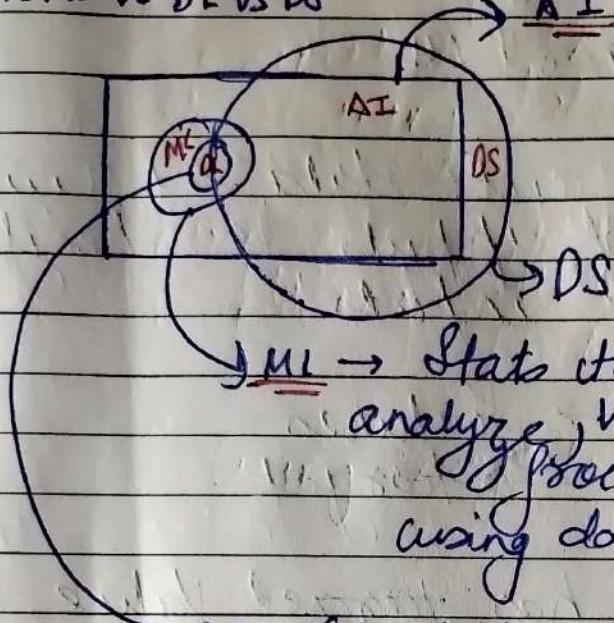


DATA SCIENCE

- 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 & often 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 is 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 Algos.

→ K-Means

→ DBSCAN

→ Hierarchical

↳ Anomaly Detection

Isolation Forest

One Class SVM

(Many more categories & algorithms.)

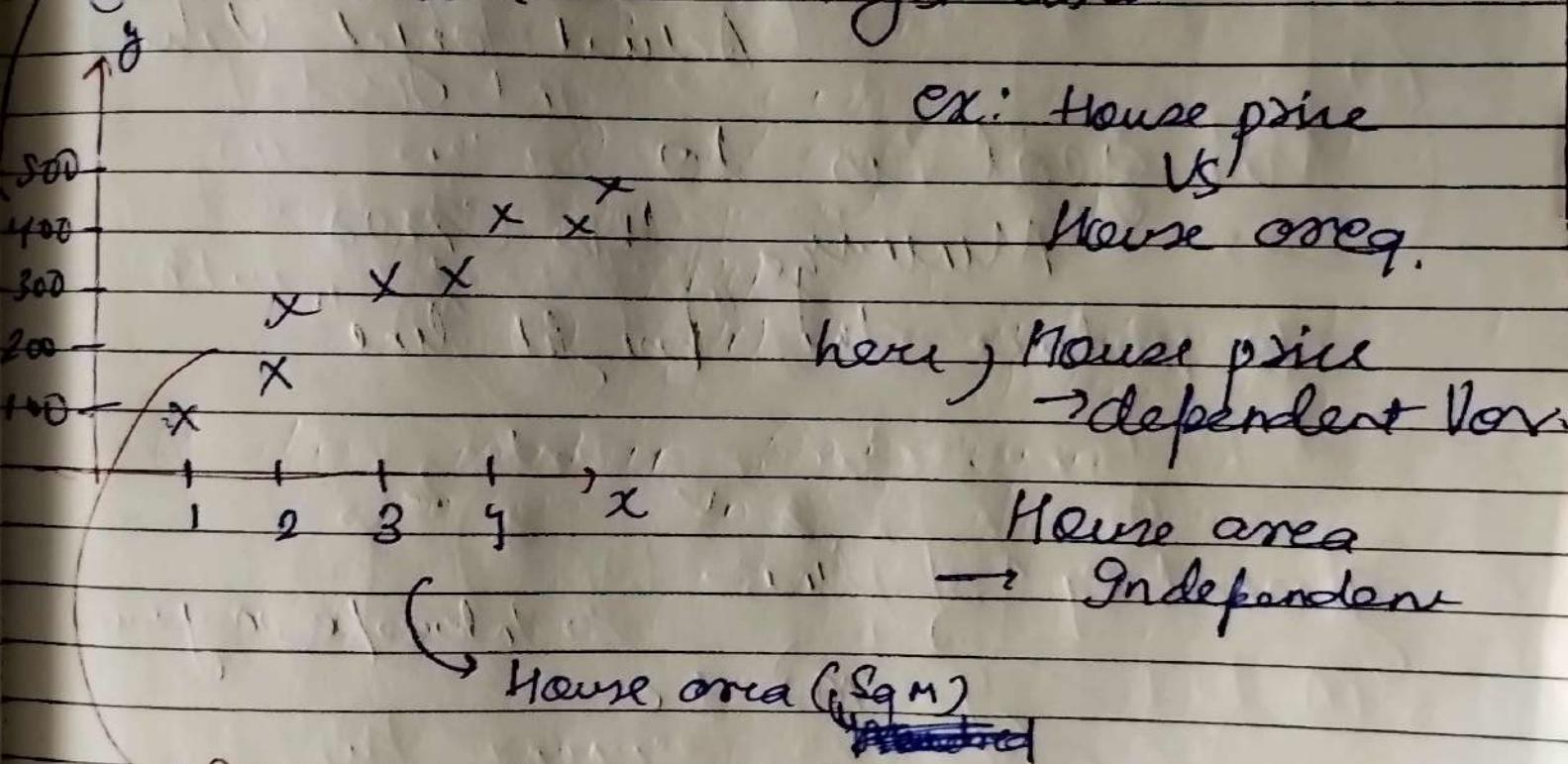
(Still learning - - -)

Supervised ML

Linear Regression

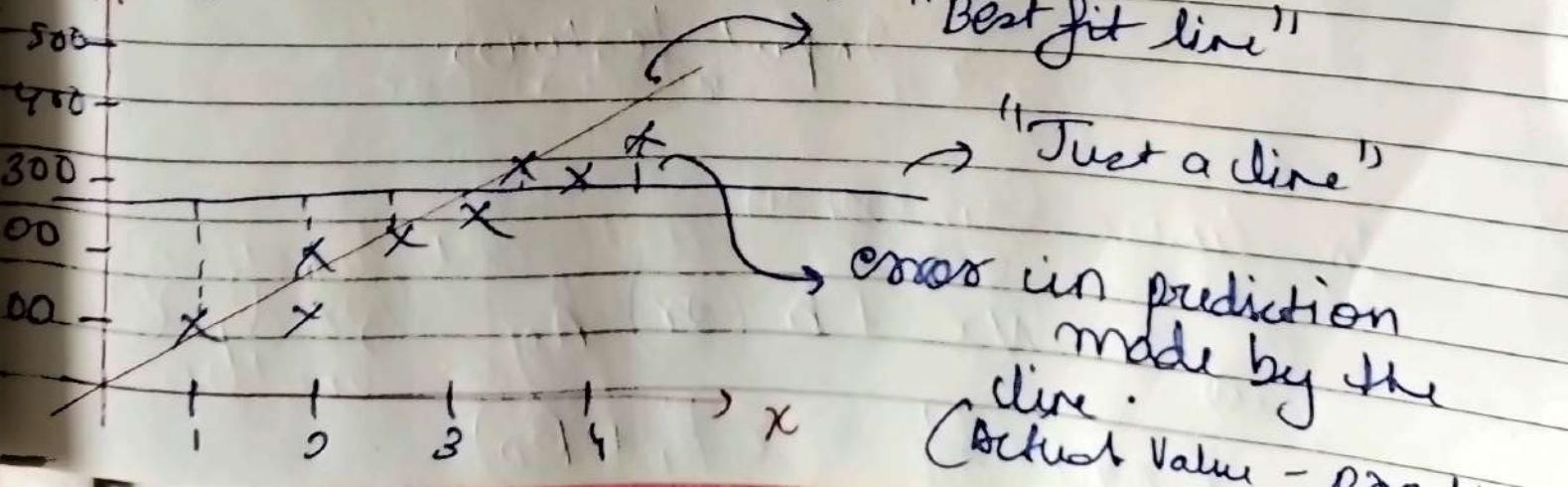
→ here we try finding an eqn
that best relates our inputs
to our output. This line is
called the "best fit line".

Monte Carlo



→ Data points

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

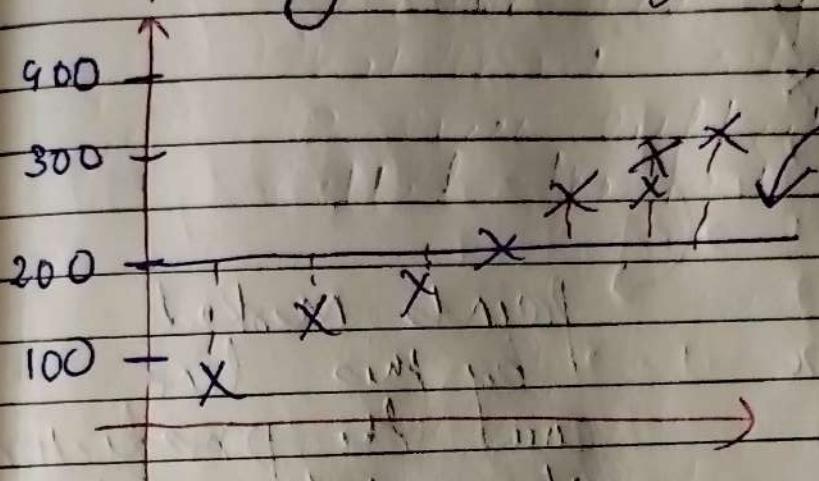
In linear Regression :-

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

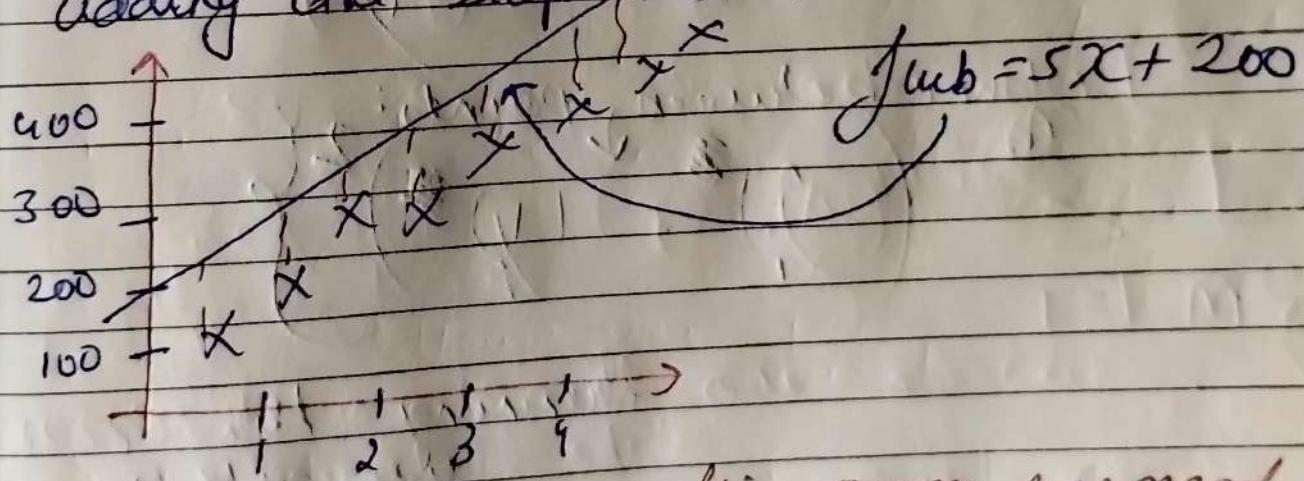
$$(\text{or}) \quad Y_i = X\beta_1 + \beta_0$$

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

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



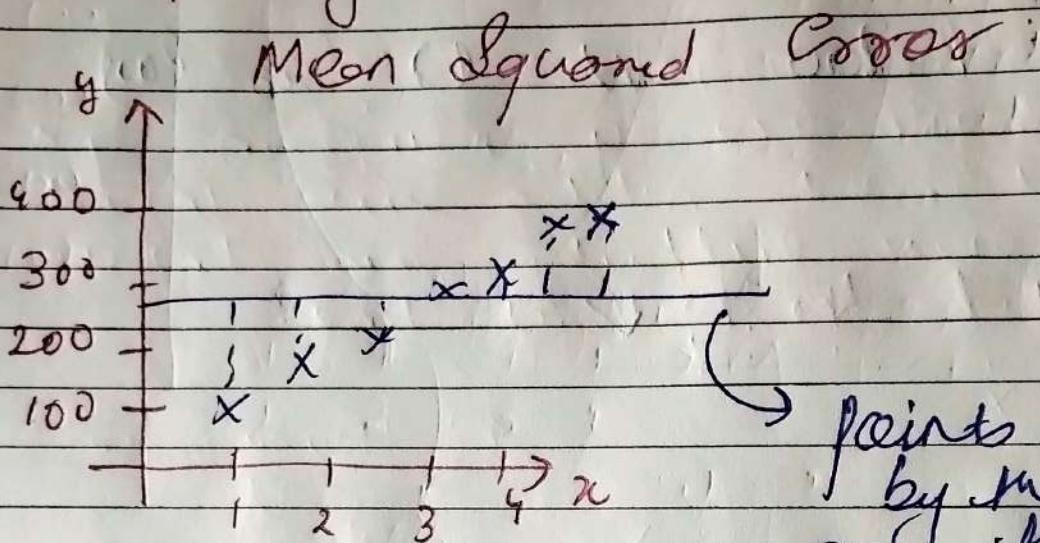
adding the slope term: $b = 200$, $W = 5$



So, to minimize this error we need
to find correct values of
 W & b .
More generally:

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

Computing the error:



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

Distance between points:

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

Running this error
over all the
data points

here for Calculation case

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

Represented as :

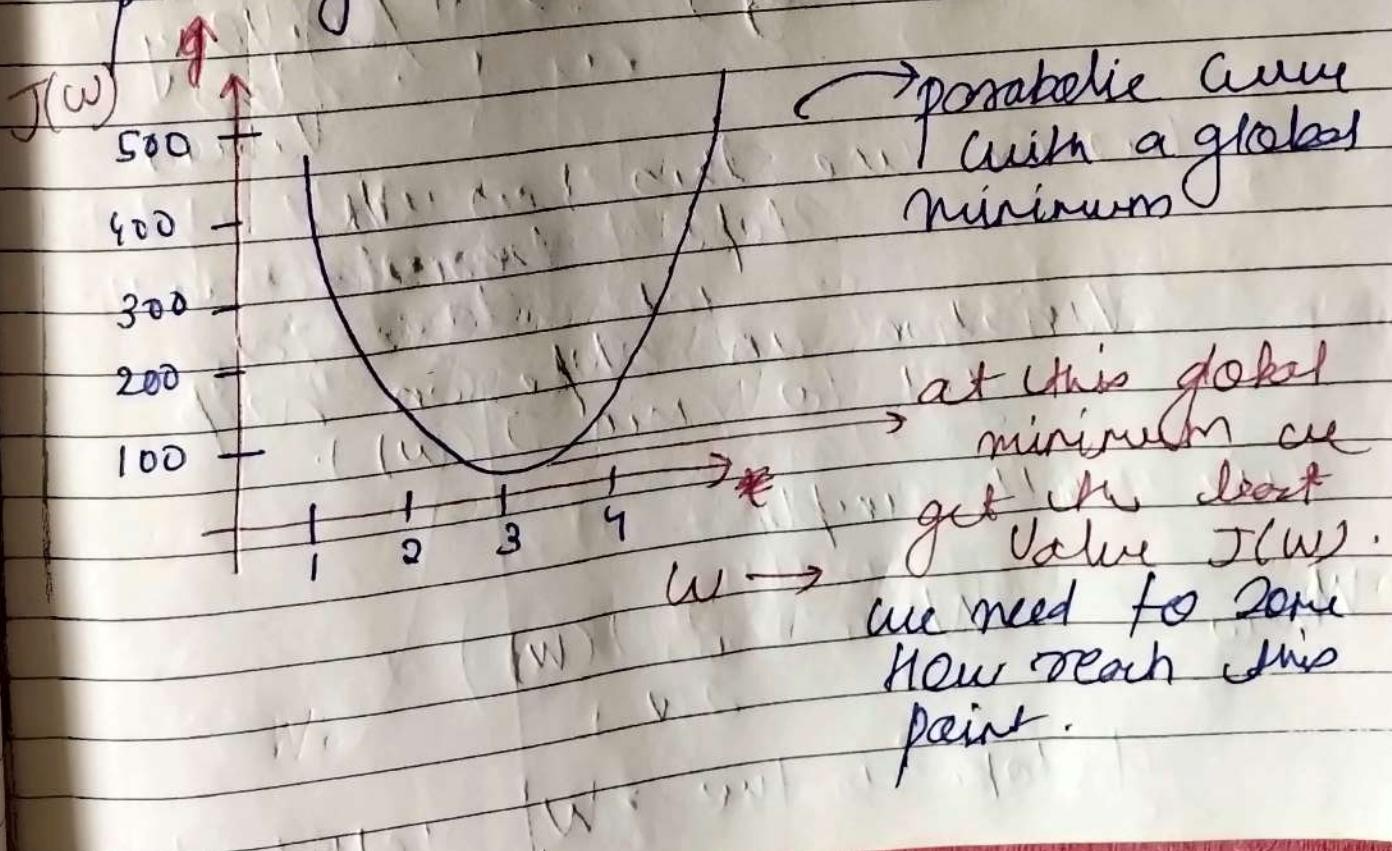
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m ((y_i) - (\hat{y}_i))^2$$

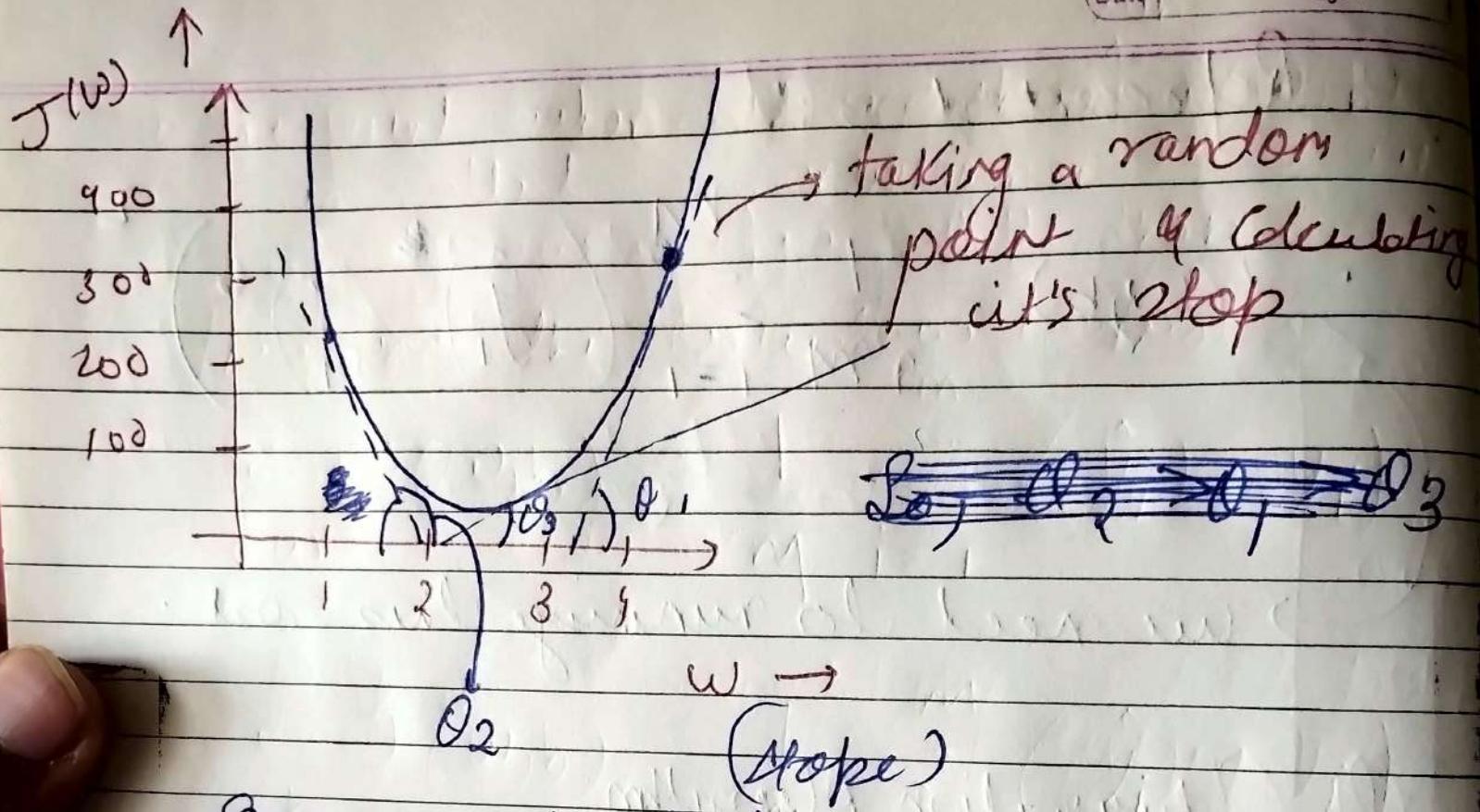
→ we need to minimize this Cost.

① Optimization Algorithm:

① "Gradient Descent Algo"

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





Now, $\tan(\theta_1) \rightarrow$ +ve

$\tan(\theta_2) \rightarrow$ -ve

$\tan(\theta_3) \rightarrow (+/-)$ (very less 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 $\rightarrow w \uparrow \rightarrow$ getting closer to min $J(w)$.
(if we are already near min $J(w)$,
we won't change the value
of desciptor will be small, so
no significant \uparrow or \downarrow)

Going from right to left $\rightarrow w \downarrow$.

↓ we encounter min $J(w)$

Similarly left to right $\rightarrow w \uparrow$

→ 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." \rightarrow "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:

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

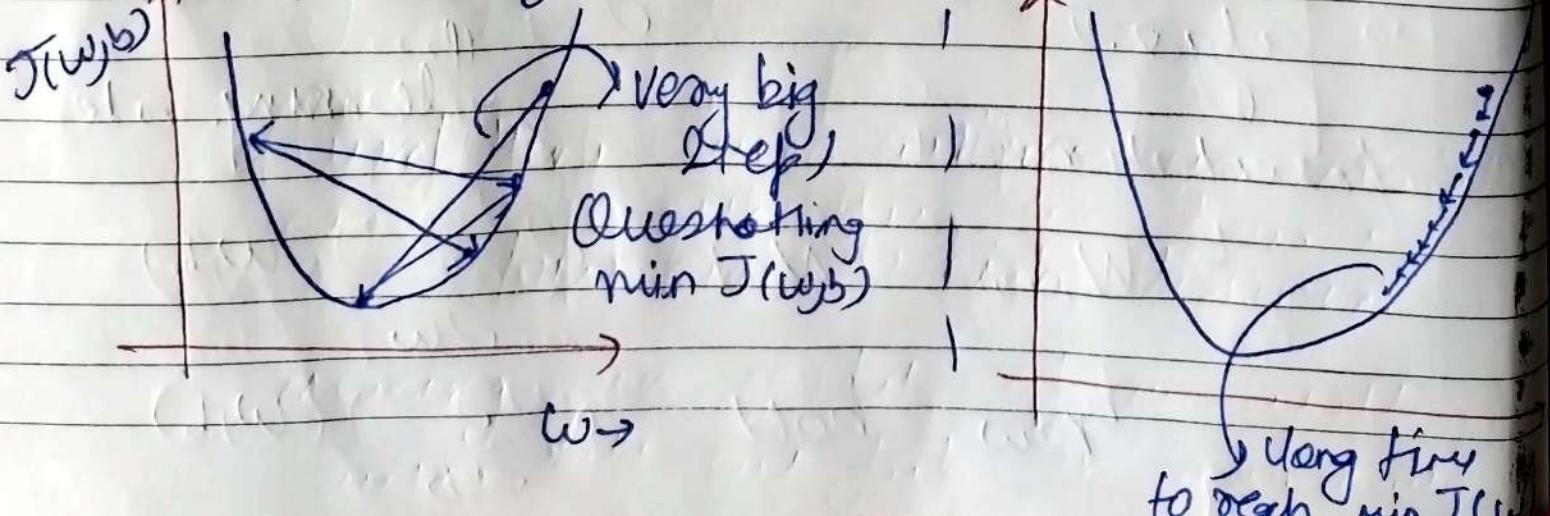
$$\frac{\partial J(w, b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m ((\hat{y}_i) - y_i)$$

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

take while updating of parameters (w, b)

$\alpha \rightarrow$ close

$\alpha \ll$ small



all, need to pick alpha Davis where it's in
the middle of this.

Matrix Notation of linear Regressions

$$y = XW + \epsilon^{\text{residuals}}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

(for pios/
Intercept: no q features
in e
term) example

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$

number of examples

$$W = \begin{bmatrix} b \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_p \end{bmatrix}$$

Intercept term
coefficients / slope

another way to find the parameters

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

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.

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

(no of features) \rightarrow Regularization parameter

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

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

- 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:-

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

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

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 → (with 1 feature)
- ② Multivariable Linear Regression
→ (with multiple features)
- ③ Multivariate
↳ predicting more than one variable on the basis of multiple input features
ex: (predicting if student passes in English & Hindi, Science, etc.)

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:

$$P(z)$$

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

$$(z = wx+b)$$

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 $p(1/2) = 1/2$)

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:

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

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 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.

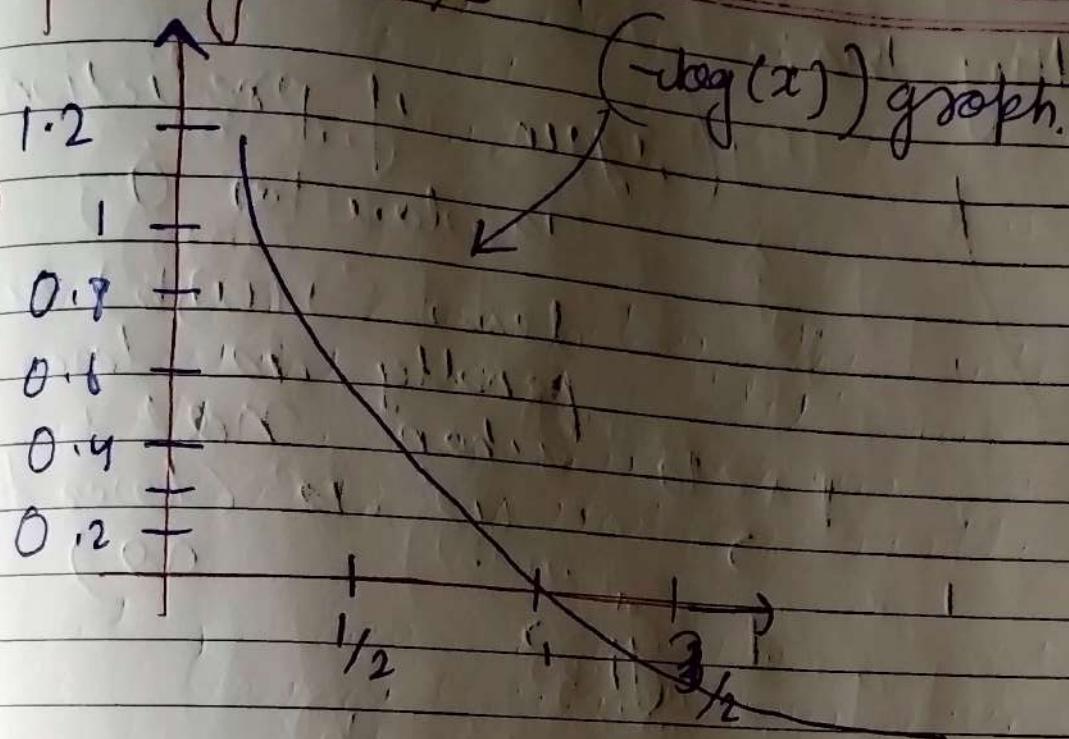
Q) this function works:-

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

plotting this:

Date:

Page:



(\hat{y}_i)

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

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

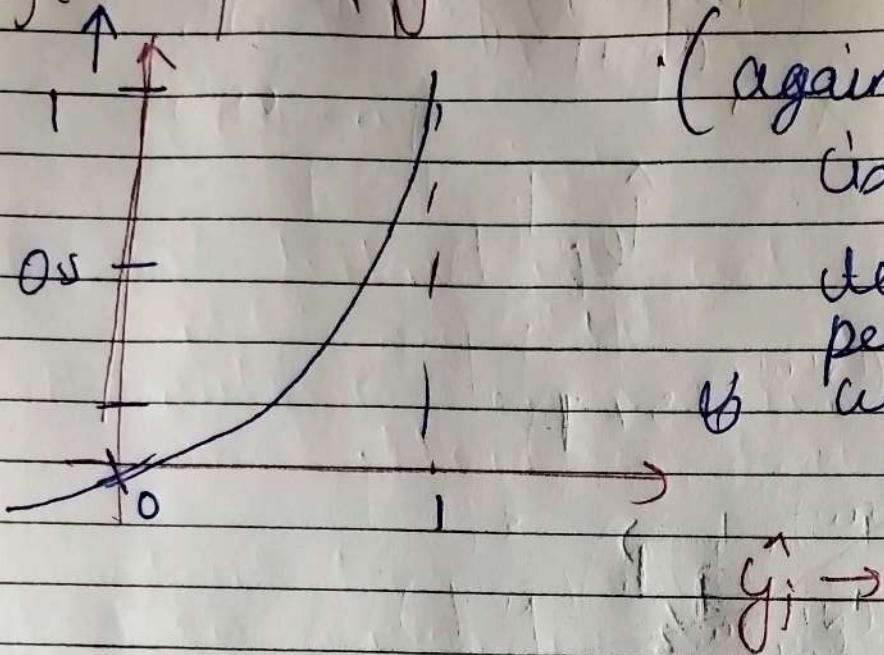
heavily penalizing the wrong predictions.

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 :-

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

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

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

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

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

or
 y

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

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

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

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

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

Now for $\frac{\partial A}{\partial w} = \frac{\partial A}{\partial w} + \frac{\partial B}{\partial w}$

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

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

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 .

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

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

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

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

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

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

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

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

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

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

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

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

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

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

another way of minimizing loss is to

maximize the inverse of the log likelihood

meaning

(the log-likelihood)

Or

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. It 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

aka
posterior

} if B happened -

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 N → Normal.
 S → Sunny C → Cool

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(~~rainy~~)

(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	$\frac{2}{5}$	$\frac{3}{5}$	$\frac{3}{5}$
	9	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	$\frac{3}{5}$	$\frac{1}{5}$
	9	5		

(Humidity)

High

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

3

4

$3/9$

$4/9$

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:

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

$$P(\text{Yes}/X) = \frac{9/14 * 3/9 * 2/5 * 3/9}{1} \Rightarrow 0.00352$$

$$P(\text{No}/X) = \frac{5/14 * 2/5 * 1/5 * 4/5 * 3/9}{1} \Rightarrow 0.02742$$

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

$$P(\text{Yes} / \bar{X}) = \frac{0.00589}{0.00589 + 0.02712}$$

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

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

$$= 92.678\%.$$

$$= 83.82\%.$$

So our prediction will be
(No)