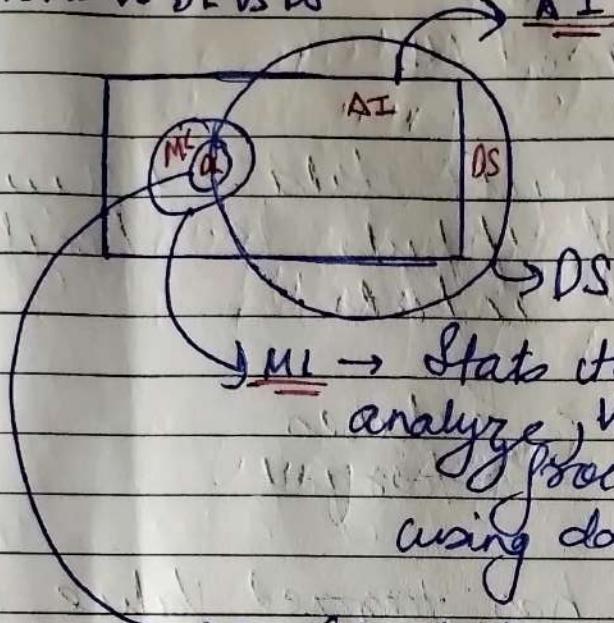


DATA SCIENCE

- Introduction
- AI vs ML vs DL vs DS



An application that can do its task without any human intervention.
Ex: (Netflix recommends system) ; (Self driving car)

ML → Stats tool to analyze, visualize, forecast, predict using data

→ DL (Deep learning)

→ trying to make a machine learn like human.
or simply
to mimic human brain.

Data Scientist DS : a person who analysis data them use any of the above to create a solution for a problem

Machine Learning

- Supervised ML
- Unsupervised ML

In machine learning we are usually given a set of inputs called "features" and it's labeled outputs called "target Values".

Supervised ML algo's are those which

We use Supervised ML algo's where we have access to both of these values but, say if we only have input value & no labeled target values, then we use Unsupervised ML algo's.

Ex: → Student pass or not
Input → Marks
Output → Passed / fail
on the basis of given inputs, we can create and train a model to predict on new data if a student passes or not.

Ex: email classification (Spam or Not)

↳ Input → Class →

Spam not-spam

↳ we need do an relation
that makes an email spam or not
spam & then predict on a new
email that we might find.

Ex: recommendation made by different
app. say Amazon:

↳ Input : Recent purchases
items, user cost, history

On basis of this input we need
to recommend a user new products
they would like to buy.

(Here we don't have target values
as we don't know what a user
might buy next).

(Supervised ML)

This consist of mainly 2 type of problems :-

Regression

Output is
a continuous

Ex House price

Classification

(Yes/No or 0/1)
or 1/2/3/4

Output is
(Categorical Value)

(Choice)

Ex: Student pass

• Regression : It is a statistical method which attempts to find a relationship between a dependent Variable & one more dependent Variable.

• Classification : It is simply the allocation of objects to certain pre-existing classes or categories. Specifically in machine learning it used to categorize a given set of inputs.

Algorithms generally used:-

Regression : ① linear Regression
② Lasso & Ridge
③ Elastic Net
④ Decision trees

- ⑤ Random Forest
- ⑥ Xgboost
- ⑦ ANN (Artificial Neural Networks)
- ⑧ RNN (Recurrent neural Network)
- ⑨ CNN

- Classification :
- ① Logistic
 - ② Lasso, Ridge
 - ③ Decision Tree Classification
 - ④ Random Forest Classification
 - ⑤ Xgboost Classification
 - ⑥ ANN
 - ⑦ CNN → (Convolutional Neural Networks)

most algos used in Regression can also be used in Classification.

(many more algos exist and are used in Supervised ML)

(still learning ...)

Unsupervised ML

↳ Clustering Algorithms

→ K-Means

→ DBSCAN

→ Hierarchical

↳ Anomaly Detection

Isolation Forest

One Class SVM

(Many more categories & algorithms.)

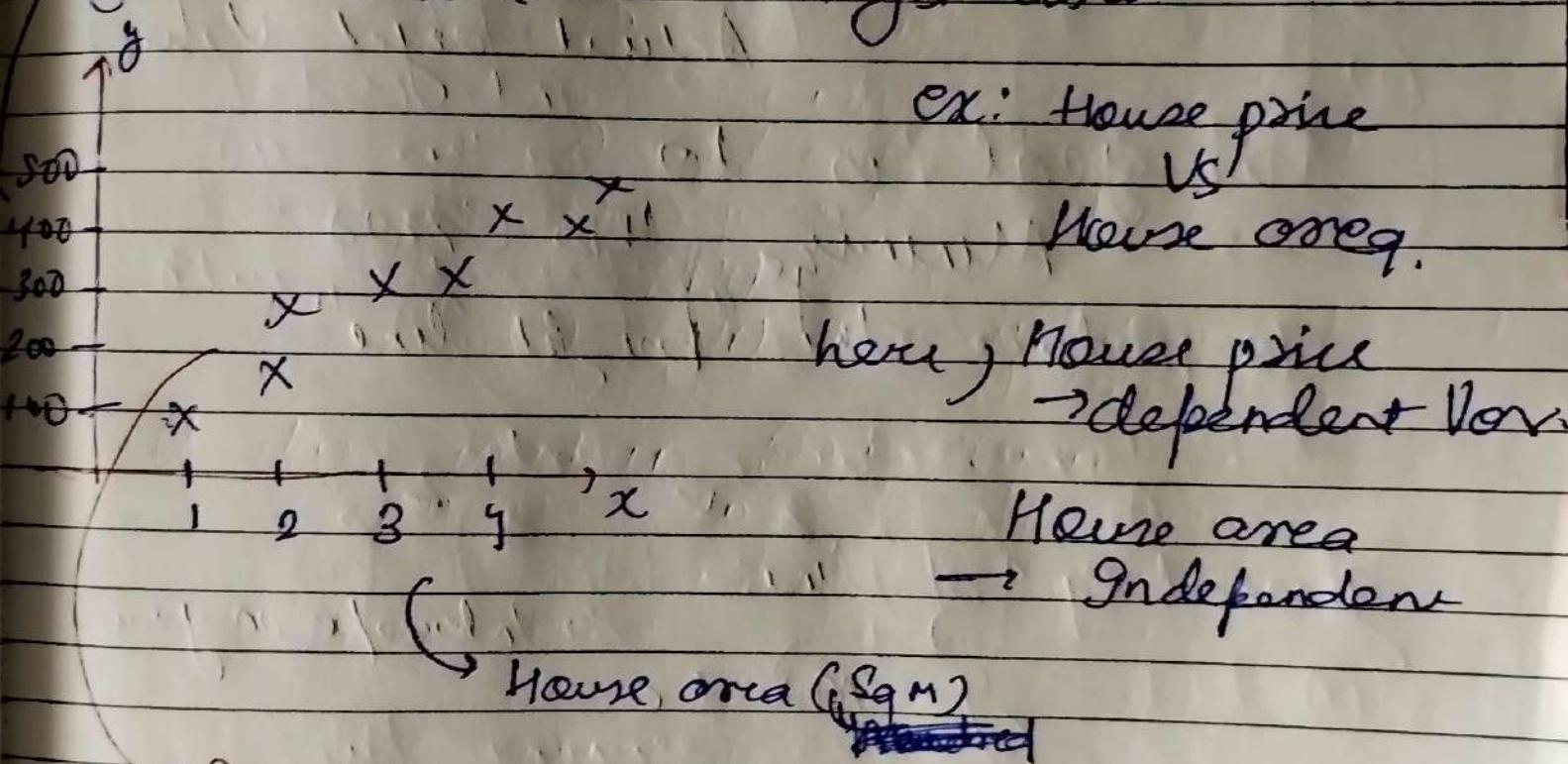
(Still learning - - -)

Supervised ML

Linear Regression

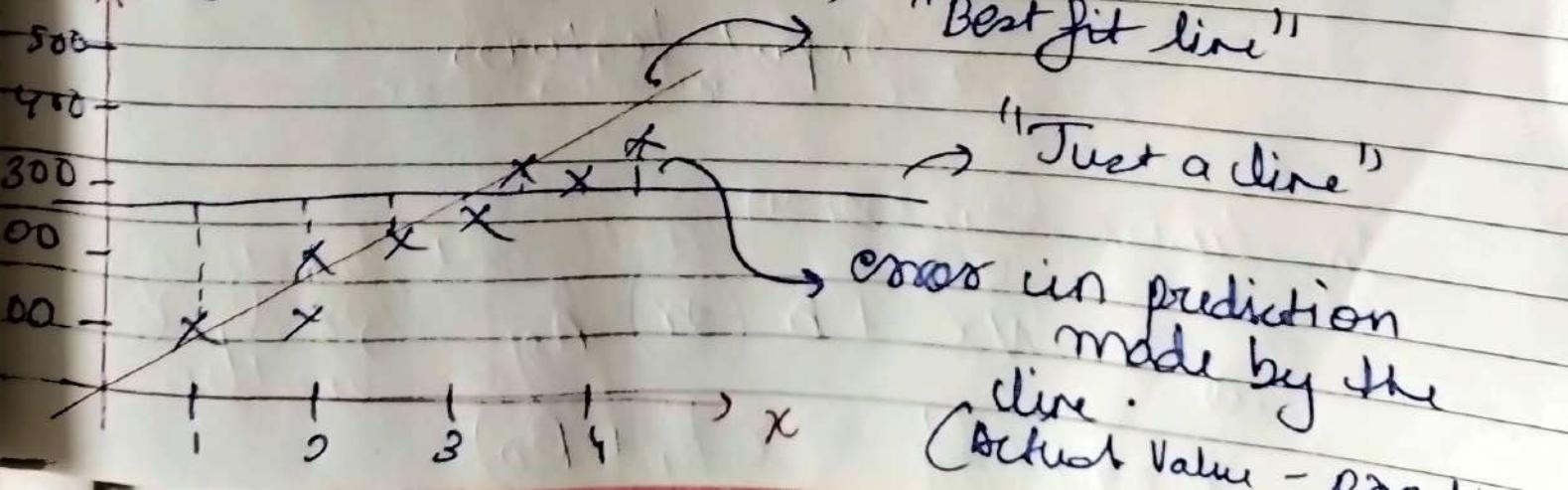
→ here we try finding an eqn
that best relates our inputs
to our output. This line is
called the "best fit line".

NOTE:



→ Data points

through this data we need to fit an ~~best~~
line: ~~the~~ "Best fit line"



"Just a line": Doesn't accurately predict a relation between input (x) and output (y). So, it's predicted to have a lot of error.

Error \rightarrow Distance between the predicted point & Actual target value

Our goal is to
 → "Minimize this error"
 to produce the "best fit line"

General eqⁿ of line:-

$$y = mx + c$$

→ Intercept on y-axis
 → slope
 → prediction \rightarrow
 → Input features

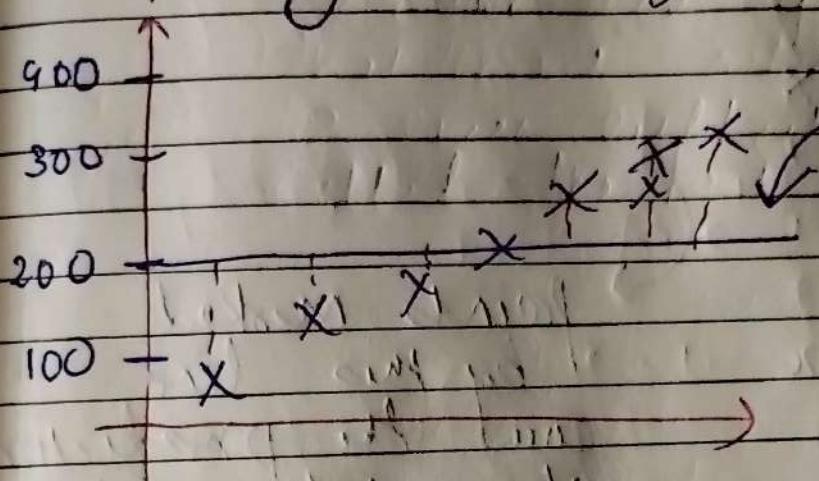
In linear Regression :-

$$f_{\text{lin}} = Wx + b \quad (\text{or}) \quad h_0 = \theta_0 + \theta_1 x$$

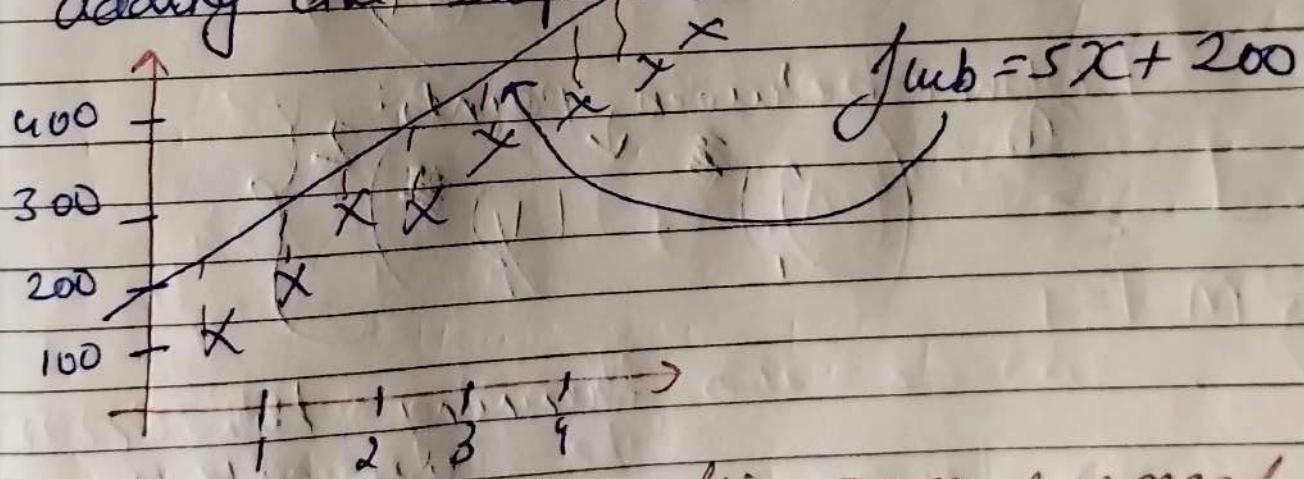
$$(\text{or}) \quad Y_i = X\beta_1 + \beta_0$$

9 use : $f_{wb} = Wx + b$

so, say $b = 200$, $W = 0$ $\{f_{wb} = 200\}$.



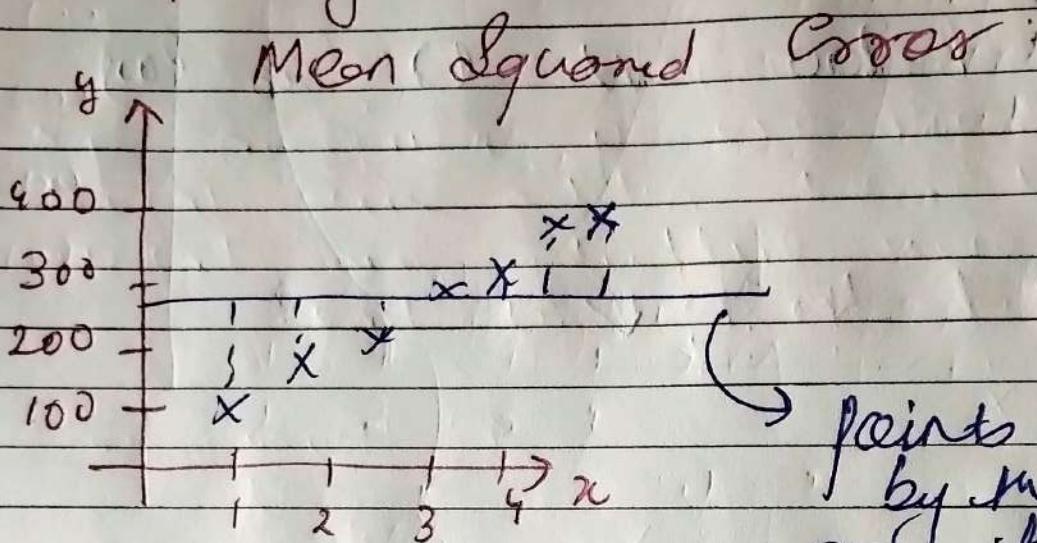
adding the slope term: $b = 200$, $W = 5$



So, to minimize this error we need
to find correct values of
more generally eqⁿ:

$$(f_{wb} = b + w_1 x_1 + w_2 x_2 + \dots + w_n x_n)$$

Computing the error:



points marked
by this line
and the prediction
represented with
(y)

Distance between points:

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

Running this error
over all the
data points

here for Calculation case

The above eqⁿ is called "Cost function"
"or
"Loss function"

Represented as :

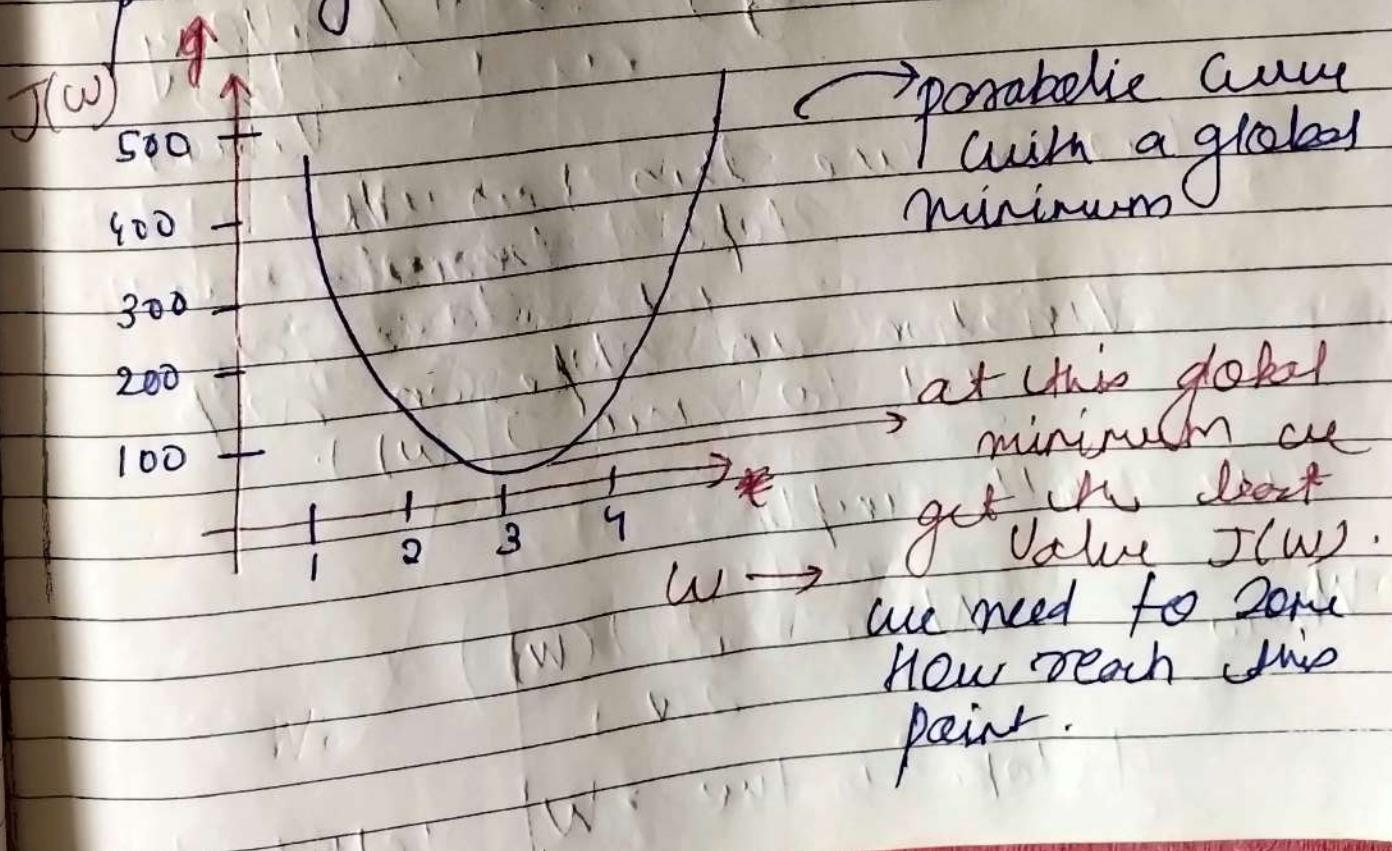
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m ((y_i) - (\hat{y}_i))^2$$

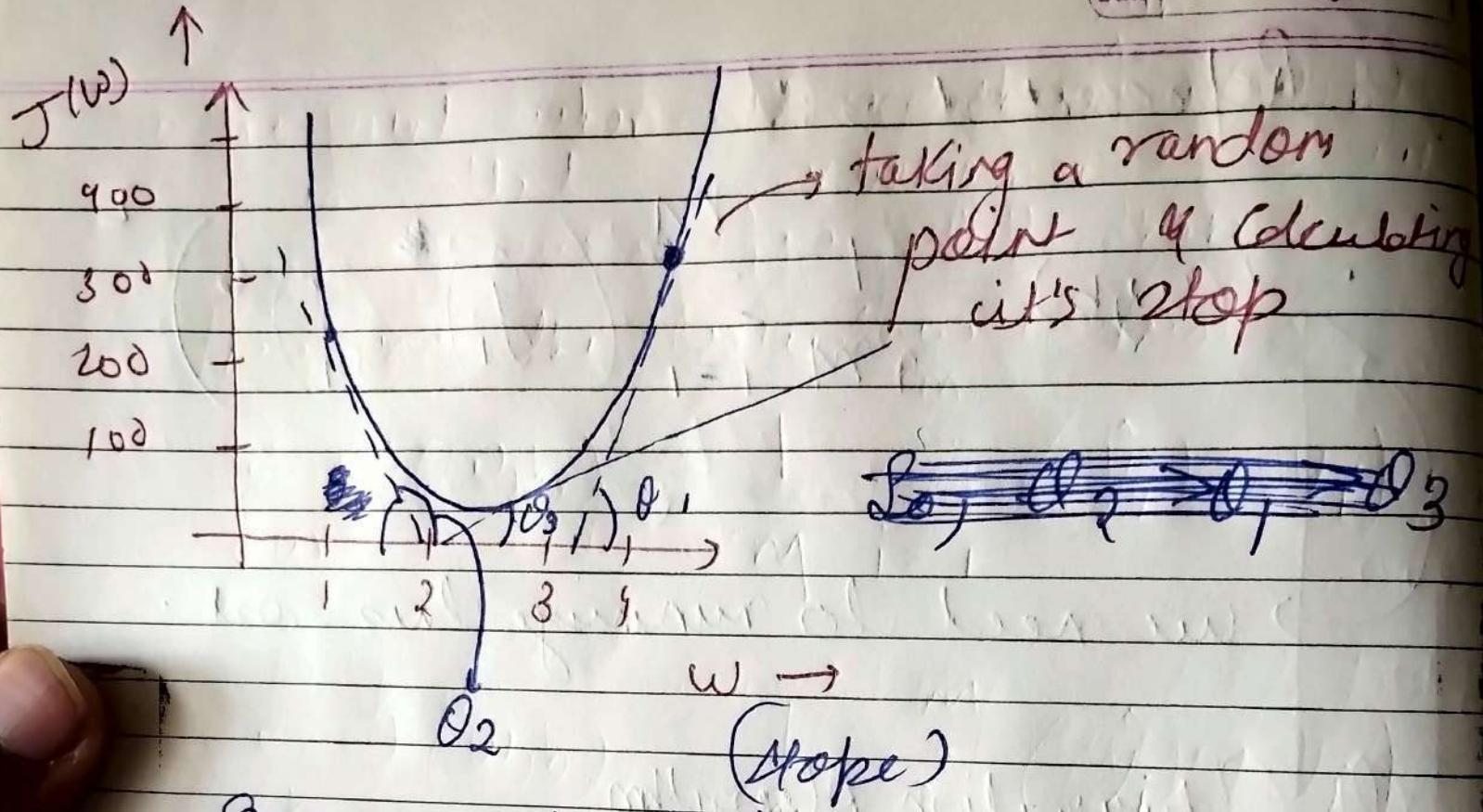
→ we need to minimize this Cost.

① Optimization Algorithm:

① "Gradient Descent Algo"

plotting this $J(w, b)$ for different w values & $b=0$.





w_{new} , $\tan(\theta_1) \rightarrow \text{true}$

$\tan(\theta_2) \rightarrow -\text{ve}$

$\tan(\theta_3) \rightarrow (+/-) \text{ very small value}$

we can use this as the slope decreases the variation in slope can help get closer to get to $(\min J(w))$.

We will simply :-

$$w_{\text{new}} = (w - \frac{\partial J(w)}{\partial w})$$

to min $J(w)$,
if slope is $\rightarrow w^+$ \rightarrow getting closer to min $J(w)$.
(if we are already near min $J(w)$,
we won't move up (the value of desciptor will be small, so no significant \uparrow or \downarrow))

Going from right to left $\rightarrow w^-$.

↓ we encounter min $J(w)$

Binary left to right $\rightarrow w^+$

→ to quantify how much to move use the value of slope.

Repeated updates to the parameters result in optimal w values

This algorithm is "gradient descent algo." "learning rate"

representation

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b_{\text{new}} = b_{\text{old}} - \alpha \frac{\partial J(w, b)}{\partial b}$$

The point when it reaches $J(w, b)$
 is called "convergence", & we
 repeat our algo till then.

These are the derivative terms:

$$\frac{\partial J(w, b)}{\partial w} = \frac{1}{m} \sum_{i=1}^m ((\hat{y}_i) - y_i) \cdot x_i$$

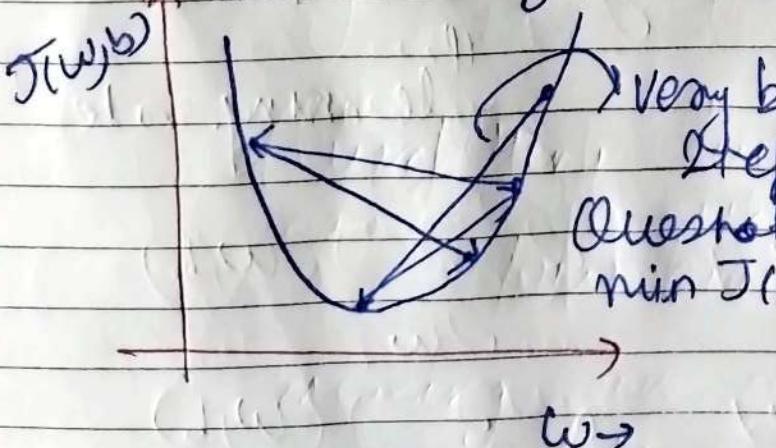
$$\frac{\partial J(w, b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m ((\hat{y}_i) - y_i)$$

Learning Rate (α): It decides how big of a step we

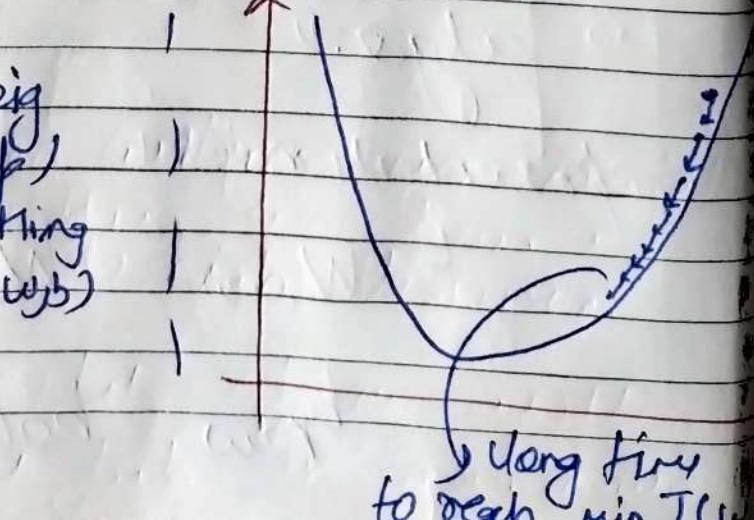
take while updating of parameters (w, b)

$\alpha \rightarrow$ close

$\alpha \ll$ small



Quenching
 $\min J(w, b)$



Long time
 to reach $\min J$

all need to pick alpha Davis where in
the middle of this

Matrix Notation of linear Regressions

$$Y = XW + \epsilon \quad \text{residuals}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \text{(for pios / intercept no q features in e)}$$

term) example

$$X = \begin{bmatrix} X^T \\ X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$

number of examples

$$W = \begin{bmatrix} b \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_p \end{bmatrix} \quad \begin{array}{l} \text{intercept term} \\ \text{coefficients / slope} \end{array}$$

another way to find the parameters

$$\star \left\{ \vec{w} = (X^T X)^{-1} X^T Y \right\} \star$$

Still learning as to how this works ??

Problem of Overfitting and Underfitting.

Bias: error in the prediction made by model.

Variance: specifies the amount our model changes in a different portion of data (or new unseen data) is used.

We want our model to have low-Bias, & low Variance (though both simultaneously are not possible)

Overfitting occurs when our model has ~~high bias & low variance~~ high bias & high variance.

meaning our model performs well on training data but fails on unseen new data.

If ours when our model is trained only on training data & the optimization algo ends up with over accurate parameters that work only to model real training data.

There are multiple way to fix this ex: (Splitting data set (into train, test, cross valid))

(using regularization methods)

Underfitting: Ours when our model has high Bias & low variance i.e. it doesn't perform well on training data if there is no change in its performance on testing data.

Cause are poor training of model among & etc.

(Ridge & Lasso Regression)

→ Solution for the problem of overfitting
This is a regularization method.

Overshooting occurs when low bias & high variance.
 So, here we intentionally add some error to our model which result's in less accuracy more errors on training data but also reduces variance.

Ridge Regression (also called L_2 norm)
 we add this term to our (MSE) or cost function.

$$\left\{ \text{Term 1} + \frac{1}{P} \left(\sum_{i=1}^P (w_i)^2 \right) \right\}$$

(no of features) \rightarrow Regularization parameter

If it's derivative w.r.t to w becomes

$$\left\{ \text{Term 1} + \frac{1}{P} \left(\sum_{i=1}^P w_i \right) \right\}$$

- In addition to adding an bit of error, to more generalize the model.

Lasso Regression also helps in reducing the dimensionality of our data.

i.e

it eliminates irrelevant features.

We add this term in Lasso: (L1 norm) (absolute)

$$\left\{ \text{original function} + \frac{\lambda}{P} \left(\sum_{i=1}^P |w_i| \right) \right\}$$

and it's derivative w.r.t w.

(newton's method doesn't have a

proper derivative use *subderivative*.

$$\left\{ \text{original function} + \lambda \underbrace{\sum_{j=1}^P \text{sign}(w_j)} \right\}$$

$$\text{sign}(w_j) = \begin{cases} 1 & \text{if } w_j > 0 \\ -1 & \text{if } w_j < 0 \\ 0 & \text{if } w_j = 0 \end{cases}$$

Difference between the Ridge

→ Ridge introduce more error

→ Lasso introduce less error than Ridge but helps in Reducing dimension.

Now, The combination of the above becomes:-

Elastic net regularization

(having good fit point of both)

this is the term added to cost function:-

$$+ \frac{1}{2} \lambda_1 \sum_{j=1}^P (w_j)^2 + \frac{1}{2} \lambda_2 \sum_{j=1}^P |w_j|$$

and derivative cost to w:-

$$+ \frac{1}{P} \lambda_1 \sum_{j=1}^P w_j + \frac{1}{P} \lambda_2 \sum_{j=1}^P \text{sign}(w_j)$$

The above is also used in this form:-

$$\cancel{\frac{1}{2} \sum_{j=1}^p w_j^2}$$

$$= + \lambda \left(\frac{1-\alpha}{2p} \sum_{j=1}^p (w_j)^2 + \frac{\alpha}{p} \sum_{j=1}^p |w_j| \right)$$

where λ is "mixing parameter".

(if $\alpha = 1$, Elastic Net becomes Lasso)

$\alpha = 0$, Elastic Net becomes Ridge

$0 < \alpha < 1$, Combination of both)

- In Ridge Regression, the penalty reduces the values of slope Coefficients but it never goes to zero.
- but, in Lasso this goes to zero resulting in some parameters reducing to 0. (It is not entirely clear as to why Lasso goes to 0.)

Linear Regression itself also has many variants:-

- ① Simple Linear Regression \rightarrow (with 1 feature)
- ② Multivariable Linear Regression \hookrightarrow (with multiple features)
- ③ Multivariate
 - \hookrightarrow predicting more than one variable on the basis of multiple input features

ex: (predicting if student passes in English & Hindi, 2 marks)

④ Polynomial Regression: Fitting a data
with higher degree (relation)
ex: $y = x^2$

many more will learning it

Logistic Regression

Here we try classifying either a continuous set of inputs or a set of binary inputs into one of 2 categories usually 1, 0 or (True, False)

We use linear model like's eqⁿ to find the value of a prediction but as we have to classify this input in terms of 0 & 1 we use

Sigmoid function

(it is a case of logistic function)

Logistic growth

~~Step function~~

$$p(x) = \frac{L}{1 + e^{-(x - x_0)/K}}$$

$L \rightarrow$ supremum of the values of function.

Location parameter which means mid point of curve where $p(x_0) = 1/2$

when the value of

Date:

Page:

$$L=1, k=1, x_0=0$$

$$P(x) = \frac{1}{1+e^{-x}}$$

↳ Sigmoid function

In our case:

$$P(z)$$

$$P(w,b) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

$$(z = w \cdot x + b)$$

This function squashes value between $(0 \leftrightarrow L=1)$.

a more accurate way of writing sigmoid function, as we won't always have ($y = 0$)

$$p(x) = \frac{1}{1 + e^{-(x-\mu)/s}}$$

Scale parameter
Controls the spread of distribution
(mid point of curve, when $p(\mu) = 1/2$)

and often written with $f(wb)$:-

$$f(wb) = p(z) = \frac{1}{1 + e^{-(b + wx)}}$$

$$\{ b = -\mu/s \} , w = 1/s$$

\curvearrowleft Intercept \curvearrowright Slope

$$\{ y = -\beta/s \}, s = 1/w$$

while implementing we mostly
care about

$$f(w_b) = \frac{1}{1 + e^{-(w^T x_b)}}$$

now, here also the initial sigmoid
might not "fit" our data
well. So we have to reduce the
error here as well.

But Linear Regression's Cost function
won't work here here. We need
to find something else:-

new one

we define

$$J(w_b) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(f(w_b), y(i))$$

$$\begin{cases} -\log(f(w_b)) & \text{if } y=1 \\ -\log(1-f(w_b)) & \text{if } y=0 \end{cases}$$

Combining we get this:

$$-y \log(f(w_b)) - (1-y) \log(1-f(w_b))$$

while implementing we mostly
care about:

$$f(w, b) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

now here also the initial sigmoid
might not "fit" our data
well. So we have to reduce the
error here as well.

But Linear Regression's Cost function
won't work here now. We need
to find something else:-

new one

we define.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(f(w, b), y(i))$$

$$\begin{cases} -\log(f(w, b)) & \text{if } y=1 \\ -\log(1-f(w, b)) & \text{if } y=0 \end{cases}$$

Combining we get this:

$$\Rightarrow -y \log(f(w, b)) - (1-y) \log(1-f(w, b))$$

also known as binary cross entropy
Date: _____
Page: _____

finally:-

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m (-y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i))$$

Reduces to

$$- \log(\hat{y}_i) \text{ for } y_i = 1$$

or to

$$\cancel{- \log(\hat{y}_i)} - \log(1-\hat{y}_i)$$

and it's total sum is ~~for~~ $y_i = 0$
also called (-ve log likelihood)

Now justify this as how does the
penalizes the input if it has a
large error.

The function is in 2 parts for a
reason:

when ($\hat{y}_i = 1$), we want

our value to be as close
to true value as possible.

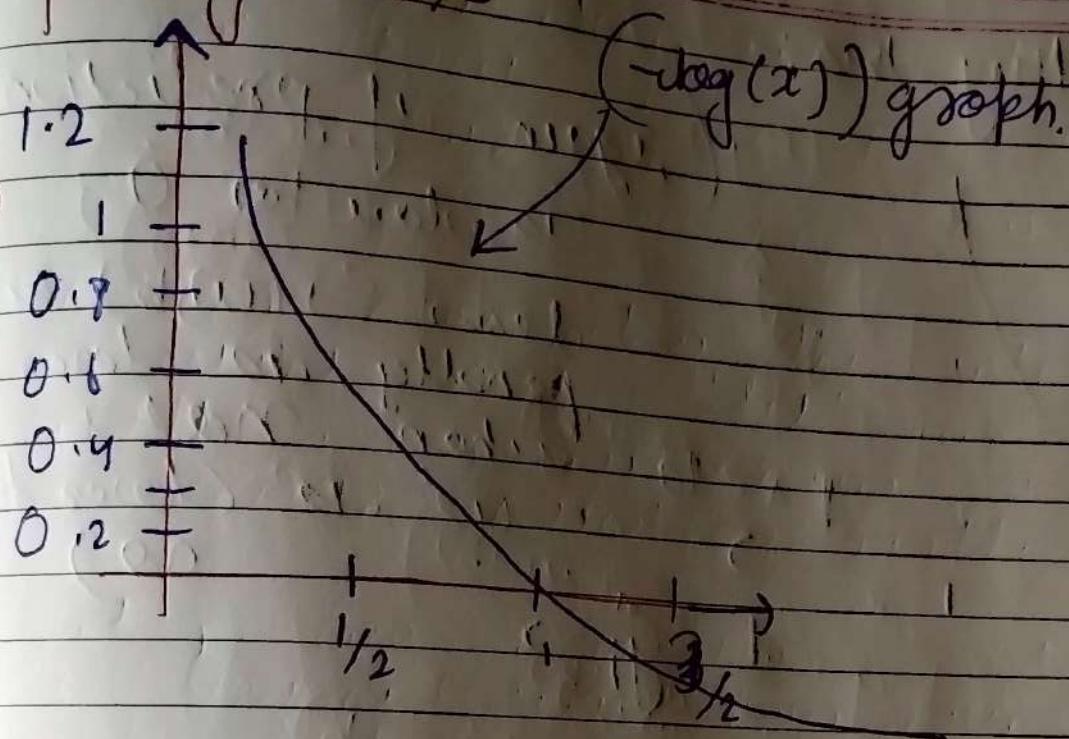
Q) this function works:-

$$-\log(\hat{y}_i)$$

plotting this:

Date:

Page:



(\hat{y}_i)

when y_i is close to 1, $J(y_i) \rightarrow 0$

if y_i is close to 0, $J(y_i) \rightarrow \infty$

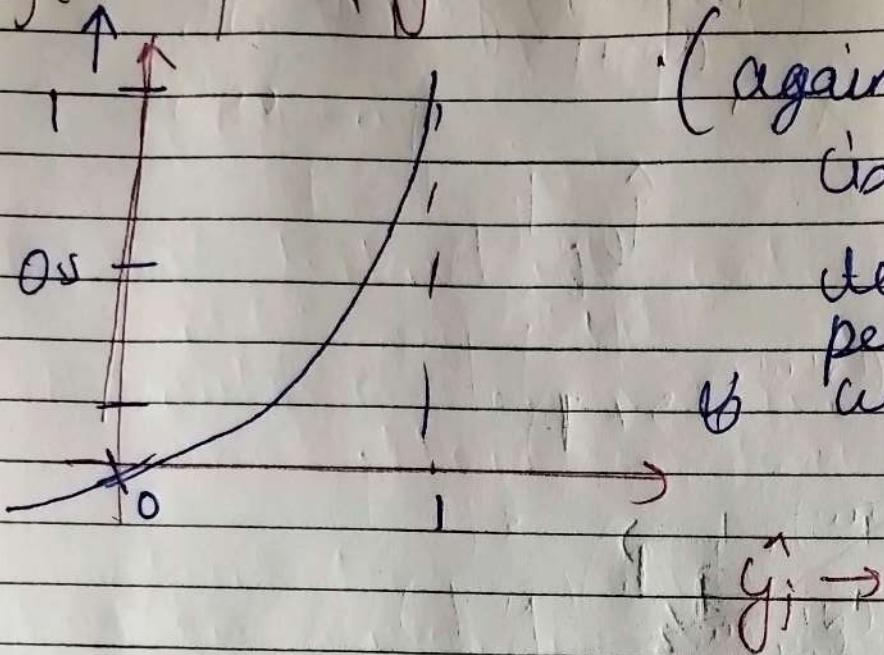
heavily penalizing the wrong predictions.

Similarly when ($\hat{y}_i = 0$),

this function works

$$-\log(1 - \hat{y}_i)$$

$J(w, b)$. Plotting this -



(again if prediction is close to 0, less $J(w, b)$, or penalty is closer to zero when next it is closer to ∞)

Gradient Descent remain same as linear Regression.

as when we take derivative of above Cost function we get the same value as, linear regression, only now definition of $J(w, b)$ is different.

• Let's find the derivative :-

$$A = -y * \log(g)$$
$$= -y * \log\left(\frac{1}{1+e^{-(wx+b)}}\right)$$
$$\frac{\partial A}{\partial w} = -y * 1 + e^{-(wx+b)} \cdot \underbrace{\frac{1}{1+e^{-(wx+b)}}}_{\frac{\partial}{\partial w}}$$

$$= -y * \left(1 + e^{-(wx+b)}\right) + \cancel{\frac{1}{1+e^{-(wx+b)}}} \cdot \cancel{\left(\frac{1}{1+e^{-(wx+b)}}\right)}$$
$$= -y * e^{-(wx+b)} \cdot x$$

$$= -y \cdot \frac{e^{-(wx+b)} \cdot x}{1+e^{-(wx+b)}}$$

$$= -y \left(\frac{1+e^{-(wx+b)} - 1}{1+e^{-(wx+b)}} \right) x$$

$$= -y x (1 - g(w,b))$$

or
 y

$$\partial B = -(1-y) \cdot \text{dlog} \left(\frac{1-y}{1+y} \right)$$

$$\frac{\partial B}{\partial w} = + \frac{(1-y)}{1-y} + \frac{1}{\partial w} \frac{1}{1+e^{-(w+b)}}$$

$$= (1-y) \cdot \frac{e^{-(w+b)} \cdot x}{(1+e^{-(w+b)})^2}$$

$$= (1-y) \cdot \cancel{(1+e^{-(w+b)})} \cdot \frac{e^{-(w+b)}}{\cancel{e^{-(w+b)}} \cdot (1+e^{-(w+b)})^2} \cdot x$$

$$= (1-y) \cdot x \cdot f(w+b)$$

Now for $\frac{\partial (\cos L)}{\partial w} = \frac{\partial A}{\partial w} + \frac{\partial B}{\partial w}$

$$= -y x (-f(w+b) + (1-y) \cdot x \cdot \cancel{f(w+b)})$$

$$= -y(x + \cancel{y}x f(w, b)) + x(f(w, b))$$

Date _____
Page _____

$$= -y x f(w, b)$$

$$= (f(w, b) - y) x$$

$$= (\hat{y} - y) \cdot x$$

Equal to the term in linear regression.

Similarly we can find b .

$$\frac{\partial A}{\partial b} = -y \cdot \log(\hat{y})$$

$$= -y \cdot \frac{2}{\hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$\frac{\partial \hat{y}}{\partial b} = e^{(wx+b)}$

$$= -y \cdot \frac{1}{\hat{y}} \cdot \frac{e}{(1+e^{-wx-b})^2}$$

$\frac{1}{\hat{y}} = \frac{1}{1+e^{-wx-b}}$

$$\frac{\partial A}{\partial b} = -y \cdot \frac{e}{(1+e^{-wx-b})^2}$$

$(1+e^{-wx-b})^2 = 1 + 2e^{-wx-b} + e^{-2wx+2b}$

$$\frac{\partial B}{\partial b} = -(1-y) \cdot \text{log}(1-y)$$

$$= + \frac{(1-y)}{1-y} \cdot \frac{\partial}{\partial b}(y)$$

$$= \frac{(1-y)}{e^{-(w_0+b)}} \cdot e^{-w_0 - b}$$

$$\frac{\partial \text{Cost}}{\partial b} = \frac{\partial A}{\partial b} + \frac{\partial B}{\partial b}$$

$$= -y \cdot \frac{e^{-(w_0+b)}}{1+e^{-(w_0+b)}} + (1-y) \cdot \frac{1}{(1+e^{-(w_0+b)})}$$

$$= \frac{1 - y(1 + e^{-(w_0+b)})}{(1 + e^{-(w_0+b)})}$$

$$= \frac{(f(w_0, b) - y)}{(y - f(w_0, b))}$$

Date: _____
Page: _____

another way of minimizing loss is to

maximize the inverse of the log likelihood

meaning

(the log-likelihood)

Considering this

Finally to convert the final Squashed Value to a probability or we call it Decision Boundary

as a threshold line above which model predict 1, below which 0. (Usually we take $y=0.5$)

Naïve Bayes Classification

Based on Bayes' Theorem:

It describes the probability of an event, based on prior knowledge of conditions that might be related to the event.

↳ "Likelihood"

$$* P(A/B) = \frac{P(B/A) * P(A)}{P(B)}$$

(PRIOR)

Conditional Probability →

probability of A happening
if B happened.
aka posterior probability.

+ Similarly for $P(B/A)$

$P(A)$ & $P(B)$ → probability of A & B respectively without any conditions.

For Machine learning :-
Say we have to predict
 (y_i) given that (x_1, x_2, \dots, x_n)
happened;

So,

$$P(y/x_1, \dots, x_n) = \frac{P(x_1, \dots, x_n/y)}{P(x_1, \dots, x_n)}$$

To understand how we use this, let's take the example of -
"Whether to play or not, on the basis of "Outlook",
"Humidity", "Temp",
"Wind".

Outlook Temp Humidity Windy play

	H	Hi	F	no
Rainy	H	Hi	T	no
R	H	Hi	F	yes
Og	M	Hi	f	yes
S	C	No	F	yes
S	C	N	T	no
OC	C	N	T	yes
R	M	Hi	F	no
R	C	N	F	yes
S	M	N	F	yes
R	M	N	T	yes
OC	M	Hi	T	yes
OC	H	N	F	yes
S	M	Hi	T	no

{ R → Rainy H → hot Hi → High }
 OC → Oceanside M → mild
 S → sunny C → cool N → Normal.

now we can further simplify the eqn.

$$P(Y/x_1, x_2, x_3, \dots, x_n) \Rightarrow$$

$$P(Y) * P(x_1, x_2, x_3, x_4, \dots, x_n / Y)$$

$$\Rightarrow P(Y) * P(x_1/Y) * P(x_2/Y) * \dots * P(x_n/Y)$$

Say $Y \rightarrow (T \text{ or } F)$,

then we can write :-

$$P(T/x_1, x_2, \dots, x_n) \Rightarrow P(T) * P(x_1/T), P(x_2/T), \dots, P(x_n/T)$$

Similarly for $P(N/x_1, x_2, \dots, x_n)$

* We can ignore the denominators as it is constant for all

So, now we compute these values

~~P(?) =~~

(PLAY)	P(Yes)	P(No)
Yes → 9	$\frac{9}{14}$	$\frac{5}{14}$
No → 5		

~~P(Outlook)~~

(Outlook)	Yes	No	P(Yes)	P(No)	or
Sunny	3	2	$\frac{3}{5}$	$\frac{2}{5}$	$\rightarrow P(\text{Sunny})$
Overcast	4	0	$\frac{4}{5}$	$\frac{1}{5}$	$\frac{1}{5}$
Rainy	2	2	$\frac{2}{5}$	$\frac{3}{5}$	$\frac{3}{5}$

(Temperature)

	Yes	No	P(Yes)	P(No)
Hot	2	2	$\frac{2}{5}$	$\frac{3}{5}$
Mild	4	2	$\frac{4}{5}$	$\frac{1}{5}$
Cool	3	1	$\frac{3}{5}$	$\frac{1}{5}$

(Humidity)

High

$P(\text{Yes})$, $P(\text{No})$, $P(\text{Yes})$, $P(\text{No})$

3

4

$3/9$

$P(\text{No})$

$4/9$

Normal

6

1

$6/9$

$1/5$

9

5

(Wind)

Yes, No $P(\text{Yes})$, $P(\text{No})$

Weak (F) 6 2 $3/9$ $2/5$

Strong (T) 3 3 $3/9$ $3/5$

9

5

$P(\text{Weak}/\text{Yes})$, $P(\text{Strong}/\text{Yes})$

Now, say we need to predict if
to play on a day with:

$X = (\text{Sunny}, \text{Hot}, \text{High}, \text{Strong})$

$$P(\text{Yes}/X) = \frac{9/14 * 3/9 * 2/9 * 3/9}{1} \Rightarrow 0.00352$$

$$P(\text{No}/X) = \frac{5/14 * 2/5 * 1/5 * 4/5 * 3/5}{1} \Rightarrow 0.02742$$

We can't really understand anything from this value so,

$$P(\text{Yes} / \bar{X}) = \frac{0.00589}{0.00589 + 0.02742}$$

$$= \frac{0.00589 \times 100}{16.529}.$$

$$P(\text{No} / \bar{X}) = \frac{0.92678 \times 100}{16.529}$$

$$= 92.678\%.$$

$$= 83.82\%.$$

So our prediction will be
(No)

Evaluating an Classification Model -

{other than logistic loss}

• Confusion Matrix.

Say, I have some data, on which I have trained a model and made some-prediction.

$$\begin{array}{l} y \rightarrow 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \\ \hat{y} \rightarrow 0 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \end{array}$$

(Predicted)

'NO(0)' YES(1)

Actual No(0)	NO(0)	2	1	3
	YES(1)	1	1	2
Actual YES(1)	NO(0)	1	1	2
	YES(1)	1	1	2

Confusion Matrix divides our model's prediction into 4 categories, which helps us understand how confused our model is.

• True Positive (TP) : predicted $\rightarrow 1$
Actual $\rightarrow 1$

• False Positive (FP) : prediction $\rightarrow 1$
Actual $\rightarrow 0$

• False Negative (FN) : prediction $\rightarrow 0$
Actual $\rightarrow 1$

• True Negative (TN) : prediction $\rightarrow 0$
Actual $\rightarrow 0$

Note: Confusion Matrix will change with
change in the threshold value -

→
"Now using this Confusion Matrix
we can determine one of five
measures of evaluation for our
model :-

* Accuracy $\rightarrow \frac{TP + TN}{TP + TN + FP + FN}$

→ It is simply how accurate
our model is predicting

(work good for balanced dataset) -

In Case of Imbalanced dataset :-

Say we have 900 True
100 False

~~Even if our model predicts True every time it will have a 90% Accuracy, but we know the model is performing poorly.~~

For, Imbalanced dataset we need to try something else.

First, we need to understand what type of prediction our ~~model~~ needs to make. i.e.

Say, Diagnosing a Cancer patient:

In this case a false negative is very dangerous as it means the person has cancer but we predicted not, then making a false positive.

So, we need our model to predict True values much more and as accurately as possible.

ex

Say, Giving Loan :

Denying a good customer loan
(False positive) is not as bad as
giving it to a bad customer.
(False negative).

Here /

precision and Recall comes in.

Precision :

$$\Rightarrow \frac{TP}{TP + FN}$$

Here we ask the question :

"Out of all the prediction, How accurately
how is predicted positive Class?"

Recall :

$$\frac{TP}{TP + FN}$$

Precision measures the accuracy of the positive prediction made by the model.

* High Precision means our model is likely to predict and actual positive (as positive).

e.g.: 10 True predicted \rightarrow 7 actual True (70%)
Precision.

Recall $\Rightarrow \frac{TP}{TP + FN}$

Here we ask the question:

"Out of all positive cases in sample, how many of those are correctly predicted by our model?"

Recall measures the ability of model to detect positive sample.

High Recall model means our model accurately predicting the positive samples.

e.g.: 10 prediction \rightarrow 3 FN (Actual Positive)

ex 10 Prediction made by our model.

5 True Positive

5 False Negatives
(meaning actual value is Positive)

So, our model has 50% Recall.

meaning it is

(predicting 50% of time

an

Positive Value as actually positive.

where in precision, out of

all the "True" prediction made, you were finding
How many of them are actually True).

(Trade off) → when we ↑ threshold
we make stricter prediction about
what is to be considered as True or
False.

i.e. ($TP \uparrow +$, $FN \uparrow +$)

($FP \downarrow -$)

i.e. (Precision $\uparrow +$, Recall $\downarrow -$)

$\{ F_1\text{-Score} \}$

~~This f1 score is kind of middle ground between the p~~

This creates a bit of middle ground between a high Recall vs high prediction model.

$$\Rightarrow 2 * \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

"Harmonic mean"

F_1 itself is actually a special case of F_β score.

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{Precision} \cdot \text{Recall}}{(\beta^2 \cdot \text{Precision}) + \text{Recall}}$$

for, $\beta > 1 \rightarrow$ recall P^+
 $\beta = 1 \rightarrow (F_1)$

$\beta < 1 \rightarrow$ Precision P^+

{ ROC & AUC → (area under ROC)

G (ROC Curve Operating Characteristics)

is Created by plotting True Positive Rate

or false positive Rate (FPR)

* → TPR, Recall, Sensitivity = $\frac{TP}{TP+FN}$

* Specificity = $\frac{TN}{TN+FP}$

* → FPR, (1-Specificity) = $\frac{FP}{FP+TN}$

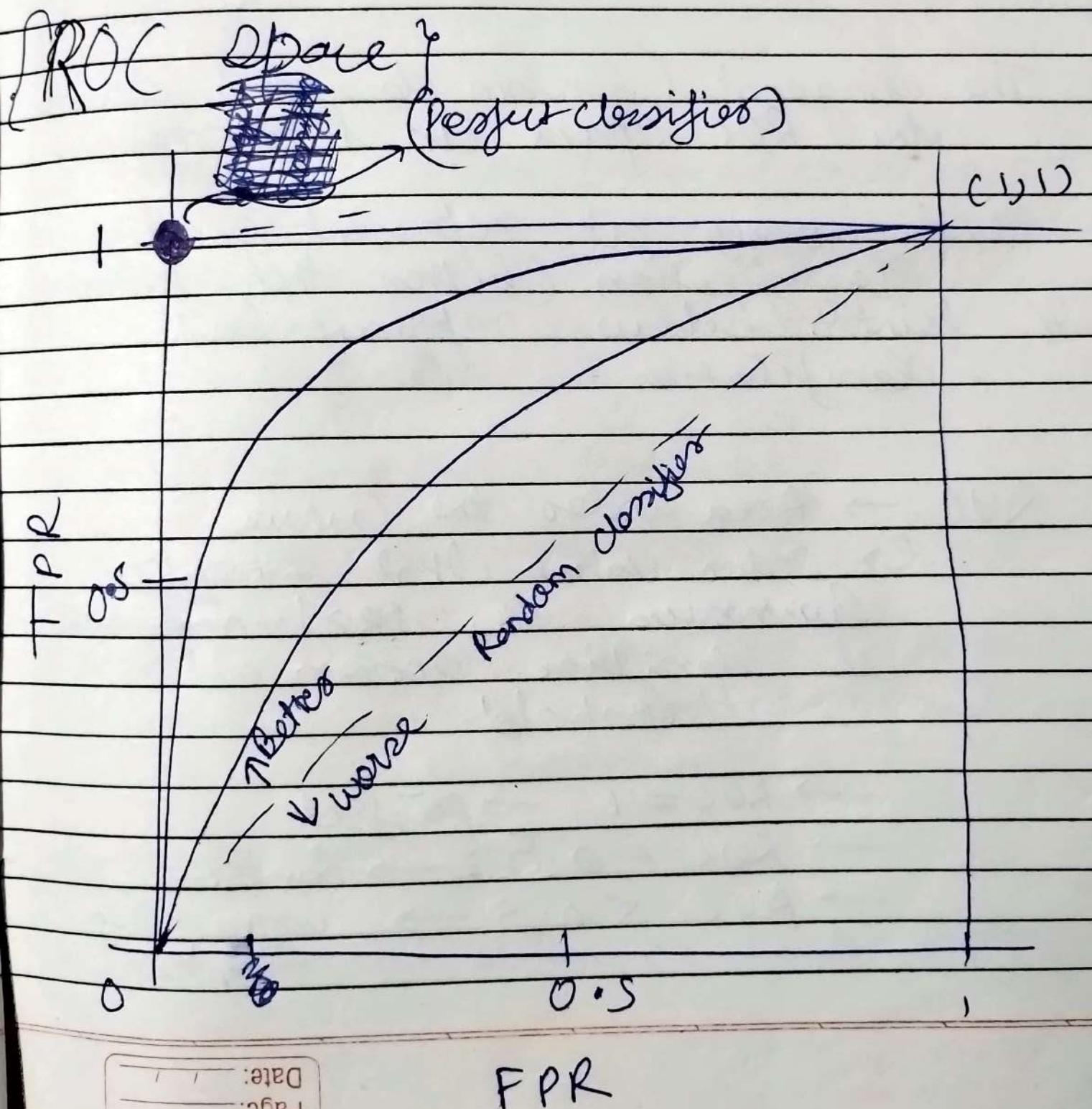
ROC → helps evaluating performance
of classification model by plotting

TPR vs FPR at various
threshold settings.

First we calculate TPR & FPR for different threshold (0...0.9).

then, we plot the TPR & FPR of (Y) (X) a corresponding threshold value.

~~3~~



In the graph best possible prediction method would be point the point "perfect classifier"

Representing

100% sensitivity (NO FN)

100% specificity (NO FP)

The diagonal (Random classifier), divides the ROC Space in two regions.

Point's above it represent good classification (better than random) & Point's below represent bad classifications.

AUC → Area under the Curve

↳ Decade value that ~~classifies~~ summarizes the performance of classifier across all threshold.

→ AUC = 1 → perfect

→ AUC = 0.5 → Random

→ AUC < 0.5 → worse than Random

{ Phi Coefficient or }

Matthews Correlation Coefficient

~~(standardized mean of regression
coefficients of problem and its class)~~

$$\cdot \text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

→ Regarded as the balanced measure
as it takes into account true & false
positive and negatives and can be used
even if the classes are of
very different size.

If Returns Value between :-

-1	0	+1
perfect Total disagreement	no better than random prediction	perfect prediction

~~for complete fit~~

{ The closer the value to 1, the better. }

{ Cohen's Kappa } → Probability of observed agreement

$$K = \frac{P_o - P_e}{1 - P_e} \rightarrow \begin{matrix} \text{Probability} \\ \text{of agreement} \\ \text{by chance.} \end{matrix}$$

$$\left. \begin{matrix} 2 \times (TP \times TN - FN \times FP) \\ (TP + FP) \times (FP + TN) + (TP + FN) \times (FN + TN) \end{matrix} \right\}$$

* It is used to measure "inter-rater reliability" for a categorical item.

{ inter-rater reliability is the degree of agreement among independent observers who rate, code or assess the same phenomenon.
 → In our case agreement between our model & Actual Values)

for example ,

Say we have 2 Judges to segregate
3 Colored balls in different bags.
Say Red, green, blue.

If both Judges Judge without bias,
they would agree on color of all balls
i.e. they will have zero disagreement.

but it's not possible , so we want to
measure the agreement between them.

{ Same they will agree actually }
Same they will agree by chance

We don't
want our ML
Model doing
this .

} by chance for our
context means, if we
start randomly guessing
we still might get an
agreement .

• We need to measure only the actual
agreement .

P_{observed} → Probability of observed agreement.

$P_{\text{by chance}}$ → Probability of ~~Obser~~ agreement by chance.

Example this is the Confusion Matrix
(JI)

	R	G	B	Total
R	44	5	1	50
G	7	20	3	30
B	9	5	6	20
Total	60	30	10	100

(Observed)

The Diagonal Values are the P_{observed} Values: ~~& Total~~

$$\textcircled{1} \quad P_{\text{Ob}} = \frac{44 + 20 + 6}{100} = 0.70$$

Or

$$0.44 + 0.2 + 0.6 = 0.70$$

new) for probability by chance

out of 100 ball how many did Judge 1
classify as red \rightarrow 60

$$P_{J_1 R} = 60/100 = 0.6$$

Similarly, $P_{J_2 R} = 50/100 = 0.5$

both classifying Red

$$P_{\text{by chance for Red}} \rightarrow P_{J_1 J_2 R} = 0.6 \times 0.5 = (0.3)$$

Similarly ~~$P_{J_1 J_2 R} = 0$~~

$$\begin{aligned} J_1 &\rightarrow \text{Red} \\ J_2 &\rightarrow \text{green} \end{aligned} \quad \left. \begin{array}{l} 0.6 \\ \times \\ 0.3 \end{array} \right\} = 0.18$$

$$\begin{aligned} J_1 &\rightarrow \text{Red} \\ J_2 &\rightarrow \text{Blue} \end{aligned} \quad \left. \begin{array}{l} 0.6 \\ \times \\ 0.2 \end{array} \right\} = 0.12$$

Similarly we construct a new confusion matrix:

		TC		
		Red	green blue	Totals
Red	Red	0.30	0.15	0.05
	green	0.18	0.09	0.03
blue	blue	0.12	0.06	0.02
Total		0.6	0.6	0.11
(by chance)				

New $P_{\text{by chance}}$ is the diagonal values

$$= 0.30 + 0.09 + 0.02 = 0.41$$

$$\text{So, } K = \frac{P_{\text{obs}} - P_{\text{by chance}}}{1 - P_{\text{by chance}}} = \frac{0.70 - 0.41}{1 - 0.41}$$

$$= \underline{0.4915}$$

A value of $K > 0.75$ is considered good.

(The one we just did is the Un-weighted)

Weighted Kappa :-

Currently our data had distinct value but say if we had "Shades of Red" (ordinal categories)

In this case we calculate "weight".
Say we have 5 shades of Red and we need to assign value from 1 to 0

In order to do this we define a "Distance" variable.

(Linear) ↘

$$R_1 \quad \text{Dist} \quad R_2 \rightarrow 1 - \frac{\text{dist}}{\max \text{dist}}$$

$$R_1 \quad 1 \quad 0$$

$$R_2 \quad 1 \quad \rightarrow 0.75$$

$$R_3 \quad 2 \quad (\text{Quadratic})$$

$$R_4 \quad 3$$

$$R_2 \rightarrow 1 - \left(\frac{\text{dist}}{\max \text{dist}} \right)$$

$$R_5 \quad 0 \quad 1$$

$$\rightarrow 0.99$$

(Linear) → weight are linearly spaced between 0 to 1)

Σ , (linearity)

(Quadratic)

	R_1	var		R_1	var
R_1	1	0	R_1	1	0
R_2	0.75	1	R_2	0.94	1
R_3	0.50	2	R_3	0.75	2
R_4	0.25	3	R_4	0.44	3
R_5	0	4	R_5	0.00	4

(Similarly we can calculate for R others)

* Similarly four are the day we need to classify,

Dark Red	Red	Light Red
Red	Light Red	Dark Red

(linear)

	Dark red	red	light red
DR	1	0.5	0
R	0.5	1	0.5
L	0	0.5	1

→ Now we can calculate Kappa like in Un-weighted) ←

Evaluating a regression Model:-
(we know MSE &)

{ Mean Absolute Error }

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

Penalizes the large errors less than MSE,
but gives a better idea of average
error magnitude.

{ Root Mean Square Error }

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

{ Mean Absolute Percentage Error }

$$\text{MAPE} \rightarrow \frac{100}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

• R-squared (R^2) → coefficient of Determination }

↳ R^2 is the proportion of Variation in the dependent Variable that is predictable from the Independent Variables.

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

$$SS_{\text{res}} \rightarrow \left(\sum_i (y_i - \hat{y}_i)^2 \right)$$

↳ residual sum of squares.

(error)

variable

$$SS_{\text{tot}} \rightarrow \left(\sum_i (y_i - \bar{y})^2 \right)$$

(total sum of squares)

mean(y_i)
actual values

best case $R^2 = 1$, worst case $R^2 = 0$
(ignorance)

R^2 value usually range between 0 to 1
but can go < 0 as well. In this case
model is performing very bad, than just
taking the mean value.

There are limitations with R^2 :-

- A high R^2 does not imply causation between independent and dependent Variable.
- R^2 measure how well the independent variables are associated with the dependent Variable, not whether they cause the change.
 \rightarrow ex: if you have a high R^2 between ice-cream sales and temperature, it shows they are related but it doesn't mean temperature directly causes people to buy ice cream.

• Omitted Variable Bias:

↳ It doesn't reveal if imp Variables
are have been left out of the Model.

• Correct Regression:

high R^2 doesn't mean you have,
chosen best model types for
your data.

• Appropriate Variable:

↳ R^2 does not indicate if the
most appropriate set of
independent variable has been
chosen

• Collinearity:

• Transformation:

↳ R^2 does not indicate if
model might be improved
by using transformed
Version of Variable.

• Sample size:
→ R^2 doesn't indicate if there are enough data points to make solid conclusions.
we can get a high R^2 with very less data points.

• Collinearity → R^2 does not reveal if collinearity is present
→ (when independent variables are highly correlated)

It is bad because:

↳ (It makes it difficult to understand effect of each variable,
if it adds redundancy)
(increasing our model complexity)

~~(To Address the issue of R² adjusted is adjusted)~~

{ Adjusted R², \bar{R}^2 }

↳ address the problem

of Overfitting to

better for model
comparison.

$$\bar{R}^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

$n \rightarrow$ number of
observation

$p \rightarrow$ is the number
of predictors
(features)

• 2 Train Test (Cross Validation) Split.

In Machine learning we split our data in 3 categories - (Train, Test, CV)

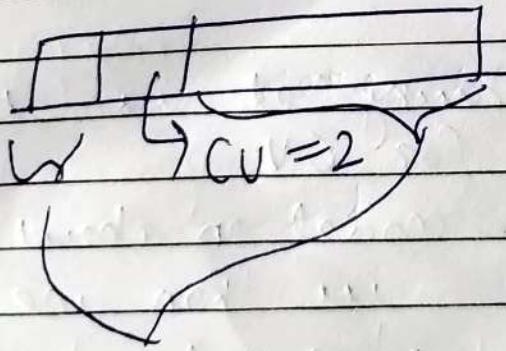
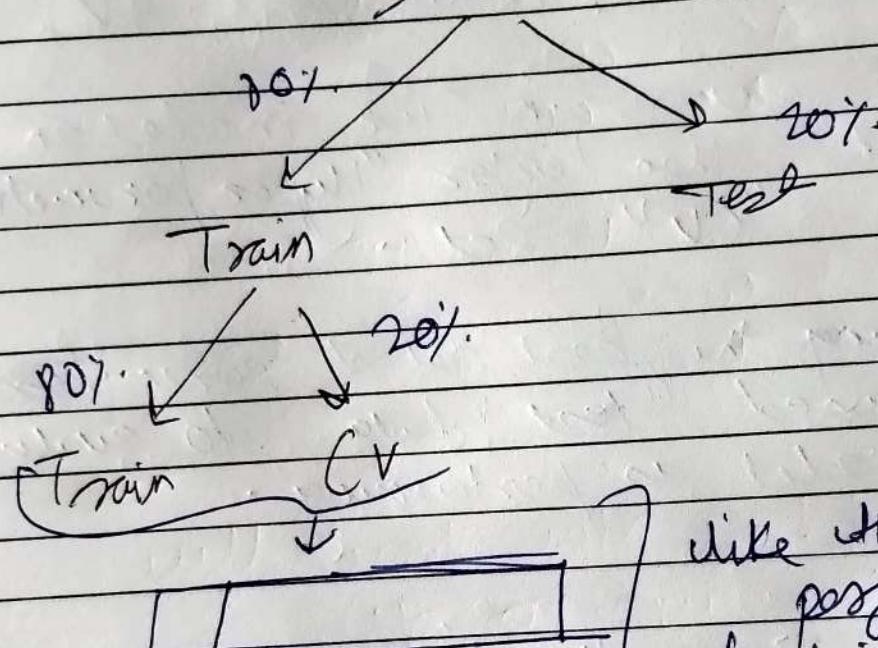
Out of this we train our model on "Train" set, then perform "Hyperparameter tuning" on "CV" set.

After finding the best set of parameters, we then used "test" data to evaluate how our model is performing.

- Training Set → usually about 70 - 80% of total Data. It is used to initially train the model.
 - From this training set we take out 20% for CV
- After training our model on about 60% of data, we then use CV data for Hyperparameters tuning.
- Test: After the above step is completed we then use the chosen set of hyperparameters & evaluate our model on test data.
(generally about 20 - 30%)

Types of taking Cross Validation

① LOOCV (Leave one out CV)
Dataset 100%.



like this we,
perform "one"
"folding" on the
data / cut the
dataset into
leaving a little
part out of
CV.

Generally not used due to being
"Time Consuming".

also this is prone to
("Overfitting").

② Hold out CV

Streaming }

"this time" we split the data base first by creating a Random-state (shuffled).

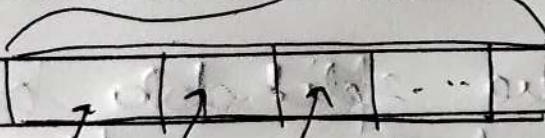
③ K Fold Cross Validation.

(Train) \rightarrow 1000 points

out $K = 10$

then $CV = \frac{1000 - 100}{10}$

10 section Rest 900 in Train.

$CV=1$  \rightarrow Acc 1 ... K)

$CV=2$ (we take Avg Acc)

$CV=3$ (rest we will use for training)

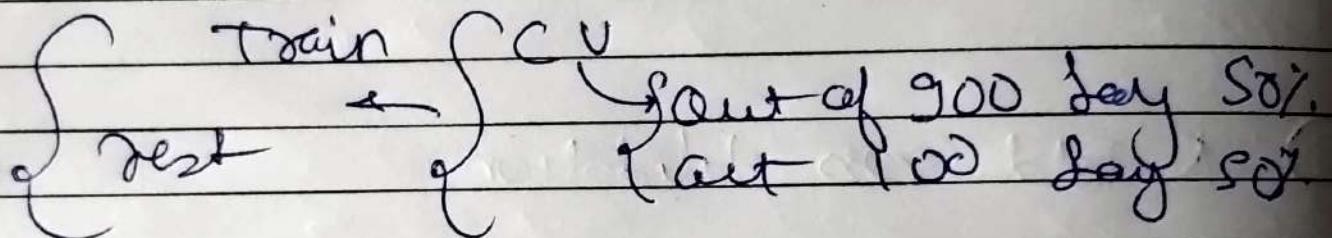
(Both of these have a problem of imbalanced dataset)

④ Stratified K fold CV

(Here we take equal proportion of data from different categories)

say

$$\begin{array}{l} 900 \rightarrow A \\ 100 \rightarrow B \end{array}$$



(will work with Tuberculosis dataset up to some extent)

⑤ Time Series Cross Val.

Here we split data based of time, older is used for training, newer is used for CV.