$$\text{SLR} \rightarrow H_0(u) = \alpha_0 + \alpha_1 u$$

$$C_F = \frac{1}{n} \sum_{i=1}^n (g_i - \hat{y})^2 + \lambda \sum_{j=1}^p \alpha_j^2$$

$n \hookrightarrow$ no. of data points \hookrightarrow no. of predicted variables

Note → The intercept α_0 is not included in final term, Regularization applies only on coefficients $\alpha_1, \alpha_2, \dots, \alpha_p$

In MLR

$$h(u) = \alpha_0 + \alpha_1 u_1 + \alpha_2 u_2 + \alpha_3 u_3$$
$$CF = \frac{1}{n} \sum_{i=1}^n (y_i - h(u_i))^2 + \lambda (\alpha_1^2 + \alpha_2^2 + \alpha_3^2)$$

Ridge advantages \rightarrow

- \rightarrow Reduces overfitting by reducing coefficient
 \rightarrow (coeff will close to zero but never be zero)
- \rightarrow effective in handling multicollinearity

$$g(y) \Rightarrow 2.3u_1 + 4.2u_2 + 0.01u_3 + 2$$

after applying ridge (ridge decrease coeff) \downarrow
 $y = 2.15u_1 + 3.98u_2 + 0.005u_3 + 2$

\downarrow
effect reduced

Disadvantage

- \rightarrow does not make less imp feature coefficient 0 as we need ~~not~~ no effect of some irrelevant feature on prediction

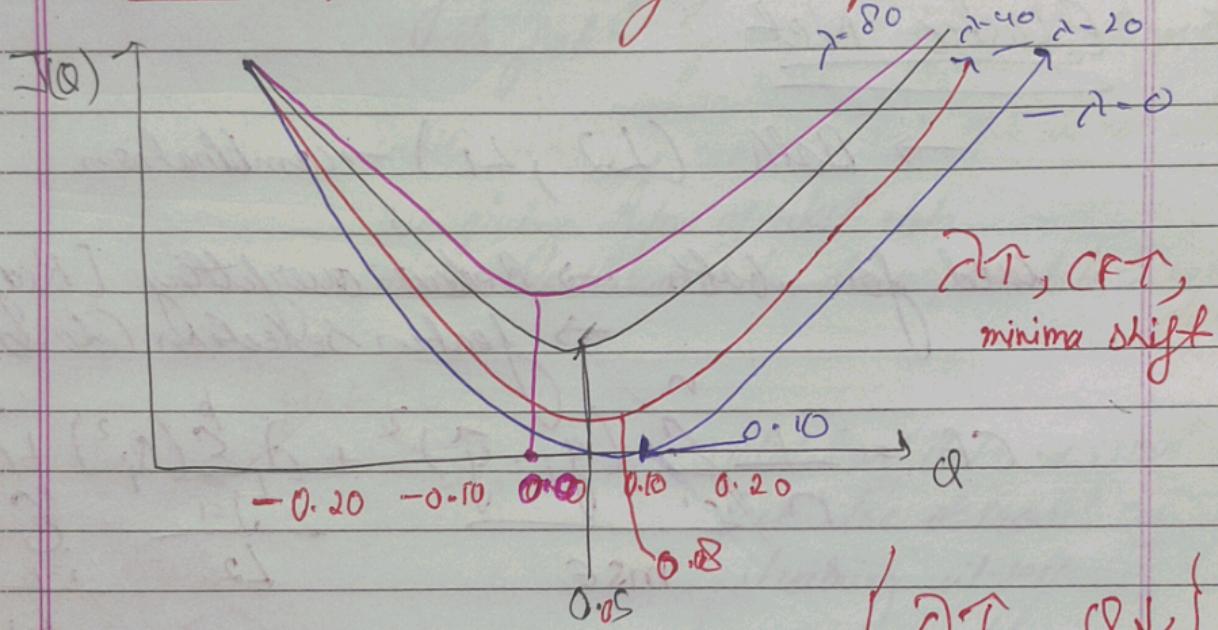
* Lasso Regression (L_1 Regularization, L_1 Norm) or path solution

$$J(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (y_i - h(\boldsymbol{x}_i))^2 + \lambda \sum_{j=1}^p |\beta_j|$$

↓

More power & Shuts out
it is $L_1 \rightarrow$ Regularization
error + $\lambda(|\beta_1| + |\beta_2| \dots |\beta_p|)$

→ similar to ridge (intercept is excluded)



How β coefficient decrease with increase in λ , and it sometimes reduced to 0 as well

↳ coefficient

irrelevant features becomes 0 in, that means no impact of feature in prediction

Advantage of Lasso

- in order to remove irrelevant features it is very useful. (Sparsity)
- Lasso → push coefficient towards 0 / exact 0
- multicollinearity & overfitting reduced.

3

Elastic Net

→ Both (L_2, L_1) → combination

Used for both → Reduce overfitting (Ridge)
 → feature selection (Lasso)

$$CF = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda_1 \sum_{j=1}^p (\beta_j)^2 + \lambda_2 \sum_{i=1}^p |\beta_i|$$

mse

$\lambda_1, \lambda_2 \rightarrow$ Hyper parameters
 range $\rightarrow (0, \infty)$

Why we using all these things.

Occam's Razor Principle → Simpler models are best model as compare to complex

Regularization does it

Cross Validation, its types and hyperparameter tuning

Data → 70 → Train Data

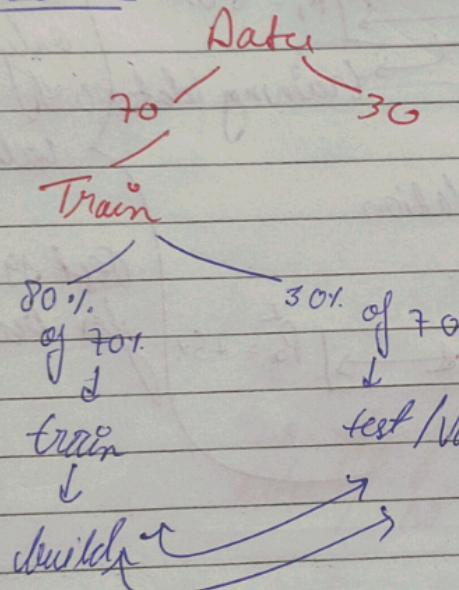
Data → 30 → Test Data → Representative of unseen data

Scenarioro → You trained a model, got 80% accuracy on train data but 40% on test data. What will you do?

→ Is you will go again to training process and do feature engineering and repeat all things?

while retraining we ← No this is not allowed
are giving info of test data
which should not be done, test data → unseen data.

⇒ Do this



Validation data is used to test the fit of model while training it self

why?

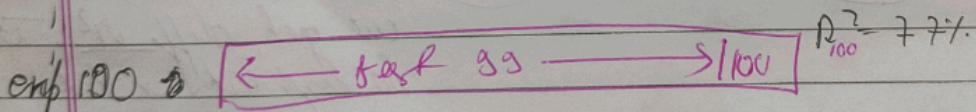
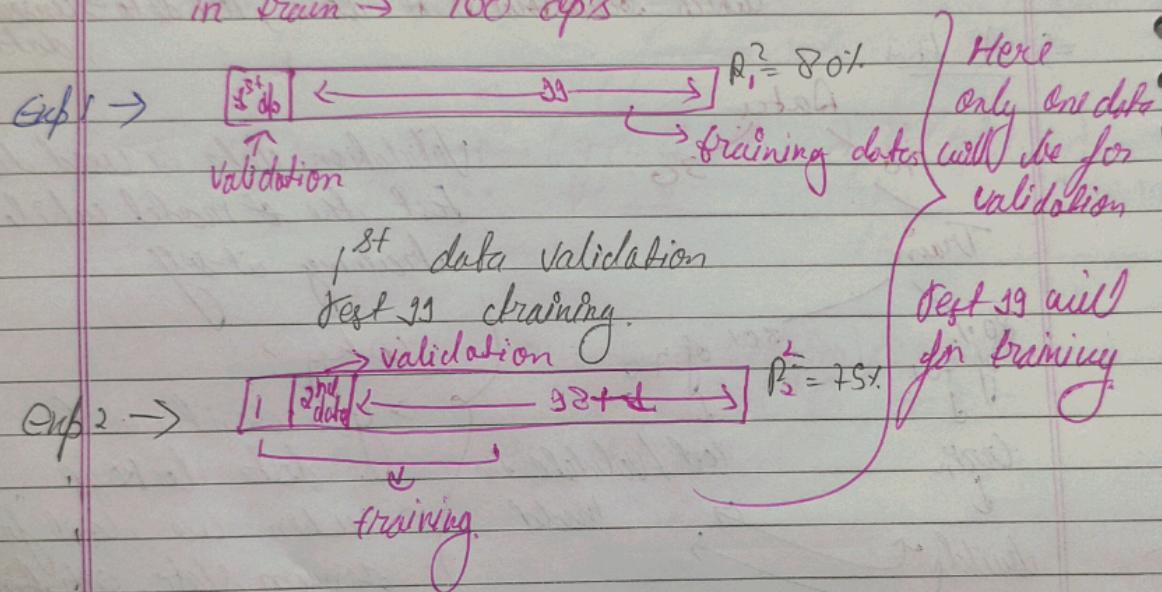
→ You don't want data to leak
→ When we don't fix random state, each time we get diff R^2 . To be confident we take average of all runs

Cross Validation → experimenting with different arrangement of some isolates to build different model of same algorithm.

train data	Exp 1	train	Validate	Accuracy 1
	Exp 2	"	"	" 2
	Exp 3	train	"	" 3
	Exp 4	"	"	" 4
				Mean of (1, 2, 3, 4)

Types of Cross Validation

① Leave One Out Cross Validation (LOOCV)
in train → 100 apx.



$$\text{avg } R^2 \Rightarrow \frac{1}{n} \sum_{i=1}^n (R_i^2)$$

Disadvantage of LOOCV

- Since only 1 df is for test, train size become very huge for large data, and time increased accordingly.
- Overfitting increases as \rightarrow training size \approx total size

② Leave ~~p outcomes~~ cross validation ($p > 1$)

Exp-1 $\boxed{S \triangleq 95 \rightarrow}$

Exp-2 $\boxed{- | S | 90+5}$

\vdots
 $\text{exp} \left(\frac{n}{p} \right) \quad \boxed{\xrightarrow{\quad} | S |}$

③ K-Fold Cross Validation \rightarrow very famous in industry

$$k = 5$$

\downarrow
 no. of group you want to divide the data / no of experiment = $\frac{100}{20} = 5$

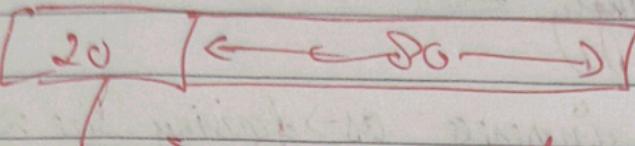
Exp-1 $\boxed{20 \leftarrow \rightarrow 20 \rightarrow} \rightarrow \text{Reg-1}$

Exp-2 $\boxed{\leftarrow \rightarrow 20 \leftarrow \rightarrow 60 + 20 \rightarrow} \rightarrow \text{Reg-2}$

Exp-3 $\boxed{\leftarrow \rightarrow 80 \rightarrow 20} \rightarrow \text{Reg-3}$

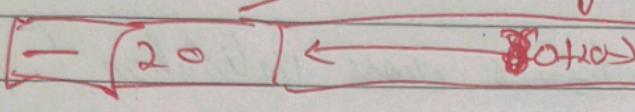
↳ Stratified k -fold cross validation.

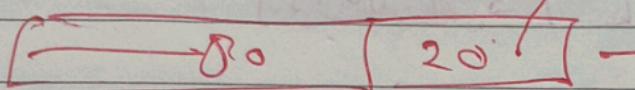
$$k=5, n=100 \quad \begin{array}{l} \text{class 0 - 60} \\ \text{class 1 - 40} \end{array}$$

Exp 1 \rightarrow  Accuracy - 1

Here 10 dfts of class 1.

10 dfts of class 2

Exp 2 \rightarrow  Accuracy - 2

Exp 3 \rightarrow  Accuracy - 3

In each validation set, there will be equal representation of each of the classes.

Hyperparameter Tuning

Hyperparameter \rightarrow external config. of model that are not learned from data but are set prior to training process.

$$\text{Ridge} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 + \frac{\lambda}{n} m^2$$

$y = mx + c$

Learned from data. \downarrow
external.

model performance is heavily influenced by
 λ \hookrightarrow Hyperparameter

$$\lambda = 0, 1, 2, \dots, 100$$

How to know best value
~~from~~ for λ from all possibilities.

Hyperparameter with cross validation

\hookrightarrow finding best hyperparameter while ~~tuning~~ training model.

1 \rightarrow Grid Search CV \rightarrow cross Validation.

2 \rightarrow Random Search CV

1 o Grid Search

$$\lambda = 1, 2, 5, 10, 15$$

$$\beta = \beta_1, 2\beta_1, 5\beta_1$$

All possible combi.
of λ, β

λ_1	λ_2	λ_3	λ_4	λ_5
β_1	$1\beta_1$	$2\beta_1$	$5\beta_1$	$10\beta_1$
β_2			$10\beta_2$	
β_3	$2\beta_3$			$15\beta_3$

Best model

\rightarrow combination of (λ, β) which gives max accuracy.

Note → if Data is small & instead of train-test
use k -fold cross validation

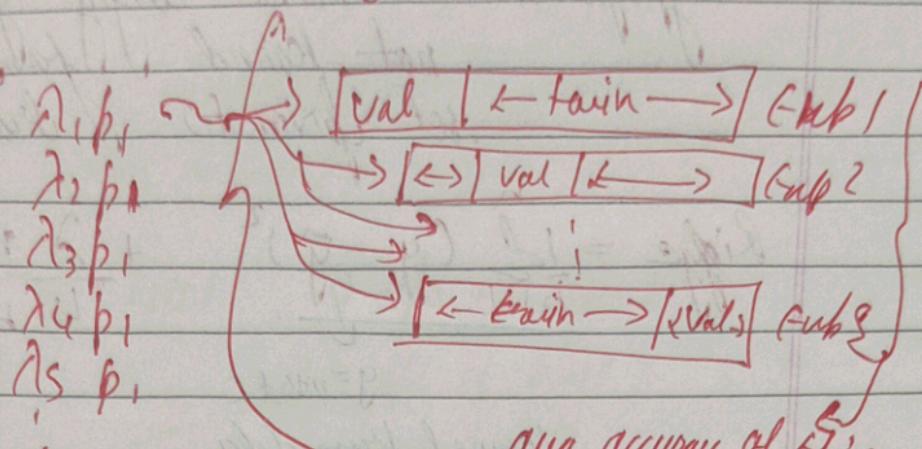
PAGE No.	
DATE	

Hyperparameters tuning with cross-validation

for each (λ, β) internally cross-validation will happen

which combi.

out of all
each combination
of (λ, β)
 λ -fold
cross-validation
will happen



avg accuracy of $\frac{1}{k}$

→ best accuracy will use for hyperparameters

Total no. of model. $\rightarrow p \times \lambda \times S$

Time Complexity is very high.

2) Randomized search CV.

→ we will not see all combi.

instead we have n iter. = S

Select random n combi.
of (λ, β)