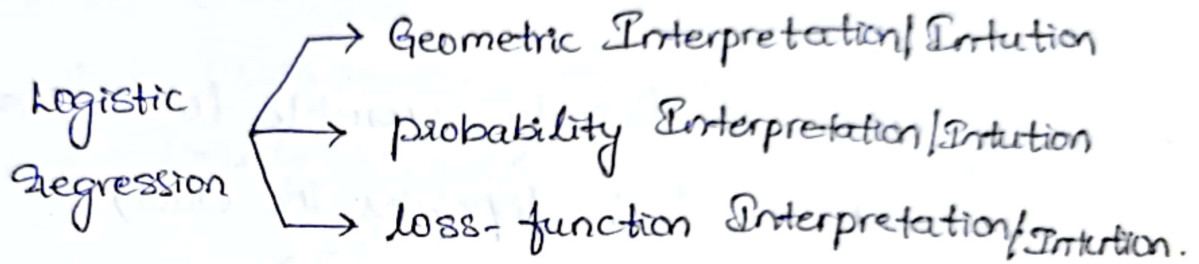# LOGISTIC REGRESSION

→ Logistic Regression is simple and elegant.

→ Even Though The name Refers to Regression It could be mainly applied for classification.

→ There are many Interpretations for logistic regression

Logistic Regression {
→ Geometric Interpretation/ Intution
→ probability Interpretation/Intution
→ loss- function Interpretation/Intution.

## * Geometric Intuition!

→ To understand this let us Consider

X: -Ve class points

X: +Ve class points.



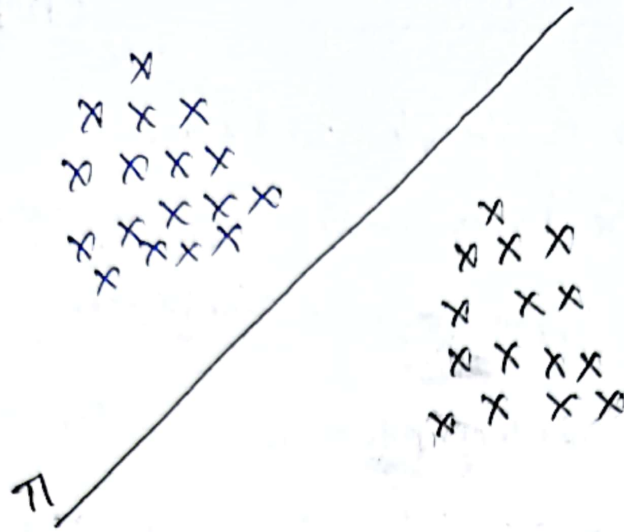Logistic Regression is a Statistical method used to predict The outcome of a dependent variable based on previous outcome.
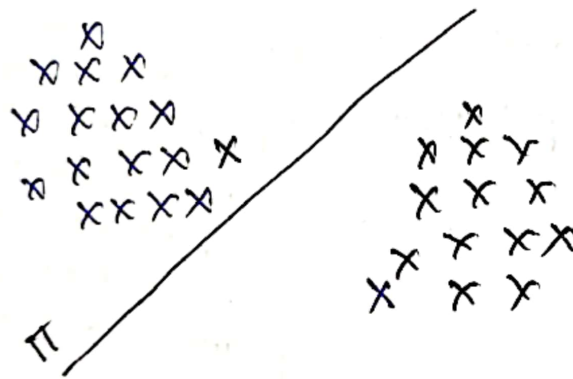
It is used to Solve binary Classification problems

It is Classification problem algorithm That predicts a binary outcome based on a Series of independent variables

→ lets use a line in 2D (or) a hyperplane:nD to seperate +ve points from The negative points.

→ If my data is linearly seperable (which means a line or a plane that seperates the data)

→ There is also a term called almost linearly seperable like the whole data got seperated except few points



→ The equation of the plane can be Represent as

$$\Pi : (\omega, b)$$

where $\omega$ - to normal to the plane $\Pi$

b: Intercept

In higher dimensions The equation of plane Represented as

$$\Pi : \omega^T x + b = 0$$

→ if the π passes Through Origin : $b=0$

$$w^T x = 0$$

π: $w^T x + b = 0$

where $x \in \mathbb{R}^d$ ; $w \in \mathbb{R}^d$ ⟹ Vectors

$b \in \mathbb{R}^1$ ⟹ Scalar

→ So- the assumption to Perform Logistic Regression
is that Class are Linearly Separable (or)
almost Linearly Separable.

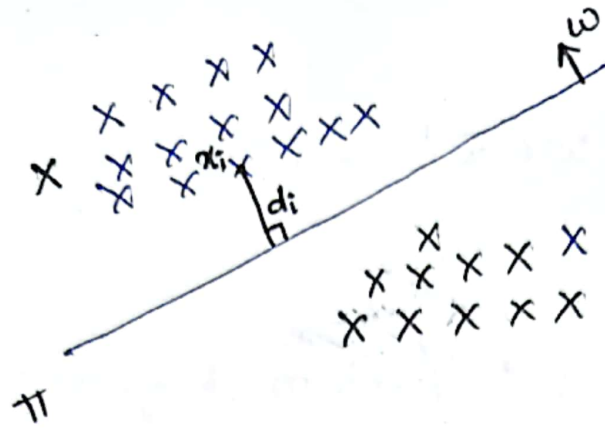

→ So from the eqn of the plane

π: $w^T x + b$

given: $Dn = \{+ve, -ve\}$

and the Task is to find $w$ and $b$

So by finding $w$ & $b$ we can find The π That
best Separates +ve points from -ve points.

**Step 1 :** finding distance of a point from the plane



where:

$x_i$ : the point

$d_i$ : distance of point from the plane.

and let $y_i$ be the class labels that Represents

$$y_i = +1 \quad : \text{+ve points}$$
$$\quad\quad -1 \quad : \text{-ve points}$$

$$y_i \in \{-1, +1\}$$

**Note :** ✗ Not +1 as +ve pts and 0 as -ve points ✗

✱ do from the basics of linear algebra distance $d_i$

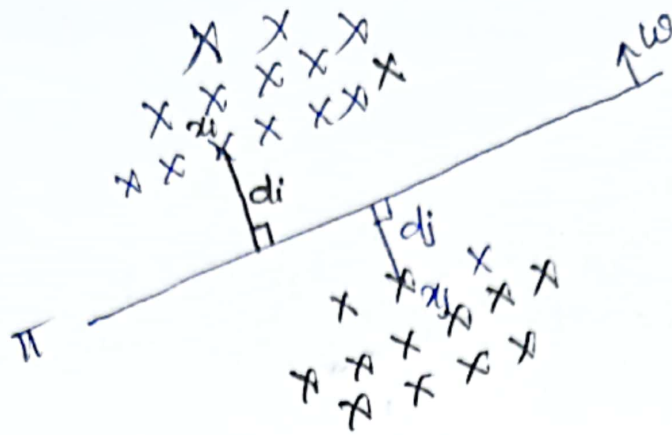$$d_i = \frac{w^T x_i}{\|w\|} \quad ; \quad$$ where $w$ is the normal to the plane

and let $\|w\|$ is a unit vector

So $\|w\| = 1$

**Step 2 :** finding other side distance of a point from the plane

→ Similarly we will find The distance $d_j$ of other side of the plane

→ as $x_i$ is towards the direction of normal $\omega$
we can say

$$d_i = \omega^T x_i > 0$$

$$d_j = \omega^T x_j < 0$$

and lly as $x_j$ is towards opp side of normal we
can say

$$d_j = \omega^T x_j < 0$$

→ And In Logistic Regression the Decision surface
will be that line/plane.

→ So our Classifier Classifies

if $\omega^T x_i > 0$ then $y_i = +1$

e₁ lly
if $\omega^T x_i < 0$ then $y_i = -1$

#Note!-
usually we need to consider plane eqn as $\omega^T x_i + b$
but here we are assuming the plane passed
Through the Origin.

Step3: Classification of points

Case 1:

let us consider the +ve point

if $y_i = +1$

and $w^T x_i > 0$ ⟹ That means if classifier is saying its +ve point

so $y_i * w^T x_i > 0$

that means $w$ is correctly classifying the point

Case 2:

let assume $y_i = -1$ : -ve point

and assume $w^T x_i < 0$ ⟹ That means LR is concluding that $x_i$ is -ve point.

⟶ Now if we take

$$\underset{-ve}{y_i} * \underset{-ve}{w^T x_i} > 0$$

So we conclude for both +ve & -ve points if $y_i * w^T x_i > 0$ ⟹ That means the LR model is correctly classifying the point
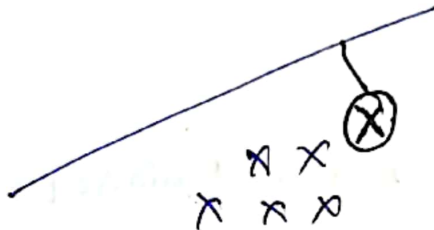
__Case 3__

if $y_i = +1$ (+ve points)

a if $\omega^T x_i < 0$ ⇒ That means LR is saying $x_i$ is −ve class.

so $y_i * \omega^T x_i < 0$

That implies

$y_i = +1$

but we got LR : −1 that means misclassified



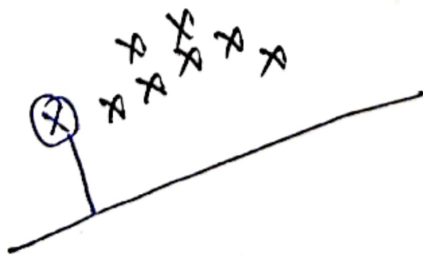__Case 4__

if $y_i = -1$

if $\omega^T x_i > 0$ ⇒ That means LR is saying $x_i$ is +ve class

$y_i * \omega^T x_i < 0$ ⇒ mis classified



* So @ The end for the classifier to be Very good
① either it has to Consider min no of misclassifi-cations

② or max no of correctly classified points.

So we have need as many points as possible to have

$$y_i \, w^T x_i > 0$$

Steps :- Mathematical Optimization.

So as we seen from previous step we need to get have

as many possible to have $y_i \, w^T x_i > 0$

so that means we need to maximize it.

$$\max_{w} \sum_{i=1}^{n} y_i \, w^T x_i$$

This means
both $x_i$ & $y_i$ are fixed in our Dataset the only
variable is $w$

↳ This Represent the term or variable That we need
to change / vary to g maximize it.

→ In python argmax is used to maximize fn

$$w^* = \underset{w}{argmax}\left(\sum_{i=1}^{n} y_i \, w^T x_i\right)$$

This means              Variable

Optimal $w$

And This $w^*$ we need to find in This math Optimiz
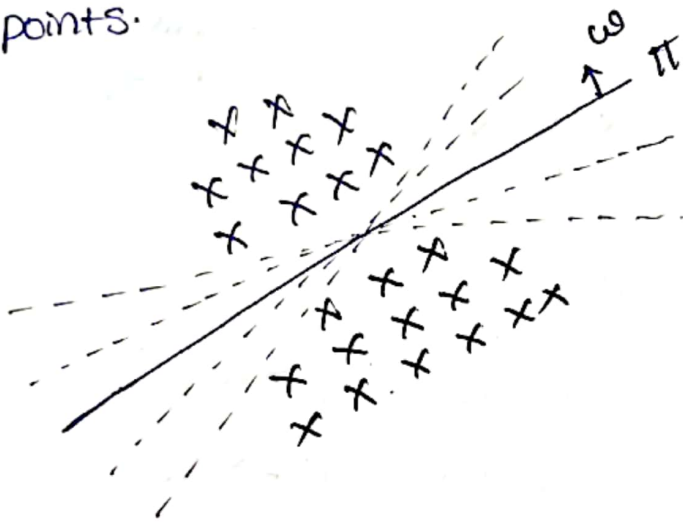
## 27.2 SIGMOID FUNCTION : SQUASHING :-

as we seen math Optimization

$$w^* = \underset{w}{argmax} \sum_{i=1}^{n} y_i \, w^T x_i$$

$\rightarrow$ optimal $w$

This means as to seperate +ve points from -ve points

There can be n no of planes could be possible

we need to find best $w_i$ corresponding to the

plane $\pi_i$ that best seperates both +ve and

-ve points.



$$\pi_i : w_i$$

$$* \quad w^* = \underset{w}{argmax} \sum_{i=1}^{n} \underbrace{y_i \, w^T x_i}$$

This portion will be called as
signed distance

Signed distance because $w^T x_i$ dist from $x_i$ to $\pi$

and

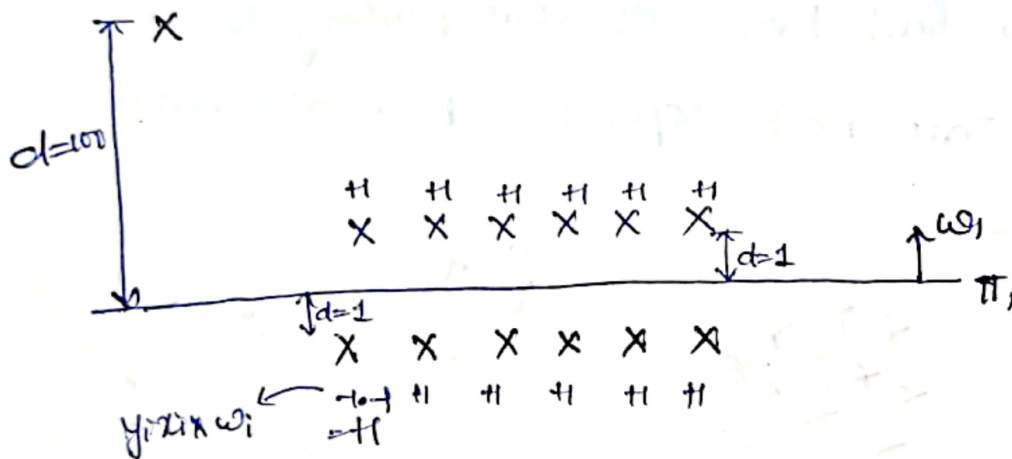$$y_i \, \omega^T x_i \; : \; +ve \Rightarrow \pi \text{ as defined by } \omega \text{ correctly}$$
$$\text{classifies } x_i.$$

$$: -ve \Rightarrow \text{incorrectly classifies } x_i.$$

* There are Cases where

* Ex:- To show where math optimization does not work well

let us assume we have +ve and -ve points



$y_i x_i \omega_i$ ←── $\boxed{-1}$ +1

Case 1:- $\pi_?$ is my Seperation

$$\sum_{i=1}^{n} y_i \, \omega_i^T x_i = \underbrace{1+1+1+1+1}_{\gg ve}$$
$$+ \underbrace{1+1+1+1+1}_{\to -ve}$$
$$-100$$

$$= -90$$

se 2

+1 ⊗─┤ +1

+1

e
pont(-1)
eside
1)
1=+1

+1 +2 +3 +4 +5
1  1  1  1  1  1
⊗ ⊗ ⊗ ⊗ ⊗ ⊗

1  ⊗ ⊗ ⊗ ⊗ ⊗ X
-1 -2 -3 -4 -5

CO₂

π2

+1 +2 +3 +4 +5 → +Ve

-1 -2 -3 -4 -5 → -Ve

+1 ⌣ outlier

⇒ +1

‿        ‿‿‿‿‿
-Ve         +Ve

l as objective is to find CO that maximizes sum of
ned distances.

$+1 > -90$

Can say π2 as our classifier for this example.

but as we see Intuitively for π1 out of 11 points
points are correctly classified whereas if we see π2
out of 10 points are correctly classified. So if we
ink logically π1 would be the best compared to π2
t as we are considering sum of signed distances
ecause of one outlier the model is choosing π2
instead of π1,

we can say sum of signed distances is easily
                        Prone to outliers.

Scanned with OKEN Scanner

→ So one single/extreme outlier point

→ So to overcome this we use a technique called Squashing.

## Squashing

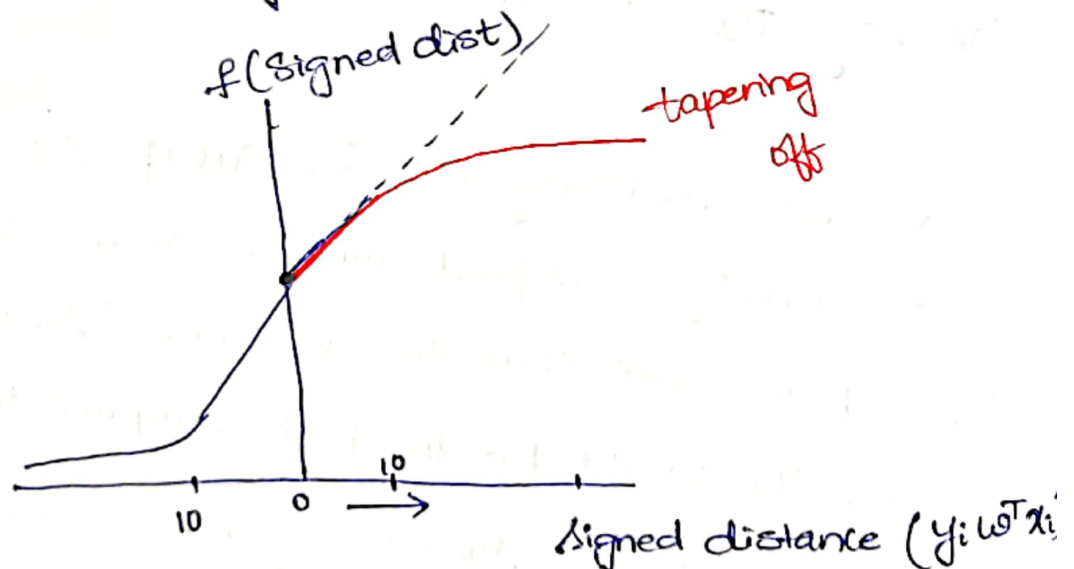* The idea is instead of using signed distance,

    We need consider

      if signed distance is small use as it is

      if signed distance is large : make it a sm
                       value.

* To do squashing



So we will design or create a function that ever the signed distance gets huge it will be tappered a particular value.
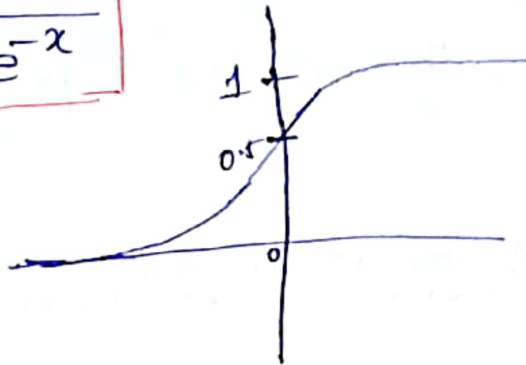
$$\to \underset{\omega}{\arg\max} \sum_{i=1}^{n} y_i \omega^T x_i \implies \underset{\omega}{\arg\max} \sum_{i=1}^{n} f(\cdot y_i \omega^T x_i)$$

and that function is called sigmoid function. $\sigma(x)$

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

The Graph of $\frac{1}{1 + e^{-x}}$



→ So by this when our data point is very far away from the plane then $\omega^T x_i$ is very large. by this fn we get $P(y_i = 1) = 1$

→ Summary

* As we are optimizing max. sum of signed distance but it got problem with Outlier

↓

So we find function called sigma $\sigma(x)$ which has properties
1. tapering behavioury linear, a probabilistic
   interpretation.
2.

↓

At last we get max sum of transformed signed distance.

$$\rightarrow \boxed{w^* = \underset{w}{\arg\max} \sum_{i=1}^{n} \sigma(y_i\, w^T x_i)}$$

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

$$\boxed{w^* = \underset{w}{\arg\max} \sum_{i=1}^{n} \frac{1}{1 + \exp(-y_i\, w^T x_i)}}$$

and this is less impacted by outliers.

## 27.3 Mathematical formulation of Objective fn :-

from Optimization problem

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

We can even Simply Optimize this further

## from Monotonic fn's :-

A fn $g(x)$ is said to be monotonic fn

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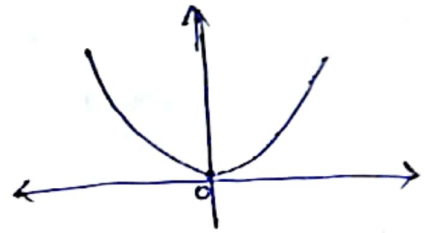
is said to be monotonically Increasing fn.

## Note

Ex of Optimization problem:

let us consider an example

$$x^* = \arg\min_x x^2$$

by using maxima and minima concepts.

we can say



so $\Rightarrow$ if $f(x) = x^2$

$$x^* = \arg\min_x x^2 \Rightarrow \arg\min_x f(x) = 0$$

## Note

If $g(x)$ is a monotonically fn.

then

$$\arg\min_x f(x) = \arg\min_x g(f(x))$$

$$\arg\max_x f(x) = \arg\max_x g(f(x))$$

Ex:- let $g(x) = \log(x)$

and let $f(x) = x^2$

$$x^* = \arg\min_x f(x) = \arg\min_x x^2 = 0 \;—\; =$$

$$x' = \arg\min \log(x^2) = 0$$

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

(α)

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

$$w^* = \arg\max_w \sum_{i=1}^{n} \log\left(\sigma(y_i w^T x_i)\right)$$

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

as $\log(1/x)$

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

* The point that minimizes a function is the same point that maximizes $-f(x)$.

$$\Rightarrow \arg\max_x f(x) = \arg\min_x -f(x) \quad ||$$

lly $\arg\max_x -f(x) = \arg\min_x f(x)$

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

and this the optimization problem of Logistic Regression.

where $y_i : +1 \; (or) -1$

Note:-

* AS we See @at last we Could get a A Optimization problem by sum of signed distance only Just by making up with log and exp

$$\arg\min_{\omega} \sum_{i=1}^{n} \log \left(1+ \exp\left(-y_i\omega^T x_i\right)\right)$$

$$\arg\min_{\omega} \sum_{i=1}^{n} \left(-y\omega^T x_i\right)$$

$$\arg\max_{\omega} \sum_{i=1}^{n} \left(y\omega^T x_i\right) \longleftarrow \text{which is sum of signed distance only}$$

* The above Optimization is of geometric Interpretation. Similarly we can also derive this Using Probabilistic method.

$$w^* = \arg\min_{\omega} \sum_{i=1}^{n} -y_i \log P_i - (1-y_i) \log (1-P_i)$$

$$\text{where } P_i = \sigma(\omega^T x_i)$$

$\rightarrow +1 \text{ or } 0$

## 2.4 WEIGHT VECTOR :-

→ As we see from optimization . problem

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

↳ and This $w^*$ is called <u>weight Vector</u> which is the best $w^T$ value.

→ and $w^*$ is a d-dimensional Vector

$$w^* = <w_1, w_2, w_3, ---w_d>$$

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

→ If we have a weighted Vector with d dimens-ions and if There are d features, for every feature there is a corresponding weight associated with it.

$$\omega = <w_1, w_2, \boxed{w_3} --- w_d>$$
$$f_1 \quad f_2 \quad \textcircled{f_3} \, - \, - \, -f_d$$

→ If we perform decision

like given a Query point $x_q$ → (we need to find Class label) $y_q$

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

where as probabilistic Interpretation. need sigmoid in

$$\sigma(\omega^T x_q) = P(y_q = +1)$$

To decide the given query point $x_q \to y_q$.

## Interpretation of $\underline{\omega}$ :-

Case 1

If $\omega_i = +ve$ for a given feature $f_i$, and if $i^{th}$ component of given Query point $x_{qi} \uparrow$ then

$$x_{qi} \uparrow \implies (\omega_i x_{qi}) \uparrow$$

$$\implies \sum_{i=1}^{d} (\omega_i x_{qi}) \uparrow$$

$$\implies \sigma(\omega^T x_q) \uparrow$$

$$P(y_q = +1) \uparrow$$

Case 2:-

if $\omega_i = -ve$

then as

$$x_{qi} \uparrow \implies (\omega_i x_{qi}) \downarrow$$

$$\implies \left( \sum_{i=1}^{a} \omega_i x_{qi} \right) \downarrow$$

$$\implies \sigma(\omega^T x_q) \downarrow$$

$$\implies P(y_q = -1) \uparrow$$

whereas probabilistic Interpretation. need sigmoid in

$$\sigma(w^T x_q) = P(y_q = +1)$$

To decide the given query point $x_q \longrightarrow y_q$.

### Interpretation of $\underline{w}$ :-

Case 1

If $w_i = +ve$ for a given feature $f_i$, and if $i^{th}$ component of given Query point $x_{qi} \uparrow$ then

$$x_{qi} \uparrow \implies (w_i x_{qi}) \uparrow$$
$$\implies \sum_{i=1}^{d} (w_i x_{qi}) \uparrow$$
$$\implies \sigma(w^T x_q) \uparrow$$
$$P(y_q = +1) \uparrow$$

Case 2 :-

if $w_i = -ve$

then as

$$x_{qi} \uparrow \implies (w_i x_{qi}) \downarrow$$
$$\implies \left(\sum_{i=1}^{a} w_i x_{qi}\right) \downarrow$$
$$\implies \sigma(w^T x_q) \downarrow$$
$$\implies P(y_q = +1) \downarrow \qquad \implies P(y_q = -1) \uparrow$$