

# LOGISTIC REGRESSION

↓ ↓ START FROM HERE ↓ ↓

PREREQUISITE For Deep learning  
Logistic Regression:-

- ① Logistic Regression can be seen as a model of a ~~uron~~ neuron
- ② Logistic Regression are combined together to give us a neural network. As such it is the foundation for deep learning.

③

Statistics VS Machine Learning:-

Machine Learning:-

- ① Linear models are a tool to introduce common ML techniques
- ② Loss function
- ③ Minimizing the loss function using calculus
- ④ Train and test sets
- ⑤ Hyperparameters, generalizations, overfitting, regularization.
- ⑥ Same techniques apply on neural networks, SVM etc.

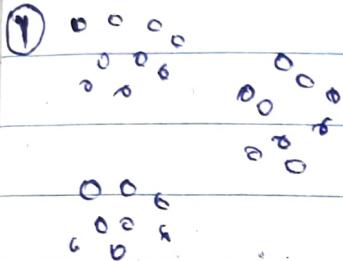
Statistical approach:-

- ① You will spend your entire degree learning about linear models.

# Machine Learning :-

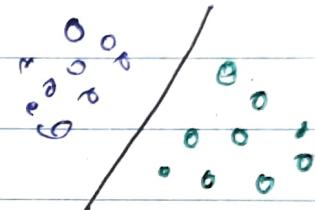
## Machine Learning

↓  
Unsupervised



①

↓  
Supervised

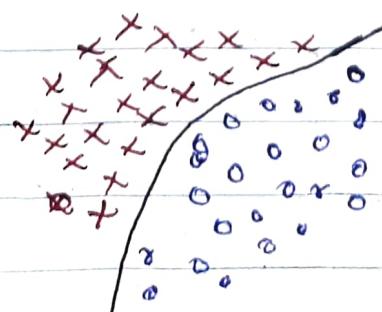


②

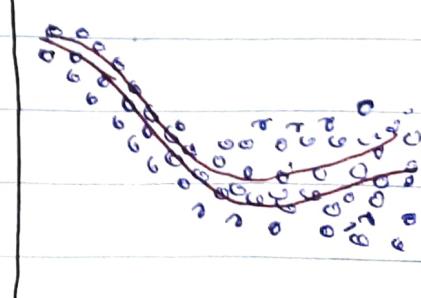
② In unsupervised learning, you usually don't have a label so you are forced to split the data into groups by looking for patterns.

② In supervised learning we already have a way to split the data

classification



Regression



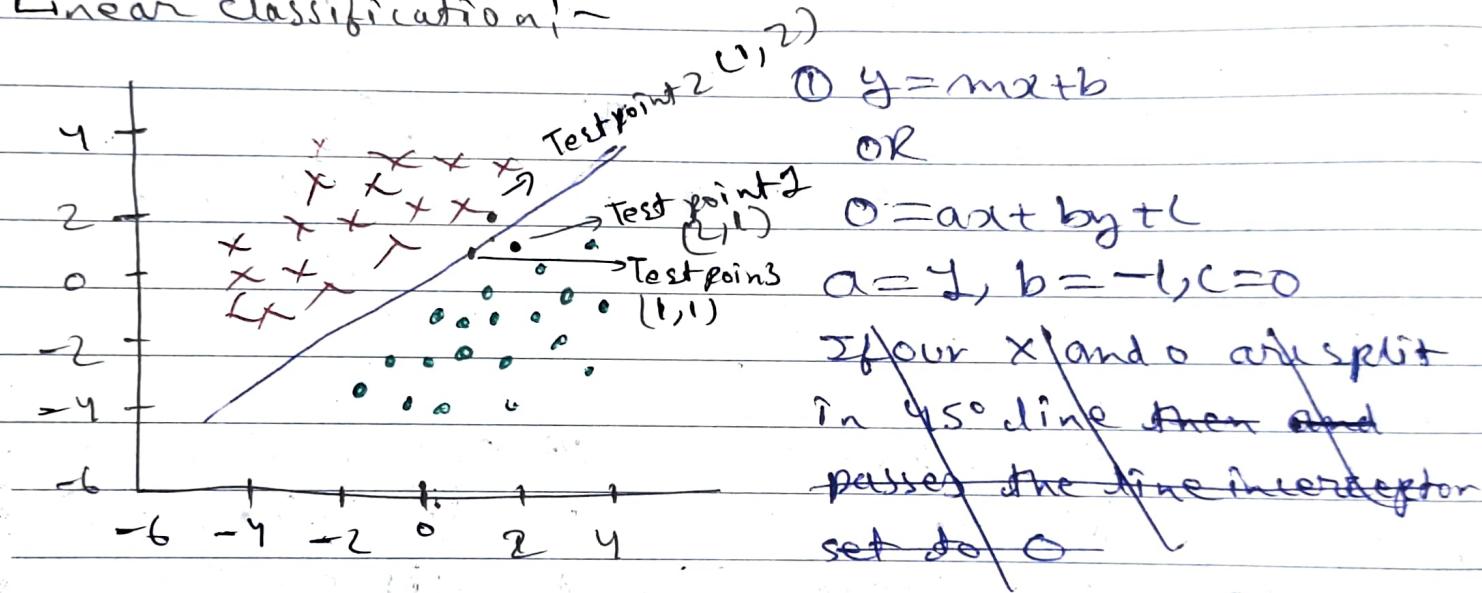
Trying to split the data into categories

Trying to predict the value of some function.

## Logistic Regression:-

- ① It is a model that is used for classification
- ② example image classification

## Linear classification:-



① You can see that If our  $x$ 's &  $o$ 's are split along a  $45^\circ$  line that crosses the line ~~the east~~ set to zero then  $A = 1$   $B = -1$  &  $C = 0$  should be the line so in other words our line is  $o = x - y$

③ So Now let's say we have a test point at  $x = 2$  &  $y = 1$  If we put plug that into

$$h(x, y) = x - y \quad h(2, 1) = 2 - 1 = 1 > 0 \rightarrow 'o'$$

Therefore we should classify this test point as an 'o'

④ Now let's say we have another test point (testpoint 2) at  $x=1, y=2$ . If we plug that into  $h(x,y) = x-y$   
 $= 1-2 = -1 < 0 \rightarrow (x)$

Therefore we should classify this test point 2 as an  $(x)$

⑤ Now let's say we have another test point (testpoint 3) at  $x=1, y=1$ . If we plug that into  $h(x,y) = x-y$   
 $= 1-1 = 0 \rightarrow \text{don't know}$

Therefore we cannot tell whether test point 3 is an  $x$  or  $o$

Understanding Linear classification in machine learning terms:-

We call  $(x,y) \rightarrow (x_1, x_2)$

We rename A, B & C constraints to  $w_i$

We call the bias term / intercept  $w_0$

$h(x) = w_0 + w_1 x_1 + w_2 x_2$

We say  $h()$  hypothesis function is a linear combination of the components of  $x$

In vector form :  $h(x) = w^T x$  (Dot product)

In 3 Dimension : line is transformed into a plane

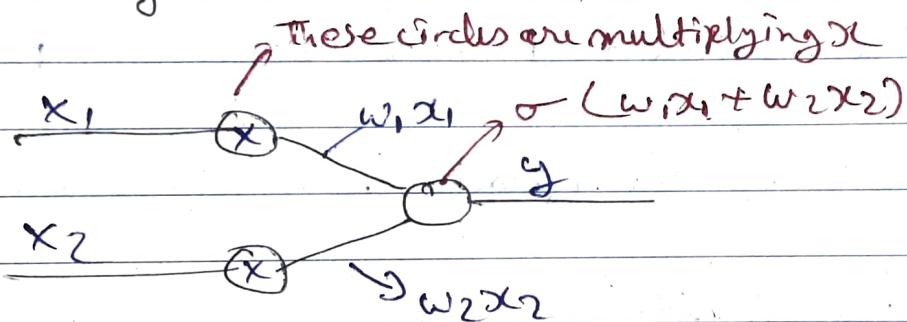
In  $\geq 3$  dimension : hyperplane

So If we have a function which is not a line function separating our data we could call that a hyper surface.

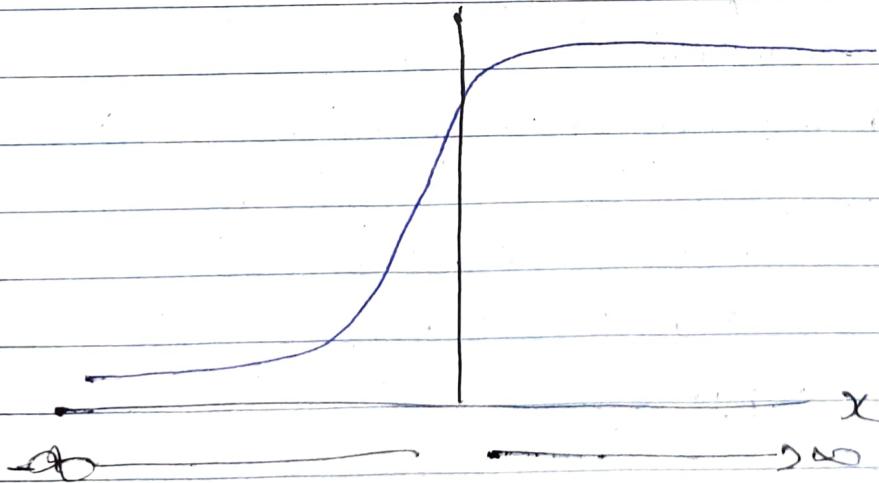
## Architecture of neurons:-

- ① neural network is a network of neurons.
- ② Many Input  $\rightarrow$  one output
- ③ neurons have an excitability threshold  
Voltage  
If a neuron passes this threshold then a spike occurs  
otherwise there is no spike.  $\text{spike} = 1$  or  $\text{no spike} = 0$
- ④ Synapses  $\rightarrow$  linear weight

Feed forward calculation: How do you get a prediction using logistic regression



logistic function = sigmoid function



a sigmoid function is an S shape. It has a finite limit as  $x$  approaches  $\infty$  in a finite limit as  $x$  approaches  $-\infty$

There are number of possible sigmoid function but the two commonly used are hyperbolic tangent

$\tanh(z) \in (-1, 1)$   $y\text{-int} = 0$  this goes from -1 to 1 and its y intercept is zero.

The other function which is the one we will be using is commonly denoted by the  $\sigma$  symbol

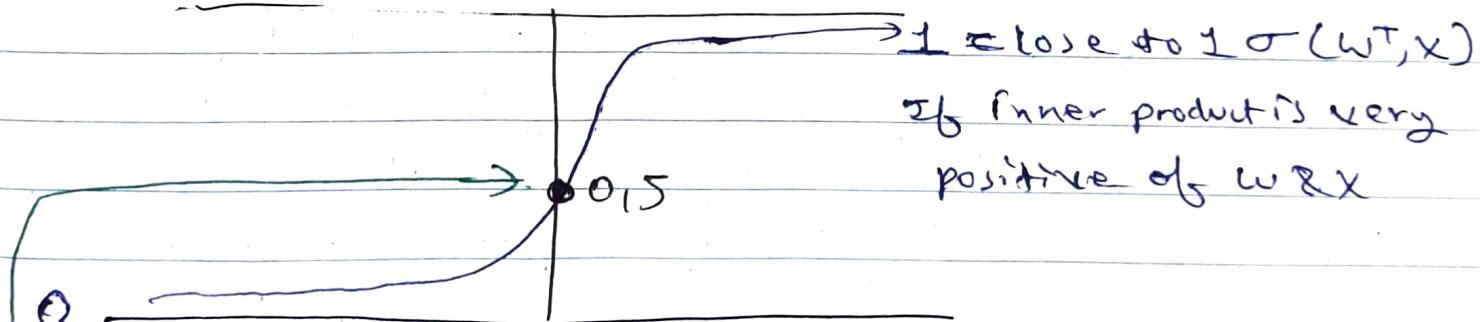
$$\sigma(z) = \frac{1}{1+e^{-z}} \in (0, 1) \quad y\text{-int} = 0.5$$

This goes from 0 to 1 & its y intercept is 0.5

So we can combine these to say the output of a logistic regression is  $\underline{\sigma(w^T x)}$



so what does this function tell us so if the inner product of  $w$  &  $x$  is very positive we will get a number very close to 1



We get a number very close to 0 the value of  $y$  equals 0.5 means the inner product is zero which means we are right in the boundary of between the two classes. So in other words in probability of belonging to the class is 50%.

Difference between logistic Regression and general Linear classifier we talked about earlier =

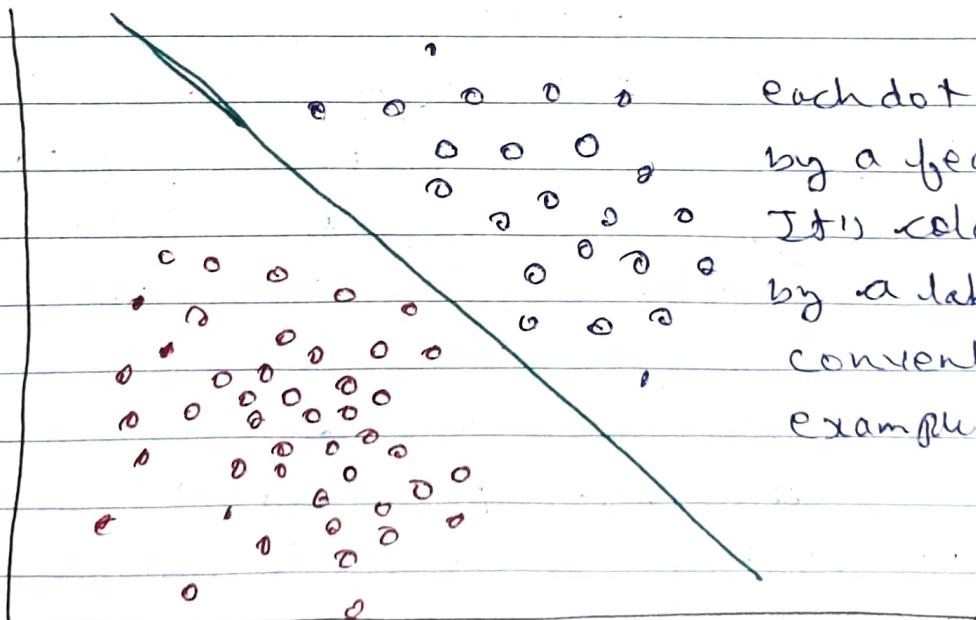
- (1) we have this logistic function at the end which gives us a number between 0 & 1 when we do our classification we can say anything above 0.5 will give us class 1 and anything else will give us class 0

What does the output mean?

$$\text{Output} = \sigma(w^T x) \in (0, 1)$$

The output of a logistic regression is ~~a~~ is a sigmoid from the sigmoid function, we are going to get a no between 0 & 1 in deep learning. This has a nice & intuitive interpretation.

Recollection of what we are trying to do in classification



- ① We have some red dots and we have blue dots & we have a line that separates them.
- ② Each dot is represented by a feature vector  $x$  and its color or in other words its label is represented by the letter  $y$ .
- ③ As per convention  $y$  is assigned a value 0 or 1
- ④ The output of logistic Regression is a number between 0 and 1 and we interpret this as
- ⑤  $P(y=1|x) = \sigma(w^T x) \in (0,1)$
- ⑥ The probability that  $y$  equals 1 given  $x$   
 Now because  $y$  can only take on two values 0 or 1 that means  $P(y=0|x)$  equals 0 given  $x$  is just  $1 - P(y=1|x)$  given  $x$  in other words the two probabilities must sum to 1
- $y$  can only have 2 values  
 $P(y=0|x) = 1 - P(y=1|x)$   
 → i.e.  
 $P(y=0|x) + P(y=1|x) = 1$   
 →  $P(y=1|x) = \sigma(w^T x) \in (0,1)$

- ⑦ So this gives us a handy way of making predictions if  $P(y=1|x) > P(y=0|x)$ :  
 predict class 1  
 else:  
 predict class 0

⑥ It's easy to show that an equivalent way of making this prediction is just to say if  $P(y)$  equals one given  $x$  is greater than 0.5 then predict class 1 otherwise predict class 0

if  $P(y=1|x) > [0.5]$ :

predict class 1

else:

predict class 0

→ Here we are using 0.5 because =

$$1 = P(y=1|x) + P(y=0|x)$$

↓  
class 1

↓  
class 2

$$0.5 + 0.5 = 1$$

So we see if we want to predict classes then

initially we were doing

if  $P(y=1|x) > [P(y=0|x)] \rightarrow$  we can replace this part with 0.5

predict class 1

else:

predict class 0

and our statement becomes as below

if  $P(y=1|x) > [0.5] \rightarrow$  0.5 here is a threshold

predict class 1

else

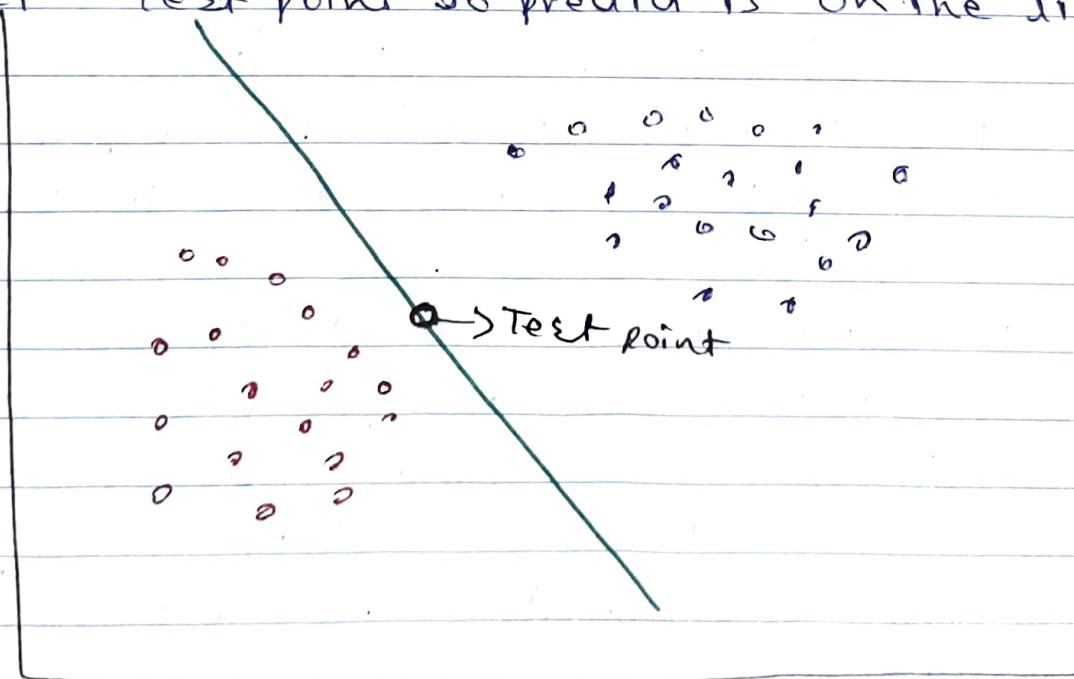
predict class 0

- ⑦ So using 0.5 as a threshold it's easy to make a prediction without using an if statement. Instead we can round the output of logistic regression
- ⑧ Rounding this number will give us either 0 or 1 which is exactly what we want. So now you can see the advantage of using 0 OR 1 as our labels.
- ⑨ Note:- Although it is possible to use numbers other than 0.5 as a threshold that usually doesn't happen too often in deep learning

⑩ Predict round( $P(y=1|x)$ ) → gives us either 0 OR 1

Another reason why ~~the~~ interpreting the output of the logistic regression as a probability makes sense →

case I:- Test point to predict is on the line

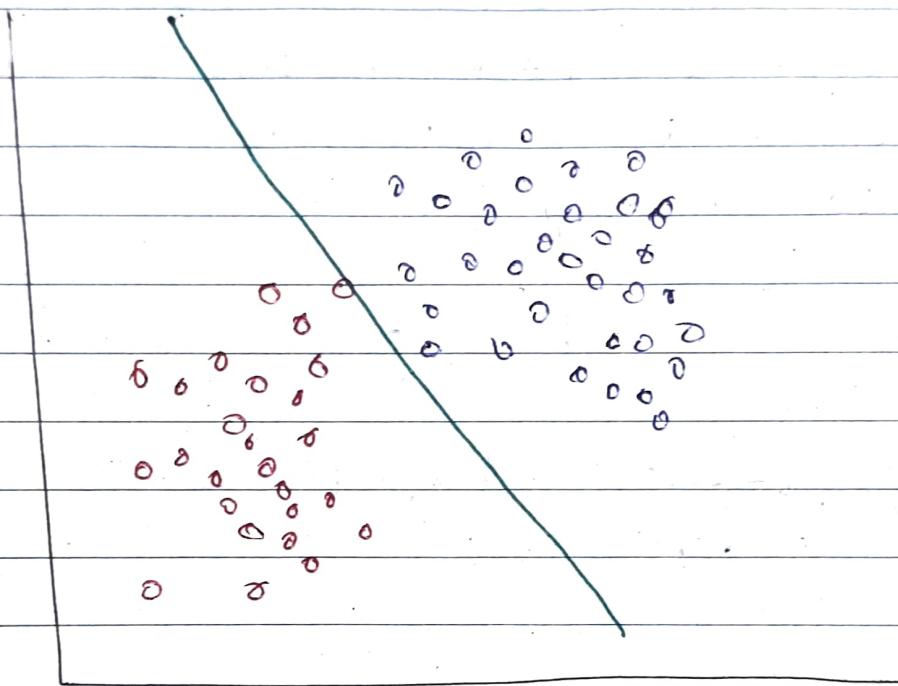


① If a point lies directly on the line the probability is exactly 0.5

② well if the point lies directly on the line we can't

really say either way whether this point should be assigned RED OR BLUE class

Case II) What happens when a point (Test point) gets further and further away from the line.



- ①  $|W^T \cdot X|$  gets larger as we move away from the line i.e. the magnitude of  $W$  transpose  $X$  gets larger and larger. In other words the arguments into the sigmoid approaches 0 or 1.
- ② The probability that  $y$  equals one given  $x$  approaches 1 or 0 depending on which way you go i.e.  $p(y=1|x)$  gets closer to 1 or 0.
- ③ The further we know that right in the middle it can means it can equally be red or blue. Thus it also makes sense that the further you go away from the line the stronger your belief gets that the target point is either Red or Blue.

## Quiz - Absence of Sigmoid

- ① To get an output predictions from a logistic regression we do the dot product between our weights & input features  $x$  and pass that to the sigmoid function

$$y = \sigma(w^T x) \rightarrow \text{dot product between } w^T x$$

$\sigma = \text{sigmoid function}$  This is used in forward function to make predictions in logistic regression

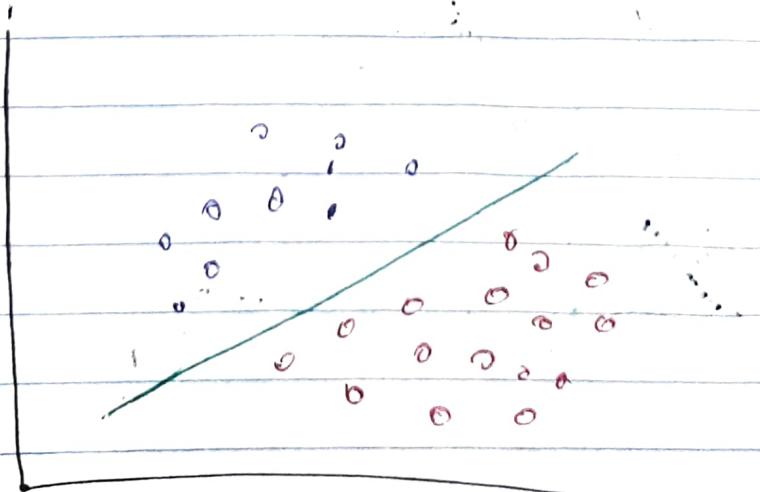
- ② The question is what kind of model do we get if we don't use the sigmoid?

$$y = w^T x$$

The answer is:- we get a model that looks like linear regression. ~~It is the only thing that makes this like linear regression is its functional form. It's a linear combination of input features weighted by the weight~~

- ③ Logistic regression makes the assumption that your data can be separated by a line / plane

$$y = \sigma(w^T x)$$



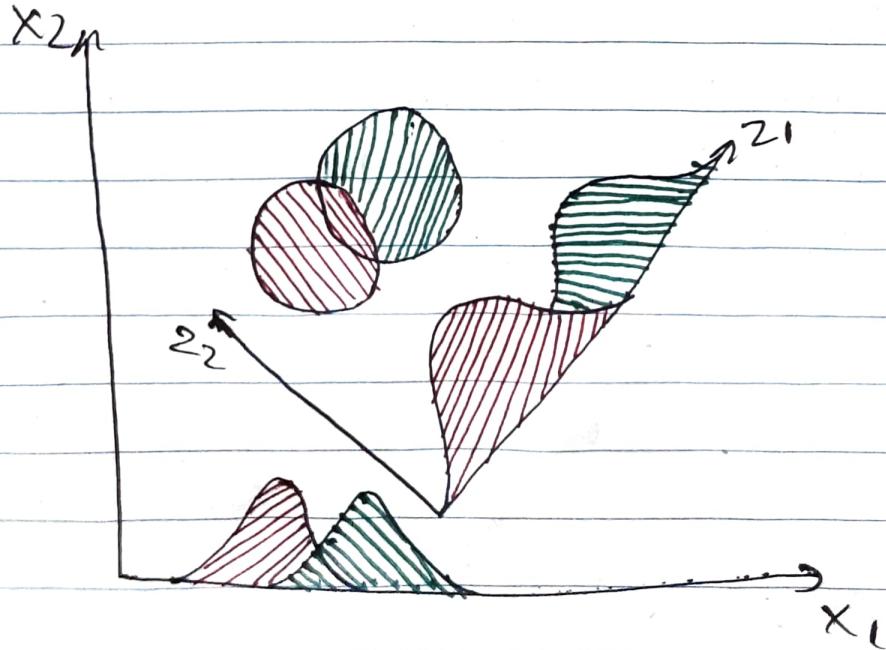
④ That's what we mean by modeling. we make some assumptions about the mathematical form and all the work we do thereafter follows from these assumptions.

⑤ we have these weights we have not found them OR choose them we just randomly select the weights using random function of python. we have to find weights using probability

### Solving for the optimal weights

- ① In machine learning finding weights  $w$  for our model is called "training"
- ② Training or learning is finding the parameters of your model that make it accurate

### Special case of logistic Regression:-



① Here we are going to look at a special case of logistic regression one where we can find a closed form solution to the problem

② In general we cannot and we will have to use gradient descent but the solution here could be useful if your data is distributed in a specific way.

③ The problem is this:-

→ Your data is from 2 classes & they are both gaussian distributed

→ They have the same same covariance but different means

→ For this example you should be familiar with multivariate Gaussian distribution

~~Bayes~~' Bayes' Rule:-

④  $p(Y|X) = p(X|Y)p(Y)/p(X)$ , This just says that the <sup>①</sup> posterior or <sup>②</sup>  $p$  of <sup>③</sup>  $y$  given <sup>④</sup>  $x$  is equal to the likelihood which is <sup>①</sup>  $p$  of <sup>②</sup>  $x$  given <sup>③</sup>  $y$  times the prior or <sup>①</sup>  $p$  of <sup>②</sup>  $y$  divided by <sup>③</sup>  $p$  of <sup>④</sup>  $x$  <sup>⑤</sup> gaussian

⑤ The likelihood part is the gaussian that we just talked about

⑥ We can calculate gaussian by taking the data from each class and getting their means and covariances the prior can just be the maximum likelihood estimate

$p(X|Y)$  is the gaussian - we calculate it over all the data that belongs to class  $Y$ .

⑦  $p(Y)$  is just the frequency estimate of  $Y \rightarrow$

example  $= p(Y=1) = 77$  times class 1 appeared / total

④ The next thing we can do is take Bayes rule for the top class and expand the bottom portion so it looks like the top portion

$$p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)} = \frac{p(x|y=1)p(y=1)}{p(x|y=1)p(y=1) + p(x|y=0)p(y=0)}$$

Next we divide the top & bottom by what's on top

$$p(y=1|x) = \frac{1}{1 + \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)}}$$

$$\textcircled{5} \quad p(y=1|x) = \frac{1}{1 + \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)}} \\ = \frac{1}{1 + \exp(-\{w^T x + b\})}$$

sigmoid used  
in program

$$-(w^T x + b) = \ln \left( \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)} \right) \rightarrow \log \text{ of this ratio from before}$$

$$\textcircled{6} \quad \ln p(y=1) = \alpha \quad \text{and} \quad p(y=0) = \alpha$$

$$\ln \left( \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)} \right) = \ln p(x|y=0) + \ln \alpha - \ln p(x|y=1) - \ln(1-\alpha)$$

$$= \ln \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} e^{-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0)} - \ln \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} e^{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)} + \ln \frac{\alpha}{1-\alpha}$$

Square root cancels out:-

$$= -\frac{1}{2} (x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x + \mu_0^T \Sigma^{-1} \mu_0) + \frac{1}{2} (x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1) + \ln \frac{\alpha}{1-\alpha}$$

The quadratic term

$$= \mu_0^T \Sigma^{-1} x - \mu_1^T \Sigma^{-1} x - (\mu_0^T + \mu_1^T) \Sigma^{-1} x = w^T x$$

$$= [(\mu_0^T - \mu_1^T) \Sigma^{-1} x] \\ = -(w^T x + b) \rightarrow \text{It depends on } x$$

If you split up the terms

$$\rightarrow w^T = (\mu_0^T - \mu_1^T)$$

$$b = \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0$$

⑦ During coding centered at (-2, -2)  
The variance is same for each dimension  
 $\mu_0 = (-2, -2)^T$ ,  
 $w = (4, 4)$ ,  $b = 0$   
each class)  
so when you calculate

The quadratic term cancels out

$$= \mu_0^T \Sigma^{-1} x - \mu_1^T \Sigma^{-1} x - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \ln \frac{\alpha}{1-\alpha}$$

$$= [(\mu_0^T - \mu_1^T) \Sigma^{-1} x] - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \ln \frac{\alpha}{1-\alpha}$$

$$= -(w^T x + b) \rightarrow \text{It does not depend on } x$$

~~w<sup>T</sup>x~~ It depends on x sigmoid used in program  $- \left( \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \ln \frac{\alpha}{1-\alpha} \right)$

If you split up the terms w & b then  $-b$

$$\rightarrow w^T = (\mu_1^T - \mu_0^T) \Sigma^{-1}$$

$$b = \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \ln \frac{\alpha}{1-\alpha}$$

② During coding we will use two gaussians one centered at  $(-2, -2)$  & one centered at  $(+2, +2)$

The variance of each dimension will be 8

each dimension will be independent

$$\mu_0 = (-2, -2)^T, \mu_1 = (+2, +2)^T, \sigma^2 = 1$$

$$w = (4, 4), b = 0 \text{ (assume equal samples of each class)}$$

so when you off diagonals covariance will be 0

③

$$\Sigma(x - \mu_1)^T + \ln \left( \frac{\alpha}{1-\alpha} \right)$$

Variables generally used in Machine learning & Deep learning:

- ①  $N$  = Number of samples we've collected
- ②  $D$  = Number of dimensions OR Features per sample
  - ↳ measure height, weight & girth to try to predict body fat %.
  - $D=3$  (height, weight & girth)
- ③  $X$  = matrix of data. It's an  $N \times D$  matrix
  - Sample  $\xrightarrow{\text{row}}$  column  $\xrightarrow{\text{Value of}}$
  - $\downarrow$  feature in each sample
- ④ Sometimes we say that we train model on inputs  $X$  and targets  $Y$
- ⑤ If we consider all  $Y$ s at the same time, then it is an  $N \times 1$  matrix (for binary classification, an  $N \times 2$  matrix of 0s & 1s)
- ⑥ Sometimes we use  $T$  for target instead
  - ↳ i.e. The cross-entropy error function! -  
 $y_{\text{test}} = \text{Targets}$   
 $y_{\text{train}} = \text{Outputs}$   
$$[- \log(Y) + (1 - Y) \log(1 - Y)] \quad \text{given } y_{\text{train}} = p_{\text{of } Y}$$
- ⑦ Now  $Y$  is used for something else; - output of logistic regression
- ⑧ Before the output was more precise  $P(Y=1 | x)$
- ⑨ Cost function = error function = objective function
- ⑩ We want to minimize the cost function or error
  - easy to see semantically
- ⑪ Objective function can mean something to maximize OR minimize
- ⑫ We will first maximize likelihood -  $P(\text{data} | \text{model})$

i.e If  $A > B$ , then  $\log(A) > \log(B)$

⑬ Sometimes we use letter  $J$

↳ maybe in form of  $J = L \rightarrow$  maximize it  
↳ OR  $J = L \rightarrow$  minimize it.

Cross Entropy error function:-

Linear Regression - squared error

$$J = \sum_n (t_n - y_n)^2$$

Assumes Gaussian error distributed error, because  $\log(\text{Gaussian}) = \text{squared function}$

Logistic Regression error can't be Gaussian distributed because:-

◦ Target is only 0/1

◦ Output is only a number between 0-1

we want: 0 if correct; > 0 if not correct,  
more wrong == bigger cost

Cross entropy error function is used in logistic.

Regression instead of squared function

Cross entropy error function  $\rightarrow$

$$J = - \left[ t \log(y) + (1-t) \log(1-y) \right] \quad (t = \text{target}, y = \text{output of logistic})$$

If  $t=1$  only first term matters, if  $t=0$  only second term matters

$\log(y) \rightarrow$  number between 0 and -infinity

Note:- we want to calculate the total error to ~~stop~~

optimize over all the training data simultaneously

$$J = - \sum_{n=1}^N t_n \log(y_n) + (1-t_n) \log(1-y_n)$$

Scalar random variable  $\boxed{x} \sim N(\mu_x, \sigma_x)$  multivariate Normal distribution

Vector of Random Variable  $\boxed{x} \sim N(\mu, \Sigma)$  Normal distribution  
mean value

$$\underline{x} \left[ \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] \sim N \left[ \begin{array}{c} \mu_{x_1} \\ \mu_{x_2} \\ \vdots \\ \mu_{x_n} \end{array} \right], \quad \begin{matrix} \text{mean vector} \\ \mu \end{matrix} \quad \begin{matrix} \text{covariance matrix} \\ \Sigma \end{matrix}$$
$$\begin{matrix} \sigma_{x_1}^2 & \sigma_{x_1 x_2} & \dots & \sigma_{x_1 x_n} \\ \sigma_{x_2 x_1} & \sigma_{x_2}^2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{x_n x_1} & \sigma_{x_n x_2} & \dots & \sigma_{x_n}^2 \end{matrix}$$

bold or underline is gonna indicate a vector OR Matrix

Closed form solution

$$\hat{\theta} = (\underline{X}^T \underline{X})^{-1} \cdot \underline{X}^T \underline{y}$$

Inverse of a matrix is one of the most computationally expensive method in computing. So if you have a very big data inverse of a matrix might not even fit in your memory. Hence we don't use closed form solution when our data is very big.

How to calculate likelihood for the biased point:-

- ① So we have a coin which have a probability of heads equal to  $p$  so it's unknown and the probability of tails is of course  $1-p$  i.e.  $P(T) = 1-p$   
 $P(H) = p$

- ② We're going to do an experiment to help us determine  $p$  we flip a coin 10 times & we get

$$N = 10 \quad 7H \cdot 3T \rightarrow 3 \text{ Tails}$$

↑  
2 Heads

- ③ How would we write the total likelihood which is the probability of seeing the result that we saw.

$$L = p^7(1-p)^3$$

we can do this because each coin tosses independent  
so we can multiply each individual probability.

- ④ Now what we want to do here is we want to maximize  $L$  with respect to  $p$  so we choose a  $p$  that maximizes  $L$  and to do this we will use calculus.

- ⑤ Now most of these problems we take the log and maximize the log likelihood which is acceptable because the log function is monotonically increasing meaning that the point at which the likelihood is maximum the log likelihood is also maximum.

$$l = \log [p^7(1-p)^3] \rightarrow \text{log likelihood}$$
$$= \log p^7 + \log (1-p)^3$$

$$= 7 \log p + 3 \log(1-p)$$

$$\frac{\partial L}{\partial p} = \frac{7}{p} + \cancel{\frac{3}{1-p}} \frac{3}{1-p} (-1) = 0$$

$$\frac{7}{p} = \frac{3}{1-p}$$

$$\frac{1-p}{p} = \frac{3}{7}$$

$$\frac{1}{p} - 1 = \frac{3}{7}$$

$$\frac{1}{p} = \frac{10}{7}$$

$$p = \frac{7}{10} = P(H)$$

② So now apply a similar idea to logistic regression.  
we have a probability that  $y$  is equal to 1 so  
we can think of it as heads.

$$p(y=1) \rightarrow p(y=1|x) = \sigma(w^T x) = y$$

$$L = \prod_{n=1}^N [y_n^{s_n} (1-y_n)^{1-s_n}]$$

$\downarrow$  target=1     $\downarrow$  target=0

taking log likelihood

$$l = \sum_n s_n \log y_n + (1-s_n) \log(1-y_n)$$

we get the same form as a cross entropy error function

⑧ So In other words maximizing the log likelihood is same as minimizing the cross entropy error.

Updating the weights using gradient descent:-

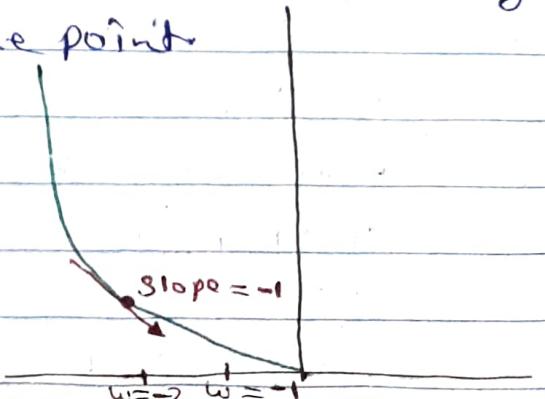
① Here we're gonna talk about how to update the weights of your logistic regression model to improve its error. So far all we know how to initialize the weights of logistic Regression Randomly.

② We also know that we can solve for the weights if we assume that the data is Gaussian distributed and has same covariance for each class, generally speaking. This isn't the case always. So we want something that will work always for all kinds of data.

③ We took the squared error calculated its derivative with respect to the ~~way it's~~ <sup>weights</sup> set it to zero and solve for the weights. With the cross entropy error we can still take the derivatives but we can't solve for the weights when we set it to zero.

④ So what can we do Instead:-

Let's say we have an error curve & we randomly initialized the weights to some point.



Suppose the derivative at this point is  $-1$  and suppose the point is at  $-2$  while the idea is we want to take small steps in the direction

Step size == learning rate (1 for this example)  
update becomes:-

$$w = -2$$

$$w = -2 - 1 * (-1) = -1$$

so our update becomes  $-1$  which bring brings us closer to the optimal point ( $w=0$ )

We just keep doing that until we get to the bottom  
Derivative at the bottom  $= 0$ , so once we get there  
the weights won't be updated

This method works for any objective function not just for logistic regression

Note:- Calculus is at the very heart of deep learning & a lot of machine learning & statistics

Q: How can we do gradient descent for logistic regression specifically.

Objective function:-  $J_n = \text{targets or } y_{\text{train}} \text{ or } y_{\text{test}}$   
 $y_n = \text{P}(\text{y given } x) \text{ or } p_{\text{train}} \text{ or } p_{\text{test}}$

$$\text{OR } J = - \sum_{n=1}^N J_n \log(y_n) + (1 - J_n) \log(1 - y_n) \rightarrow \begin{matrix} \text{logistic regression} \\ \text{train \& test error function} \end{matrix}$$

We have to find the derivative of this with respect to each weight.

Split into 3 derivatives

$$\frac{\partial J}{\partial w_i} = \sum_{n=1}^N \frac{\partial J}{\partial y_n} \frac{\partial y_n}{\partial a_n} \frac{\partial a_n}{\partial w_i}$$

Input of sigmoid OR  
Sometimes we call it activation

$$a_n = w^T \cdot z_n$$

Solution  $\Rightarrow$

$$\textcircled{a} J = -\sum_{n=1}^N t_n \log(y_n) + (1-t_n) \log(1-y_n)$$

The derivative of  $J$  with respect to  $y$  is just  $\frac{d}{dy} J$  of  $N$  over  $y$  of  $n$  minus one minus the event over

$$\boxed{\frac{\partial J}{\partial y_n} = -t_n \frac{1}{y_n} + (1-t_n) \frac{1}{1-y_n} (-1)}$$

~~a~~  $\textcircled{b}$  Next we need to find the derivative of sigmoid  $y_n = \sigma(a_n) = \frac{1}{1+e^{-a_n}}$   $\rightarrow$  sigmoid function

use the division rule to differentiate  $y_n$  wrt  $a_n$

$$\frac{\partial y_n}{\partial a_n} = \frac{\partial}{\partial a_n} \left( \frac{1}{1+e^{-a_n}} \right)$$

By division Rule:-

$$\frac{\partial y}{\partial x} = \frac{v \frac{du}{dx} - u \frac{dv}{dx}}{v^2}$$

v2

OR

$$\text{we can write } \frac{1}{1+e^{-a_n}} \text{ as } (1+e^{-a_n})^{-1}$$

then simply differentiate it wrt  $a_n$

$$-1 (1+e^{-a_n})^{-2} \frac{d}{da_n} (1+e^{-a_n})$$

$$\left| \begin{array}{l} \frac{d}{da_n} 1 = 0 \\ \frac{d}{da_n} e^{-a_n} = e^{-a_n} \cdot (-1) \end{array} \right.$$

$$\frac{-1}{(1+e^{-a_n})^2} (0 + e^{-a_n} \cdot (-1))$$

$$\boxed{\frac{\partial y_n}{\partial a_n} = \frac{-1}{(1+e^{-a_n})^2} (e^{-a_n}) (-1)}$$

$$\frac{\partial y_n}{\partial a_n} = \frac{-1}{(1+e^{-a_n})^2} (e^{-a_n})(-1)$$

$$\frac{\partial y_n}{\partial a_n} = \frac{e^{-a_n}}{(1+e^{-a_n})^2} = \frac{1}{(1+e^{-a_n})} \cdot \frac{e^{-a_n}}{(1+e^{-a_n})}$$

$y_n \quad (1-y_n)$

$$1 - \frac{1}{1+e^{-a_n}}$$

$$\cancel{x} e^{-a_n} \cancel{x}$$

$$(1-y_n) - \frac{e^{-a_n}}{1+e^{-a_n}}$$

So we can say that

$$\boxed{\frac{\partial y_n}{\partial a_n} = y_n(1-y_n)}$$

(c) Take the derivative of a wrt w.

$$a_n = w^T \cdot x_n$$

$$a_n = w_0 x_{n0} + w_1 x_{n1} + w_2 x_{n2} + \dots$$

$$\boxed{\frac{\partial a_n}{\partial w_i} = x_{ni}}$$

Finally we can put all this together

$$\frac{\partial J}{\partial w_i} = \sum_{n=1}^N \frac{\partial J}{\partial y_n} \frac{\partial y_n}{\partial a_n} \frac{\partial a_n}{\partial w_i}$$

$$= \sum_{n=1}^N \frac{\partial J}{\partial y_n} y_n(1-y_n)x_{ni} - \frac{1-a_n}{1+y_n} y_n(1-y_n)x_{ni}$$

$$\frac{\partial J}{\partial w_i} = - \sum_{n=1}^N \frac{t_n y_n (1-y_n)^2 x_{ni}}{y_n} - \frac{1-t_n y_n (1-y_n) x_{ni}}{1-y_n}$$

taking  $(1-y_n)$  common

$$\frac{\partial J}{\partial w_i} = \frac{1-y_n}{y_n} \sum_{n=1}^N t_n y_n (1-y_n) x_{ni} - (1-t_n y_n) x_{ni}$$

$$\frac{\partial J}{\partial w_i} = - \sum_{n=1}^N t_n (1-y_n) x_{ni} - (1-t_n y_n) x_{ni}$$

taking  $x_{ni}$  common

$$\frac{\partial J}{\partial w_i} = - \sum_{n=1}^N \left[ t_n (1-y_n) - (1-t_n y_n) \right] x_{ni}$$

$$\frac{\partial J}{\partial w_i} = - \sum_{n=1}^N \left[ t_n - t_n y_n - y_n + t_n y_n \right] x_{ni}$$

$$\frac{\partial J}{\partial w_i} = - \sum_{n=1}^N [t_n - y_n] x_{ni}$$

$$\boxed{\frac{\partial J}{\partial w_i} = \sum_{n=1}^N [y_n - t_n] x_{ni}}$$

$y_n$  = prob given x  
by train or predict  
for targets (Y train or Y test)  
 $x_{ni} = (x_{train})_i$

since each of  $i^{th}$  component on the left side only depends on the  $i^{th}$  component on the right side we can write this in vector form.

$$\frac{\partial J}{\partial w} = \sum_{n=1}^N (y_n - t_n) x_n \text{ where } x \text{ is a vector}$$

you can vectorize gradient descent formula for finding weights

$$\frac{\partial J}{\partial w} = \sum_{n=1}^N (y_n - \hat{y}_n) x_n \quad \text{for further weight finding formula}$$

$$\frac{\partial J}{\partial w} = X^T \cdot (Y - T) \quad \text{weights formula}$$

$$a^T b = \sum_{n=1}^N a_n b_n$$

$X^T = (X_{\text{train}})^T$   
Transpose

$Y = y_{\text{train}}$   
 $P$  of  $y_{\text{given } x}$

$T = \text{targets}$   
 $X_{\text{train}}$

$$X = N \times D \quad (X \text{ matrix})$$

$$Y, T \text{ are } N \times 1 \quad (\text{matrices } Y, T)$$

$$w = D \times 1 \quad (w \text{ matrix})$$

$$X^T \cdot (Y - T) \quad \text{shape is } (D \times N) \cdot (N \times 1) \rightarrow (D \times 1)$$

which is the correct shape of  $w$

$N$  gets summed over

Explanation of above →

(a) A dot product is summation over an index

(b) we can achieve  $w$  as  $D \times 1$  matrix by writing

$X^T \cdot (Y - T)$  This is a  $D \times \boxed{N}$  matrix,  $\boxed{N}$  times  $\boxed{N \times 1}$  matrix

since  $N$  is the inner dimension so it goes away  
and you are left with the outer dimension

(c) what happens with the biased term in logistic regression it's easy to incorporate biased term into  $x$  itself just by adding a column of ones

(d) let's say you want to find the derivative of the biased term explicitly for that let's return to the original expression

$$\frac{\partial J}{\partial w_0} = \sum_{n=1}^N (y_n - \hat{y}_n) \boxed{x_{0n}} = \sum_{n=1}^N (y_n - \hat{y}_n) \quad \xrightarrow{\text{bias term}}$$

Here all the  $x_0$  are 1

## Summary of Section 3:-

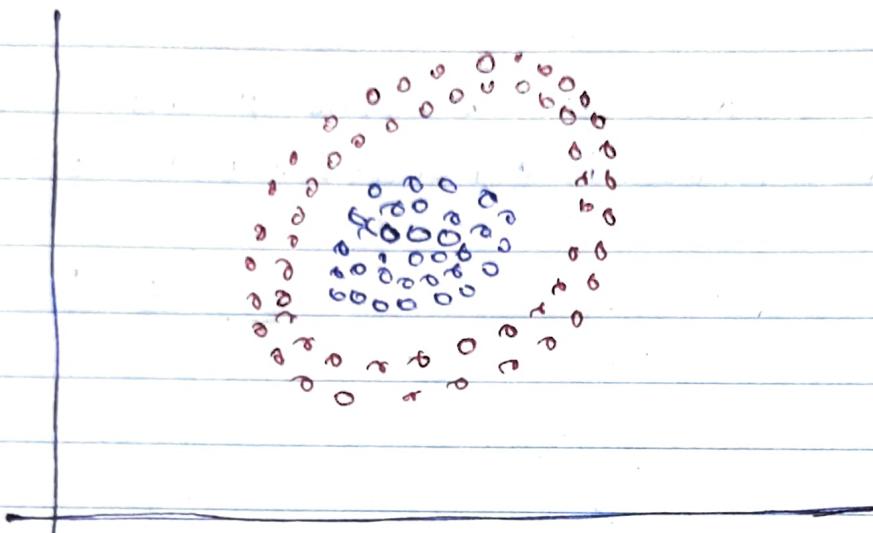
- ① How to find weights of a logistic regression model given some training data
- ② we started by making an assumption that our data is gaussian distributed with equal covariances , we ~~saw~~ saw that this leads to linear classifier exactly like logistic regression (use Bayes Rule)
- ③ General solution:- Use gradient descent on cross entropy (negative log likelihood) function and thought about how to maximize it or rather how to maximize the log likely hood. There is no closed form solution to this problem. So we need a more general optimization method gradient descent
- ④ we saw that the hard part of gradient descent is finding the gradient. To find the gradient we have to use the rules of calculus
- ⑤ We also used our of vectors & matrices to vectorize the solution. Since we know from our work in numpy that vectorized expression can be calculated faster than for loops

## Practical issues surrounding Logistic regression

- ① For any supervised ML model, we MUST know how to do these two (② things)
  - (a) Prediction (Inputs  $\rightarrow$  Outputs)
  - (b) Training (Finding  $w$ )
- ② Regularization & overfitting (what happens when your model performs too well on the training data?)
- ③ Overfitting happens when your model performs too well on the training data but training data doesn't matter really on the real world
- ④ where machine learning model becomes powerful when they predict things for us in the future
- ⑤ It doesn't matter how good our model is at predicting stock returns over the past year we want to be able to predict stock returns tomorrow
- ⑥ overfitting can happen for any number of reasons one of which is having irrelevant inputs in your model we will see how to bring the weights for irrelevant inputs down to zero
- ⑦ a related concept in Logistic Regression is that there is a subtle problem with our cost function the way that our cost function is constructed what can happen is that the ideal weights are actually infinite. In this section we are gonna see how it happens & how we can fix it

## Practical limitation of logistic Regression -

- ① One practical limitation of logistic regression is that the model is a straight line or a hyper plane in other words it cannot curve. This could be problematic in many situations.
- ② For example:-



In the above data we can see that the two classes are clearly separable

If a new test point is given then a human can easily classify whether the test point is Red or Blue

→ So why can't logistic Regression do this?

↳ It's because logistic Regression model is constrained to be straight

P

Q

R

S

T

U

V

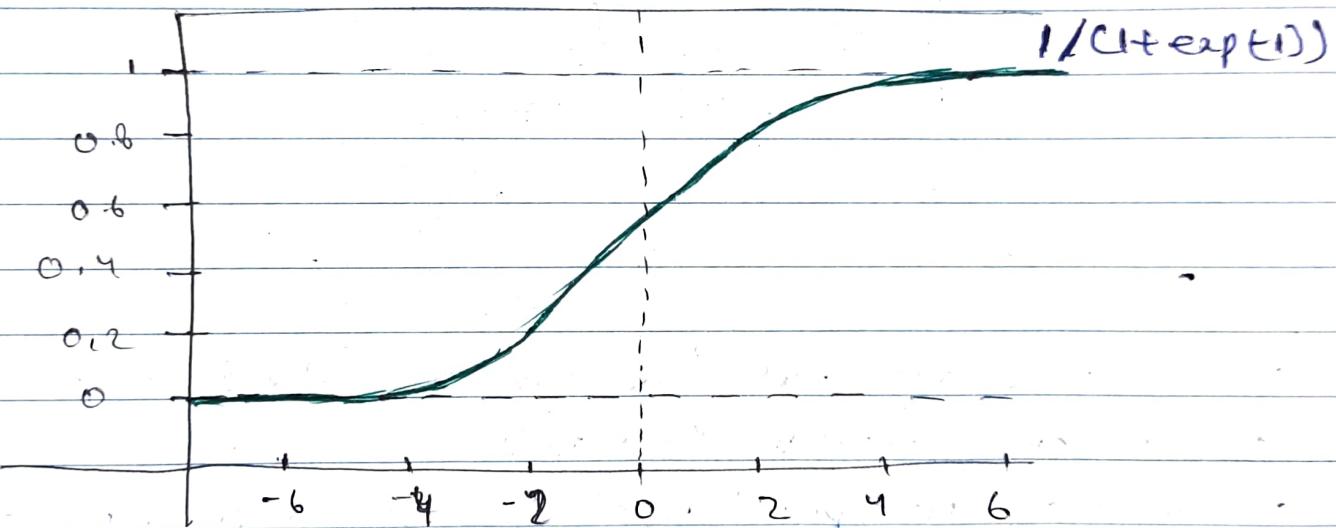
W

X

Y

How to interpret the weights learnt by logistic regression:-

- ①  $w_i$  is the amount of  $y$  will increase if  $x_i$  is increased by 1 and all other  $x_i$ 's remain constant. i.e. The weight  $w_i$  can be interpreted as the amount of that  $y$  will increase if  $x_i$  increases by 1 and all other  $x_i$ 's remain constant.
- ② In the binary scenario the output prediction from a logistic regression model is going to be either 0 or 1



- ③ So the weight is either going to either bring the output closer to 1 or 0:
- ④ A bigger weight means bigger effect. If weight  $w_i$  is big and positive than a small then a small increase in  $x_i$  assuming  $x_i$  is also positive will push the output closer to 1 too.
- ⑤ If  $w_i$  is big and negative than a small increase in  $x_i$  assuming  $x_i$  is positive will push the output closer to zero

⑥ In other words larger magnitude means larger impacts on the output.

⑦ If we see the interpretation of the output of the logistic model

$$P(y=1|x) = \sigma(w^T x) = \frac{1}{1 + \exp(-w^T x)}$$

Therefore if we subtract 2 we get

$$P(y=0|x) = 1 - \sigma(w^T x) = 1 - \frac{1}{1 + \exp(-w^T x)}$$

$$= \frac{1 + \exp(-w^T x) - 1}{1 + \exp(-w^T x)} = \frac{\exp(-w^T x)}{1 + \exp(-w^T x)}$$

⑧ The next step is to calculate the odds.

Odds =  $\frac{P(\text{Event})}{P(\text{not Event})}$ . Probability that event happens / Probability that event doesn't happen

$$\text{In our case} = \frac{P(y=1|x)}{P(y=0|x)} = \text{odds}$$

⑨ If we write this down in terms of  $w^T x$  we get exponential of  $w^T x$

$$\frac{P(y=1|x)}{P(y=0|x)} = \frac{1}{1 + \exp(-w^T x)} \times \frac{\exp(-w^T x)}{1 + \exp(-w^T x)}$$

$$= \frac{1}{1 + \exp(-w^T x)} \times \frac{1 + \exp(-w^T x)}{\exp(-w^T x)} = \frac{1}{\exp(-w^T x)}$$
$$\approx \exp(w^T x)$$

$$\boxed{\frac{P(y=1|x)}{P(y=0|x)} \approx \exp(w^T x)}$$

(10) If we take log of both sides :-

$$\log\left(\frac{P(y=1|x)}{P(y=0|x)}\right) = \underbrace{w^T x}_T + w_0 + \text{two terms}$$

It's exactly same as linear regression.

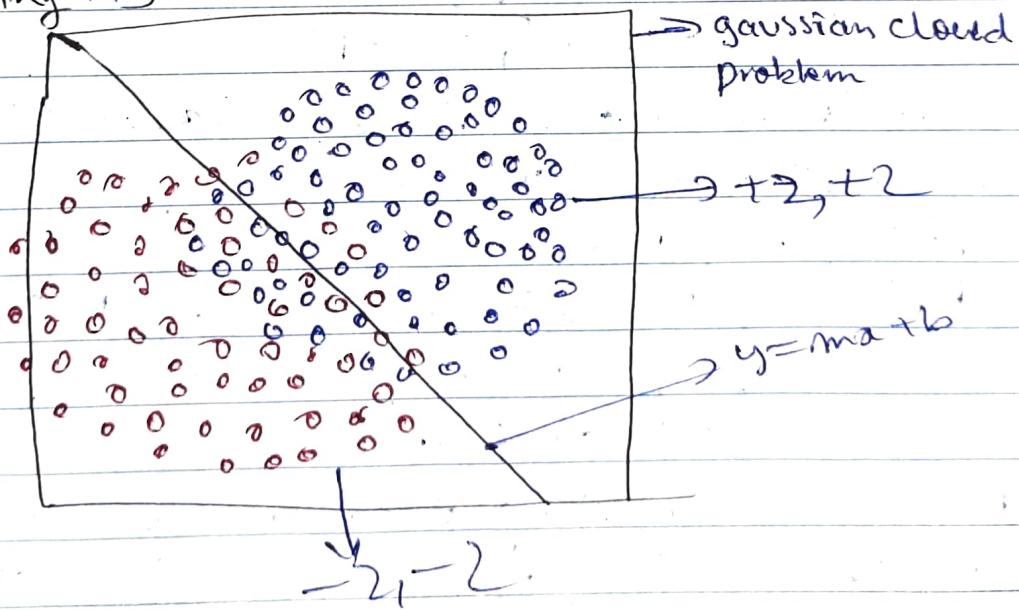
Since on the left we took the log of the odds. This quantity is appropriately known as log odds. Therefore we are doing linear regression on the log odds.

(11) This means we can interpret the weights like in linear regression.

The weight  $w_i$  can be interpreted as the amount that a log odds will increase. If  $x_i$  increases by 1 and all other  $x_j$ 's remain constant this interpretation is something that shows.

Regularization:- It is a generalization method

(A) Overfitting  $\rightarrow$



① Again we are gonna look at Gaussian cloud problem with 1 gaussian centered at  $+2, +2$  & one gaussian centered at  $-2, -2$ . This the same problem that we found out that exact Bayesine solution. It's easy to visualize because it's just two circles & you draw a line to separate them.

② The Bayesine solution that we found was  $w = [0, 1]$  where  $\omega$  is the bias

③ So what this line corresponds to  $= y = mx + b$  is the representation of line

④  $y = mx + b$

$y$  coordinate of the  $xy$  plane & not  $y$  as the output of Logistic Regression

⑤ So we can say  $0 + 1x + 1y = 0$

This gives us  $y = -x$

so that means we get a slope of  $-1$  & y-intercept  $0$

⑥ Why is the Bayesine solution  $(0, 0)$ ? why not  $1, 1$  or  $10, 10$ ?  
In fact these all represent the same line. This is why we need Regularization here

⑦ Consider the objective/error function so now this  $y$  is the output of logistic Regression.

→ Let's take a test point  $(x_1, x_2) = (1, 1)$

→ We know this should be in the same class as all the points surrounding  $(2, 2)$  so the class should be  $1$

→ So if we plug this into logistic Regression model

with the current setting of w

$$[\sigma(0+4+4) = \sigma(8) = 0.99966]$$

→ so with sigmoid of eight we get this objective

which is approximately ~~bad~~ • 0.003

$J = \sigma(8)$ ,  $J = -0.0003$  (good)  
output Sigmoid      objective  
of  
logistic Regression

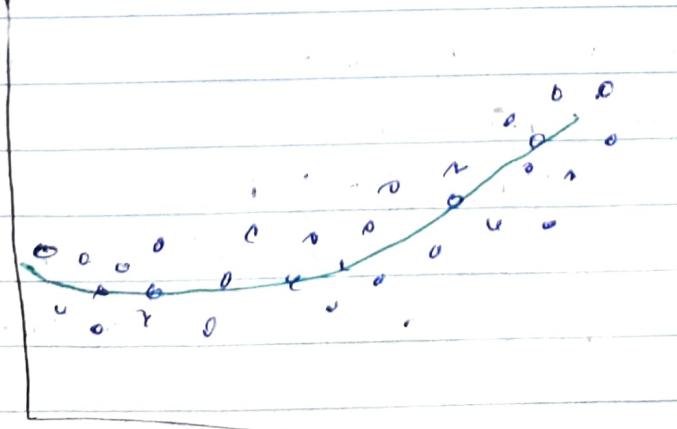
→ Now what if my weights were (0, 1, 1) 0 is the bias  
 $J$  would become = -0.1269 (Not good)

→ what if my weights were (0, 10, 10) 0 is the bias  
 $J$  would become =  $-2.06115 \times 10^{-9}$  OR  
=  $2.0 \times 10^{-9}$  (Not good)

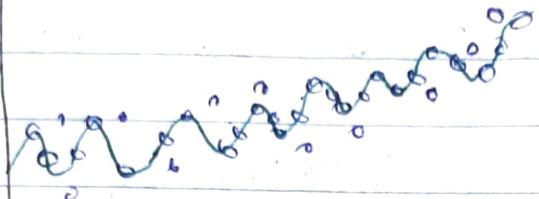
→ Moral of this is that the best weight of this model is (0,  $\infty$ ,  $\infty$ ) but of course computers can't represent  $\infty$ . So eventually we will just get an error

### B) Regularization:

① a lot of the times people explain regularization in terms of regression which I've already done



Good fit



## over fitting

- ② In fact your data completely fills up all possible inputs you shouldn't overfit even if your model is very complex
- ③ That's why we like having lots of data your model over fits when it has to guess what the output should be in a space that that it is never seen before
- ④ But if your training data is nicely spread out and covers all the different possibilities you can train to do well on that data
- ⑤ If your test data looks exactly like a training data and you do well on your training data then you'll also do well on your test data
- ⑥ Even though our data is nicely shaped and represents all the possible inputs the logistic Regression model still tries to give us a solution that the computer can't calculate (the solution to  $w=(0, \infty, \infty)$ ). The solution to this is Regularization
- ⑦ Regularization penalizes very large weight
- ⑧ So we have a cost function which is our cross entropy
 
$$J = -[t \log(y) + (1-t) \log(1-y)]$$

we just add another term to this cost so that it makes our cost grow larger If any weight grow bigger.

So the thing we add is the  $\frac{\lambda}{2}$  times the L2 norm of the weights

$$J_{\text{reg}} = J + \left(\frac{\lambda}{2}\right) \|W\|^2 = J + \left(\frac{\lambda}{2}\right) w^T w$$

what it does it encourages the weights to be close to zero. So now you will not get a weight like (zero, 10/10) because that will make your cost very large.

$\lambda$  is what we call a smoothing parameter

$\lambda$  balances the cross entropy cost & regularization penalty.

If  $\lambda$  is larger the weights will go closer to zero

If  $\lambda$  is smaller the weights will try to minimize the cross entropy

Typical values of  $\lambda = 0.1$  or  $1$  or something in between

How to choose a cross  $\lambda$  value.

observe the behaviour of the cost function and look at the results

There is no universal way to choose the best LAMBDA

How does  $\lambda$  effect our solution for  $w$ :-

- ① To perform gradient descent all we need to do is take the gradient of the objective function & move in that direction since addition doesn't effect the gradient calculation

$dJ_{reg}/dw$  instead of  $dJ/dw$

Add doesn't effect other terms in derivative

$$\text{reg. cost} = \left(\frac{\lambda}{2}\right) (w_0^2 + w_1^2 + w_2^2 + \dots)$$

You just take the gradient separately. So we just need to focus on the gradient of a regularization penalty.

$$\frac{d(\text{reg-cost})}{dw_i} = \lambda w_i \rightarrow \text{The derivative wrt any particular } w_i \text{ is just Lambda times } w_i$$

So now we can add this full vectorized form to our original gradient

$$\text{In vector form} \leftarrow \frac{d(\text{reg-cost})}{dw} = \lambda w$$

$$dJ_{reg}/dw = X^T (X - T) + \lambda w \rightarrow \text{in program we set } \lambda = 0.01$$

learning rate & npidot

L2 Regularization component.

Regularization in Probabilistic prospective:-

① We know that within the cross entropy what we are really trying to do is maximizing the likelihood.

$$J = -[t \log y + (1-t) \log(1-y)] \geq t \log y + (1-t) \log(1-y)$$

$$J = -[t \log y + (1-t) \log(1-y)] - \left(\frac{\lambda}{2}\right) \|w\|^2$$

$$\exp(J) = \frac{y^t (1-y)^{1-t}}{\exp(-\lambda \|w\|^2 / 2)} \downarrow \text{When I exponentiate a -ve squared term}$$

Binomial distribution of the likelihood

We call this gaussian distribution over

the weights prior i.e. It represents our

Prior beliefs about the weights

That's just a gaussian distribution

Impartial that they should be small & centered around zero

The variance of this gaussian distribution is  $\frac{1}{\lambda}$

$$w \sim N(0, \frac{1}{\lambda}) \rightarrow (\text{prior})$$

$t \sim \text{Bernoulli}(y) \rightarrow (\text{likelihood})$

posterior  $\propto$  likelihood  $\times$  prior

$\lambda$  is also called the precision which is the inverse variance of  $w$

This will keep occurring in Bayesian machine learning  
The basic rule is posterior  $\propto$  likelihood  $\times$  prior

Bayes Rule is -

$$P(w|x, y) \propto P(x, y|w) P(w)$$

$$\text{Consequence of: } P(B|A) = \frac{P(A|B) P(B)}{P(A)} = \frac{P(A|B)}{P(A)}$$

$$\text{bottom part is just } \sum_B P(A|B) = \sum_B P(A|B) P(B)$$

with regularization we maximize the posterior  
without regularization we maximize the likelihood  
This is called maximum a posteriori or MAP estimation.

## L1 Regularization

- ① In L2 Regularization even when we add a column of completely random noise to our data we can still improve our training score
- ② In general we want the dimensionality of  $\mathbf{z}$  to be much smaller than the number of samples  $N$ 

D < N (ideal)	D > N Not ideal
------------------	--------------------
- ③ In certain data sets the dimensionality is much greater than the number of samples. One way to visualize this via the  $\mathbf{X}$  matrix, we want this  $\mathbf{X}$  matrix to be skinny i.e.  $\boxed{N}$  number of samples features per sample do be large &  $D$  to be small. If we have opposite situation where  $\mathbf{X}$  is fat where  $N$  is small and these large that's not a good situation & we need to take steps to avoid potential problems.
- ④ what we would like to be able to do is select the small number of important features that produce that trend & remove all the features that are just noise. we call it sparsity since the end result will be that most of the weights will be zero and only a few of the weights will be non-zero
- ⑤ L1 Regularization which helps us achieve sparsity
- ⑥ L1 Regularization has a pretty simple setup similar to L2 we just add a penalty term to the original cost

⑦ As you might assume since we use a multiple of the L<sub>2</sub> norm of the weights in L<sub>2</sub> regularization.

$$L_{\text{RIDGE}} = \sum_{n=1}^N (\text{t}_n \log(y_n) + (1-\text{t}_n) \log(1-y_n)) + \lambda \|w\|_2^2$$

$$L_{\text{LASSO}} = \sum_{n=1}^N (\text{t}_n \log(y_n) + (1-\text{t}_n) \log(1-y_n)) + \lambda \|w\|_1$$

we will similarly use L<sub>1</sub> norm of the weights in L<sub>1</sub> regularization.

→ L<sub>2</sub> Regularization is called Ridge Regression

→ L<sub>1</sub> Regularization is called Lasso Regression

⑧ Let's look at the probability distribution involved

$$p(w) \propto \exp(-\frac{1}{2}\lambda\|w\|^2),$$

Gaussian distribution

⑨ So for L<sub>2</sub> regularization we had Gaussian likelihood and a Gaussian prior on w in this case we no longer have a gaussian prior on w

Gaussian likelihood = Bernoulli likelihood

⑩ So what distribution has a negative absolute value in the exponent. This would be the Laplace distribution.

⑪ So with L<sub>1</sub> regularization we are putting a Laplace distributed prior on the weights and we are solving for the posterior of w which applies

$$\text{⑫ } J = -\sum_{n=1}^N (\text{t}_n \log(y_n) + (1-\text{t}_n) \log(1-y_n) + \lambda \|w\|_1) \quad \text{Cost function}$$



महाराष्ट्र  
N.Rly.

gradient descent equation

finding weights with L1 Regularization

L1 Regularization constant OR L1 Penalty

$$\frac{\partial J}{\partial w} = X^T(Y - T) + \lambda \text{sign}(w), \quad (\|w\|_1)$$

For logistic regression we just calculate the gradient & move in the direction of 0.

For L1 regularisation you already know to take the gradient of the first part. The 2nd part is just lambda times the sign of w.

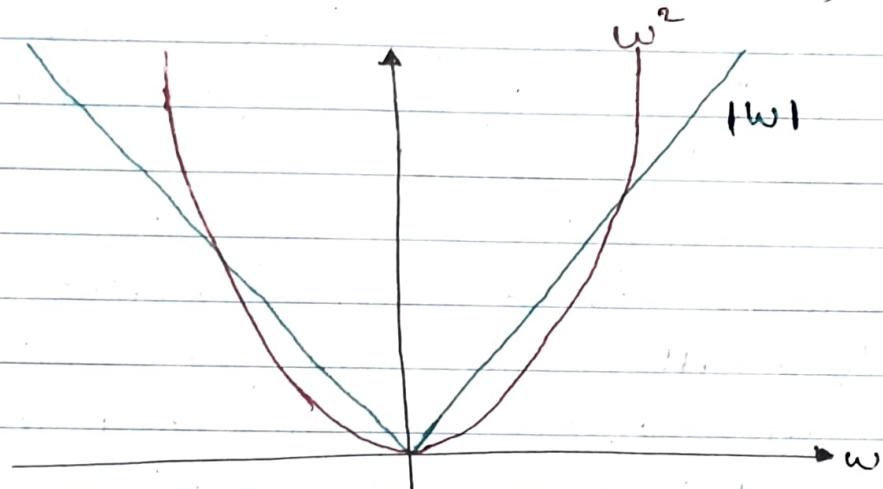
The sign function returns 1 if its argument is greater than zero negative 1 if its argument is less than zero & zero if its argument is = 0

### L1 vs L2 Regularization:-

- ① L1 Regularization encourages feature weights to be sparse ie only a few will be non-zero to regularisation
- ② L2 Regularization encourages feature weights to be small but will not encourage weights to go exactly zero
- ③ Both L1 & L2 Help improve generalization & test error because they prevent you from overfitting to the noise in the data L1 does this by choosing the most important features (the ones that have the most impact with respect to the output) you can imagine that the features with a smaller impact might in fact just help you predict that the noise in the training set. L2 does prevent you from overfitting by making the assertion that none of the weights are disproportionately large

④ Think of a one dimensional weight for L2 Regularization we get the -logprior -ve log prior of  $w$  to be quadratic. For L1 regularization we get the negative log prior of  $w$  to be the absolute function

⑤ The important thing to consider here is the derivative of this functions because gradient descent goes in the direction of the derivative, for a quadratic the closer you get to zero the smaller your derivative gets



Since it also approaches to zero.

⑥ So for L2 Regularization where your weight is already small, The further gradient descent won't change it that much.  $w \rightarrow 0$ , derivative  $\rightarrow 0$

⑦ For the absolute function the gradient is constant it always has a magnitude of 1 ( $+/-1$ ). (~~at  $w=0$~~ )

The derivative at 0 is undefined but we treat it as 0.

So for L1 regularization gradient descent will go towards zero at a constant rate. Then when it arrives at 0 it stays there forever. So

that is why L2 regularization encourages weights to be small while L1 regularization encourages

weights to be zero therefore encouraging sparsity.

Note:- It is possible to include both L<sub>1</sub> & L<sub>2</sub> regularization in your model. This also has a special name called elastic net. Here we will add L<sub>1</sub> & L<sub>2</sub> penalty to the cost function.

$$\rightarrow J_{\text{RIDGE}} = J + \lambda_2 \|w\|^2$$

$$J_{\text{LASSO}} = J + \lambda_1 |w|$$

$$J_{\text{ELASTICNET}} = J + \lambda_1 |w| + \lambda_2 \|w\|^2$$

Why we decided  $w = \frac{\text{np.random.randn}(D)}{\text{np.sqrt}(D)}$

$w = \frac{\text{random } D}{\sqrt{D}}$  where D is number of features per sample  
OR Dimensions

- ① We did  $w = \frac{\text{np.random.randn}(D)}{\text{np.sqrt}(D)}$  in L<sub>1</sub> regularization program.
- ② poor weight initialization can lead to poor convergence of your loss per iteration
- ③ sometimes your loss might even explode to  $\infty$ . If you've never encountered this then ~~likely~~ it is likely you just haven't practiced enough
- ④ Once you've seen enough datasets and done enough machine learning work you will encounter a scenario where weight initialization matters. That's not to say it ~~always~~ always matters and that's not to say there won't be some instances where these recommendation will appear to perform poorly

## (5) How it works:-

# weight initialization

$$w = np.random.randn(D) / np.sqrt(D),$$

→ we know that the randn function draws samples from the standard normal and that's usually denoted as  $N(0, 1)$

→ what happens when we multiply these samples by some number of  $C$ . Well you can show that if you multiply random samples from the standard normal by a number  $C$  they will ~~not~~ have mean zero and standard deviation  $C$

$$\text{randn}() * C == \text{Samples from } N(0, C^2)$$

→ In other words it's equivalent to sampling from the distribution ~~and a general zero~~  $C^2$

→ Infact this is just the opposite of standardization or normalization

→ When we want to standardize a set of samples we subtract their mean and divide by their standard deviation sigma

$$\text{To standardize } \Rightarrow z = (x - \mu) / \sigma$$

$z$  has mean 0, standard deviation 1

$$\text{To inverse transform } \Rightarrow x = z \sigma + \mu,$$

$x$  has mean  $\mu$ , standard deviation  $\sigma$

this makes it so that the standard samples have mean zero and standard deviation 1.

→ Thus we can perform the inverse operation as well  
If we multiply it by a ~~mean~~ number sigma and add a ~~number you can't~~ ~~mean~~ ~~a standard dev~~

a number  $\mu$  on a standardized sample the result will be a sample that has mean  $\mu$  and standard deviation sigma  
Inverse transform:  $x = \sigma + \mu / \sqrt{D}$  //  $x$  has mean  $\mu$ , s.d.  $\sigma$

⑥ Let's go back to our initialization of weights

$$w = np.random.randn(D) / np.sqrt(D)$$

→ If we divide by the square root of D, This is equivalent to multiplying by 1 over square root of D

$$\frac{1}{\sqrt{D}}$$

→ Therefore we are saying that we want the initial weight vector  $w$  to have mean zero and variance one over D.

→ So why do we want to be that

If you don't do this →

$$w = np.random.randn(D) / np.sqrt(D)$$

then our loss per iteration might explode and go to infinity.

→ The reason it's exploding is because those weights are too large and if they are too large that means their variance is too large. Which means we may be able to avoid this problem if we make the variance smaller

→ In fact in some code there are people who will just multiply by a constant like 0.01. This is also perfectly ok. It comes down to doing experiment and finding the optimal value for your own dataset. This is always the ultimate foolproof method.

P

Q

R

S

T

U

V

W

X

Y

Z

⑦ Okay so we know that we want weights with small variance but why exactly? Variance  
→ Let's go back to multiple linear regression where we have  $D$  input features and 1 output; since we have  $D$  input features our weight vector will be a vector of size  $D$ . As you know

→ As you know we like to normalize our data before passing it through any kind of model whether it's a linear regression or a neural network. Again it's all about experimentation.

$$y = w_1x_1 + w_2x_2 + \dots + w_Dx_D$$

$$y = w_1x_1 + w_2x_2 + \dots + w_Dx_D$$

→ If we normalize ~~out~~ our data then that means each input feature will have mean zero & variance 1

→ Suppose we want to know the variance of the output  $y$ . Why do we care about the variance of  $y$

→ Well in neural networks: After the first layer you simply have more layers of the same type. Therefore the output of the first layer will be the input of the second layer and so on. But since those will be the inputs to the next layer you want those to be normalized as well.

→ If it's just linear regression your target will probably be standardized so there's still a ~~reason~~ reason to care about the variance of  $y$

→ we know that the expression of  $y \Rightarrow$

$$y = w_1x_1 + w_2x_2 + \dots + w_Dx_D$$

So how do we calculate the variance of  $Y$

→ Since all the  $w_i$ 's &  $X_i$ 's are independent this expression simplifies to

$$\text{Var}(Y) = \underbrace{\text{Var}(w_1)}_{\text{Variance of } w_i} \underbrace{\text{Var}(X_1)}_{\text{Variance of } X_i} + \text{Var}(w_2) \text{Var}(X_2) + \dots + \text{Var}(w_D) \text{Var}(X_D)$$
$$= \text{Var}(w_1) + \text{Var}(w_2) + \dots + \text{Var}(w_D) / \text{since } \text{Var}(X_i) = 1$$
$$= D \text{Var}(w)$$

→ Variance of all the  $X_i$ 's = 1 [" $\text{Var}X_i = 1$ "] because we have standardized our data & we also know that variance of all the  $w_i$ 's will be the same.

at last we get that variance of  $Y = D * \underline{\text{Variance of } w}$

→ It means if we want our variance of  $Y$  to be  $1$  which we do then the variance of  $w$  must be  ~~$\frac{1}{D}$~~  or standard deviation of  $w = \frac{1}{\sqrt{D}}$

## NOTE:-

\* → This isn't the only way to initialize the weights

\* → You want the variance of the weights to be initialized to small values

\* → Other methods include He-normal, Glorot-normal etc.

To initialize the weights

Typically these are used in Deep neural networks & not linear models like Linear Regression

\* → It's normal to just multiply by a constant value like 0.01

Q

R

S

T

U

V

W

X

Y

Z

## Sentiment analyzer (theory):-

① What is sentiment analyzer?

→ Sentiment = how positive ~~or~~ or negative some text  
is

amazon reviews, Yelp reviews, hotel reviews etc.

② Our data is in form of XML file so you will need an XML parser

③ Outline of the sentiment analyzer:-

a) we are gonna use the reviews on the electronics section but this code should work on others as well

b) XML parser used here is Beautiful Soup

c) only look at the key "review\_text"

d) we will need 2 passes, one to determine vocabulary size and which index corresponds to which word, and one to create data vectors

e) after that we could just use SKLearn classifier as we did previously

f) But we will use logistic regression so we can interpret the weights.