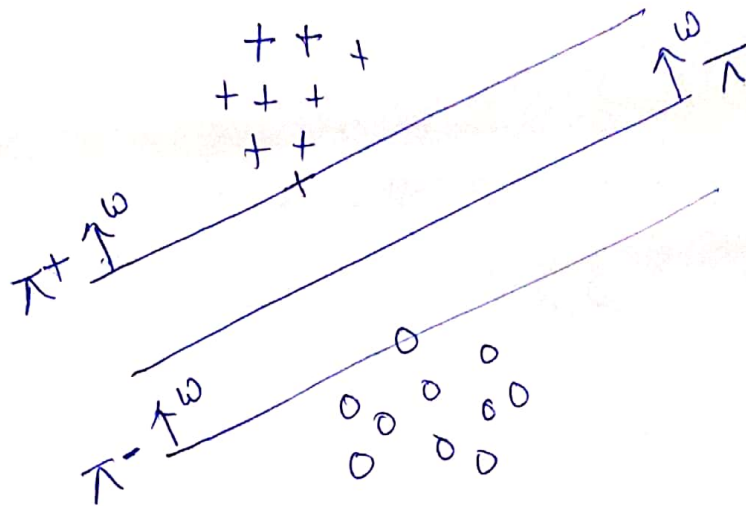→ "Alternate Mathematical formulation of SVM"

We already know from geometric intuition of SVM, that, we want to find a hyperplane 'π' that does ["margin maximization"]



Suppose if 'π' is our best hyperplane, let's write hyperplane. 'π' as $w^T x + b$

$$\pi: w^T x + b = 0$$

Here "w" is $\perp$ to the hyperplane"

One thing that we quickly realize is, if $\pi^+$ is my positive hyperplane & $\pi^-$ is my negative hyperplane & since 'w' is $\perp$ to 'π' & $\pi^+, \pi^-$ & π are parallel to each other then 'w' will also be $\perp$ to { $\pi^+$ & $\pi^-$ }

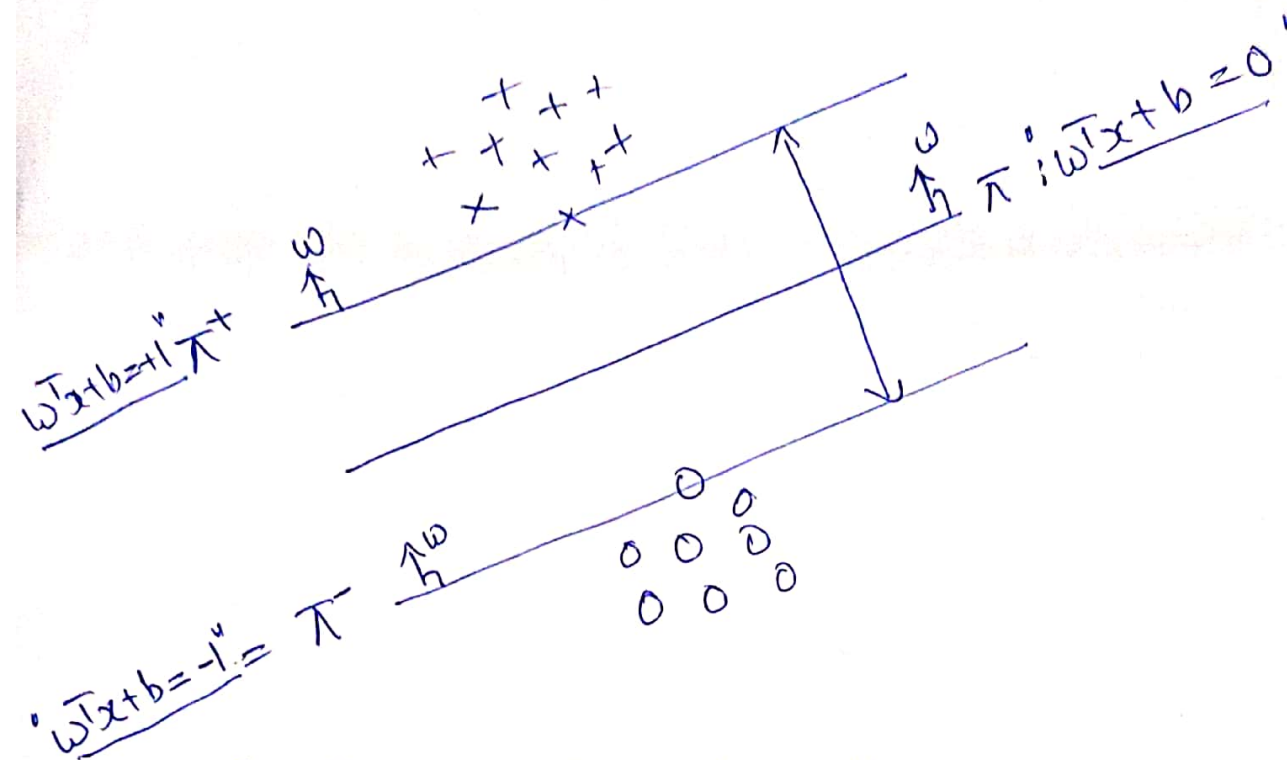let's say "$\pi^+$" has this form.

$$\pi^+: w^T x + b = 1$$

& $\pi^-$ has this form

$$\pi^-: w^T x + b = -1$$

Note :> "$w^T w \neq 1$"

['w' is not a 'unit vector"]

let's assume that "$w$" is some vector & not necessarily a "unit vector". & it is $\perp$ to "$\pi$", "$\pi^+$" & "$\pi^-$".

$\pi^+ : w^T x + b = +1$

$\pi : w^T x + b = 0$

$\pi^- : w^T x + b = -1$

In SVM we care about this _margin_.

The margin is $d = \left\{ \dfrac{2}{||w||} \right\}$ & "$w$" is some vector which $\left\{ \text{is} \perp \text{to } \text{"}\pi\text{", "}\pi^+\text{" } \& \text{ "}\pi^-\text{"} \right\}$

We want to find a "$w^*$" & "$b^*$" in such a way that the _margin is maximized_.

$$\left\{ (w^*, b^*) = \arg\max_{w,b} \dfrac{2}{||w||} \right\}$$

$\Rightarrow$ Constraint

_Such that_

All the +ive points are on the side of "$\pi^+$" & all the -ive points are on the side of "$\pi^-$"

This is what we want to find

There are some constraint.

$$(w^*, b^*) = \underset{w,b}{\arg\max} \; \frac{2}{||w||} \longrightarrow \text{margin.}$$

such that any of the +ve point is on the side of "$\pi^+$" & all of the -ve points are on the side of "$\pi^-$"

$y_i(w^T x_i + b) > 1$ $\longleftarrow$

$\pi^+ = w^T x + b = +1$

$\pi = w^T x + b = 0$

$\pi^- = w^T x + b = -1$

$y_i(w^T x_i + b) = 1$

'Support vectors'

$y_i(w^T x_i + b) = 1$

for +ve points
class label $y_i = +1$
& for -ve points
class label $y_i = -1$

$y_i(w^T x_i + b) > 1$

Now the constraint that we have is +
& our optimization problem will eventually look like.

$$(w^*, b^*) = \underset{w,b}{\arg\max} \; \frac{2}{||w||} \longrightarrow \text{margin}$$

such that $y_i(w^T x_i + b) \geq 1$ for all $x_i$'s

Note :- It is exactly equal to '1' for support vectors &
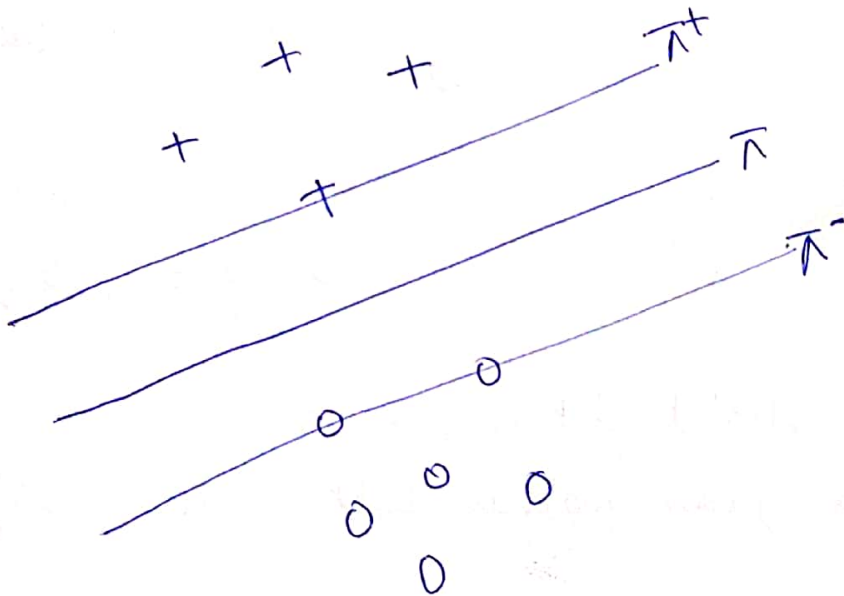for Non-support vectors, it is greater than '1'

So, literally it is not one constraint, it is 'n' constraints becoz we have 'n' points in our training data.

So, the final problem that we have is :-

$$w^*, b^* = \underset{w, b}{\arg\max} \frac{2}{||w||}$$

$$\text{such that } \forall i, \; y_i(\overline{w}^T x_i + b) \geq 1$$

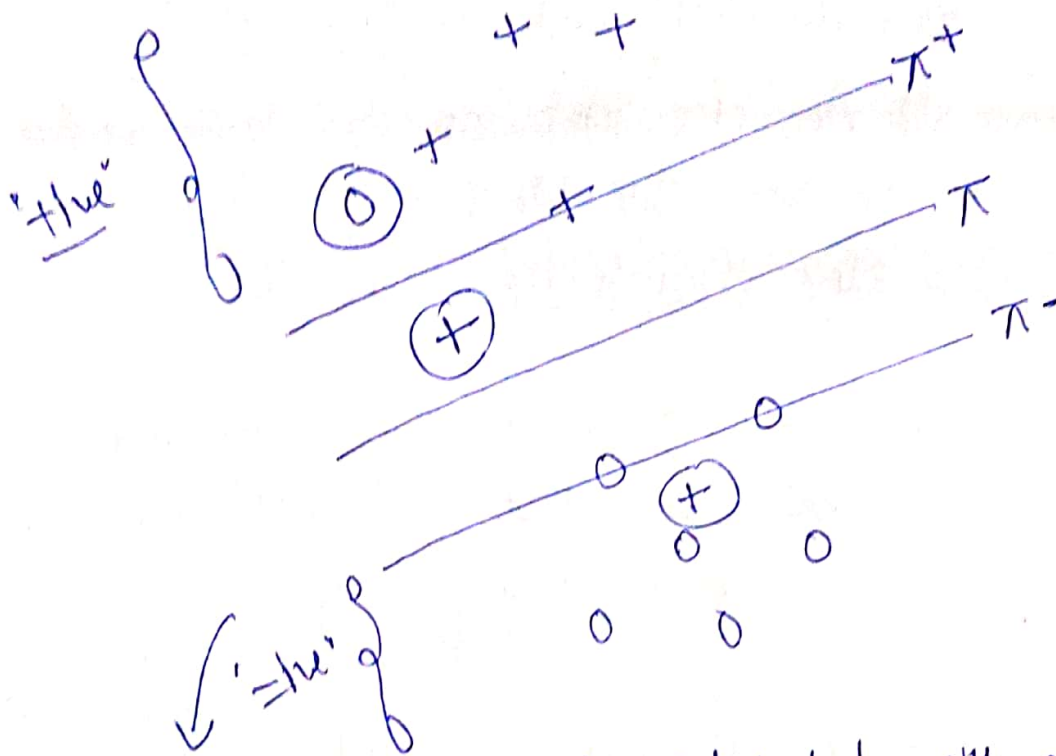↳ This is Constraint optimization problem of SVM

There is one fundamental problem with this formulation.
What if we have our data like shown below :-



↳ This woofee when our data is linearly separable. becoz the constraint that we have is, that every +ive point should be in its upper region of $\pi^+$ & every negative point should be in its lower region of $\pi^-$

There should be no +ve & -ve points in the opposite directions. Actually there should be no point between the planes also (in the margin area).

Now what if we have a "-ve" point on the upper side of "$\pi^+$" or a "+ve" point on the lower side of "$\pi^-$"



This dataset 'cant be separated with a hyperplane

These three encircled points in above case will never satisfy the constraints.

So, if we try to solve the problem (optimization problem) for a dataset like shown above which is not linearly separable but it is almost linearly separable except for just a few points, most of the points are ok.

In such a case you may never be able to find a "w & b" that satisfy these conditions.
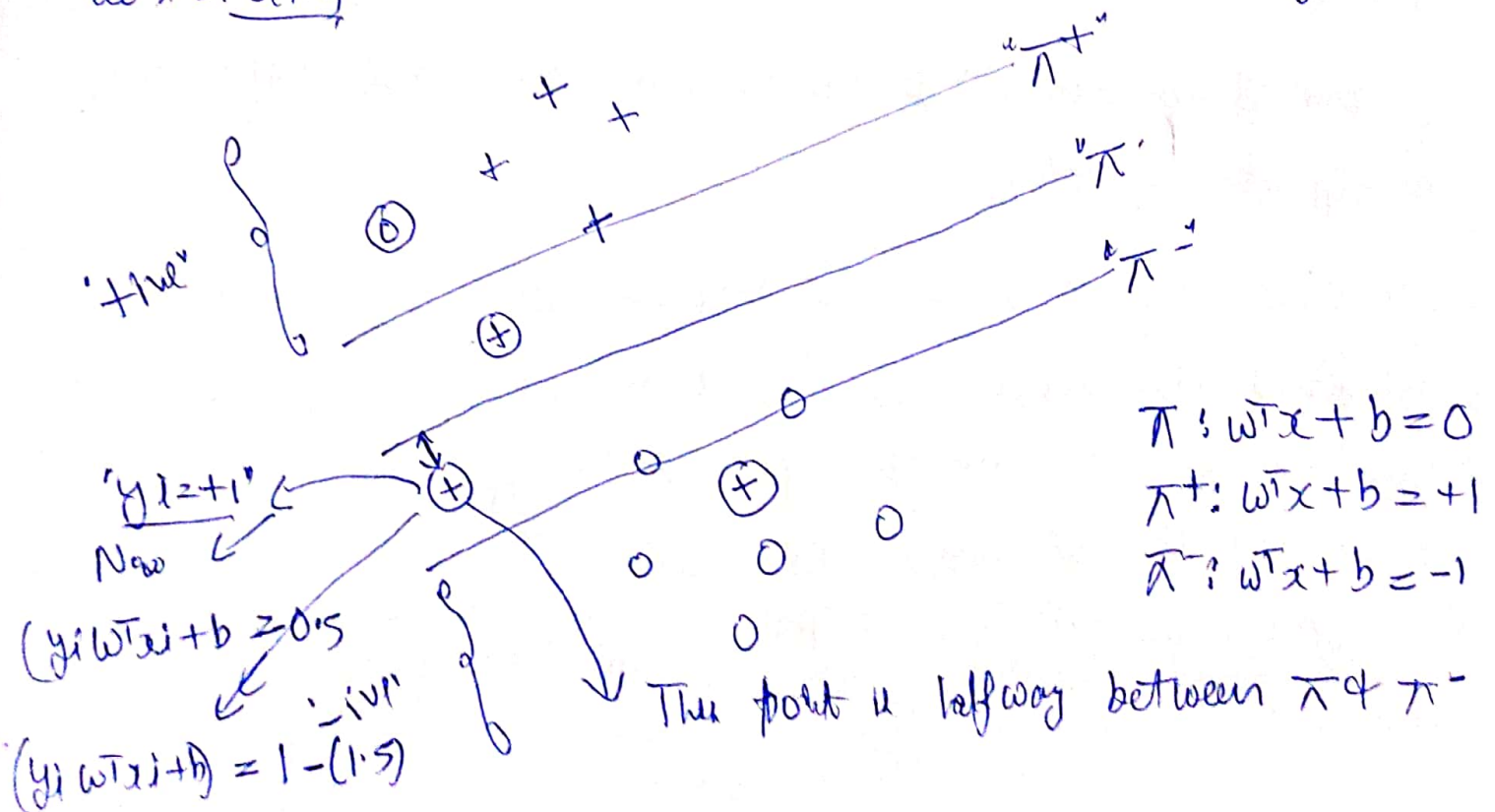
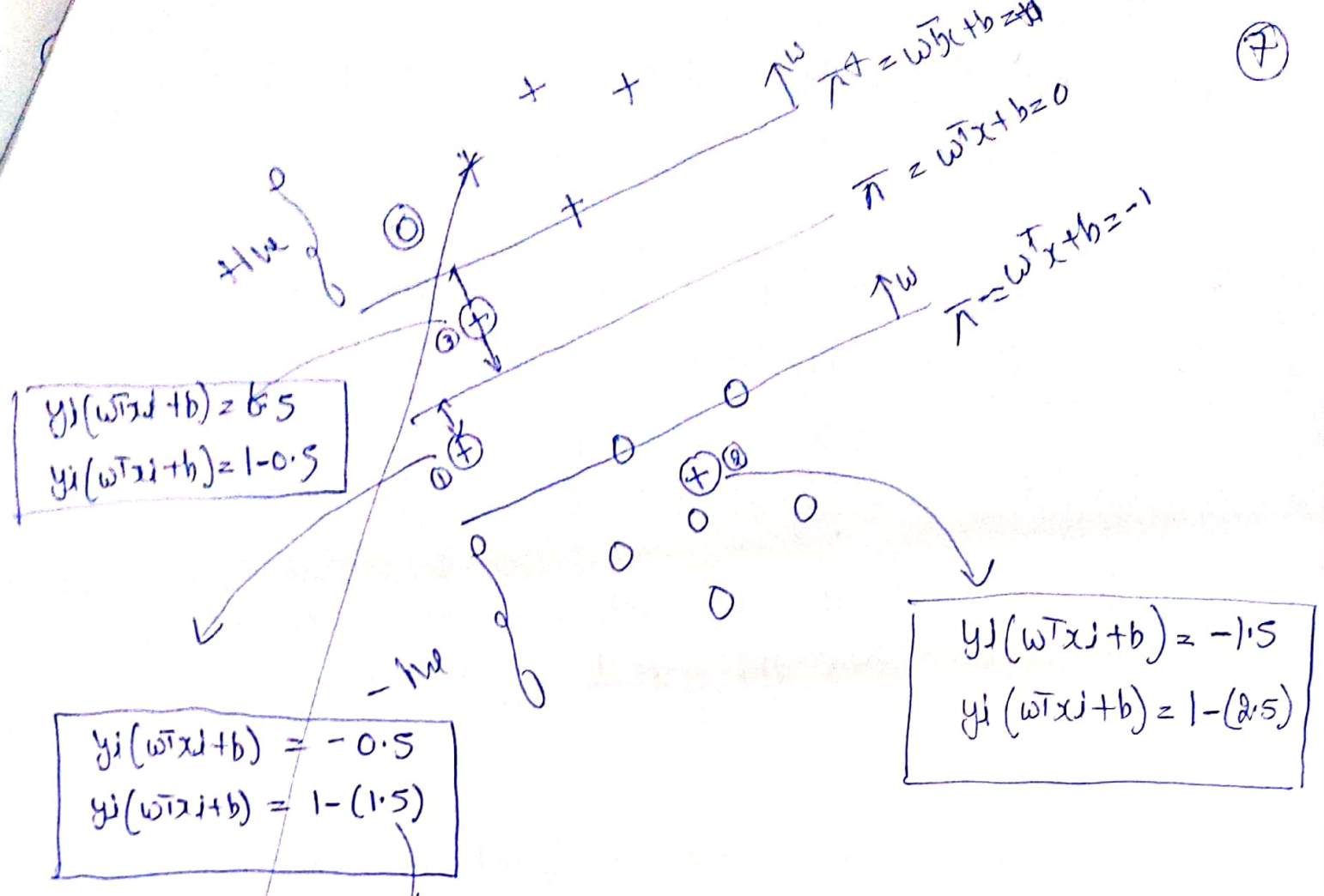It is impossible to find something like that, which satisfy the constraint.

This optimization problem is the diagrammatic representation, where all the +ve points are on the upper side of "$\pi^+$" & all the -ve points are on the lower side of "$\pi^-$". This formulation is called [hard-margin SVM"].

becoz we are saying all the +ve points are on one side & all the -ve points are on other side & nothing else. & we are imposing this thing through the constraint.

Now for almost linearly separable data, let's see can we somehow modify this formulation of ["SVM"]. slightly to find a hyperplane for almost linearly separable cases also.

We take each of these misclassified points & try to do something



'$y1 = +1$' ←

Now ←

$(y_i \omega^T x_i + b \geqslant 0.5$

←— '$i.v.$'

$(y_i \omega^T x_i + b) = 1 - (1.5)$

This point is halfway between $\pi$ & $\pi^-$

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

$\pi^+ : \omega^T x + b = +1$

$\pi^- : \omega^T x + b = -1$

$$\pi^+ = w^Tx + b = +1$$

$$\pi = w^Tx + b = 0$$

$$\pi^- = w^Tx + b = -1$$

$$y_j(w^Tx_j + b) = 1.5$$
$$y_i(w^Tx_i + b) = 1 - 0.5$$

$$y_j(w^Tx_j + b) = -1.5$$
$$y_i(w^Tx_i + b) = 1 - (2.5)$$

$$y_j(w^Tx_j + b) = -0.5$$
$$y_i(w^Tx_i + b) = 1 - (1.5)$$

let us call that term at $\xi_i$

for ① $\xi_i$ is 1.5      for ② $\xi_i$ is 2.5   & for ③ $\xi_i$ is 0.5.

& for correct point $\xi_j$ will be "0" becz it is already greater than '1'.

So, what we do is, we create a new variable called "$\xi_i$"

If a +ve point lies in correct region, its $\xi_i$ will be '0' same is with -ve point, if it lies in correct region. i.e below $\pi^-$ plane.

If a +ve point lies anywhere else other than the region specified. Then its $\xi_i$ will be +ve

Let's assume '$\xi_i$ to be $\geq 0$"

So, if $\xi_i \uparrow$ then the point is farther away from the correct hyperplane in _incorrect direction_.

So, for every point '$x_i$' we are creating '$\xi_i$' such that $\xi_i = 0$ if $y_i(\omega^T x_i + b) \geq 1$

which means they are correctly classified, as per not '$\pi$' but as per '$\underline{\pi^+}$' & '$\underline{\pi^-}$'

But $\xi_i > 0$ & it is equal to some units of distance away from the correct hyperplane either $\pi^+$ or $\pi^-$ in the _incorrect direction_.

$\xi_i$'s are telling us whether a point is correctly classified or not & how far it is away from the correct hyperplane in its incorrect direction..

Now let us formulate our _optimization function_

Our initial formulation is, we have to find $(\omega^*, b^*)$ such that $\frac{2}{||\omega||}$ (marginal distance) gets maximized

$$(\omega^*, b^*) \ \underset{\omega, b}{argmax} \ \frac{2}{||\omega||}$$

& maximizing that is same as minimizing $\frac{-||\omega||}{2}$

$$(\omega^*, b^*) \ \underset{\omega, b}{argmin} \ \frac{||\omega||}{2}$$

becoz

$\max f(x) == \min \frac{1}{f(x)}$

Now let us see why & how can we write our whole optimization problem ⟶ ½margin

$$\{ (w^*, b^*) = \underset{w,b}{\arg\min} \left( \frac{\|w\|}{2} \right) + C \cdot \frac{1}{n} \sum_{i=1}^{n} \xi_i \}$$

such that  $y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i$

& $\xi_i > 0$

we have seen that for every misclassified point . we can write    $y_i(w^T x_i + b) = 1 - \xi_i$    where $\xi_i$ is +ve

so for all misclassified points, we can write it as

$y_i(w^T x_i + b) \geq 1 - \xi_i$    where $\xi_i$ is +ve

Here what it happenly is, for correctly classified points '$\xi_i$' will be equal to '0'.

Now what do we want to minimize.,

we want to minimize the errors or we want to minimize misclassifications.

Minimizing misclassification means that, since $\xi_i > 0$ for all misclassified points. This means we want to minimize the sum of '$\xi_i$'

so, in our objective function, earlier we have only margin

Now along with margin if I say that I want to minimize the average distance (becoz, $\varepsilon_i$ represent the distance of incorrectly classified points for correct hyperplane in opposite direction).

$$\frac{1}{n}\sum_{i=1}^{n}\varepsilon_i \rightarrow$$ This is the avg. distance of misclassified points, becoz for all the correctly classified points $\varepsilon_i$ will be equal to '0'.

& we want to minimize that

'C' here is <u>hyperparameter</u>.

Here $\dfrac{\|w\|}{2}$ is $\dfrac{1}{margin}$

$$\frac{1}{n}\sum_{i=1}^{n}\varepsilon_i \rightarrow$$ avg distance for misclassified points

We can think of this $\rightarrow \left(\dfrac{1}{n}\sum_{j=1}^{n}\varepsilon_i\right)$ as a loss becoz we want to minimize it. We want to minimize the no. of misclassified points. So whenever there is a misclassified point '$\varepsilon_i$' is greater than '0'. This is basically a loss to the model that we want to <u>minimize</u>

'C here is <u>hyperparameter</u>        C is the

as C↑; we are giving more importance to not making errors. As C↑ we are saying, we don't want to make mistakes.

C↑ ; tendency to make mistakes on Dtrain reduces

<s>High variance model</s> → - This means we are going to <u>overfit</u>.

C↓ ; $\underline{\|w\|}$ will get more importance, & we have a
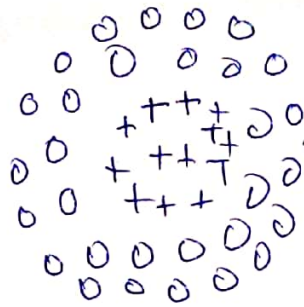↳ tendency to <u>underfit</u>.
↳ high bias

The formulation of SVM is called 'soft-margin SVM'

Hard-margin SVM's donot allow errors.

but a 'soft-margin SVM' says make errors but minimize them

→ 'Polynomial Kernel"

let's take an example where we have a bunch of 'true points,'
surrounded by a bunch of "negative points" & the datasets looks like
two concentric circles.



In logistic regression, we can separate these points by feature transfor-
mation such.

$$(f_1, f_2) \xrightarrow{\text{'FT'}} (f_1^2, f_2^2)$$

↳ & finally we can separate them with
a line —

Now let's look a polynomial kernel :→

The general definition of a polynomial kernel is
given two datapoints $(x_1 \& x_2)$ the general polynomial kernel is
$(x_1^T x_2 + c)^d$ .

$$K(x_1, x_2) = (x_1^T x_2 + c)^d$$

where 'c' & 'd' are constants.

let's take an example of a quadratic kernel.

eg:→ $K(x_1, x_2) = (1 + x_1^T x_2)^2$

↓ quadratic kernel

Here $c = 1$ & $d = 2$

If we apply this, let's see what is $K(x_1, x_2)$

$K_\phi(x_1 + x_2) = (1 + x_1^T x_2)^2$

$= (1 + [x_{11}, x_{12}] \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix})^2$

$= (1 + x_{11} x_{21} + x_{12} x_{22})^2$

$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 + x_{22}^2 + 2 x_{11} x_{21} + 2 x_{12} x_{22} + 2 x_{11} x_{21} x_{12} x_{22}$

Let's assume

$x_1 = \langle x_{11}, x_{12} \rangle$
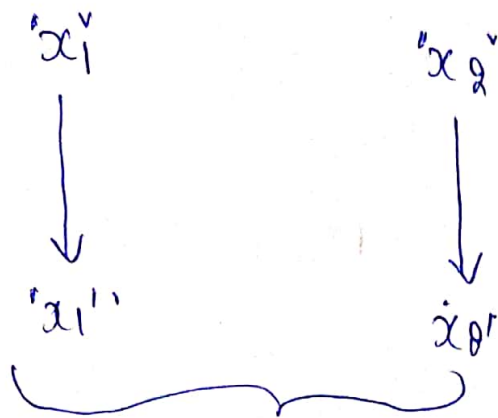
↓
vector of two points

$x_2 = \langle x_{21}, x_{22} \rangle$

↓
This can be represented as a product of vectors.

let $\# [1, ; x_{11}^2, x_{12}^2, \sqrt{2} x_{11}, \sqrt{2} x_{12}, \sqrt{2} x_{11} x_{12}] : (x_1')^T$

$[1, x_{21}^2, x_{22}^2, \sqrt{2} x_{21}, \sqrt{2} x_{22}, \sqrt{2} x_{21} x_{22}] : x_2'$ vectors

Now we can show that this product we have written above

is equivalent to

$\{ (x_1')^T (x_2') \}$

If we look at it carefully, we will find out $x_1$ & $x_2$ are

in 2d & now for $x_1'$ we have from only $x_{11}$ & $x_{12}$

& for $x_2'$ we have only $x_{21}$ & $x_{22}$ terms

(14) So imagine if we are given "$x_i$" & "$x_2$"
we have transformed them into $x_1'$ & $x_2'$

"$x_i$"        "$x_2$"

↓        ↓

'$x_1$''        $\dot{x}_{2'}$

Now instead of doing

$$\boxed{x_1^T x_2}$$ we can do

$$\boxed{x_1'^T x_2'}$$

⤷ This is equivalent to "feature transform"

So, what 'kernalization' is doing internally is exactly equal to 'feature transformation'

Kernelization takes "d –dimensional data" & does a feature transformation internally & implicitly

Kernelization :-  $d \xrightarrow[\substack{\text{internally} \\ \& \\ \text{implicitly}}]{FT} d'$

Feature Transformation :→  $d \xrightarrow[\text{Explicitly}]{FT} d'$

⤷ very powerful trick, we are converting d dimensional points into d' dimensional points, where d' > d typically

$x_1'$ is 6-d data

$x_2'$ is 6-d data

So using keonel trick we went from 2d data to '6-d data'

$$x_1' = [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}]$$

$$x_2' = [1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}]$$

Since in there we have squared terms $(x_{11}^2, x_{12}^2, x_{21}^2, x_{22}^2)$ & these terms are very similar to $(f_1^2, f_2^2)$ terms the data will become separable.

———————