

## HOW TO SOLVE / <sup>N</sup>HANDLE OVERFITTING &

### UNDERFITTING ?

#### <sup>N</sup>HANDLING OVERFITTING :

- ↳ Reduce the network's capacity by removing layers or reducing the number of elements in the hidden layers.
- ↳ Apply regularization, which comes down to adding a cost to the loss function for large weights / values.
- ↳ Use Dropout layers, which will randomly remove certain features by setting them to zero.

## HANDLING UNDERFITTING:

- ↳ Get more training data.
- ↳ Increase the size or number of parameters in the model.
- ↳ Increase the complexity of the model.
- ↳ Increase the training time, until the cost function is minimised.

We know that,

→ A line eq<sup>n</sup> in 2D is  $ax_1 + bx_2 + c = 0$

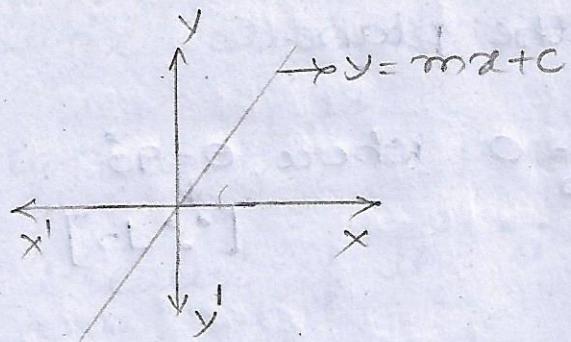
$$\Rightarrow w_1x_1 + w_2x_2 + w_0 = 0 \rightarrow \textcircled{1}$$

→ A line eq<sup>n</sup> in 4D is

$$\Rightarrow w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 = 0 \rightarrow \textcircled{2}$$

⇒  $\textcircled{2}$  can also be written as

$$\sum_{i=1}^d w_i x_i + w_0 = 0$$



From the graph  
we can say that  
a line is passing  
through  $(0,0)$

$$\Rightarrow y = mx + c$$

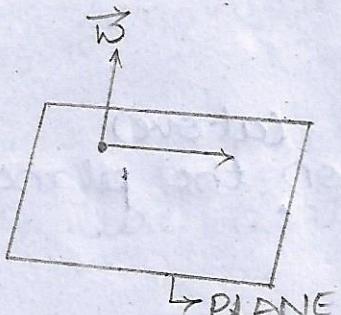
$\rightarrow * c$  tends to be zero.

$\Rightarrow w_0$  is zero.

let us consider a plane is passing  
through  $(0,0,0)$ ; then

$$\Rightarrow \textcircled{2} \Rightarrow x_2 = -\frac{w_1}{w_2} x_1 - \frac{w_3}{w_2} x_3 - \frac{w_0}{w_2} \stackrel{0}{\nearrow}$$

$$\Rightarrow x_2 = -\frac{w_1}{w_2} x_1 - \frac{w_3}{w_2} x_3$$



$w_1 x_1 + w_2 x_2 + w_3 x_3 = 0$  can  
also be written as  $\vec{w} \cdot \vec{x} = 0$

where  $\vec{w}$  is perpendicular to the plane  
and also shows the planes facing

(68)

$\vec{x}$  is any point on the plane.

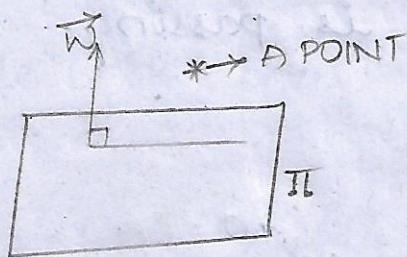
then,  $\|\mathbf{w}\| \|\mathbf{x}\| \cos \theta_{\mathbf{x}, \mathbf{w}} = 0$  where  $\theta = 90^\circ$   
 $\downarrow$   
 $[0^\circ \text{ or } 180^\circ]$

If we have a eqn of  $2x + 3y - z = 4$  then

$$\vec{\mathbf{w}} = \begin{bmatrix} 2 \\ 3 \\ -1 \end{bmatrix} \quad \text{where,}$$

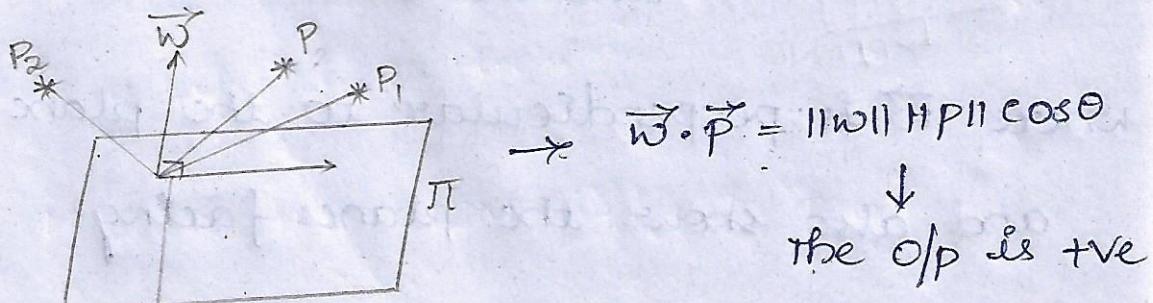
$\vec{\mathbf{w}}$  — NORMAL / NORM OF THE PLANE (or)

HYPERPLANE  $\rightarrow$  if it is in higher dimension.



From the above figure, where does the point lie i.e., above the plane or below the plane?

ANS: The point is present on the plane. (above)



$$\vec{\mathbf{w}} \cdot \vec{\mathbf{p}} = \|\mathbf{w}\| \|\mathbf{p}\| \cos \theta$$

$\downarrow$   
 The o/p is +ve

$\rightarrow$  The angle is between  $0^\circ$  &  $90^\circ$

(69)

↳ the values of  $w \cdot p$  never be "-ve".

↳ But the dot product will be "-ve"

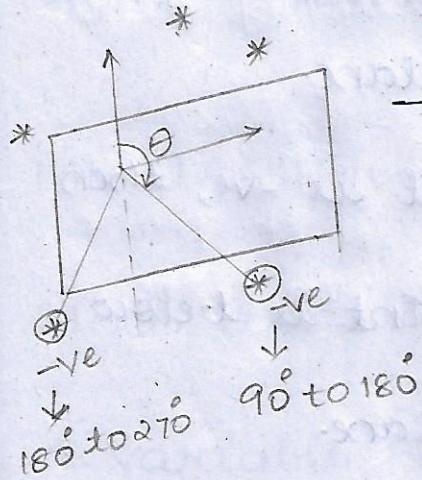
because of the " $\theta$  value"

$\Rightarrow \vec{w} \cdot \vec{p}_1 \rightarrow$  the O/P is +ve

$\Rightarrow \vec{w} \cdot \vec{p}_2 \rightarrow$  the O/P is +ve.

NOTE :

$\rightarrow$  the normal to the plane and the dot product is "+ve" if the point is above the plane.



$\rightarrow$  In this case, the value of  $\theta$  lies between  $90^\circ$  to  $180^\circ$ .

$\rightarrow$  This clearly states that the dot product is "-ve".

We use a line/hyperplane/plane where

In linear Regression & logistic Regression

there is a bit change.

LINEAR  
REGRESSION

LOGISTIC  
REGRESSION

→ Here,

$$y_{\text{pred}} = mx + c$$

where  $mx$  can be written as  $w^T x$ .

$$\rightarrow \text{error} = y_{\text{act}} - y_{\text{pred}}$$

→ Best fitting means minimizing the errors

$$\Rightarrow \min \sum e^2$$

→ Here,

$$y_{\text{pred}} = \text{sign}(mx + c)$$

where  $mx$  can be written as  $w^T x$ .

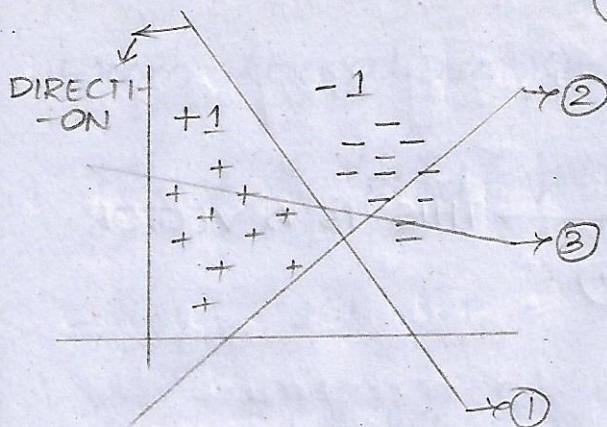
$$\rightarrow \text{sign} \leftrightarrow '+'/-'$$

→ If it is +ve, then the point is above the plane.

→ If it is -ve, then the point is below the plane.

Over here, we need to understand the misclassification i.e.,

$$\text{MISCLASSIFICATION} = y_{\text{act}} * y_{\text{pred}}$$



In this case line ① is the best separator of a hyper plane.

$\Rightarrow Y_{act} \rightarrow$  positive = +1 & negative = -1

NOTE:

↳ The separation will never be clear.

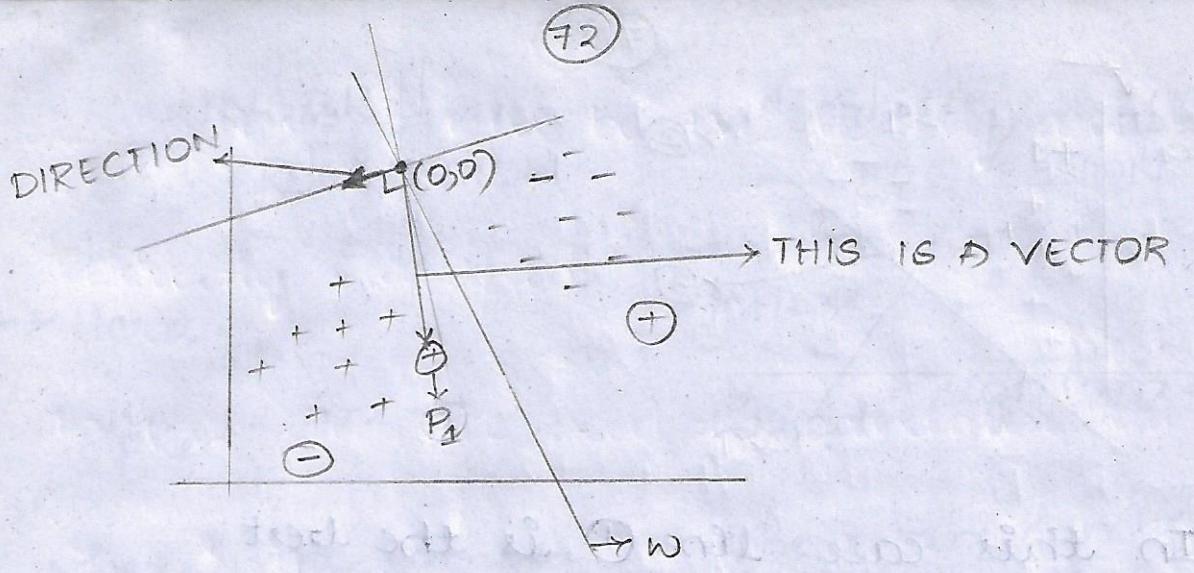
### MISCLASSIFICATION:

When all classes, groups or categories of a variable have the same error rate or

Probability of being misclassified then

it is said to be misclassification.

↳ SVM algorithm can be used for the analysis of misclassification.



In this case, there are 2 misclassifiers.

If we talk about the point  $P_1$ ,

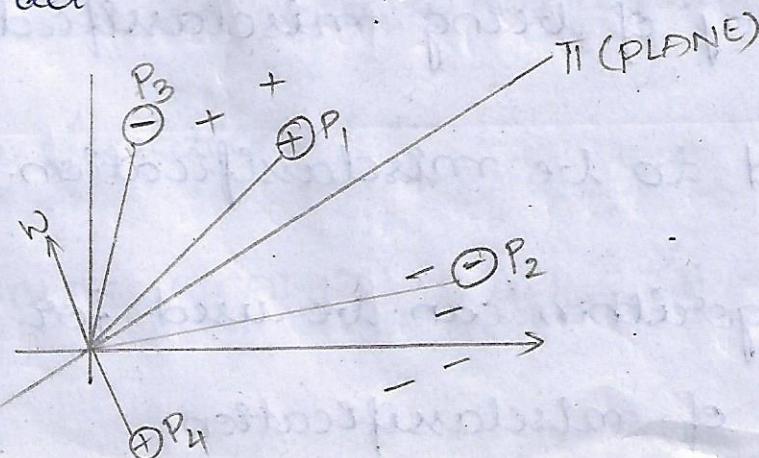
$\hookrightarrow \vec{w} \perp$  the plane are at  $90^\circ$ .

$\hookrightarrow$  For  $P_1 \rightarrow \vec{w} \cdot \vec{P}_1 \Rightarrow$  the O/p is +ve because  
 $\Rightarrow \|\vec{w}\| \cdot \|\vec{P}_1\| \cdot \cos \theta_{w, P_1}$   $\theta$  lies b/w  $0^\circ \text{ to } 90^\circ$

Here  $P_1$  is  $y_{pred}$ :

so, we can say that the Norm and the  
point has the " +ve " value i.e.,

$$\Rightarrow y_{act} = 1 \quad \& \quad y_{pred} = \vec{w} \cdot \vec{P}_1$$



(73)

From the above graph, lets check the  $y_{act} * y_{pred}$  for each point.

$$\Rightarrow P_1 = 1 * w \cdot P_1 = +ve \text{ value}$$

$$\Rightarrow P_2 = -1 * -w \cdot P_2 = +ve \text{ value.}$$

In this case,  $\cos\theta$  value ~~is~~ <sup>is</sup> +ve since  $\theta$  value is  $> 90^\circ \approx 180^\circ$ .

$$\Rightarrow P_3 = -1 * w \cdot P_3 = -ve \text{ value.}$$

In this case,  $\cos\theta$  value is +ve since  $\theta$  lies between  $0^\circ$  and  $90^\circ$  i.e.,  $> 0^\circ \approx 90^\circ$ .

$$\Rightarrow P_4 = +1 * -w \cdot P_4 \rightarrow -ve \text{ value.}$$

In this case,  $\cos\theta$  value is -ve since  $\theta$  value is  $> 180^\circ \approx 270^\circ$ .

$\therefore$  From  $P_1, P_2, P_3, P_4$  we can say that

$\hookrightarrow P_1, P_2$  are in correct positions  $\rightarrow +ve$  value

$\hookrightarrow P_3, P_4$  are misclassified  $\rightarrow -ve$  value.

$\rightarrow y_{act} * y_{pred}$  will be "tve" if they are correctly classified.

In this case, we need to maximise the

$$y_{act} * y_{pred}$$

$$\Rightarrow \max \sum y_{act} * y_{pred}$$

$$\Rightarrow m^*, c^* = \arg \max_{m, c} \sum (y_{act} * y_{pred})$$

$$\Rightarrow m^*, c^* = \arg \max_{m, c} \sum (y_{act} * \{mx + c\})$$

The above eq<sup>n</sup> is used in classification.

$$\Rightarrow m^*, c^* = \arg \max_{m, c} \sum (y_{act} * \{mx_c + c\})$$

can also be written as

$$\Rightarrow w^*, w_0^* = \arg \max_{w, w_0} \sum (y_{act} * \{w^T \cdot x\})$$

[ $\because$  From line eq<sup>n</sup>  $w_0$  tend to be '0'  $\infty$

$w_0$  is a scalar quantity]

$\rightarrow w, w_0$  will be same i.e., scalar quantity if it is in 2D.

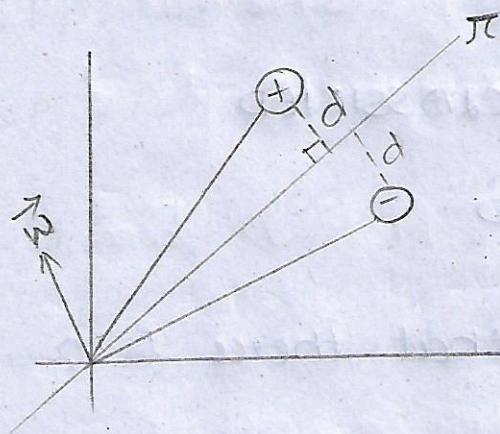
$\rightarrow$  If  $w$  - vector quantity  $\neq w_0$  - scalar quantity then it is said to be in higher dimensions.

-ons.

NOTE:

$w$  - vector quantity  $\neq w_0$  - scalar quantity

applies same in LINEAR REGRESSION.



$\rightarrow$  In this case,

$$y_{\text{pred}} = w \cdot p \text{ L.e.}$$

the "DISTANCE OF A POINT FROM THE PLANE".

$$\Rightarrow d = \frac{w \cdot p}{\|w\|}$$

where  $\|w\| = 1 \rightarrow$  since unit vector.

$\Rightarrow y_{\text{act}} * y_{\text{pred}}$  is called as the "SIGNED DISTANCE".

## SIGNED DISTANCE :

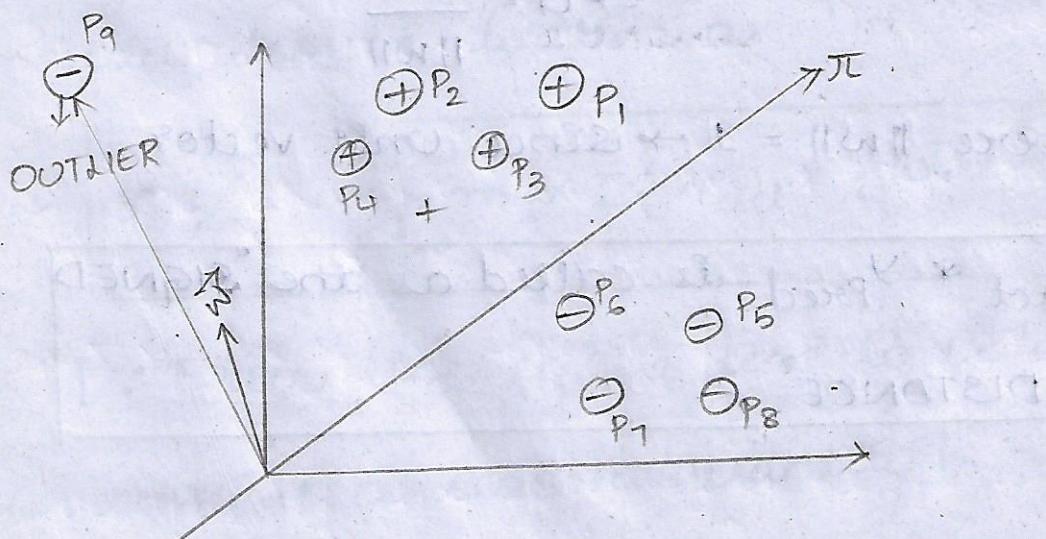
When passed the coordinates of a point in space, return the shortest distance between that point and some surface.

→ The sign of return value indicates whether the point is inside or outside of that surface (or) outside of that surface.

↳ If +ve - CORRECTLY CLASSIFIED

↳ If -ve - MISCLASSIFIED

Now let us consider that there is a misclassified point



(77)

From the above graph, lets calculate

$y_{act} * y_{pred}$

$$\Rightarrow P_1 = +1 * +2 = +2$$

$[y_{pred}]$  = HOW FAR IS THE POINT]

$$\Rightarrow P_2 = +1 * +2 = +2$$

$$\Rightarrow P_3 = +1 * +1 = +1$$

$$\Rightarrow P_4 = +1 * +1 = +1$$

$$\Rightarrow P_5 = -1 * -1 = +1$$

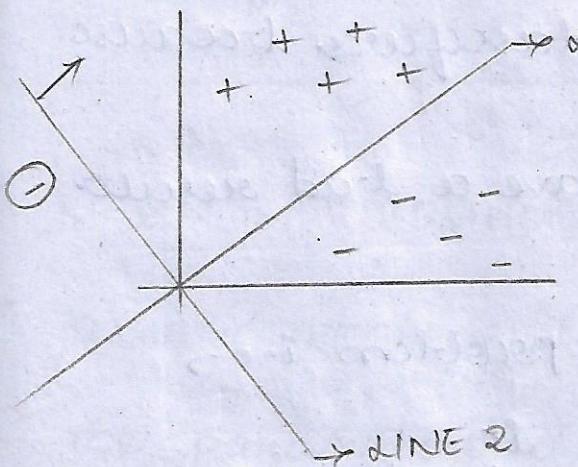
$$\Rightarrow P_6 = -1 * -1 = +1$$

$$\Rightarrow P_7 = -1 * -2 = +2$$

$$\Rightarrow P_8 = -1 * -2 = +2$$

$$\Rightarrow P_9 = -1 * 100 = -100$$

$$\Rightarrow \sum y_{act} * y_{pred} = -88$$



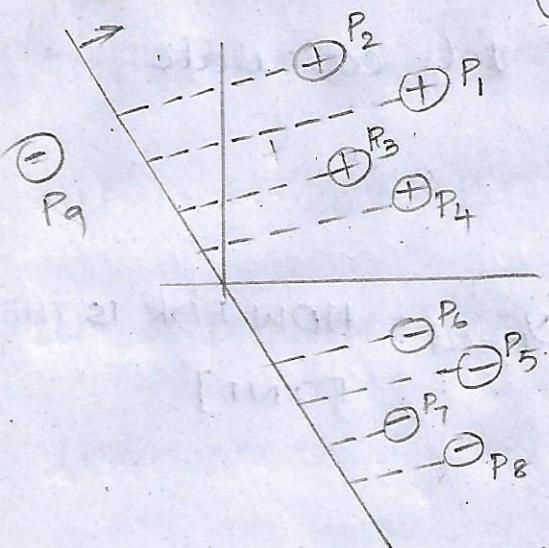
$\rightarrow$  In this case,

Misclassification

done by

$\rightarrow$  LINE 1 - 1

$\rightarrow$  LINE 2 - 4



(78)

In this case,

$$y_{\text{act}} * y_{\text{pred}}$$

$$\Rightarrow P_1 = +1 * +3 = +3$$

$$\Rightarrow P_2 = +1 * +2 = +2$$

$$\Rightarrow P_3 = +1 * +2 = +2$$

$$\Rightarrow P_4 = +1 * +1 = +1$$

$$\Rightarrow P_5 = -1 * +3 = -3$$

$$\Rightarrow P_6 = -1 * +2 = -2$$

$$\Rightarrow P_7 = -1 * +2 = -2$$

$$\Rightarrow P_8 = -1 * +4 = -4$$

$$\Rightarrow P_9 = -1 * -1 = 1$$

$\therefore \sum y_{\text{act}} * y_{\text{pred}} = -1$  which is remaining

$\therefore$  From LINE 2 we can say that we

have maximum misclassifiers because

of one outlier, we have a bad result.

$\Rightarrow \sum y_{\text{act}} * y_{\text{pred}}$  has a problem i.e.,

because of the outlier data, this

(79)

optimization eq<sup>n</sup> is giving an incorrect result.

→ The reason is, LINE 1 has 1 misclassifier,  
i.e., outlier  
and LINE 2 has 4 misclassifiers.

→ In order to solve the  $\sum Y_{act} * Y_{pred}$  eq<sup>n</sup>  
we use Gradient Descent.

→ In this case we will minimise the  
eq<sup>n</sup> to maximise state.

i.e., in linear Regression  $\rightarrow \min \sum e^2$  &  
also called as CONVEX OPTIMIZATION

PROBLEM

→ In logistic Regression  $\rightarrow \max \sum \text{signed distance}$

(MLE)

MAXIMUM LIKELIHOOD ESTIMATION OF

LOGISTIC REGRESSION:

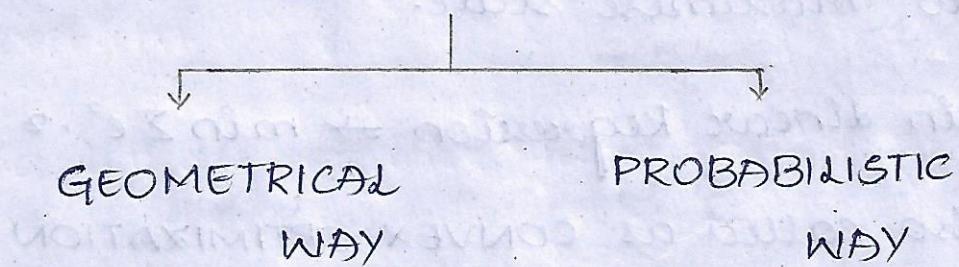
The parameters of a logistic regression model can be estimated by the

80

Probabilistic framework called maximum likelihood estimation.

→ the parameters of the model can be estimated by maximizing a likelihood function that predicts the mean of a Bernoulli distribution for each example.

## MAXIMUM LIKELIHOOD ESTIMATION (MLE)



As probabilistic way is a bit difficult we go with Geometrical way.

## WHY DO WE USE MLE?

For logistic regression, least squares estimation is not capable of producing

(81)

minimum variance unbiased estimators  
for the actual parameters.

- ↳ In its place, maximum likelihood estimation is used to solve for the parameters that best fits the data.
- ↳ In other words, it is a technique used for estimating the parameters of a given distribution, using some observed data.