

LOGISTIC REGRESSION

[GEOMETRIC INTUITION](#)

[SIGMOID FUNCTION SQUASHING](#)

[MATHEMATICAL FORMULATION OF OBJECTIVE FUNCTION](#)

[WEIGHT VECTOR](#)

[L2 REGULARIZATION OVERRFITTING AND UNDERFITTING](#)

[L1 REGULARIZATION AND SPARSITY](#)

[PROBABILISTIC DERIVATION OF LOGISTIC REGRESSION](#)

[LOSS MINIMIZATION INTERPRETATION](#)

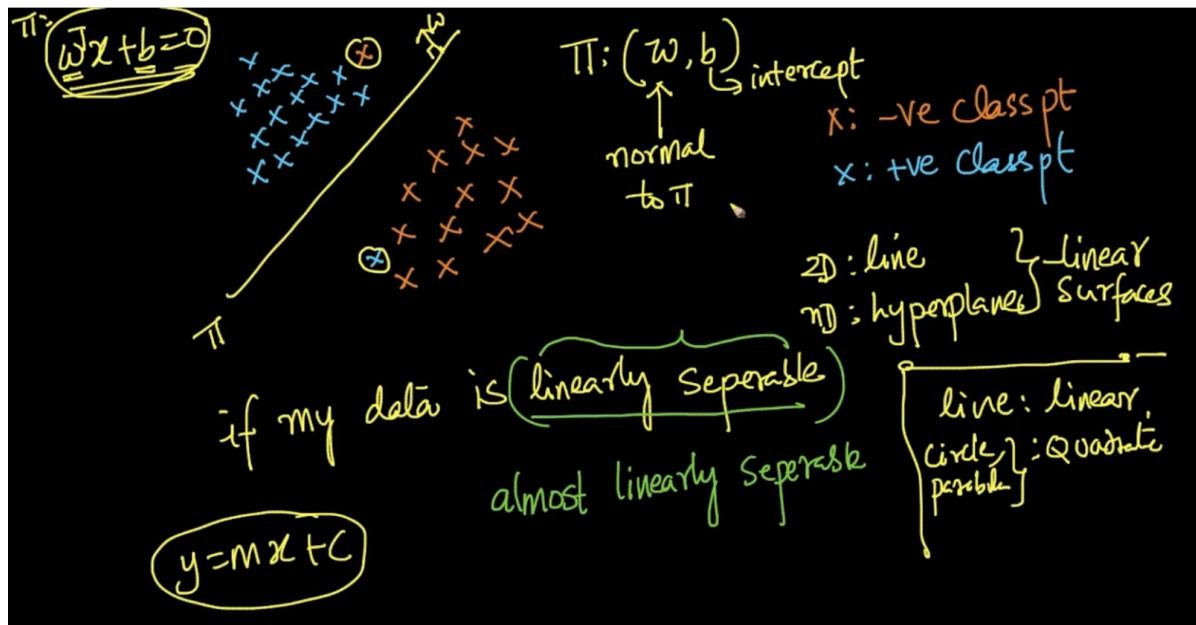
[HYPERPARAMETER SEARCH/OPTIMIZATION](#)

[COLUMN / FEATURE STANDARDIZATION](#)

[FEATURE IMPORTANCE AND FEATURE INTERPRETABILITY](#)

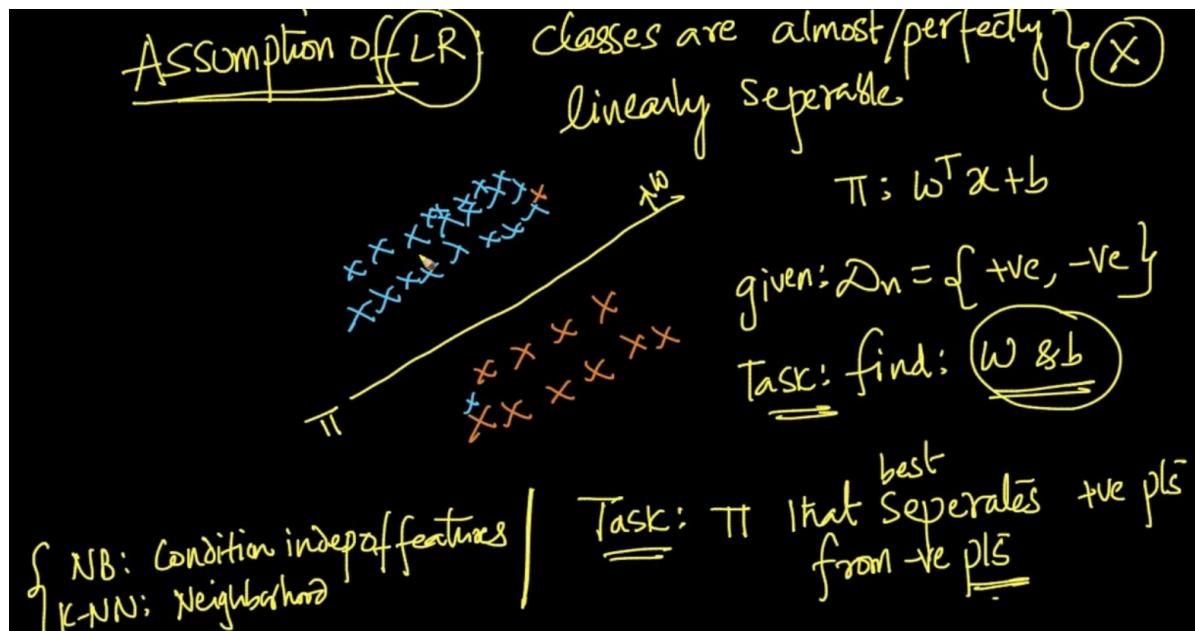
[COLLINEARITY OF FEATURES](#)

GEOMETRIC INTUITION

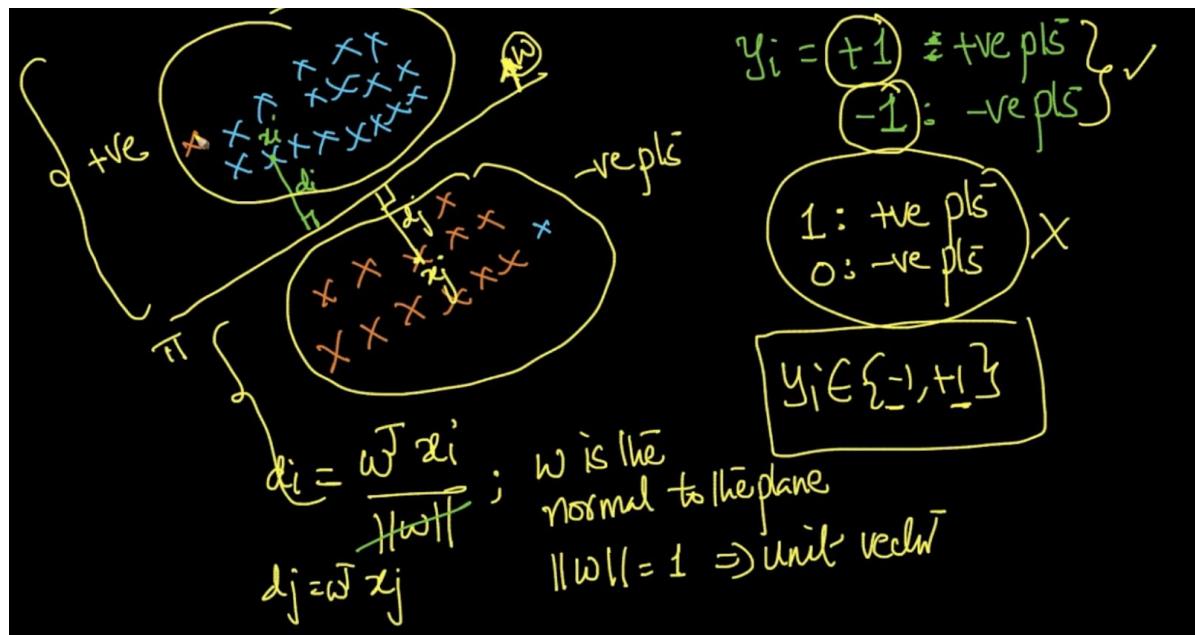


We need to find the hyperplane which makes the data linearly separable.

Equation of line/ hyperplane $\pi: w^T x + b = 0$



Assumption of Logistic Regression is that our data is perfectly linearly separable. So, our task is to find the hyper-plane π that separate both the +ve and -ve points.



We are calculating distance d_i of a point x_i which is $d_i = w^T x_i$. We've labeled all the blue points as +ve points and the orange ones are -ve points.

$$\left. \begin{array}{l} d_i = \omega^T x_i > 0 \\ d_j = \omega^T x_j < 0 \end{array} \right\} \text{classify} \\
 \left\{ \begin{array}{l} \text{if } \omega^T x_i > 0 \text{ then } y_i = +1 \\ \text{if } \omega^T x_i < 0 \text{ then } y_i = -1 \end{array} \right. \\
 \text{decision surface in LR} \\
 \underline{\text{; Line / plane}}$$

Point x_i lies in the +ve region and the direction of normal w is also at +ve points direction

So, $w^T x_i > 0$. Hence we built a classifier saying that if $w^T x_i > 0$ then $y_i = +1$ elif $w^T x_i < 0$ then $y_i = -1$. That's why we wanted an ideal plane π which is decision surface.

$$\text{Case 1: } \begin{array}{l} (+ve pt) \\ y_i * w^T x_i > 0 \\ \quad \quad \quad \boxed{y_i = +1} \end{array} \rightarrow \begin{array}{l} \text{if } y_i \neq +1 \rightarrow +ve pt \\ (\omega^T x_i > 0 \Rightarrow \text{classifier is saying it's the pt}) \\ \text{① is correctly classifying the pt} \end{array}$$

$$\text{Case 2: } \begin{array}{l} (-ve pt) \\ y_i = -1 : -ve pt \\ w^T x_i < 0 \Rightarrow \text{LR concluding that } x_i \text{ is a -ve pt} \\ \boxed{y_i | w^T x_i > 0} \end{array}$$

Case 1 : if $y_i = +1$ and $w^T x_i > 0$ then the classifier is saying x_i is a +ve point. So, if

$y_i * w^T x_i > 0$ and $y_i = +1$ then we can say that our normal w is correctly classifying the point

Case 2 : Same as case 1 except that case 2 is for -ve point . Here also $y_i * w^T x_i > 0$.

NOTE : Both the cases are of w correctly classifying data points and in both $y_i * w^T x_i > 0$
So, we want $y_i * w^T x_i$ greater than 0 to make sure that our classifier is working correctly

both +ve & -ve pt
 $y_i w^T x_i > 0 \Rightarrow$ the LR model is correctly classifying the pt x_i

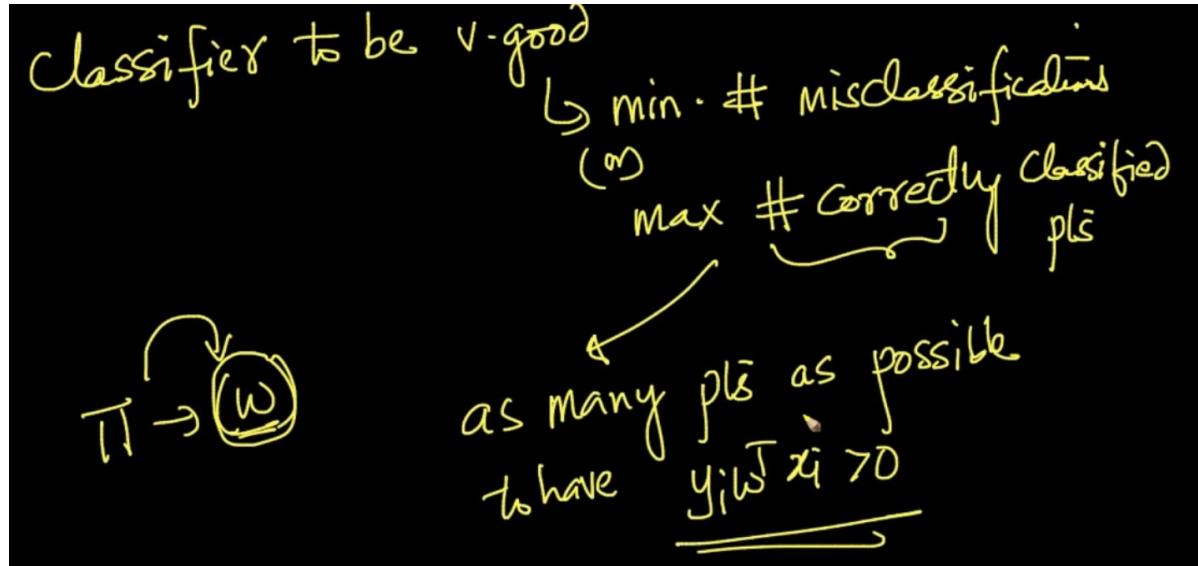
Case 3: $y_i = \pm 1$ (+ve pt)
 $w^T x_i \leq 0 \Rightarrow$ LR is saying x_i is -ve class
 $y_i w^T x_i < 0 \Rightarrow$ $y_i = +1$ misclassified $LR = -1$

Case 3 : Here $y_i = +ve$ point but $w^T x_i < 0$ so $y_i * w^T x_i < 0$ that means our Logistic Regression is misclassifying the data as LR : -1. Ex: Blue point in the orange points.

Case 4: $\begin{cases} y_i = -1 \Rightarrow -ve \text{ class} \\ w^T x_i > 0 \Rightarrow \text{LR is saying } x_i \text{ is +ve pt} \end{cases}$
misclassified $y_i w^T x_i < 0$

Case 4: Same as case 3. Here also $y_i * w^T x_i < 0$.

NOTE : Case 3 and Case 4 both the cases are of w incorrectly classifying data points and in both $y_i * w^T x_i < 0$. So, we want $y_i * w^T x_i < 0$ points to be less than 0 to make sure that our classifier is working correctly.



So in order for our classifier to be good we want as many points as possible to have

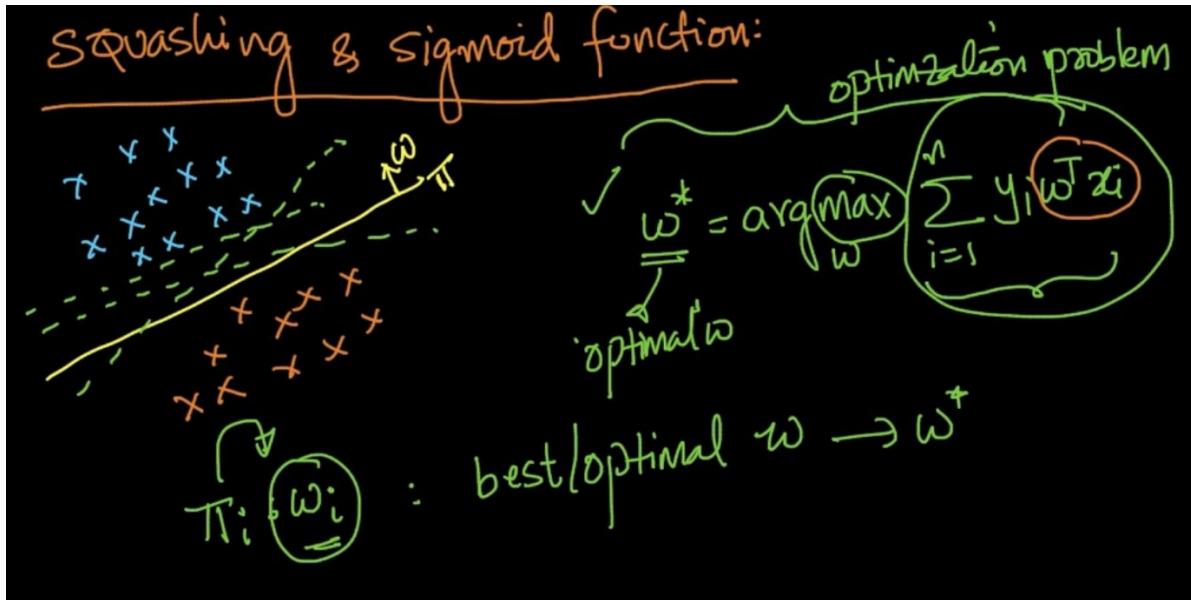
$$y_i * w^T x_i > 0$$

$$\begin{aligned}
 & \left\{ \max_{\underline{w}} \sum_{i=1}^n (y_i \underline{w}^T \underline{x}_i) \right. \\
 & \quad \text{variable} \quad \left. \begin{array}{l} n \text{ datapts} \\ \text{train } D_n \end{array} \right\} \\
 & \hookrightarrow \text{Change/Vary } \underline{w} \\
 & \left\{ \begin{array}{l} \text{optimal } \underline{w}^* \\ \underline{w}^* = \arg \max_{\underline{w}} \left(\sum_{i=1}^n y_i \underline{w}^T \underline{x}_i \right) \end{array} \right. \\
 & \quad \text{variable} \quad \left. \begin{array}{l} \text{Math.} \\ \text{optimization} \\ \text{problem} \end{array} \right\}
 \end{aligned}$$

In this our x_i and y_i are fixed so the only thing that can be varied in $y_i * w^T x_i$ is w^T .

So, our goal is to find the optimal $w \Rightarrow w^*$ which gives us maximum in summation of $y_i * w^T x_i$ as more $y_i * w^T x_i > 0$ means more correctly classified data. Our maximum of summation gives that. So finding this optimal w is the crux of this algorithm called Logistic Regression

SIGMOID FUNCTION SQUASHING



We need to find the optimal $w \rightarrow w^*$

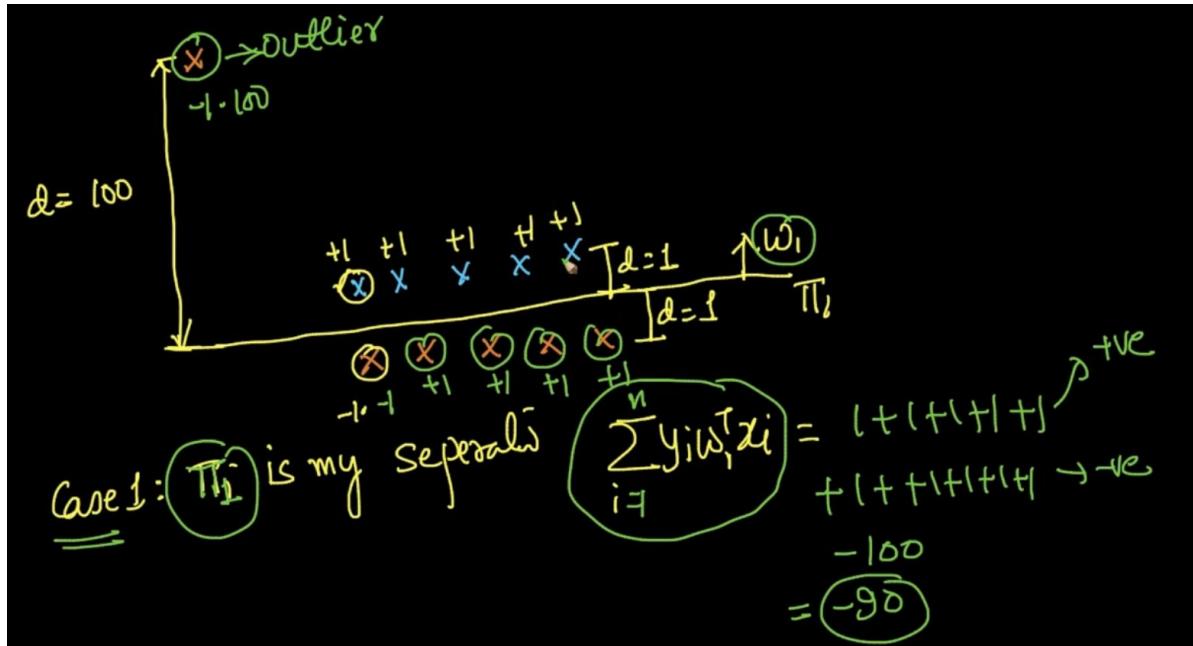
$$\arg \max_w \sum_{i=1}^n y_i w^T x_i$$

↑ signed distance
dist from x_i to π (w is a unit vector)

$y_i w^T x_i$: +ve $\Rightarrow \pi$ as defined by w correctly classifies x_i
 -ve \Rightarrow incorrectly classifies x_i

y gives us sign of the point +1 or -1. $w^T x_i$ is the distance of data point x_i from plane π

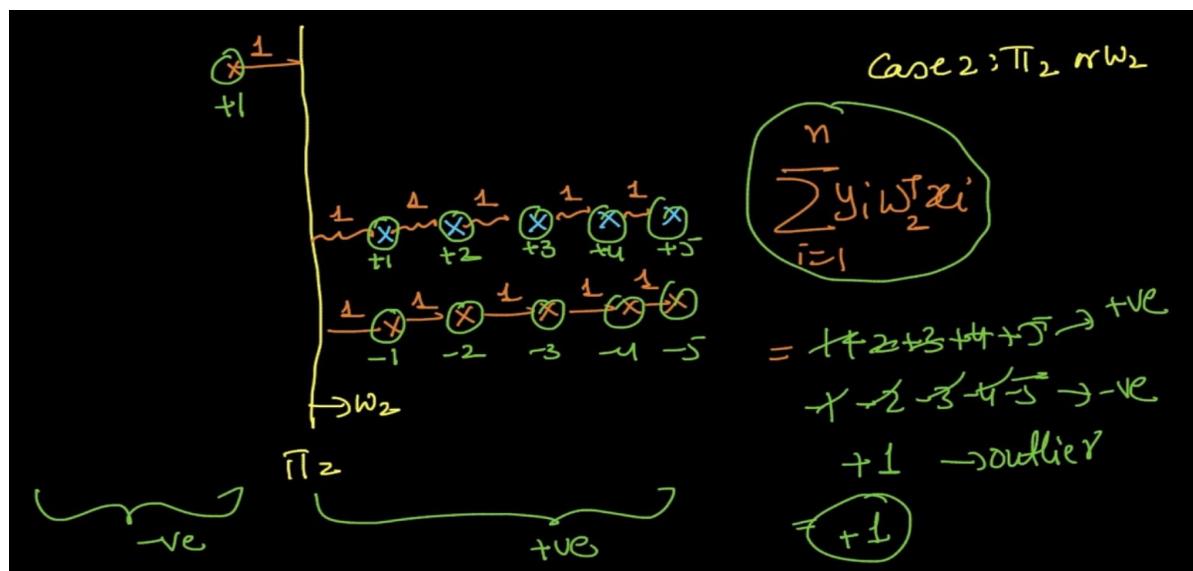
If $y_i w^T x_i$: +ve then our π defined by w correctly classifies x_i . So we want it to be maxed



Case 1: Here for +ve points(blue) distance $d \Rightarrow w^T x = 1$ and $y = +1$ since positive

Here for -ve points(orange) distance $d \Rightarrow w^T x = -1$ since opposite to direction of w and $y = -1$ since negative. For the outlier $d = +100$ since in the direction of w and -ve point

$$\text{Therefore } y = -1. \sum_{i=1}^n y_i w^T x_i = -90$$



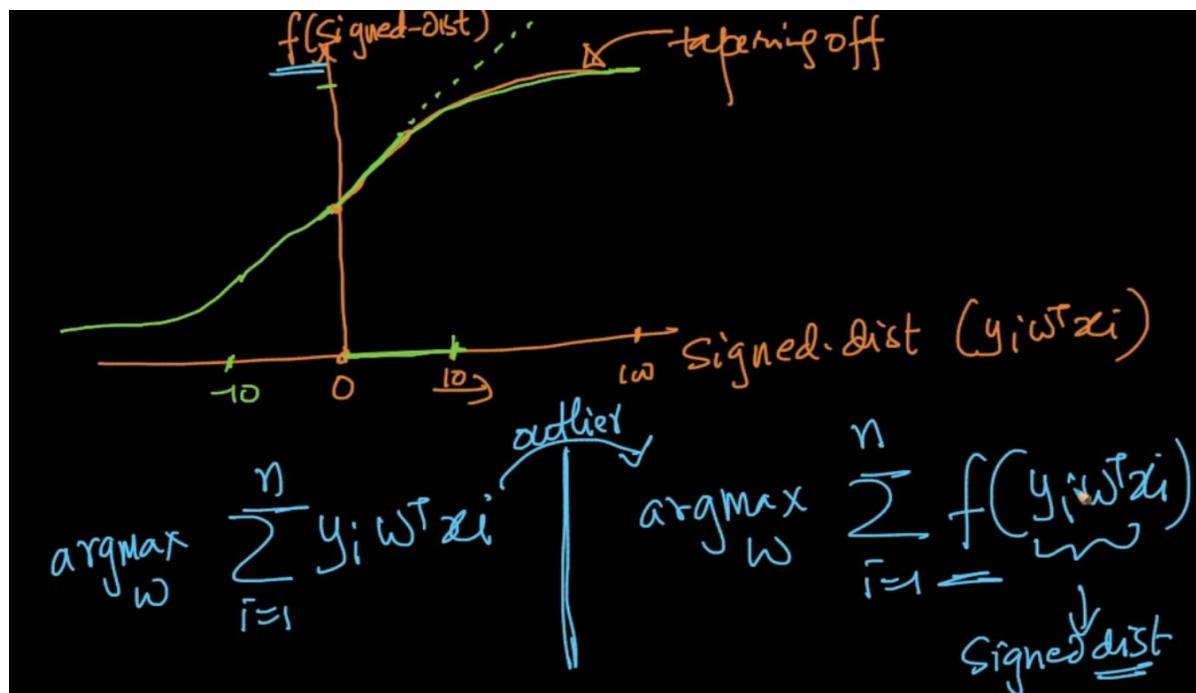
Case 2: Here the plane is separating the outlier and points. Here $\sum_{i=1}^n y_i w^T x_i = +1$

We want $\sum_{i=1}^n y_i w^T x_i$ to be maxed so our model will say Case 2 but we know intuitively that the plane in Case 2 is doing a shit job as it's misclassifying more data. It's happening because of our outlier. So our max $\sum_{i=1}^n y_i w^T x_i$ isn't prone to outliers

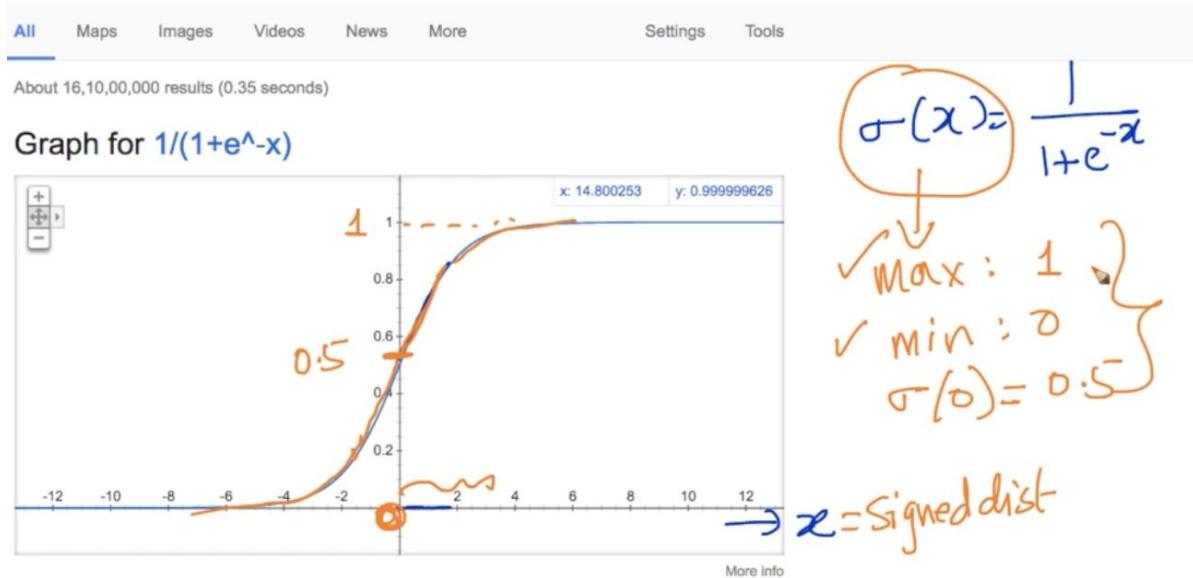
Squashing:

- Idea:- instead of using signed-dist
- if signed-dist is small :- use it as is
- " " " is large: make it a smaller value

So to deal with it we do squashing. It means if the distance is small use it as it is but if it's large make it a smaller value



To get rid of the outlier we need a function such that after a certain distance ,10 here it'll taper off. That's max $f(\sum_{i=1}^n y_i w^T x_i)$. In the above example our outlier with distance 100 is changing our whole model so it should be tapered down . So we need to find that function

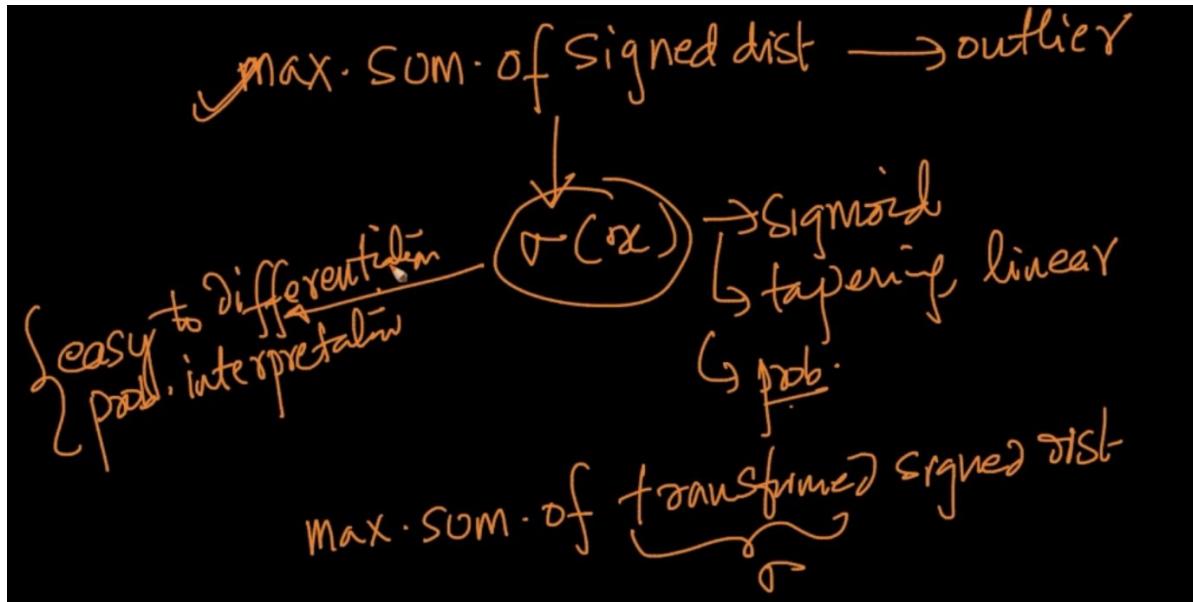


f

One such function is sigmoid function. We can see that till a certain value it's increasing linearly but after that it's getting tapered down and it's maxed value is 1 so if our x: signed distance > 3 here it's y will be 1

This sigmoid function has nice probabilistic interpretability properties . For ex: If the point is at plane then $w^T x_i = 0$ therefore $P(y_i = 1) = 0.5$ because we don't know whether it's +ve or -ve . If the distance $w^T x_i$ is large say 12 here our $P(y_i = 1) = 1$.





We wanted max sum of signed distance but it's problem of outlier so we are using a function $\sigma(x)$ which has some nice properties mentioned above. So we want max sum of transformed signed distance with σ .

$$\omega^* = \arg \max_{\omega} \sum_{i=1}^n \sigma(y_i w^T x_i)$$

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

$$\omega^* = \arg \max_{\omega} \sum_{i=1}^n \frac{1}{1+e^{-y_i w^T x_i}}$$

↔ less impacted by outliers

So our goal is to find optimal w : $w^* = \operatorname{argmax}_{\omega} \sum_{i=1}^n \frac{1}{1+e^{-y_i w^T x_i}}$

MATHEMATICAL FORMULATION OF OBJECTIVE FUNCTION

monotonic fn:- $g(x)$

$x \uparrow; g(x) \uparrow \rightarrow$ monotonically incr-fn

if $x_1 > x_2$ then $g(x_1) > g(x_2)$
then $g(x)$ is said to be mon. incr fn

A **monotonic function** is a function which gets increased with x

~~opt. prob:~~

$x^* = \arg \min_x x^2$ Minima & Maxima

$f(x) = x^2$

$x^* = \arg \min_x f(x)$

We are finding the minima of x^2 which is we can see from the plot is 0

mon. incv $g(x) = \underline{\log(x)}$ optimization

$\textcircled{1} \rightarrow \boxed{0} = \underline{x^*} = \arg \min_x f(x) \quad ; f(x) = \underline{x^L}$

$\textcircled{2} \rightarrow x' = \arg \min_x g(f(x))$
 $\log(x^2)$

claim: $\underline{x^*} = x'$ $\because g(x)$ is a monotonic fn

$g(x) = \log(x)$ which is a monotonically increasing function. We say here that $x^* = x'$

If $g(x)$ is a monotonic fn. $\underline{x^*} = \arg \min_x g(f(x))$
 $\underline{x^*} = \arg \max_x g(f(x))$

Here's the condition here.

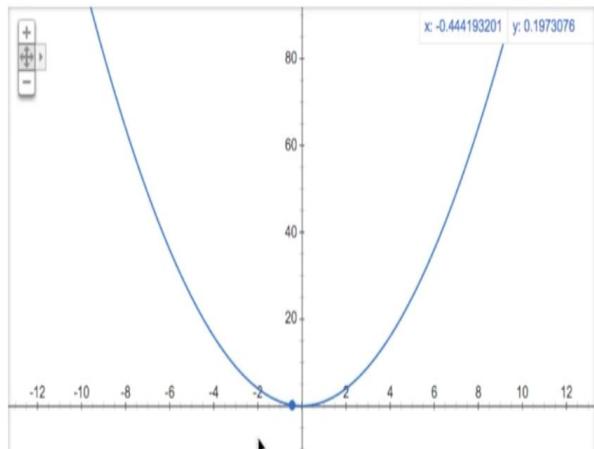
$$\omega^* = \operatorname{argmax}_{\omega} \sum_{i=1}^n \log \left\{ \frac{1}{1 + \exp(-y_i \omega^T x_i)} \right\}$$

$\log(1/x) = -\log(x)$

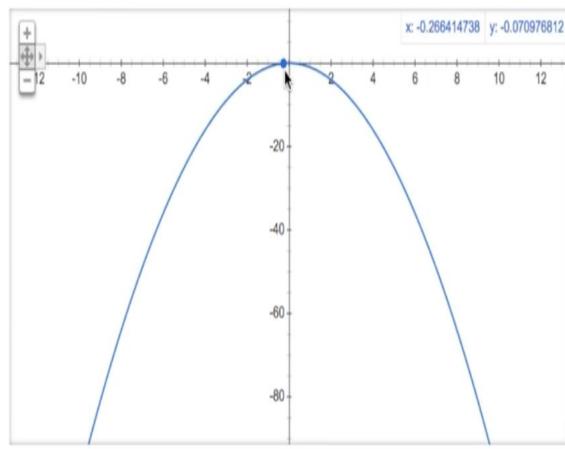
$$\omega^* = \operatorname{argmax}_{\omega} \sum_{i=1}^n -\log \left(1 + \exp(-y_i \omega^T x_i) \right)$$

We have applied log to the sigmoid function and getting that equation

Graph for x^2



Graph for $(-1)^*x^2$



Minimum of x^2 = Maximum of $-x^2$. So we do the same to the above function as well

$$\boxed{\arg \max_x f(x) = \arg \min_x -f(x)}$$

$$\boxed{\arg \max_x f(x) = \arg \min_x f(x) \text{ } \xrightarrow{\text{+1 or -1}}}$$

$$\boxed{w^* = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))}$$

We want to find $\min(\log(1 + \exp(-y_i w^T x_i)))$

$$\arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) \xrightarrow{\text{opt-LR}}$$

$$\arg \min_w \sum y_i w^T x_i$$

$$\cancel{\arg \max_w \sum y_i w^T x_i} \xrightarrow{\text{huge outlier problem}}$$

Sum of signed dist

If we remove 1 from that we'll only get $y_i w^T x_i$ which is sum of signed distance but we are not solving that problem since it gives us a huge outlier problem

WEIGHT VECTOR

Weight -Vector

$$\check{w} = \arg \min_w \sum_{i=1}^n \log (1 + \exp (-y_i(w^T x_i)))$$

$w \in \mathbb{R}^d$ $x_i \in \mathbb{R}^d$

$w = \langle w_1, w_2, w_3, w_4, \dots, w_d \rangle$

We want to get the optimal w which is also known as weight vector which has d -dimensions as well as x has d -dimensions

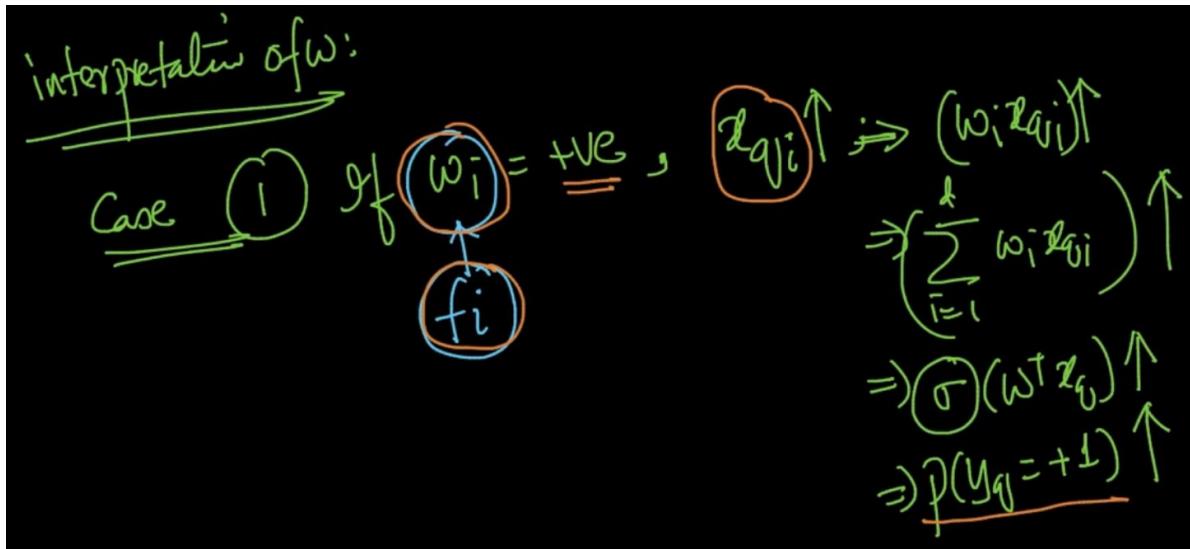
$$w = \langle w_1, w_2, \dots, w_d \rangle$$

$$f_1, f_2, \dots, f_d$$

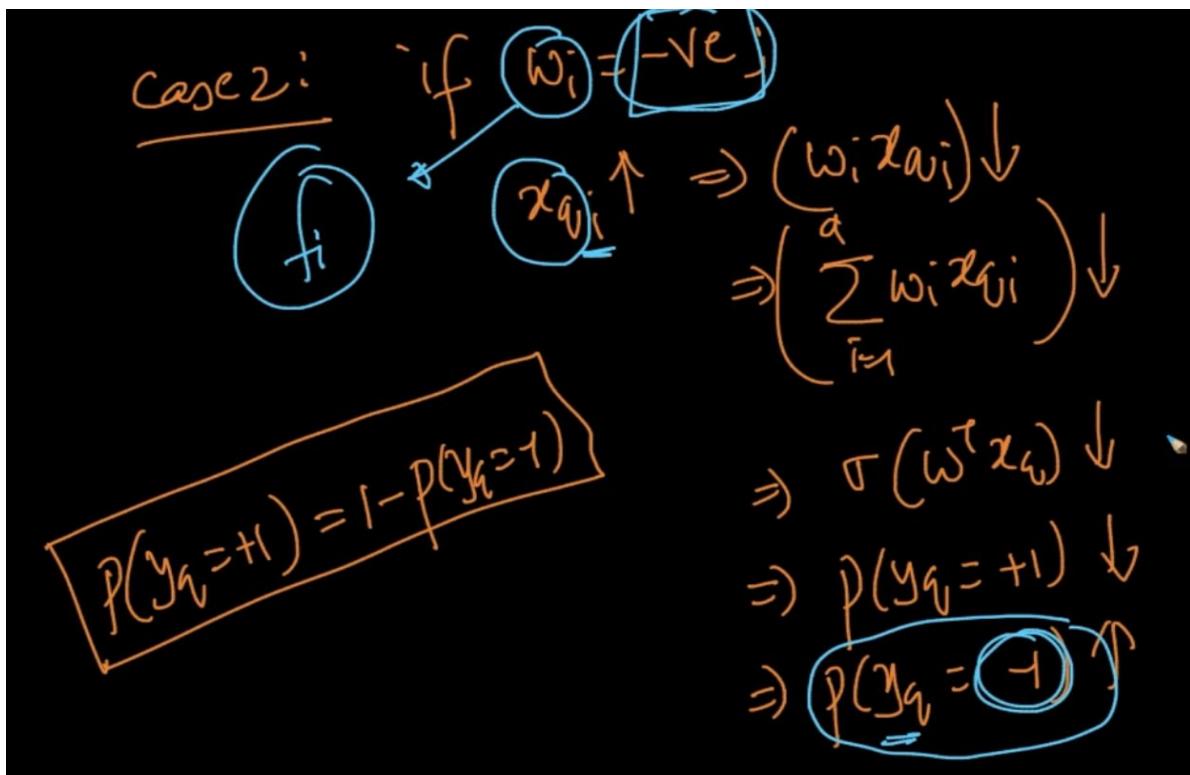
decision: $x_q \rightarrow y_q$ $\begin{cases} \text{if } w^T x_q > 0 \text{ then } y_q = +1 \\ \text{if } w^T x_q < 0 \text{ then } y_q = -1 \end{cases}$

prob.interp: $P(y_q = +1) = \sigma(w^T x_q)$

This is how we predict the class label y . We can also get probabilistic interpretation using sigmoid function



Case 1: If w_i corresponding to a feature $f_i = +ve$ and $x_{qi} \uparrow$ then all the other things get increased as well. As $x_{qi} \uparrow \Rightarrow P(y_q = +1) \uparrow$



Case 2: If $w_i = -ve$ and $x_{qi} \uparrow$ then all the other things will get decreased and $P(y_q = +1) \downarrow$ but $P(y_q = -1) \uparrow$

L2 REGULARIZATION OVERFITTING AND UNDERFITTING

$$\sum_{i=1}^n \log \left(1 + \underbrace{\exp(-z_i)}_{>0} \right)$$

$\exp(-z_i) > 0$

$$\log \left(1 + \underbrace{\exp(-z_i)}_{>0} \right) > 0$$

$\log(1) = 0$
 $\log(2) > \log(1)$
 $\left\{ \begin{array}{l} \log(1+\delta) > \log(1) \\ \delta > 0 \end{array} \right.$

$\exp(-z) \gg 0$ and the value of our whole function > 0 since $\log(1 + \delta) \gg \log(1)$ and
 $\delta = \exp(-z)$

$$w^* = \arg \min_w \sum_{i=1}^n \log \left(1 + \underbrace{\exp(-z_i)}_s \right)$$

if $z_i = +ve$, $z_i \rightarrow +\infty$
then $\exp(-z) \rightarrow 0$
 $\log \left(1 + \underbrace{\exp(-z)} \right) \rightarrow 0$

We want our z to be +ve for correct classification and if $z \rightarrow \infty$ then $\exp(-z) \rightarrow 0$

$$w^* = \arg \min_w \left[\sum_{i=1}^n \log(1 + \exp(-z_i)) \right] > 0$$

minimal value of $\sum_{i=1}^n \log(1 + \exp(-z_i))$ is "0"
 { it occurs when $z_i \rightarrow \infty$ for all i

We want to minimize it and its minimal value is 0 when $z_i \rightarrow \infty$ for all i

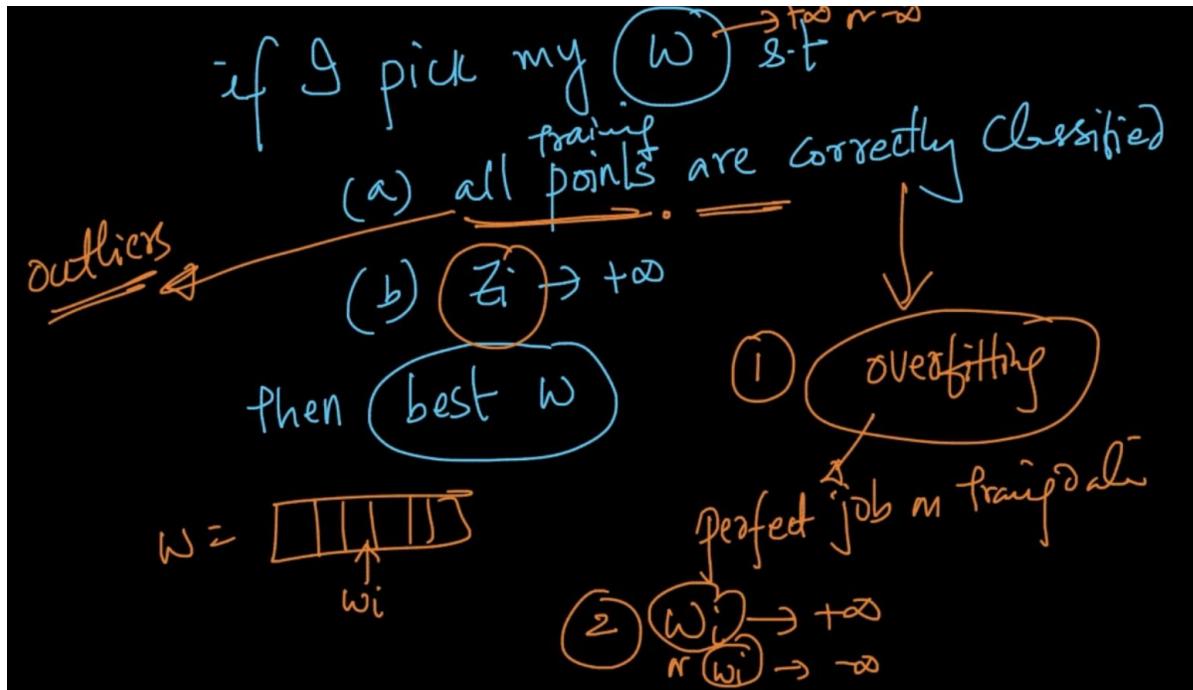
$$z_i = y_i (w^T x_i) + v$$

$D_{\text{train}} = \{(x_i, y_i); i=1:n\}$

$z_i \rightarrow \infty$ $\hookrightarrow w$ modify my w in such a way that each $z_i \rightarrow \infty$

- ① $z_i = +ve$; x_i is correctly classified by w
- ② $z_i \rightarrow \infty$;

For $z_i \rightarrow +\infty$ we need to modify w such that $z_i \rightarrow +\infty$. Since x and y are fixed w is the only variable here



Since, our whole function is becoming 0 we are perfectly fitting the data . It means that it can be overfitting and since $z \rightarrow +\infty$, $w_i \rightarrow +\infty$ so it's becoming too large

regularization

$$w^* = \arg \min_w \left[\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda w^T w \right]$$

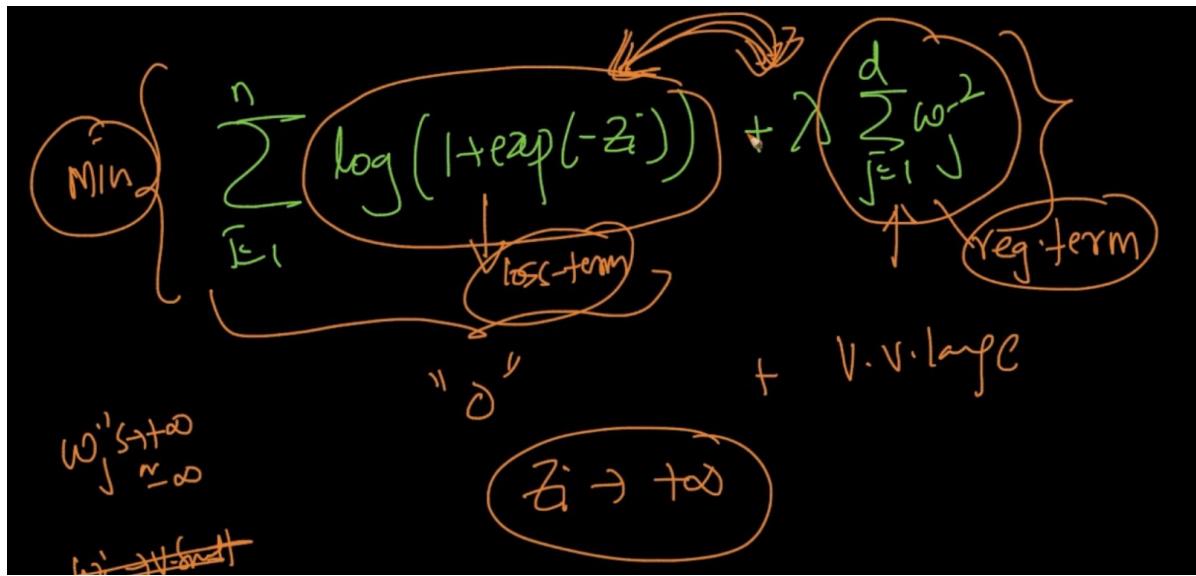
$w_j \rightarrow \infty$
 $w_j \rightarrow -\infty$

loss-term
 $\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$

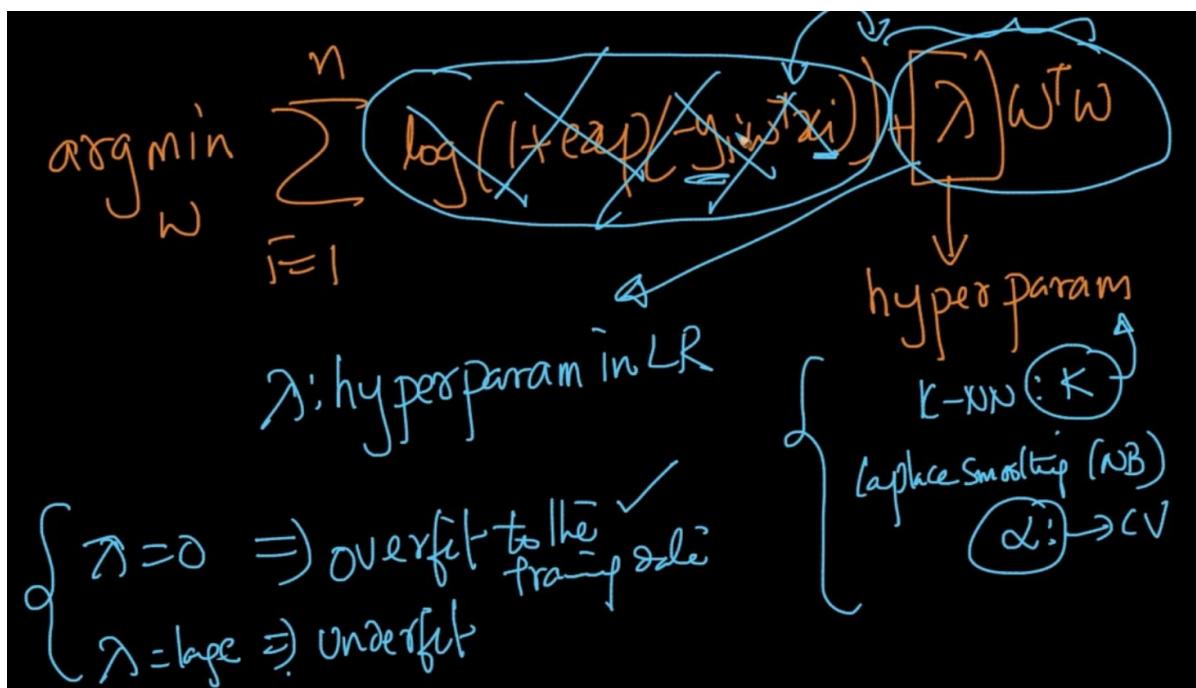
$\lambda w^T w \rightarrow \|w\|_2^2$

$\lambda \sum_{j=1}^d w_j^2$
regularization term

So we are adding $\lambda w^T w$ to our optimal w , known as regularization term. As $w \rightarrow \infty$ Then our regularization term will also be ∞ but we need to minimize it.



As $w_j \rightarrow \infty$ our log function will be 0 as $\exp(-z_i)$ but our regularization term will be very large and since we want to minimize the whole it won't happen as log function $\rightarrow 0$ means we are overfitting our data which we don't want. So both loss term and regular term are fighting with each other to minimize the whole



λ is hyper-parameter as $\lambda = 0 \Rightarrow$ It'll mean overfit to Training data and if $\lambda = \text{large} \Rightarrow$ Underfit since we are not using the training data i.e log thing and only w

L1 REGULARIZATION AND SPARSITY

L_1 regularization and Sparsity

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left[\underbrace{\sum_{i=1}^n \log(1 + \exp(-y_i w^\top x_i))}_{\text{logistic loss}} + \lambda \|w\|_1 \right]$$

$\begin{cases} z_i \rightarrow +\infty \\ w_i \rightarrow +\infty \\ \text{or} \\ w_i \rightarrow -\infty \end{cases}$
 overfit
 $\underbrace{\quad}_{L_1-\text{reg.}}$

Question: alternatives to L_2 -Reg.

popular:- L_1 -Reg

We'd seen L2 regularization now we see the other alternatives

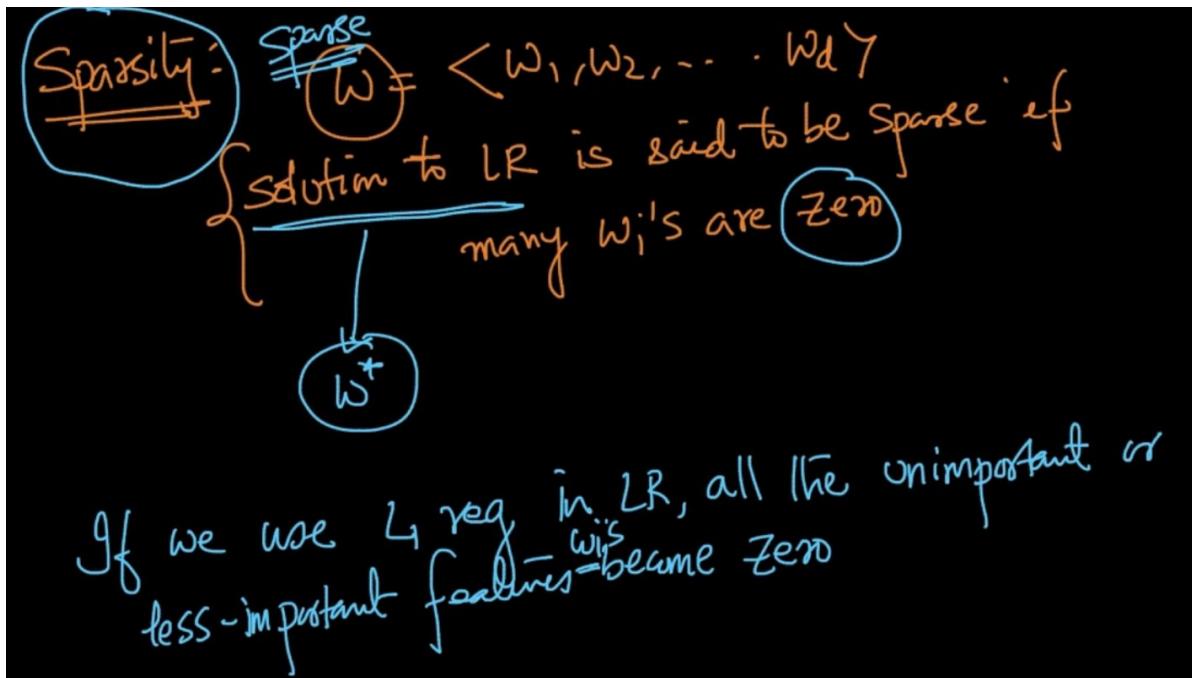
$$\|w\|_2^2 \text{ for reg.} \quad \|w\|_1 = \sum_{i=1}^d |w_i| \text{ abs}$$

$$\|w\|_1 \text{ for reg.:-}$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left[\underbrace{\text{(logistic loss for training data)}}_{\text{hyperparam}} + \lambda \|w\|_1 \right]$$

$\cancel{\text{will avoid } w_i \rightarrow +\infty \text{ or } -\infty}$

This is L1-Regularizer. It has absolute value of w. It's different advantage as compared to L2 which we'll see below



If we use L1 regularization then all the less important features become zero

$$f_1, f_2, \dots, \overset{\text{less-important}}{f_i}, \dots, f_d$$

$$w = \langle w_1, w_2, \dots, \overset{\text{highlighted}}{w_i}, \dots, w_d \rangle$$

Zero - if L_1 is used

if L_2 reg is used; w_i becomes a small values but not necessarily zero

If we are using L1 regularization and feature f_1 is less important then its corresponding w_i will be zero as compared to L2 where w_i becomes a small value

Elastic-net:

either L₁ or L₂

$$\left\{ \begin{array}{l} w^t = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-z_i)) + \\ + \lambda_1 ||w||_1 \\ + \lambda_2 ||w||_2 \end{array} \right.$$

Two hyper-param:- λ_1 & λ_2

Elastic-Net is a combination of L1 and L2 regularization. It has two Hyper-parameters λ_1 and λ_2 which can be found out through cross-validation

PROBABILISTIC DERIVATION OF LOGISTIC REGRESSION

prob. derivation of (LR)

Naive Bayes

→ ① features are real valued:-
↪ Gaussian dist

→ ② $y_i = +1 \text{ or } 0$
Bernoulli r.v (coin-toss)

If these two points are true the Logistic Regression can be derived using probability

tic Regression and summarized in equations (16) and (17). In particular, consider
a GNB based on the following modeling assumptions:

- ① GNB + Bernoulli
- Y is boolean, governed by a Bernoulli distribution, with parameter $\pi = P(Y=1)$
 - $X = \langle X_1 \dots X_n \rangle$, where each X_i is a continuous random variable
 - For each X_i , $P(X_i|Y=y_k)$ is a Gaussian distribution of the form $N(\mu_{ik}, \sigma_i)$
 - For all i and $j \neq i$, X_i and X_j are conditionally independent given Y

If the above assumptions are true we can easily derive i.e $P(y=1|X) = \frac{1}{1+exp(\sum_{i=1}^n w_i X_i)}$

geom: $\omega^t = \arg \min_{\omega} \sum_{i=1}^n \log(1 + \exp(-y_i \omega^T x_i)) + \gamma_{\text{reg}}$

$+1 \text{ or } -1$

prob: $\omega^t = \arg \min_{\omega} \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log(1-p_i) + \gamma_{\text{reg}}$

where $p_i = \sigma(\omega^T x_i)$

$+1 \text{ or } 0$

Both the formulae are same but in Geometric , y has +1 or -1 value and Probabilistic y is +1 or 0

Case 1: $y_i : +ve$

<u>geom:-</u>	$y_i = +1$
<u>prob:-</u>	$y_i = +1$

<u>geom:-</u>	$\log(1 + \exp(-\omega^T x_i))$	$p_i = \sigma(\omega^T x_i)$
<u>prob:-</u>	$\log\left(\frac{1}{1 + \exp(-\omega^T x_i)}\right)$ $= \log\left(1 + \exp(-\omega^T x_i)\right)$	$- \log(x) = \log(1/x)$

Case 1: Here, $y = +1$ and if we do the math both will be same

<u>Case 2:</u>	$y_i = -ve;$	$y_i = -1 \xleftarrow{\text{geom}}$
		$y_i = 0 \xleftarrow{\text{prob}}$
<u>geom:</u>	$\log(1 + \exp(w^T x_i))$	$\begin{aligned} -\log(x) \\ = \log(1/x) \end{aligned}$
<u>prob:</u>	$-1 \cdot \log\left(1 - \frac{1}{1 + \exp(-w^T x_i)}\right)$	$-1 \cdot \log\left(\frac{\exp(-w^T x_i)}{1 + \exp(-w^T x_i)}\right)$

Case 2: Here, $y = -1$ in Geometry and $y = 0$ in probability. We are inverting the numerator and denominator using property $-\log(x) = \log(1/x)$

$$\begin{aligned}
 &= \log \left(\frac{1 + \exp(-\omega^T x_i)}{\exp(-\omega^T x_i)} \right) \\
 &\quad \text{divide numerator & den by } (\exp(-\omega^T x_i)) \\
 &= \log \left(1 + \frac{1}{\exp(\omega^T x_i)} \right) \\
 &= \log \left(1 + \exp(-\omega^T x_i) \right)
 \end{aligned}$$

$\boxed{\frac{1}{e^{-x}} = e^x}$

As we can see the Probabilistic measure is deriving to the same thing

LOSS MINIMIZATION INTERPRETATION

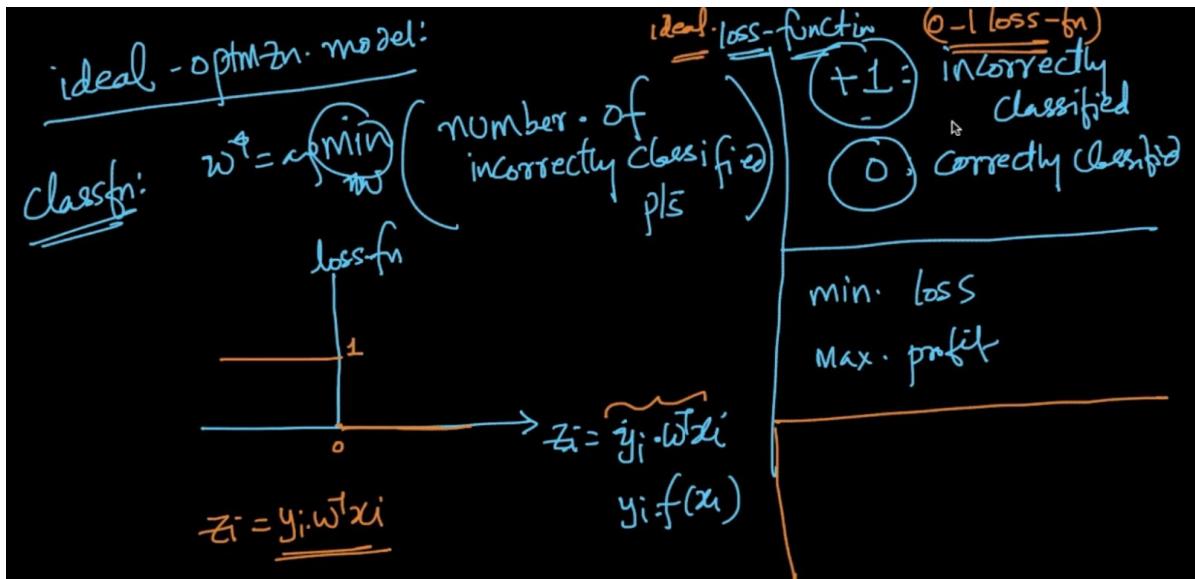
Loss minimization interpretation of LR

✓ geometric & probabilistic

✓ loss-minimization

$$\left\{ \begin{array}{l} f(x_i) = \underline{\omega^T x_i} \\ \omega^T = \arg \min_{\omega} \sum_{i=1}^n \log \left(1 + \exp(-y_i \omega^T x_i) \right) \\ z_i = y_i \omega^T x_i = \underline{y_i \cdot f(x_i)} \end{array} \right.$$

We'd seen Geometric and Probabilistic interpretation now we'll see loss minimization



In our ideal optimization model we want to minimize the number of incorrectly classified points. Our loss function is saying that if +1: Then incorrect classification and 0 if correctly classified. So in the graph if $z < 0 \Rightarrow 1$ and if $z > 0$ then 0

ideal

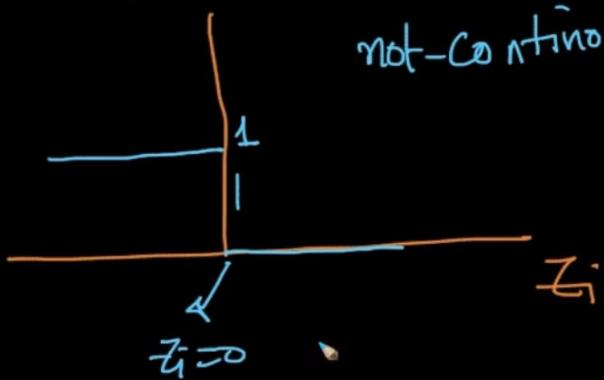
$$\hat{\omega}^* = \operatorname{argmin}_{\omega} \sum_{i=1}^n \underbrace{\text{o-1 loss}(x_i, y_i, \omega)}_{\text{loss}}$$

$$\text{o-1 loss}(z_i) = \begin{cases} 1 & \text{if } z_i < 0 \\ 0 & \text{if } z_i \geq 0 \end{cases}$$

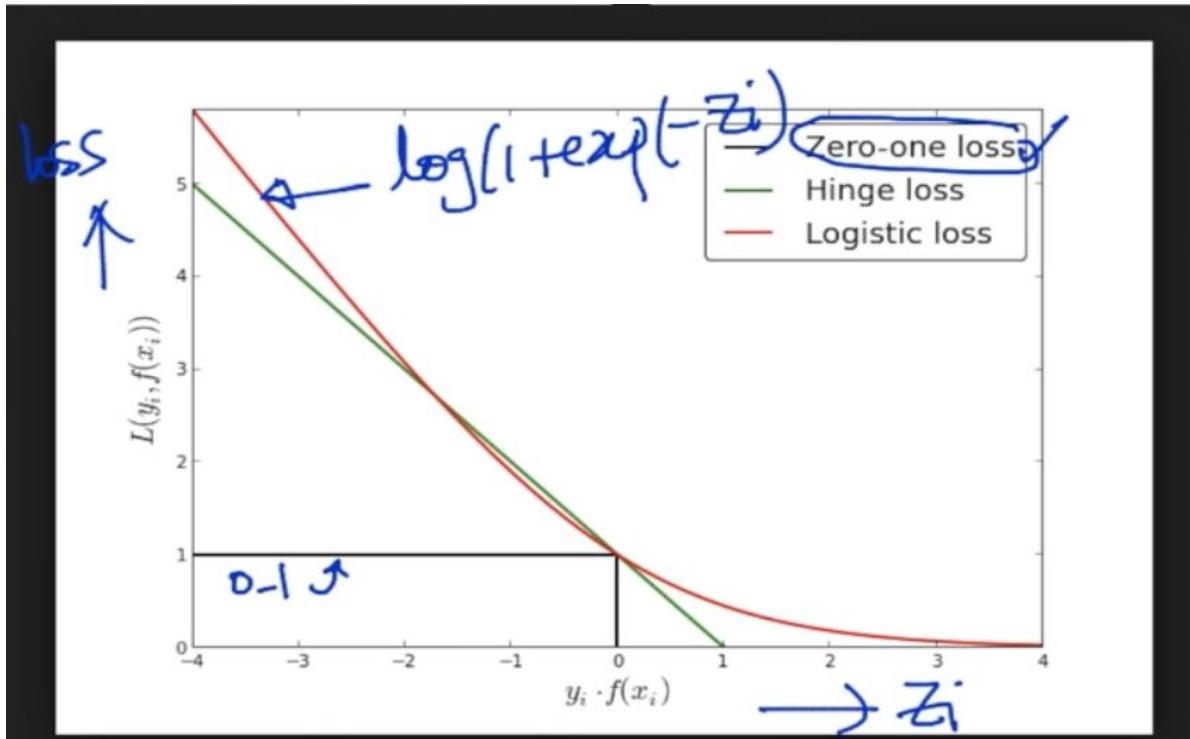
$f(x)$ is differentiable

\hookrightarrow continuous

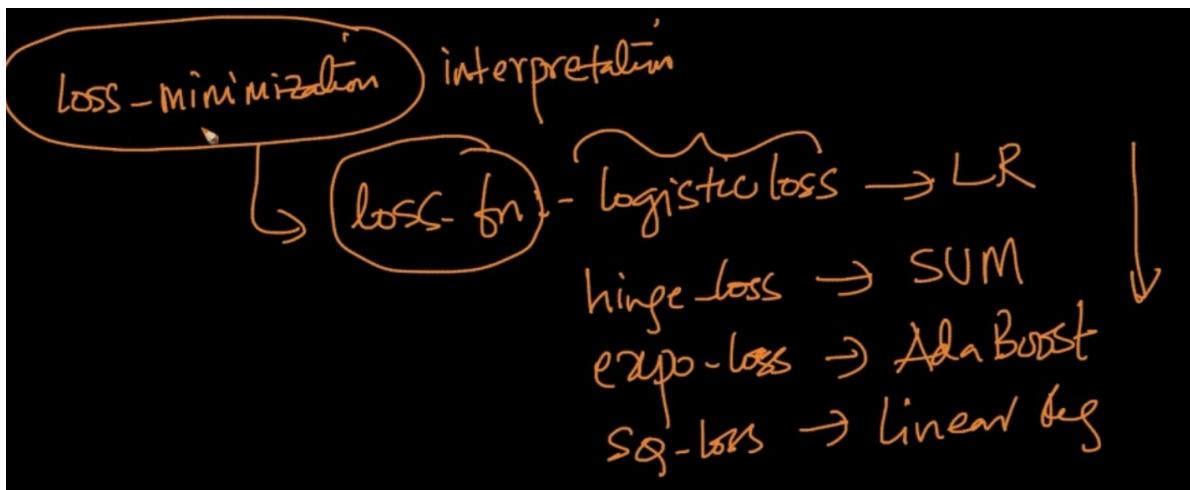
not-continuous \Rightarrow not-differentiable



But to solve the optimal model we need to differentiate the function $f(x)$ and for function $f(x)$ to be differentiable $f(x)$ should be continuous

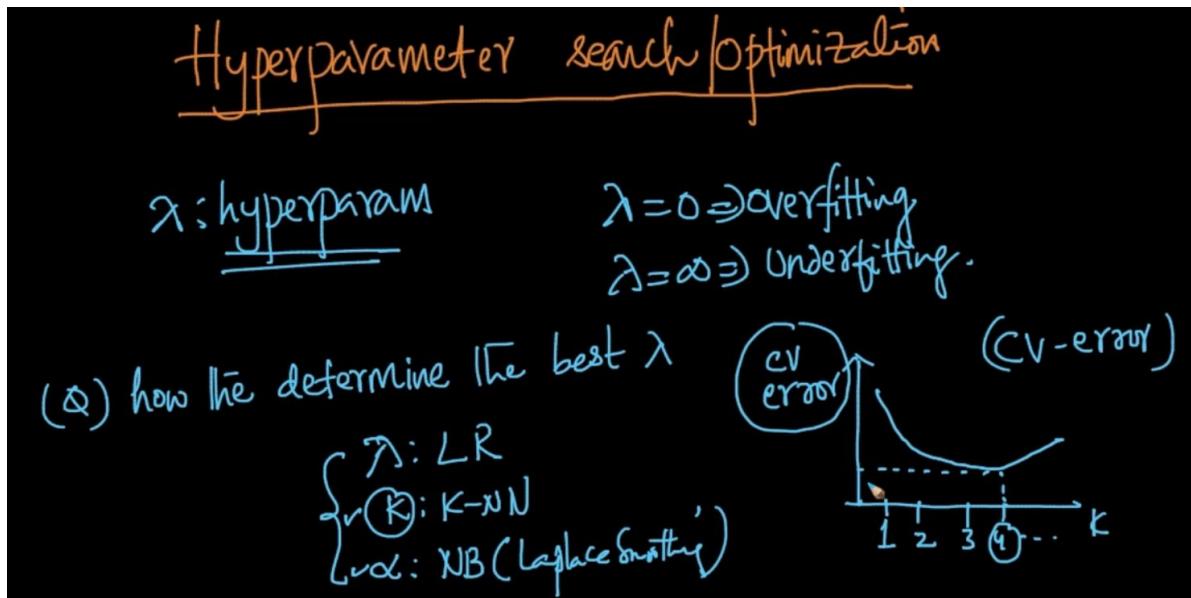


0-1 is not continuous hence not differentiable. So we'll approximate it by using continuous functions. We are using log-loss where $y = \log(1 + \exp(-z))$ where $z = y \cdot f(x)$ and $f(x) = w^T x$. We could've drawn any other loss function as well but for logistic regression Logistic loss is used. Here, if $z > 0$ it's getting close to 0 after certain value but when $z < 0$ it's increasing sharply.

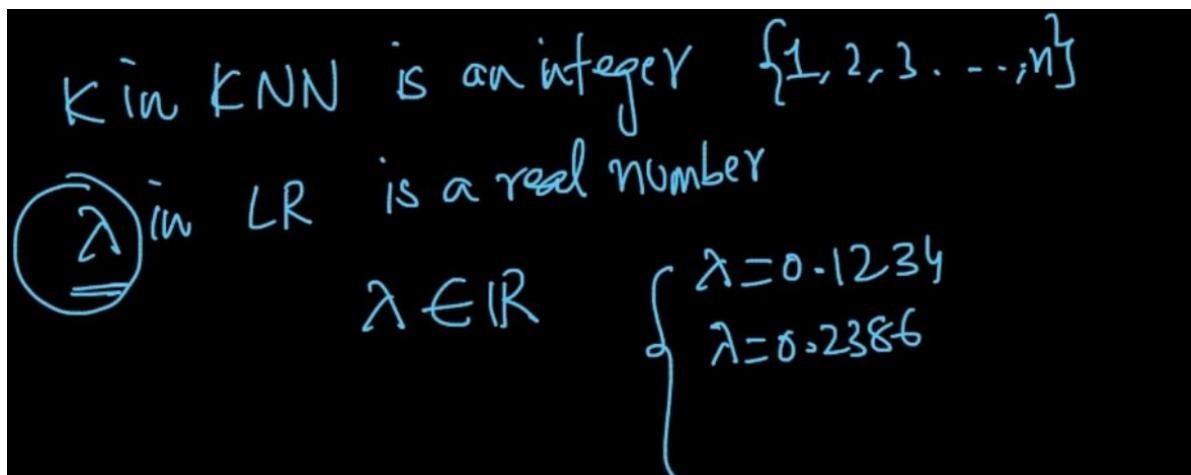


Different loss function for different ML algorithms

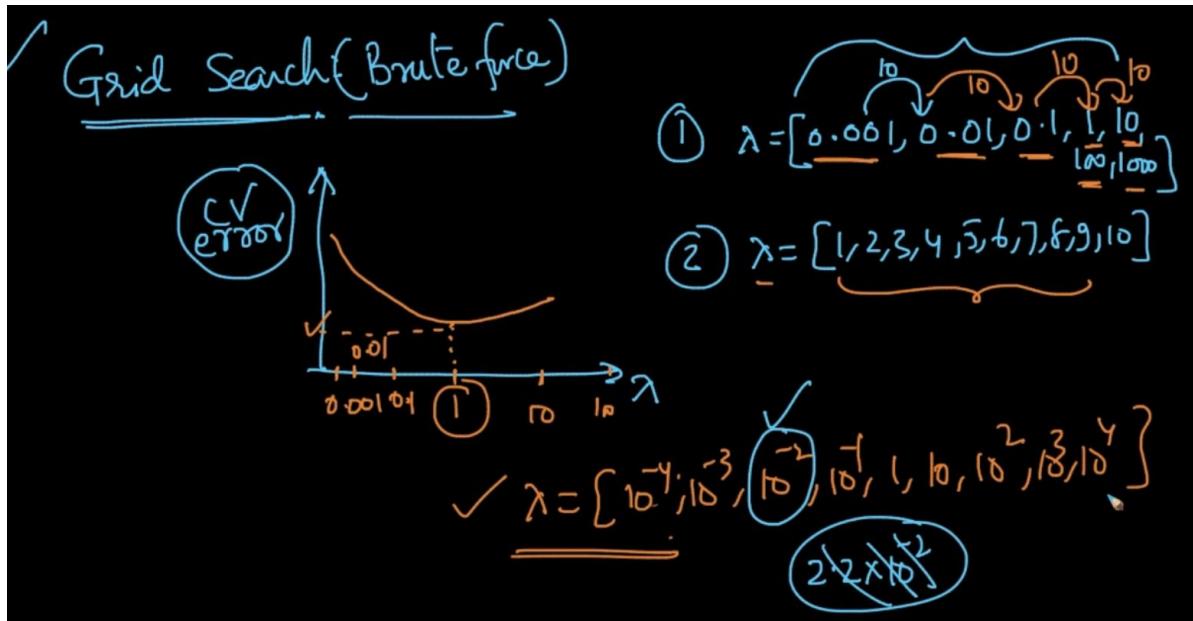
HYPERPARAMETER SEARCH/OPTIMIZATION



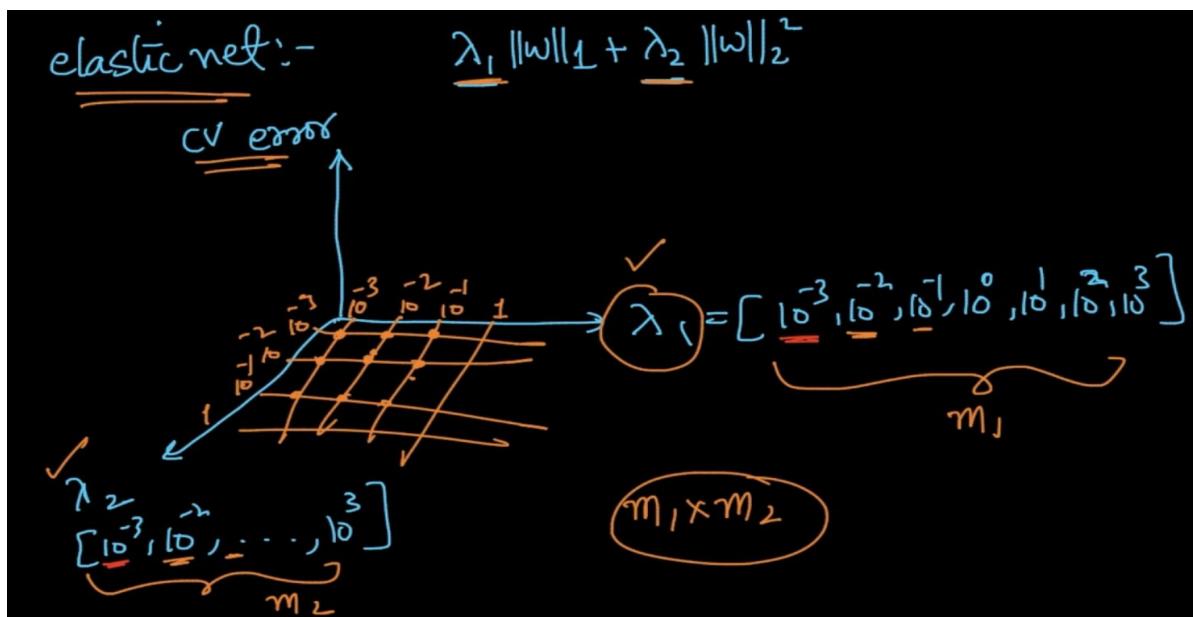
We want to identify the best λ for Logistic Regression. In k-NN, we found the best k by plotting the CV-error vs k graph and checking whichever k gives lowest CV error is our k



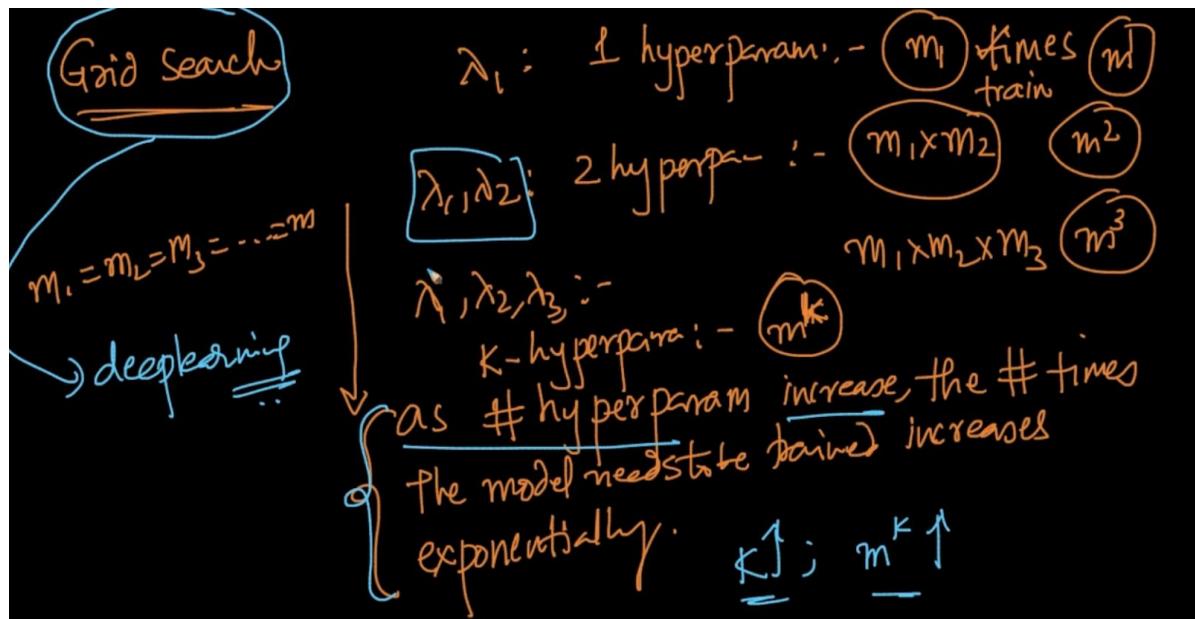
But our k in k-NN is an integer but our λ in LR is a real number



In order to tackle that real value situation of λ we'll do a Brute Force technique called Grid Search where we create a window from $\lambda = [0.001, 0.01, \dots, 100, 1000]$ and we plot the CV error v/s λ graph and select the best λ which has lowest CV error

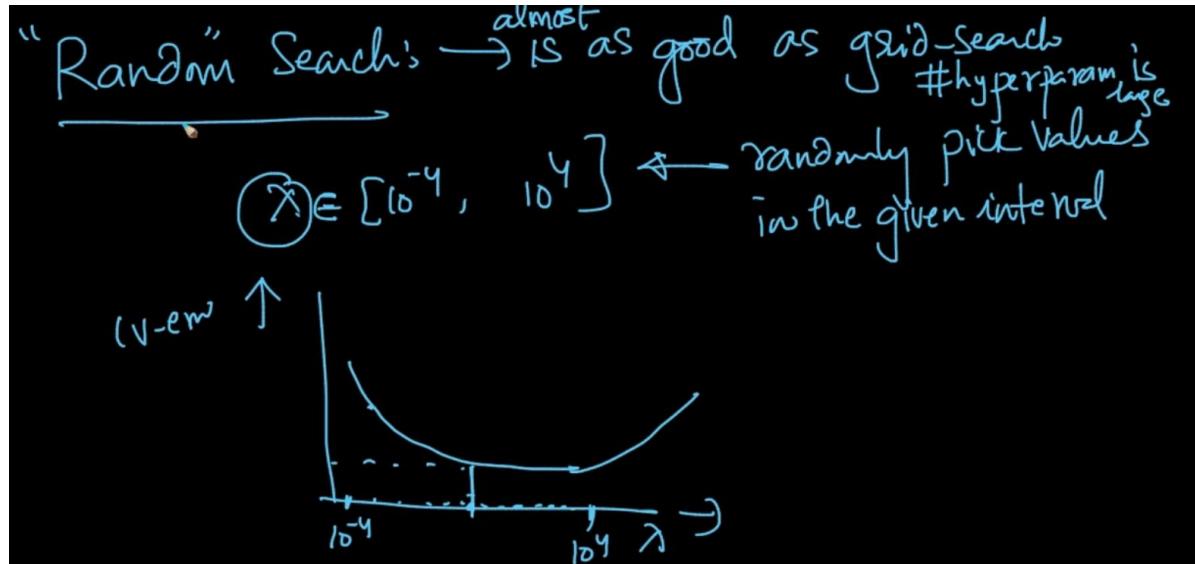


If we've more than one parameters, lets say above. We'll plot the graph of λ_1 v/s λ_2 v/s CV error to select the best λ_1, λ_2 which has lowest CV error. If λ_1 has m_1 dimensions and λ_2 has m_2 dimensions so total calculations will be $m_1 * m_2$



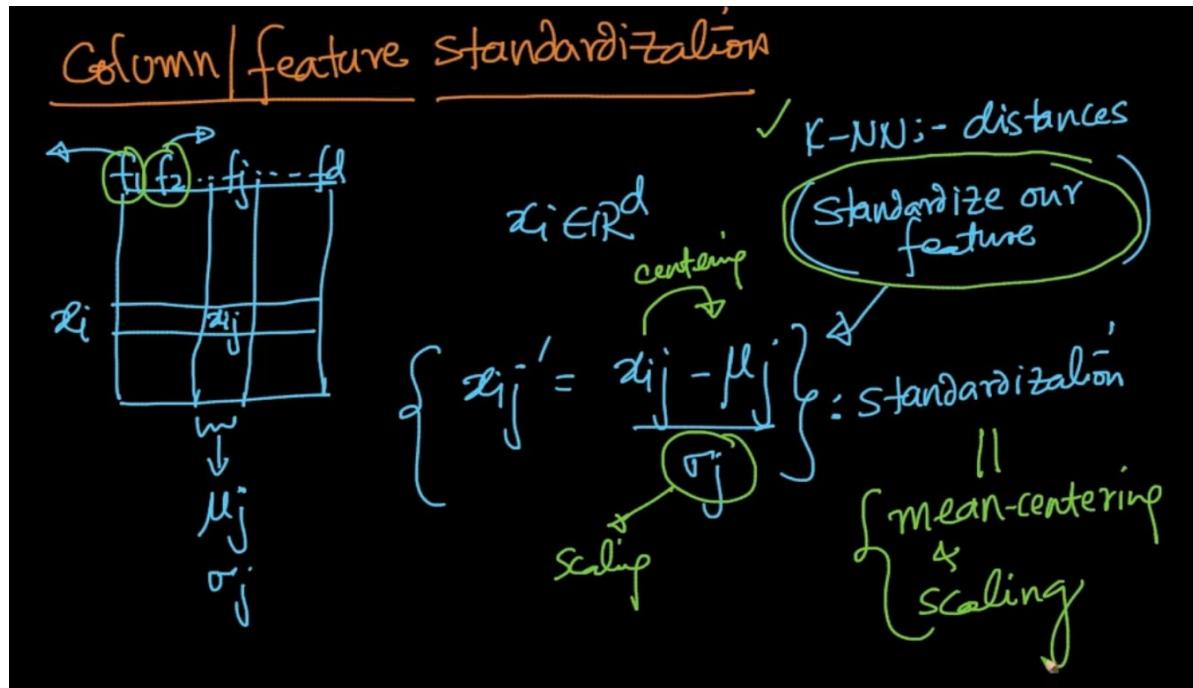
There's one problem in Grid Search. It's that as number of hyperparameters increase The no. of times the model needs to be trained is increased exponentially

Ex : For 3 hyperparameters with m -dimensions has m^3 calculations.

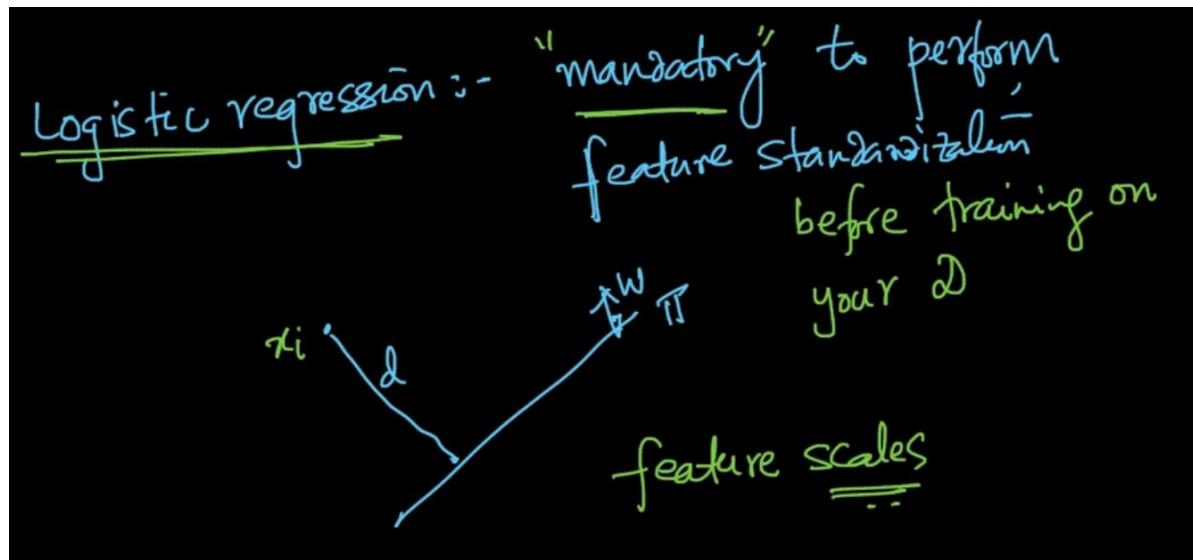


There's another technique called Random Search in which we select an interval and from that interval we randomly pick values and get the value of λ with the lowest CV error. It's said that it's as good as Grid search . Used when hyperparameters are large

COLUMN / FEATURE STANDARDIZATION



We are performing Column standardization by mean-centering and scaling



In Logistic Regression, it's mandatory to perform Column Standardization because if two features are at different scales then it can create problems as distances can be affected by scales

FEATURE IMPORTANCE AND FEATURE INTERPRETABILITY

Feature importance & Model Interpretability

$f_1, f_2, \dots, f_i, \dots, f_d$

$w \rightarrow w_1, w_2, \dots, w_i, \dots, w_d$

$\text{assume: if all features are independent (Naive Bayes)}$

feature-imp: w_j^j

$LR: GNB + \text{Bernoulli}$

We can decide the important features by using their corresponding w's

Note- We assume here that all features are independent just like in Naive Bayes

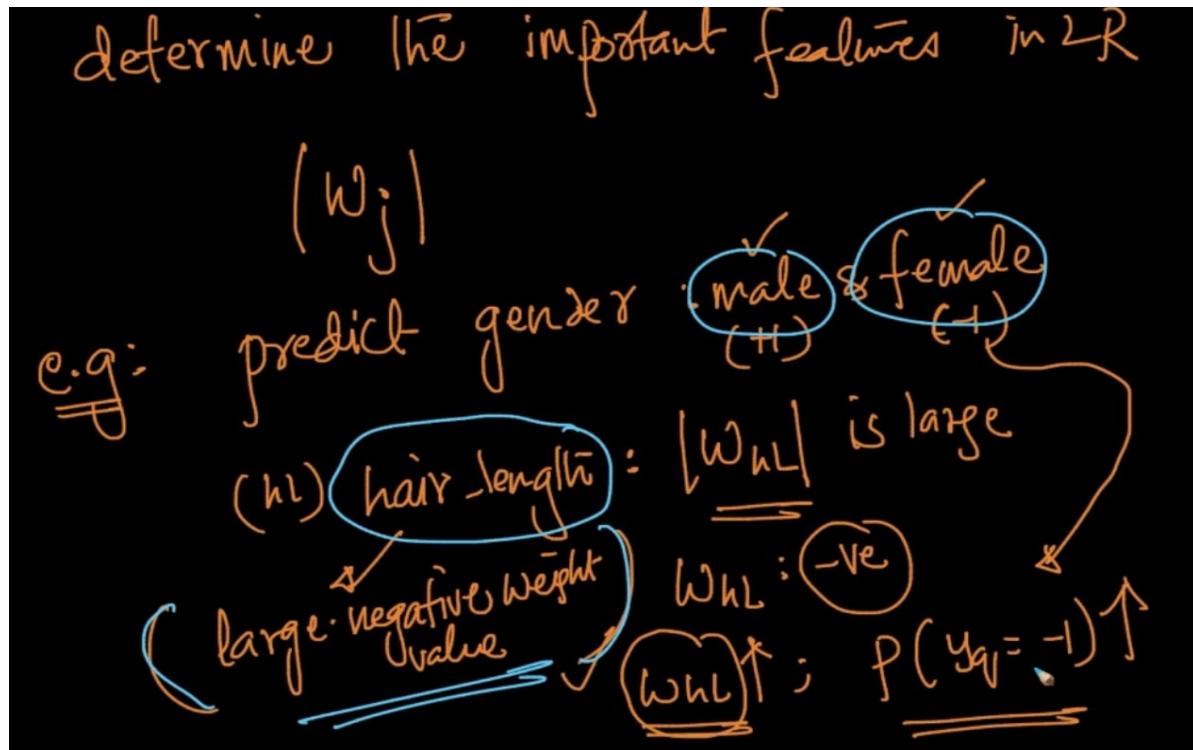
$|w_j| = \text{absolute value of weight corresponding to } f_j$

$|w_j| \uparrow ; (w^T x_{qj}) \uparrow$

Case 1: $w_j = \text{+ve \& large} ; \sum_{j=1}^d w_j \cdot x_{qj} = w^T x_q$

Case 2: $w_j = \text{-ve \& large} ; - \sum_{j=1}^d w_j x_{qj} = \underline{\underline{w^T x_q}}$

Feature importance of feature f_j is dependent on absolute value corresponding weight vector w_j . If $|w_j|$ is large and positive .It'll also contribute to $w^T x_q$ therefore $P(y_q = +1)$ is also increased. Same analogy for case 2



E.x: If we are predicting gender : Male = +1 and Female = -1 and we've a feature Hair Length(h_L) and if $|W_{hL}|$ is large and hair_length has a large negative weight value because as W_{hL} increases $P(y_q = -1)$ also increases so it's an important feature to classify male and female

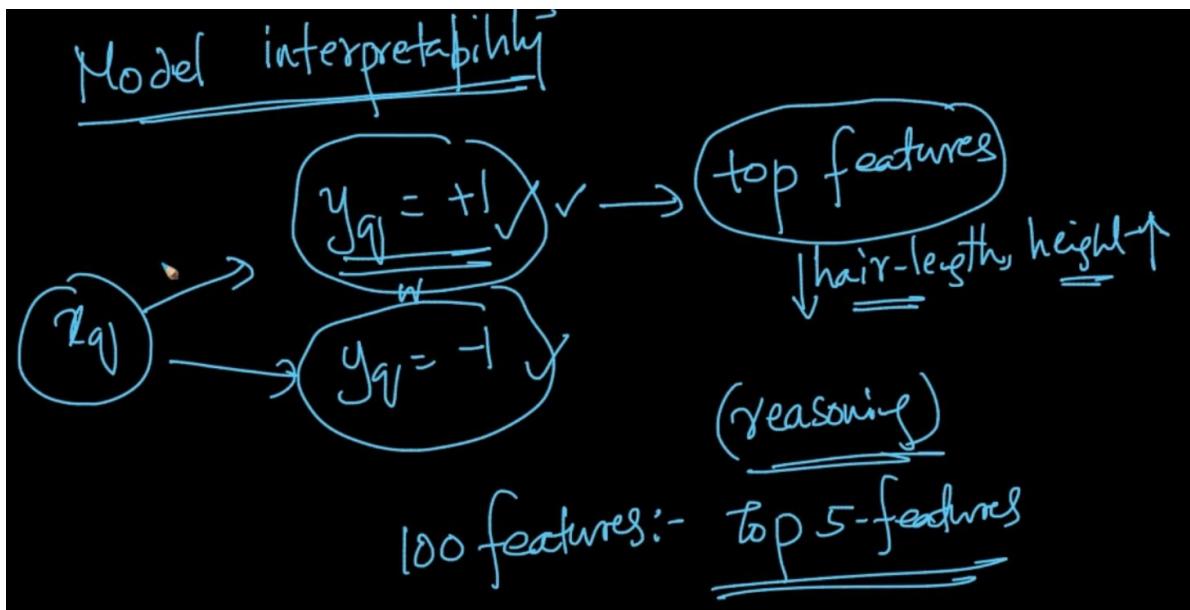
$\{h\}$ height \uparrow ; $P(y_q = +1) \uparrow$
~~height~~ male

w_h : +ve

↳ medium positive weight



If we are taking height as feature , weight w_h : =ve but medium positive not large positive because if we plot the distributions between male and female hL has greater seperability
 But height doesn't have that



It also gives us Model Interpretability because if hair-length is low and height is high then we can say that $y_q = +1$ saying that it's a male. So we select the top features from the absolute value of the corresponding w. If they are large then that feature is important

COLLINEARITY OF FEATURES

"Collinearity" (or) Multicollinearity

f.I :- features are indep; $|w_j|$ as f.I values

Collinearity: \tilde{f}_i, \tilde{f}_j
s.t if, $\tilde{f}_i = \alpha \tilde{f}_j + \beta$
then \tilde{f}_i & \tilde{f}_j are collinear

Features f_i, f_j are collinear if $f_i = \alpha f_j + \beta$ where α & β are constants

Multicollinearity: if f_1, f_2, f_3 & f_4 s.t
 $f_1 = \alpha_1 + \alpha_2 f_2 + \alpha_3 f_3 + \alpha_4 f_4$
 then f_1, f_2, f_3 & f_4 are said to be multicollinear

(Q) Why does $|w_j|$ not be useful as f.I if features are collinear

We'll see that $|w_j|$ for a feature f_j is not as useful for feature importance if features are collinear

$$\mathcal{D} = \langle x_i, y_i \rangle_{i=1}^n$$

$$w^* = \langle 1, 2, 3 \rangle ; x_q = \langle x_{q1}, x_{q2}, x_{q3} \rangle$$

f_1 f_2 f_3

$$w^T x_q = x_{q1} + 2x_{q2} + 3x_{q3}$$

if $f_2 = 1.5 f_1 \Rightarrow f_1 \text{ & } f_2 \text{ are collinear}$

$$w^T x_q = \underline{x_{q1} + 3x_{q1}} + 3x_{q3} = 4x_{q1} + 3x_{q3}$$

$$= \langle 4, 0, 3 \rangle$$

Our weight - vector $w^* = \langle 1, 2, 3 \rangle$ here but $f_2 = 1.5 f_1$ therefore they are collinear so our weight vector will get converted to $\bar{w} = \langle 4, 0, 3 \rangle$

$$\checkmark w^* = \langle 1, 2, 3 \rangle$$

f_3 is the most imp

$$\checkmark \bar{w} = \langle 4, 0, 3 \rangle$$

f_1 is the most imp

not all imp

Same classifier

$\therefore f_2 = 1.5 f_1$

$x_q \rightarrow (w^T x_q)$

In w^* , f_3 is the most important feature as it has the highest absolute weight vector value
 But in \bar{w} f_1 is the most important feature and both w's give the same classifier since features are collinear

NOTE - If features are collinear then weight vectors can change arbitrarily so $|w|$ cannot be used for Feature Importance as seen from the example above

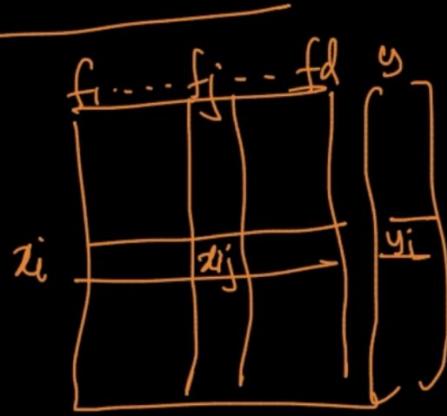
$|w_j|$ as f. I

determine. If features are multicollinear:

Perturbation technique:

$$x_{ij} + \epsilon$$

Small noise
 $N(0, D - \frac{1}{D})$



To determine the features are multicollinear we use perturbation technique in which we add a small noise ϵ to the data point x_{ij} which has small variance

before perturbation: $\underline{w} = \langle w_1, w_2, \dots, w_j, \dots, w_d \rangle$

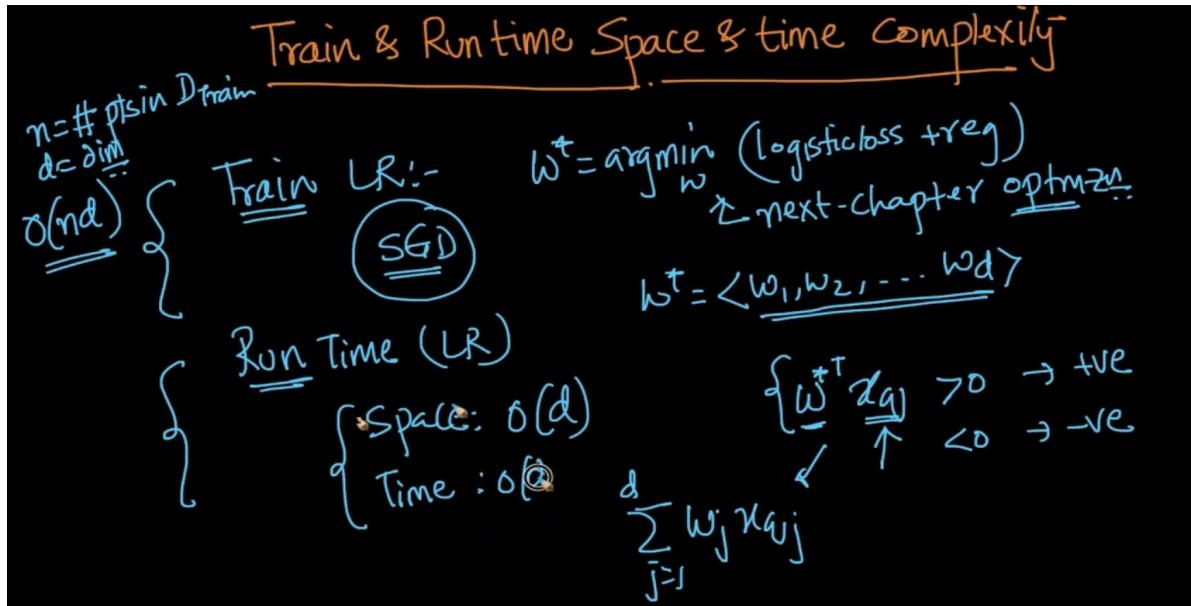
after perturbation: $\tilde{w} = \langle \tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_j, \dots, \tilde{w}_d \rangle$

If $w_i \neq \tilde{w}_i$ differ significantly then your features are collinear

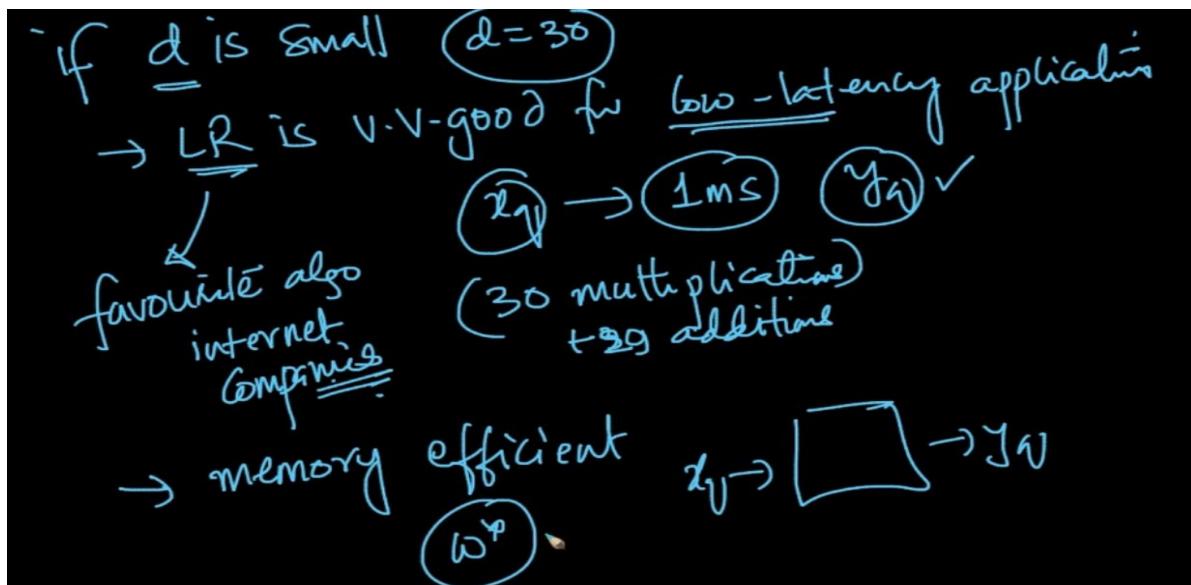
Cannot use $|w_j|$'s as f. I

So we compare the weight vectors before and after the perturbation technique. If the difference is very significant then the features are collinear. So $|w_j|$ cannot be used

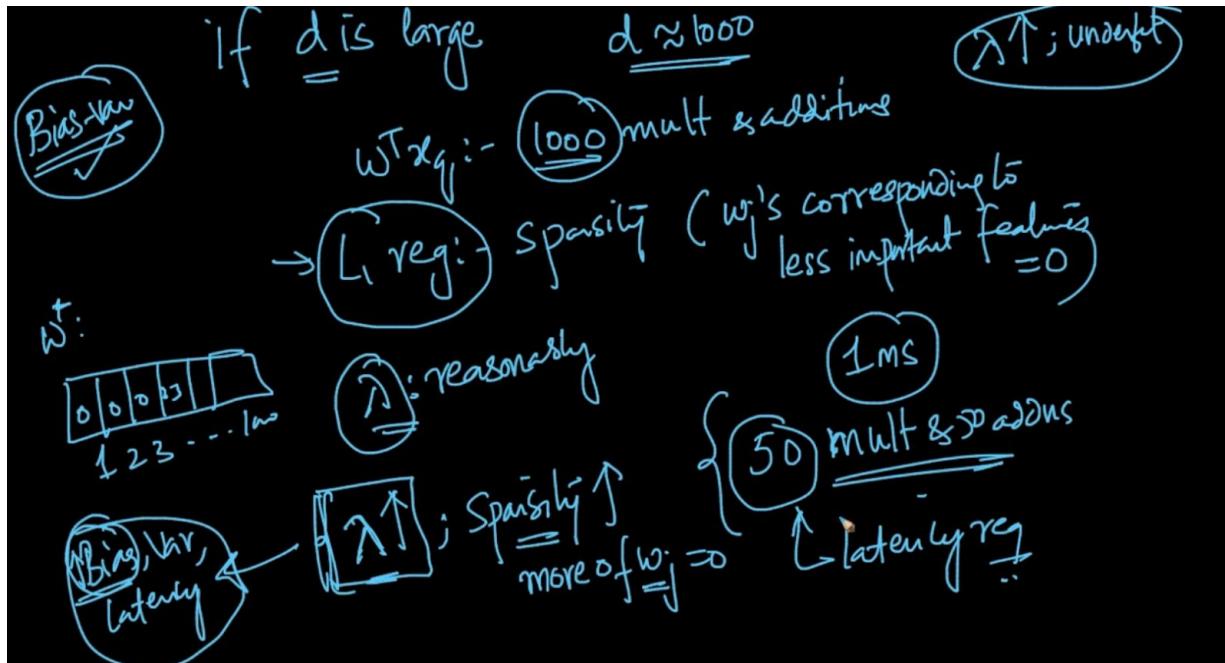
TRAIN & RUNTIME SPACE COMPLEXITY



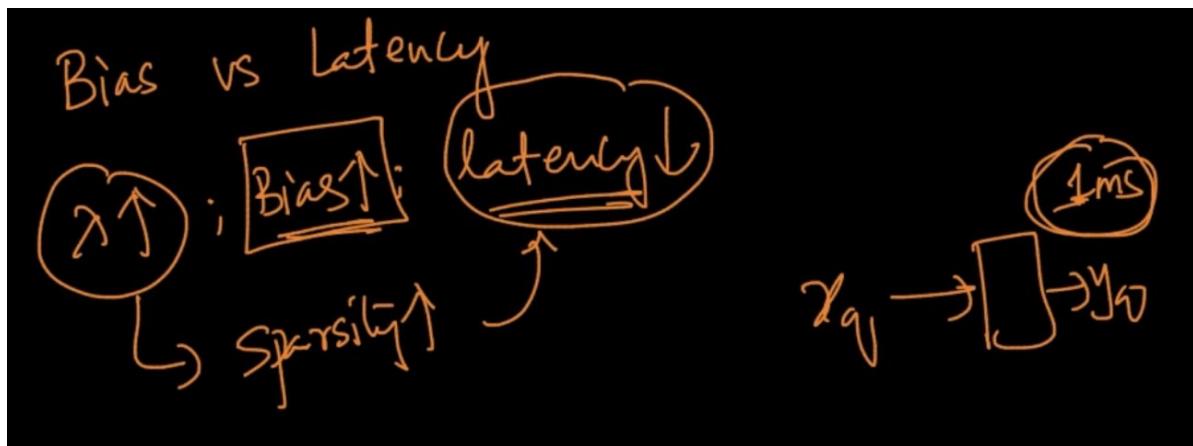
In Runtime, we need to only store the optimal w^* , if it has d - dimensions then it's Space complexity will be $O(d)$ because x_q is given at run time. It's time complexity will also be $O(d)$ as only d - multiplications and $d-1$ addition



If d is small, Logistic Regression is very good for low latency. Ex: If $d = 30$ then 30 multiplications and 29 additions. It's also memory efficient as only need to store w^*



If d is large then there'll be slightly more time as 1000 multiplications and additions but we know that L_1 regularization we can create sparsity i.e w_j corresponding to less important features = 0. If the required latency is only achieved by 50 multiplications and 50 additions per microsecond then we increase the λ as increasing the λ increases sparsity but it means also increase in bias means underfitting . So, we need to find the reasonable λ



So there's a tradeoff we do between Bias and Latency that we've to do to achieve our goals.Logistic Regression one of the most popular algorithm since it has low Latency

REAL WORLD CASES

Cases:-

decision surface:- Linear / hyperplane

assumption :- data is linearly Separable (or almost lin. sep)

feature imp & Interpretability

$|w_j|$ if features are not multicollinear

$x_{qj} \rightarrow y_{qj} = +$

Assumption :- Our data is linearly separable. Feature Importance: $|w_j|$ i.e it depends on the absolute value of the weight vector i'th dimension but our features are not multi-collinear

Imbalanced data:- upsampling or down sampling

Outliers:- → less impact ∵ of $\sigma(x)$

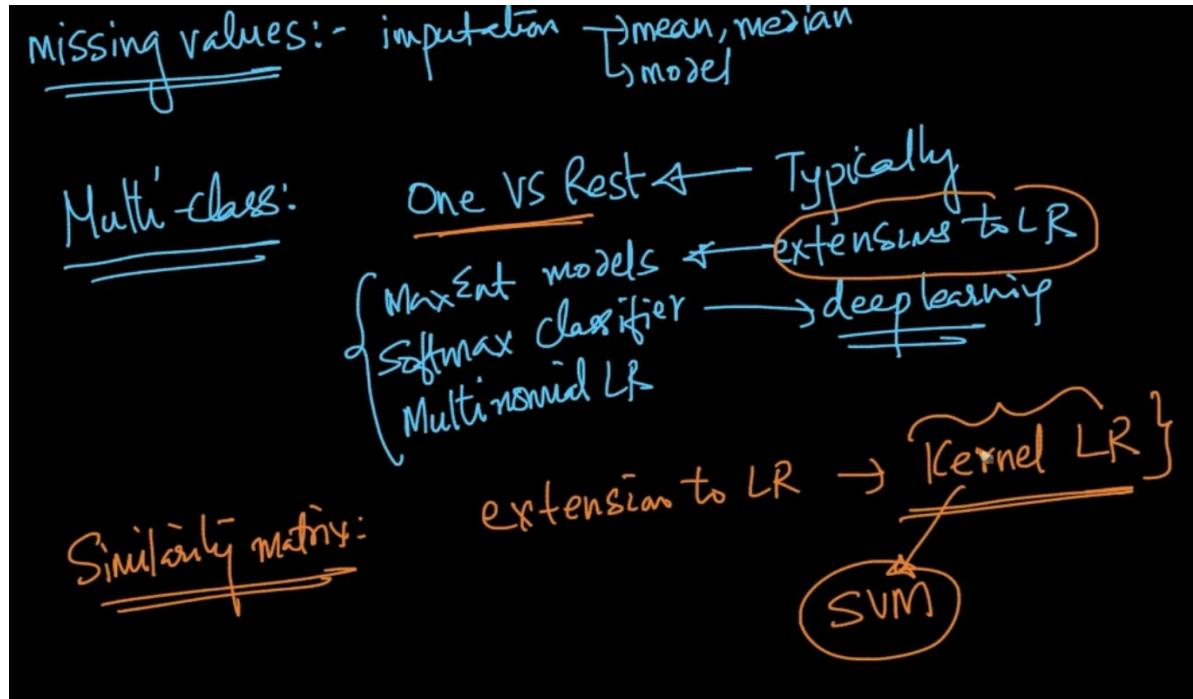
$\rightarrow D_{train}$ → w^*

$\rightarrow x_i \rightarrow w^T x_i$: dist fr π to x_i

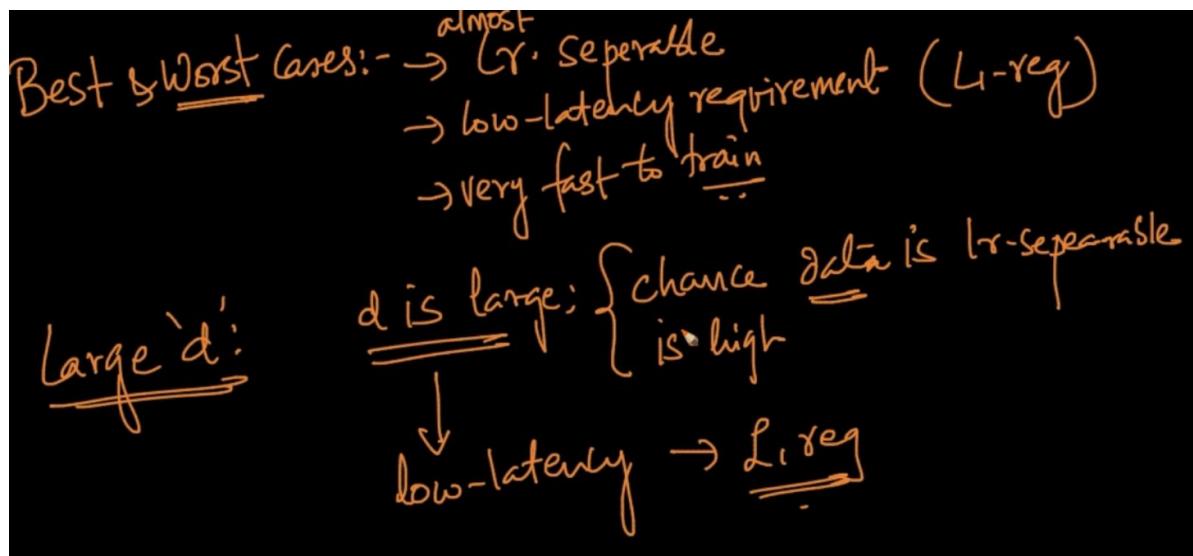
remove pts which are very far away from π from $D_{train} \rightarrow D_{train}'$

$\rightarrow D_{train}' \rightarrow \tilde{w}^*$ final soln

If we've outliers then also it has less impact because of sigmoid function. We'll get the optimal w^* . If the data point x_i has great distance from hyper-plane π then we remove that data point and then train again. Our new \bar{w}^* is our final solution

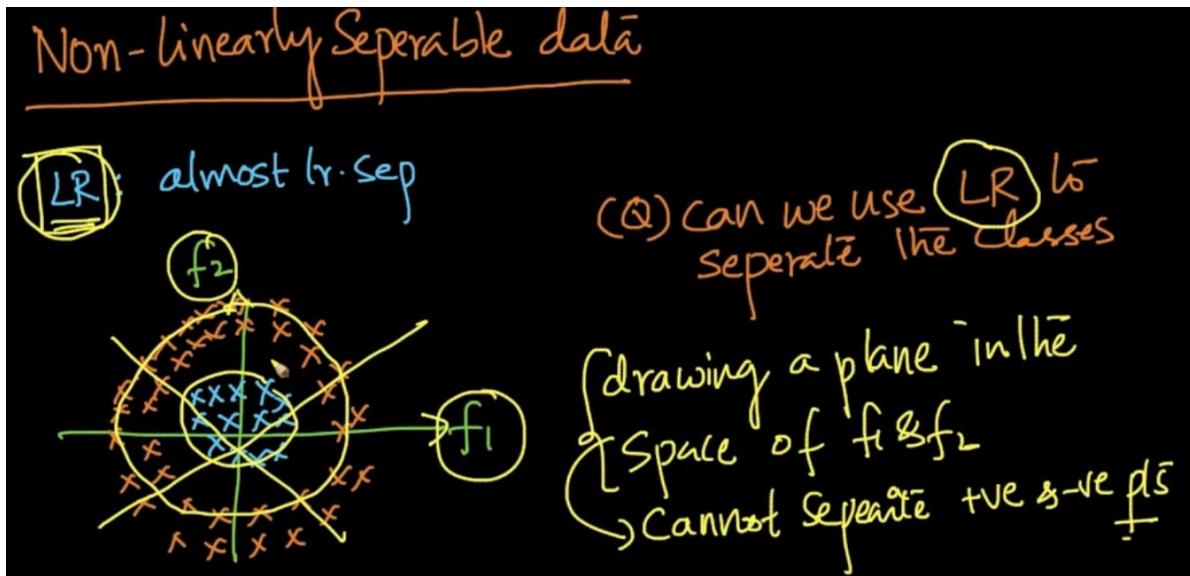


Other cases where Logistic Regression can be applied

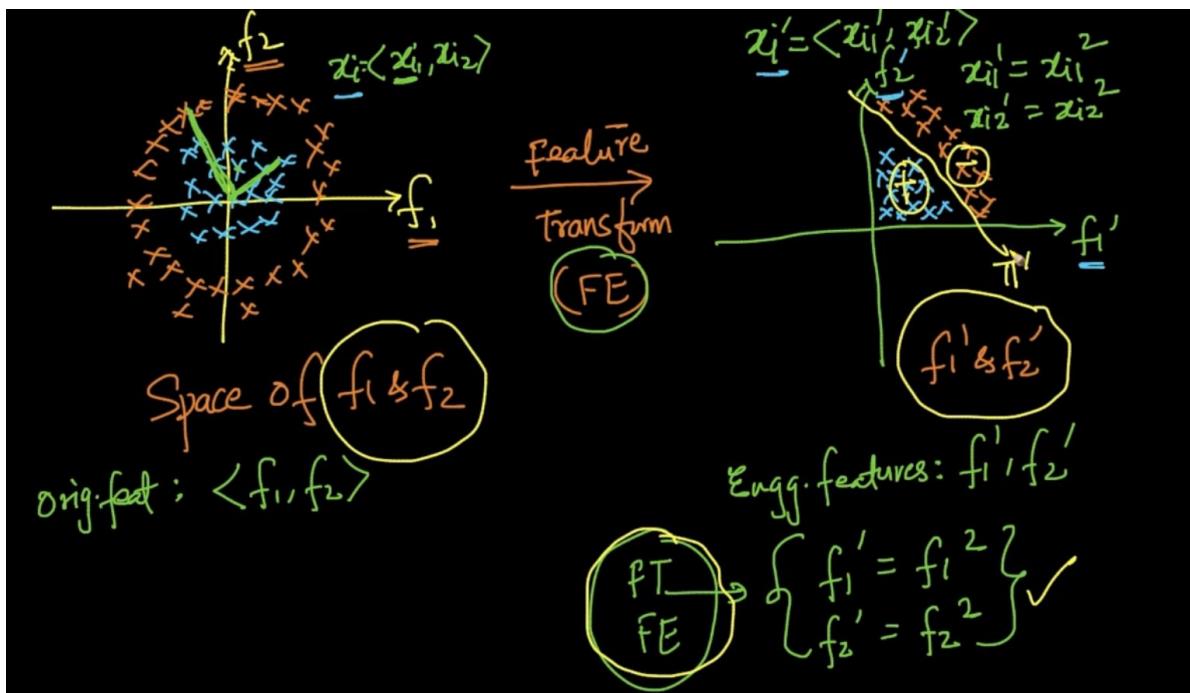


We've seen the best cases above but even if 'd' is large then chance that data is linearly separable is high and if requirement is low latency then L1-regularization

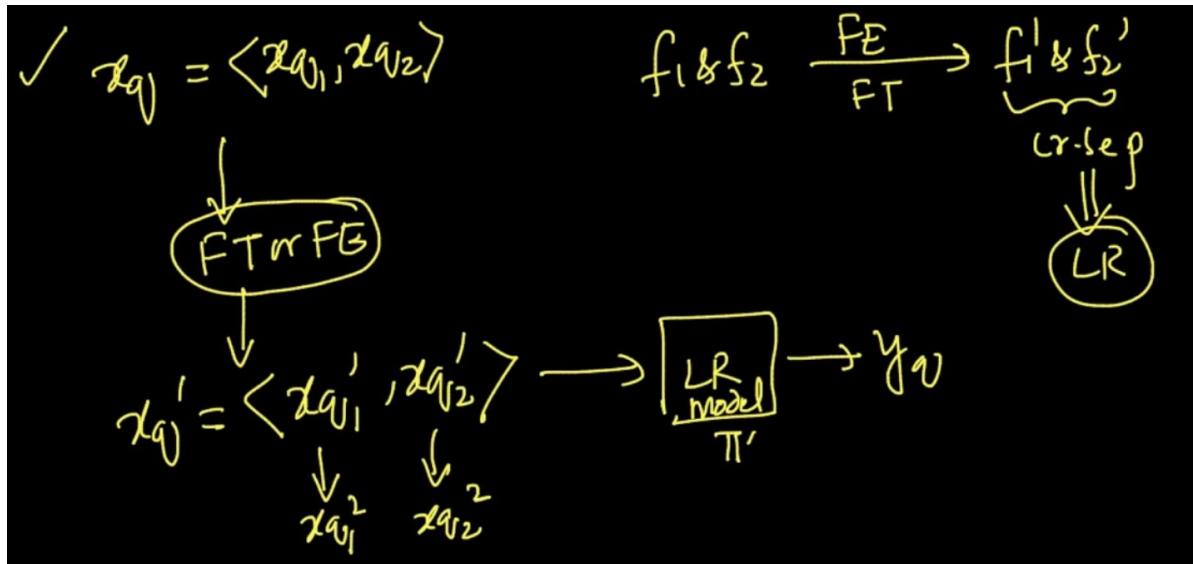
FEATURE ENGINEERING



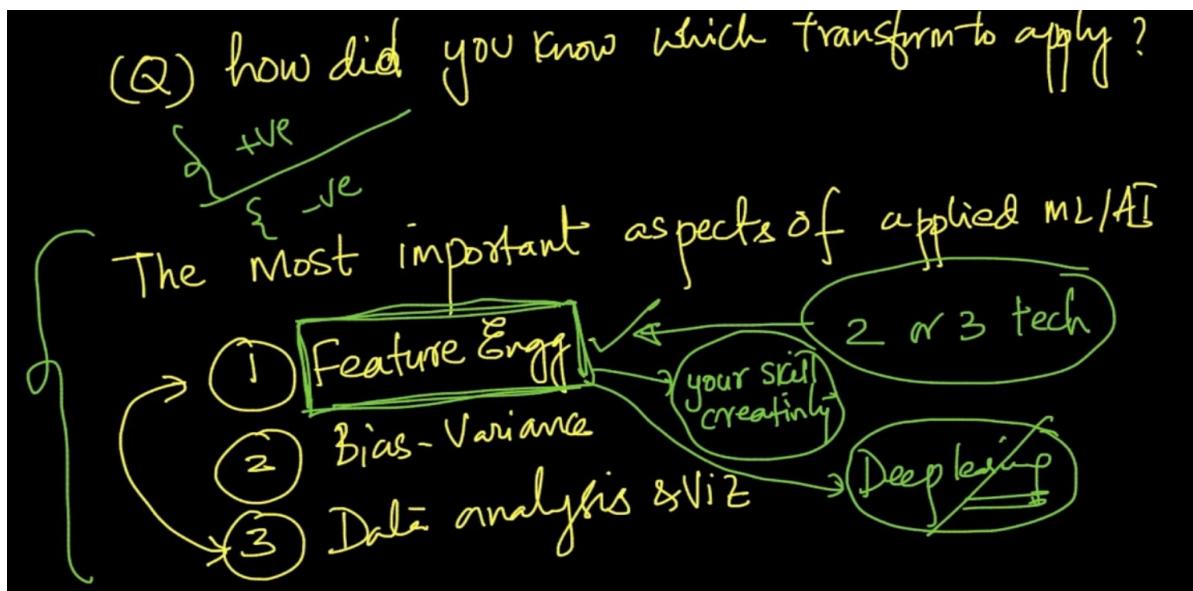
We apply Logistic to the data which is Linearly separable but in the above data Drawing a plane/line cannot separate +ve and -ve points. So, the question is can we apply LR to separate the classes?



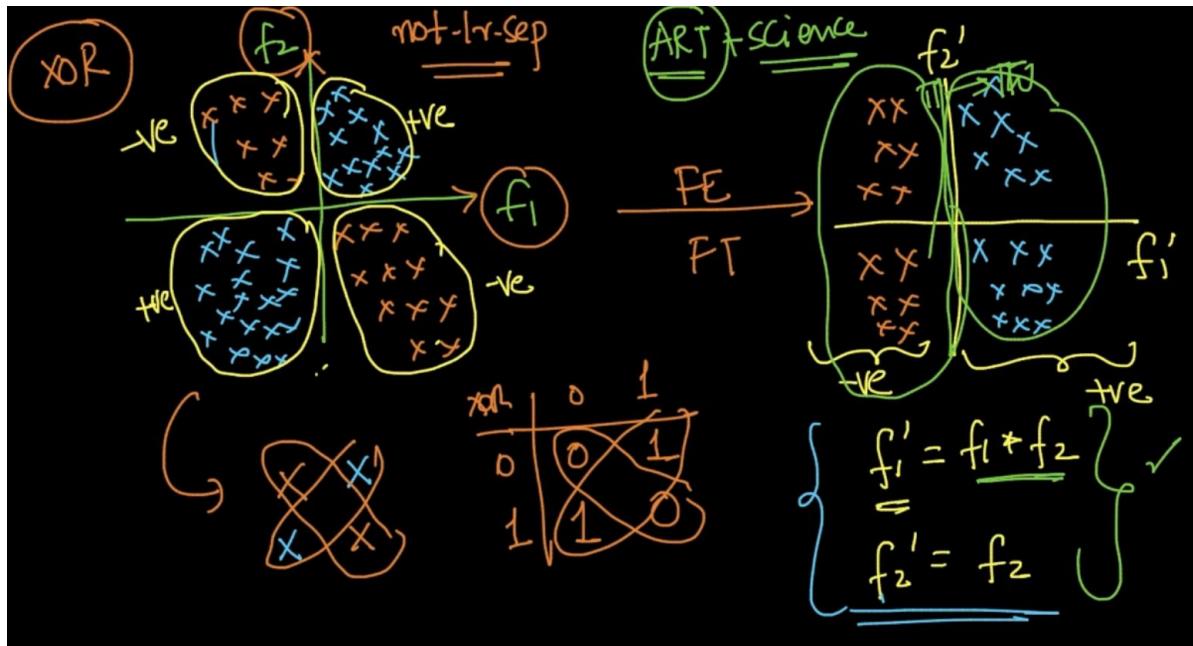
We are applying Feature Transformation by taking the features to a different space f'_1 and f'_2 . Here, $f'_1 = f_1^2$ and $f'_2 = f_2^2$. As the distance of -ve point is more than +ve points they'll be clustered far away as seen but now they can be separated with a hyper-plane π . So Logistic Regression definitely works here



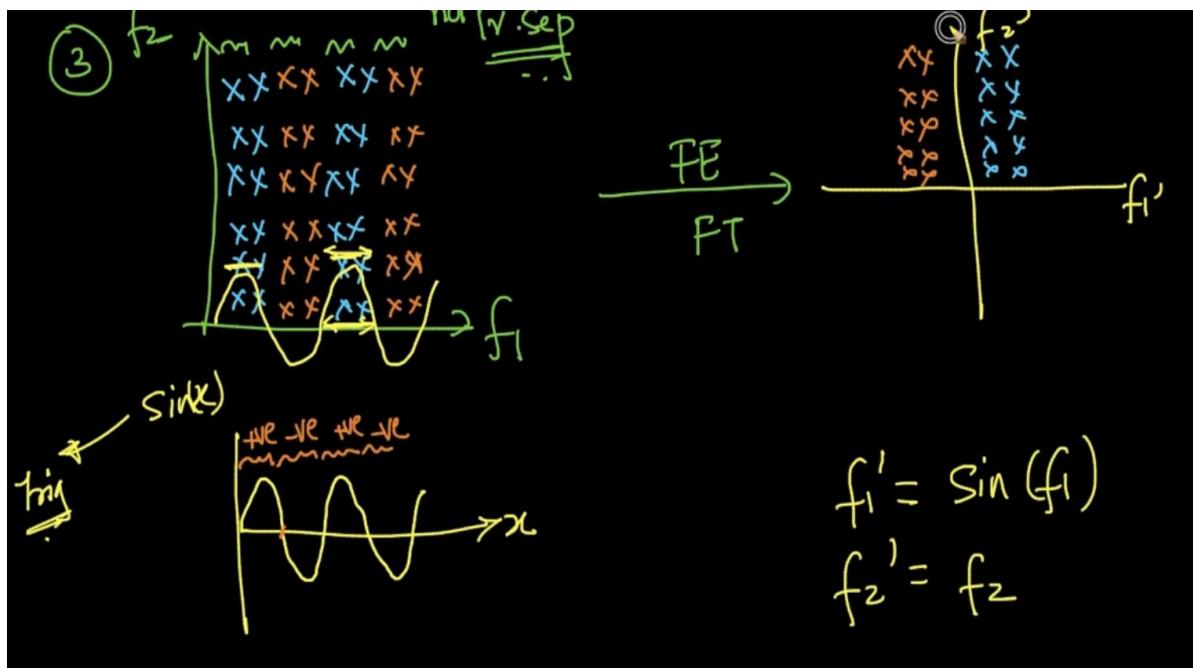
Feature Engineering is done on the features taking them to a different space which are linearly separable and Logistic Regression can be applied there



Feature Engineering is the important aspect of applied Machine Learning



This is XOR dataset. To apply feature engineering , $f_1' = f_1 * f_2$, $f_2' = f_2$. When we do this all the blue points become +ve and the orange points become -ve so can be linearly separated



Here, feature engineering is done by applying $\sin()$ function as we can see all the blue points are in positive region of sinusoidal wave and orange points in negative region . So, it can be separated after that as seen above

Typical transforms for real-valued features:

✓ ① $f_1 * f_2, f_1^2, f_2^2, f_2^3, f_1^3, f_1^2 f_2$,
polynomial features

Math:
 $\log(f_1),$
 $e^{f_1},$

② Trigonometric :- $\sin(f_1); \cos(f_1)$
 $\underline{\sin(f_1)} * \underline{\cos(f_2)} (\sin(f_1))^2$

③ Boolean features :- OR, AND, XOR

These are the typical transformation methods applied in Feature Engineering.