LOSS MINIMIZATION INTERPRETATION :-

\* From optimization problem

$$w^* = \underset{w}{\arg\min} \sum_{i=1}^{n} \log\left(1 + \exp\left(-y_i\, w^T x_i\right)\right)$$

where $z_i = y_i\, w^T x_i = y_i \cdot f(x_i)$

→ If we want "to build an Ideal optimization model

Then we try to minimize the no of incorrectly classified points."

$$w^* = \underset{w}{\arg\min} \left(\begin{array}{c} \text{no of incorrectly classified} \\ \text{points} \end{array}\right)$$

→ As the whole "classification" is all about find the correct parameter that minimizes no of incorrectly classified points.·
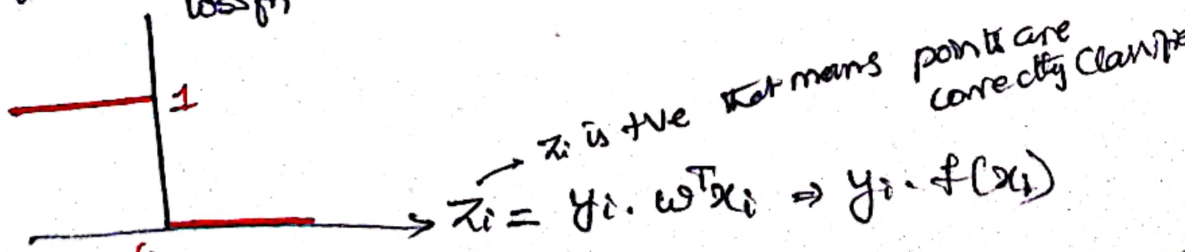
→ If we define a loss function for ex:

+1 : for incorrectly classified

0 : correctly classified.

as we always try to min the loss

So if we take ideal loss fn

loss fn



$z_i$ is +ve that means point are correctly classified

$z_i = y_i \cdot w^T x_i \Rightarrow y_i \cdot f(x_i)$

,This is also Called as 0-1 loss fn.

So if we have a loss fn like This

So     0-1 loss fn $(Z_i) = \begin{cases} 1 & \text{if } Z_i < 0 \\ 0 & \text{if } Z_i > 0 \end{cases}$

So if we have a loss fn like This we are trying to minimize $\omega$ Such that

$$\omega^* = \arg\min_{\omega} \sum_{i=1}^{n} 0\text{-}1 \text{ loss}(x_i, y_i, \omega)$$

So if given $x_i, y_i$ & $\omega$ we Can compute ~~0-1 loss~~ $Z_i$ if we Can Compute $Z_i$, given the 0-1 loss fn. Such that

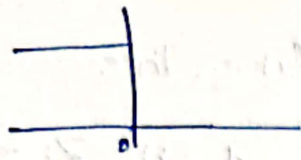$$0\text{-}1 \text{loss fn }(Z_i) = 1 \text{ if } Z_i < 0$$
$$0 \text{ if } Z_i > 0$$

→ One of the problem with 0-1 loss fn is (as we need to solve Optimization problems in ML The functions needs to be differentiable Then only We Can do some much operations on optimization Problem)

\#from basics we know The function Can differen- tiable if it is on Continous only.

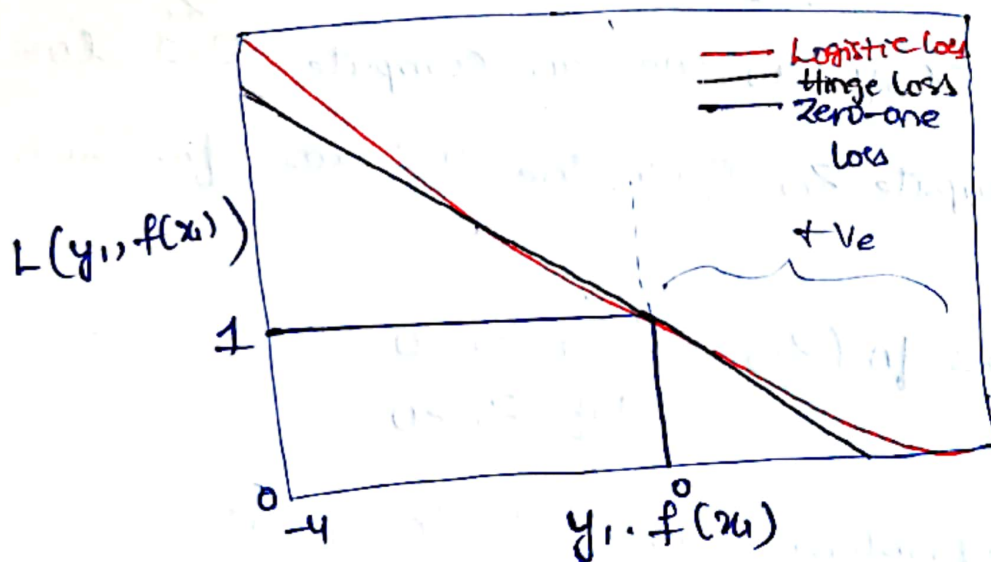\# as we See 0.1 loss function is discontinous

@ $z_i = 0$ ⌐ So it is not differentiable.

This is the problem.

\# So to avoid This problem, as it is not differentiable we will try to approximate it·

we can try many approximations. one such approxi-mation is called Logistic loss



$$L(y_i, f(x_i))$$

—— Logistic loss
—— Hinge loss
—— Zero-one loss

+ve

1

0
-4

$$y_i \cdot \overset{o}{f}(x_i)$$

→ So when we use Logistic loss as an approximation to zero-one loss we get Logistic Regression.

→ lly when we use Hinge Loss as an approximation to zero-one loss we⁻ get other algorithm called Support Vector Machines·

od loss
→ lly exponential loss ——→ Ada boost

→ `Asquared loss` $\xrightarrow{0-1 \text{ loss}}$ linear Regression

→ So we need to check every possibility of $\lambda_1$ & $\lambda_2$ values to get the best hyperparameter.

→ So If there are $m_1$ values $[<m_1>]$ in the when evaluating $\lambda_1$ : 1 hyperparameter. We need to perform $m_1$ times of searching across all the values

ly if there are 2 hyperparameters

     $\lambda_1, \lambda_2$ : 2 hyperparameters : $m_1 \times m_2$    $m^2$

ly if 3 hyperparameters

     $\lambda_1, \lambda_2, \lambda_3$ : 3 hyperparameters : $m_1 \times m_2 \times m_3$    $m^3$

       :
       :
       :

   if K hyperparameters    $m^K$

So as There are " no of hyperparameter increase The # no of times the model needs to be trained increases exponentially."
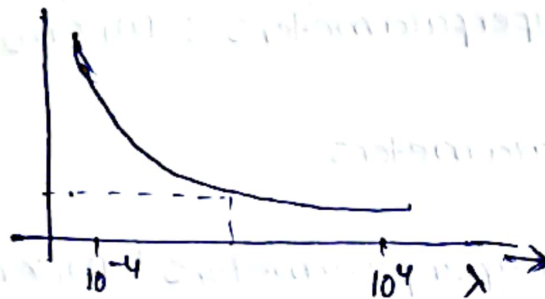
→ Even though. there are we are only 2 hyperparameters in Logistic Regressions there are Cases in deep learning where more than 2 hyperparameters would be needed That's why Gridsearch is not best technique

→ we use another technique which is as good as
  for Grid search Called Random Search

Random Search

\# In "Random Search we Randomly pick values
  in the given Interval"
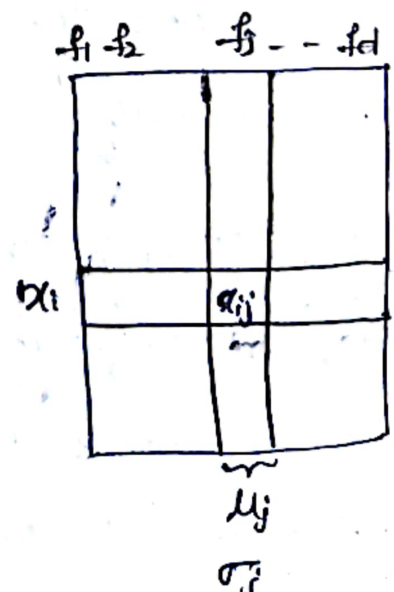
$$\lambda \in \left[10^{-4}, 10^{4}\right]$$



## 27.10 COLUMN STANDARDIZATION

→ for a dataset given with $x_i$ as a d-dimensional
  point and with d features.

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

→ In Column Standardization
  we transform



$$\boxed{x_{ij}' = \frac{x_{ij} - \mu_j}{\sigma_j}}$$

**HYPERPARAMETER SEARCH: GRID SEARCH AND RANDOM SEARCH:**

As for the logistic regression $\lambda$ is the hyper para-meter

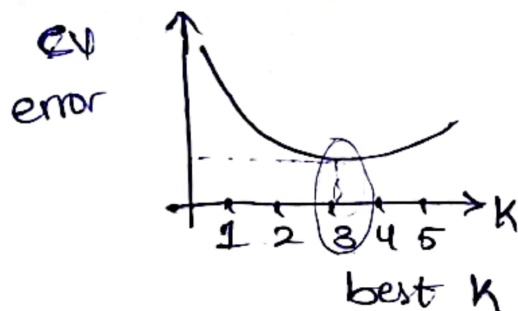and as we seen if $\lambda = 0 \Rightarrow$ over-fitting

$\lambda = \infty \Rightarrow$ under-fitting.

as in case of

K: for KNN

and $\alpha$: for NB (Laplace smoothing)

we use cross validation to find the best k



CV error

best K

lly for Logistic regression.
But there is a small problem that in b.

≠ k in KNN is an integer $\{1, 2, 3, ---N\}$

But is the $\lambda$ in logistic regression is a real number

$$\lambda \in \mathbb{R}$$

so the values that can accomodate is infinity so
to avoid such problem as it is difficult to get a
value of $\lambda$ using General techniques there is

Called Grid Search technique.

→ **Grid Search :-**

# Gridsearch is a brute force technique where people typically uses a grid or bunch of values to search for the best $\lambda$ value.

ex:- $\lambda = \begin{bmatrix} 0.001, & 0.01, & 0.1, & 1, & 10, 100, & 1000 \end{bmatrix}$

             10     10     10     10   10

(or) $\lambda = \begin{bmatrix} 1, 2, 3, 4, 5, 6, 7, 8, 9, 10 \end{bmatrix}$

(or)

$$\lambda = \begin{bmatrix} 10^{-4}, & 10^{-3}, & 10^{-2}, & 10^{-1}, & 1, 10, 10^{2}, 10^{3} \end{bmatrix}$$

# Instead of one $\lambda$ value if we considered

elastic net!

Then     $\lambda_1 \|\omega\|_1 + \lambda_2 \|\omega\|_2^2$



let

$\lambda_1 = \begin{bmatrix} 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, \\ 10^{1}, 10^{2}, 10^{3} \end{bmatrix}$

$\lambda_2 = \begin{bmatrix} 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, \\ 10^{1}, 10^{2}, 10^{3} \end{bmatrix}$

→ In KNN since we use distances b/w data points
The features could be in different scale.

so whereever we are using distance we need to
standardize the features.

→ Similarly in Logistic Regression also It is ~~made~~
mandatory to perform feature Standardization.
before training on Your data.



as Logistic regression also deals with distances

→ column / feature Standardization also called
   mean Centering and scaling.

because

$$x_{ij}' = \frac{x_{ij} - \mu_j}{\sigma_j}$$

It is like Centering the mean by substracting

→ scaling.