/→ "Support Vector Machines" (SVM)
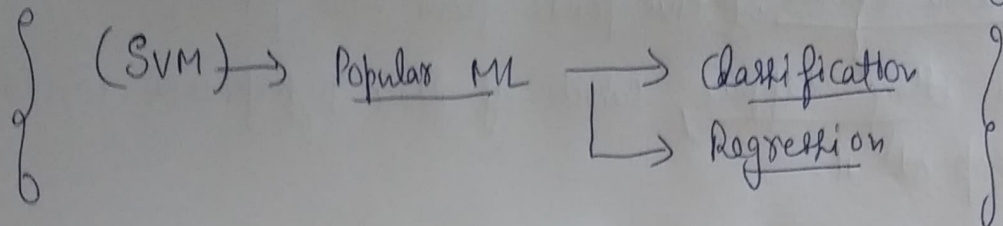
Support vector machines often referred to as SVM's are very popular machine learning techniques.

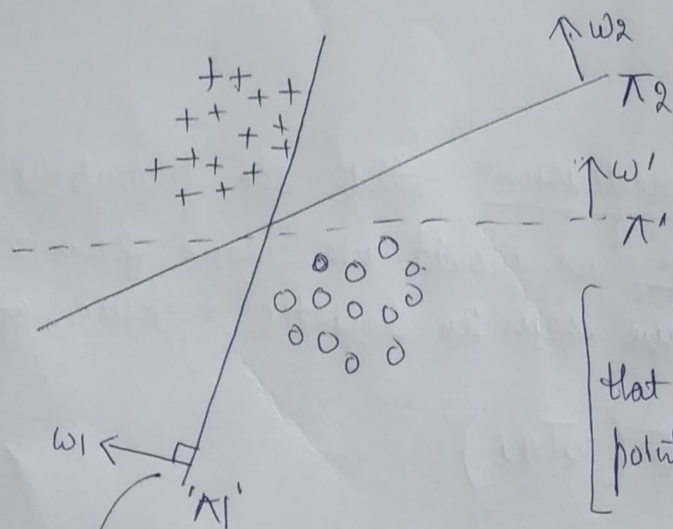They can be used for both "classification as well as Regression purposes"

$$\Big\{ \quad (SVM) \to \text{Popular ML} \quad \overset{\to \text{Classification}}{\underset{\to \text{Regression}}{\boxed{\phantom{x}}}} \quad \Big\}$$

It become extremely popular in the late 1990's.

They are very very elegant & simple algorithms.

Before we go into mathematical details, let's understand the geometric intuition behind "SVM's".

Imagine, we have a dataset comprising of [+tive] & [-tue) points.

$$\Big\{ \quad \text{+ive points are Represented using "+"} \atop \text{& -1ve points are represented using "o"} \quad \Big\}$$



Let's assume for simplicity that this data is linearly separable.

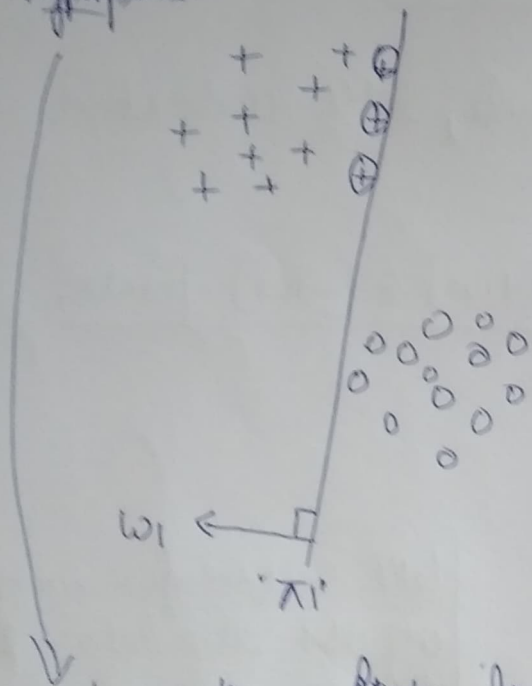Now I have multiple hyperplanes that I can draw to separate these points from one another.

↓ I can draw this hyperplane, let met call this hyperplane as "$\pi_1$". Correspondingly "$w_1$" is a vector normal or perpendicular to it. This hyperplane separate these two classes of points.

Similarly if you look at it, there is another hyperplane that separates these two clusters. let's call this hyperplane as $\pi_2$ & $w_2$ is a normal to it.

like these two there are many many hyperplanes possible which separates these two classes (+ve points from -ve points)

When such a thing happens, which of these hyperplanes would you prefer. let's look at it.

let's take $\pi_1$ as separating hyperplane. Now if you look closely there are lots of points which are lying very close to the separating hyperplane.
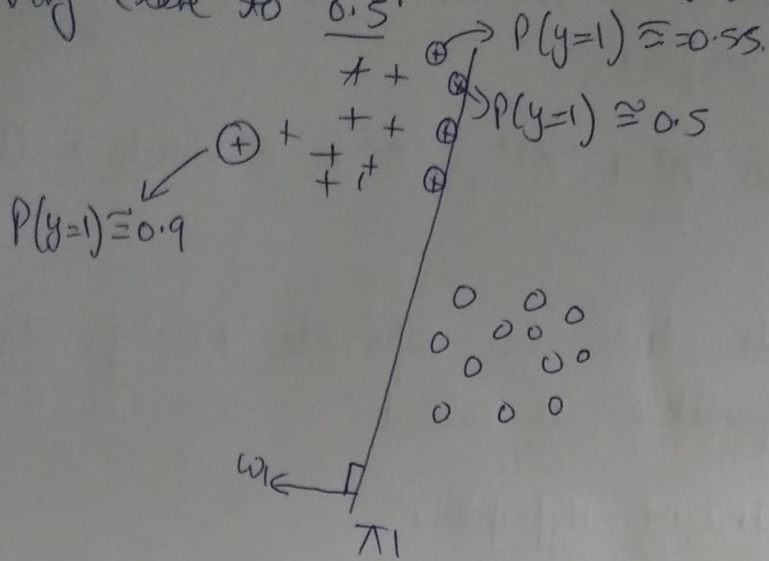


if we know from 'logistic Regression' that the probability of a point belonging to 'the class', we would use some function like sigmoid function as we have seen in logistic regressor of the distance

$$P(y_i = 1) = \sigma(w^T x_i)$$

So, if a point is very close for eg. the encircled point you see above, since they are very very close to hyperplane, if the hyperplane changes slightly, they could get easily misclassified.

So, if a point is lying very close to the hyperplane, then ③ the probability of the point belonging to 'class 1' let's say will be very close to '0.5'

$P(y=1) \cong 0.55$

$P(y=1) \cong 0.5$

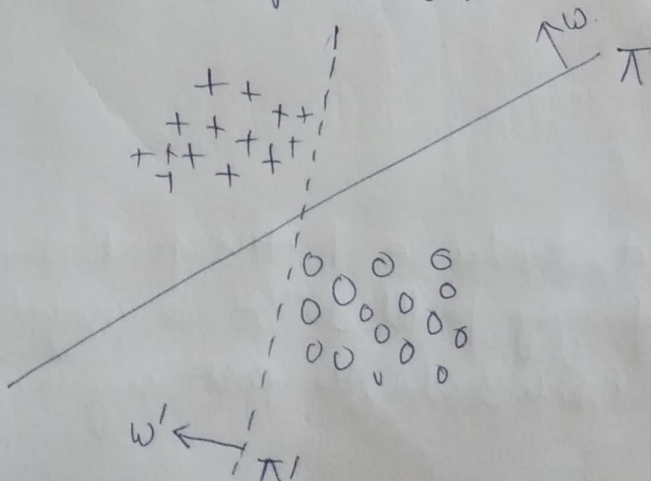$P(y=1) \cong 0.9$

$w \longleftarrow \quad \pi_1$

So, what we want is a hyperplane that separates the positive points (+ve) & -ve points as far away as possible.

↓ This is the key geometric idea of SVM's.

So, what SVM's are trying to do, is it is trying to find a hyperplane that separates the +ve & -ve points as widely as possible.

Let's understand what does it mean actually from a geometric perspective :-

Let's have a bunch of the (+) & -ve points (O) as shown.

$\pi$

$w' \longleftarrow \quad \pi_1$

If you choose a hyper plane "$\pi$" with its normal "$w$". You will notice that +ve points & -ve points are as far away as possible to this hyper plane as compared to "$\pi_1$" with normal $w'$

So, we see that $\boxed{``\pi"}$ is better than $\boxed{``\pi'"}$.

becoz in case of "$\pi'$" there are some "+ve" & "-ve" points which are very very close to the 'hyperplane'.
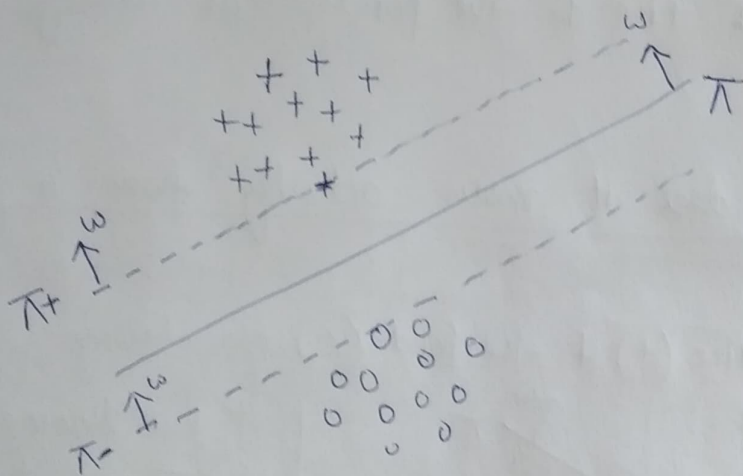
$\therefore$ If we have to choose b/w "$\pi$" & "$\pi'$", "$\pi$" is a much better choice.

& Such a hyper plane which tries to separate '+ve' points from '-ve' points as much as possible is called

"Margin Maximizing Hyperplane"

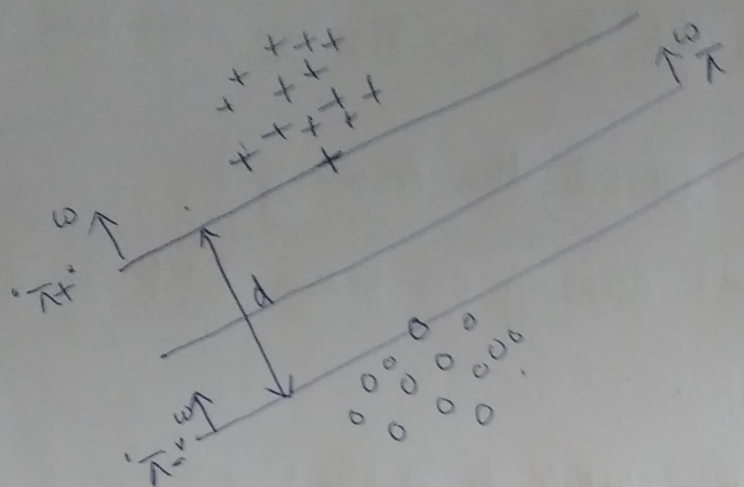[ $\pi$: Margin maximizing hyperplane ]

let's understand margin maximizing hyperplane, So, if we choose "$\pi$" as a hyperplane which separates the '+ve' & '-ve' points & if we go parallel to this hyperplane which means if we keep on drawing the hyperplanes in +ve direction then at some point we will intersect the first positive point. let's call this hyperplane as $\pi+$.



$\pi+$ is parallel to $\pi$ & it touches a positive point & it is parallel to $\pi$. Similarly if we keep going parallel to "$\pi$" in negative direction we will get a hyperplane which touches the first negative point. let's call that as "$\pi-$".

'π+' is parallel to "π" and 'π-' is also parallel to 'π' except ⑤
that it passes through the -ive point.

Now the gap between then there two, let's denote it by d.



Now since ["π+"] & ["π-"] are parallel hyperplanes, they have a constant distance between them.

$$d = distance (\overline{\pi}+, \pi-)$$

↓This distance is also called as __margin__, because it is a margin or a gap b/w '+ive points' & '-ive points'.

So, we want to find a hyperplane "π" such that, if we go parallel to 'π' & if we touch a positive point with a hyper plane 'π+' & if we go parallel to 'π' & we have a hyperplane 'π-' which touches the first '-ive' points. Then the margin that we get, which is basically the distance b/w 'π+ & π-', we want to maximize it.

∴ when this distance 'd' is maximized, then the '+ive points' & '-ive points' are quite far away from the actual separating hyperplane, and the farther they are, the wider the gap, the better it is.
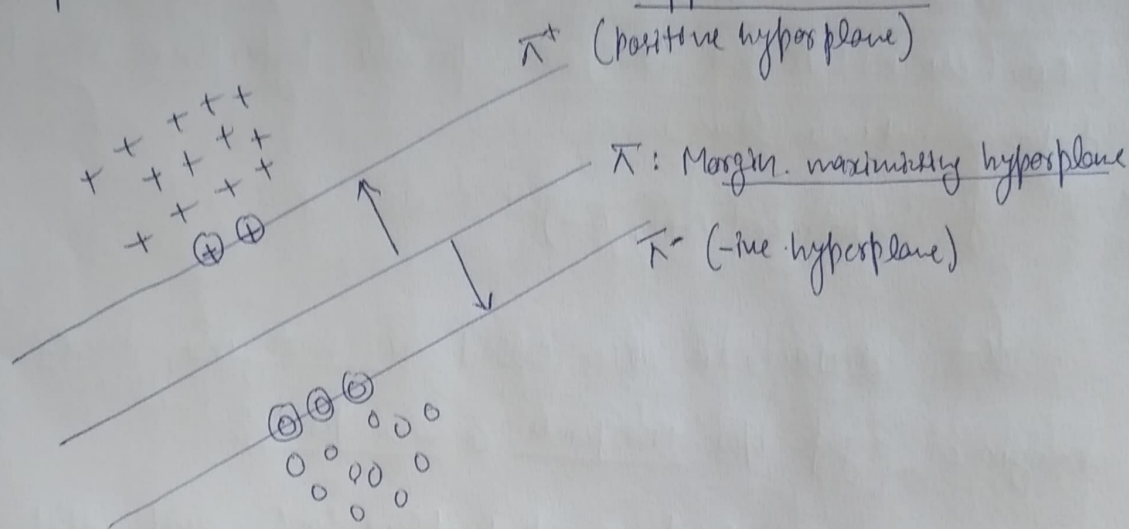
So, what SVM's try to do it, they attempt or try to find a hyperplane $\pi$ that maximizes the [margin = dist $(\pi^+, \pi^-)$].

Note :⇒ If the margin is high, the chance that we will misclassify will decrease, which means that the generalization error will improve.

as, Margin ↑; generalization Accuracy ↑;

generalization accuracy means. accuracy on unseen datapoints in future. It will also improve.

Another important idea in SVM's is "support vector".



$\pi^+$ (positive hyperplane)

$\pi$ : Margin. maximizing hyperplane

$\pi^-$ (-ive hyperplane)

"$\pi^+$" goes through ["2 positive points"] + "$\pi^-$" goes through ["3 negative points".]

Note :⇒ The points through with "$\pi^+$" + "$\pi^-$" passes through are called support vectors...

Given a query point, we will use original "$\pi$", the actual separating hyperplane to to classify the query point into "+ives or -ives but we maximize the margin b/w "$\pi^+$" + "$\pi^-$". So the points which are on either "$\pi^+$" or "$\pi^-$" are called as "support vectors".
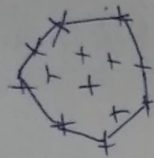
Alternative geometric intuition behind SUM's :->

It is as follows :->

Suppose we have a bunch of "+ve" & "-ve" points. we can find the margin maximising hyperplane by drawing something known as a "convex hull".
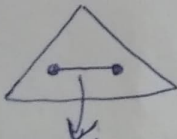
↳ let's understand a "convex hull". Suppose we have a bunch of points. A convex hull is defined as a polygon that
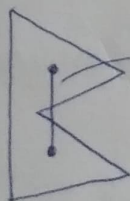


we can draw, such that, it is a smallest polygon, such that all the points are either inside the polygon or on the polygon. & this polygon has to be convex polygon.

Note :- A shape is said to convex if the line joining any two points inside that shape is also present inside the shape/region.

Eg :-



The line joining these two points is lying inside the triangle, hence it a convex polygon.

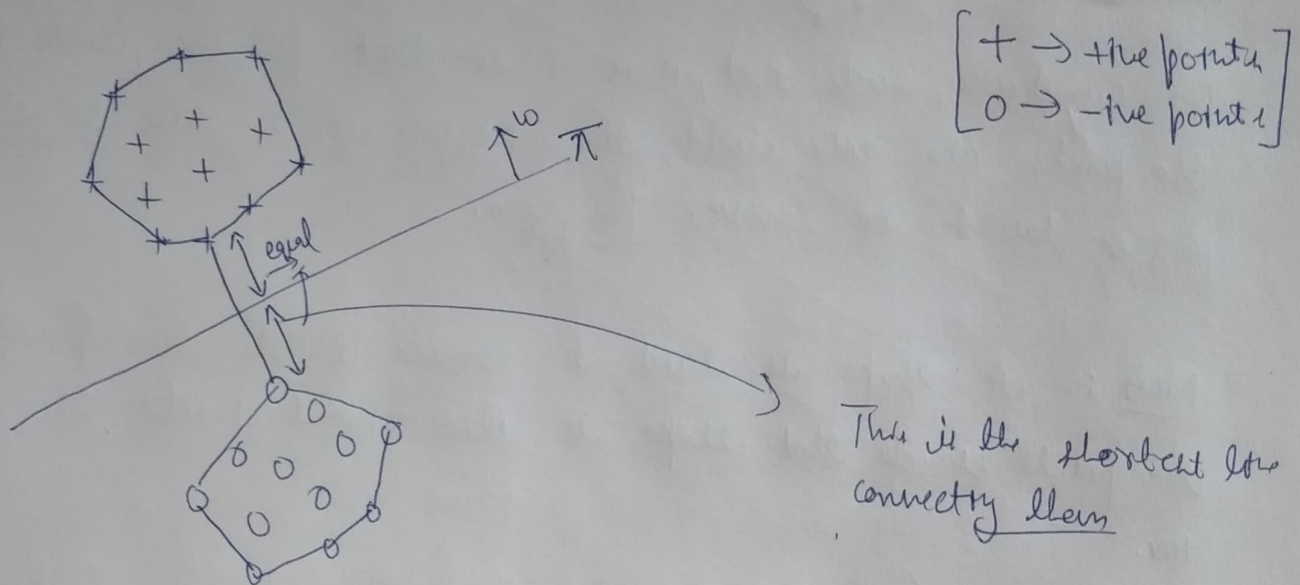 → This shape is not convex, because the line joining these points, some part of it is lying outside the region.

↓ This is called → non-convex polygon.

'Convex-polygon' :→ It is a polygon, such that if we went to go from one point to another point, the line connecting both of them is lying inside the polygon.

'Convex-hull" :→ Given a bunch of points, if we can build a smallest convex polygon, that has every point either inside the polygon or on the polygon.

Let's understand the alternative _geometric intuition_ behind behind these SVM's.

Suppose we are given a bunch of _positive_ & _negative_ points.



$$\left[ \begin{array}{l} + \rightarrow +ve\ points \\ 0 \rightarrow -ve\ points \end{array} \right]$$

This is the shortest line connecting them
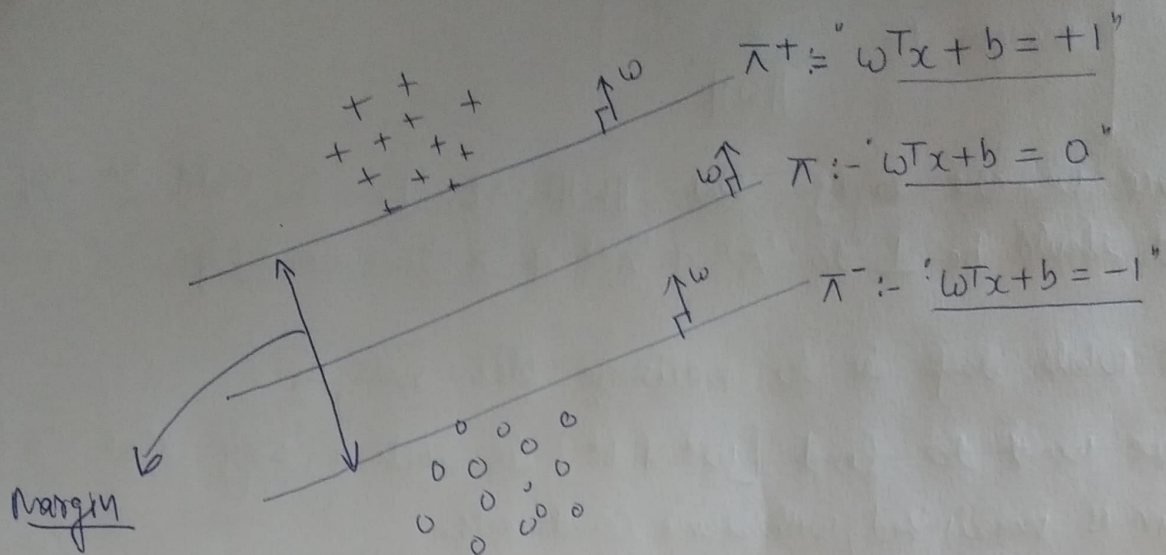
Step ① :→ Construct a convex hull for the +ve points.
& Construct a convex hull for -ve points separately.

Step ② :→ Find the shortest line connecting these convex hulls.

Step ③ :→ Bisect the line, (Break it into two equal parts).
The plane that bisects this line, is the margin maximizing hyper plane.

→ "Mathematical derivation of SVM"



$$\bar{\pi}^+ = \text{"}\omega^T x + b = +1\text{"}$$

$$\omega \hat{1} \quad \pi :- \text{"}\omega^T x + b = 0\text{"}$$

$$\bar{\pi}^- :- \text{"}\omega^T x + b = -1\text{"}$$

margin

In the formulation for __SVM's__ the equations that we use to represent the planes are.

$$
\begin{bmatrix}
\pi : & \omega^T x + b = 0 \\
\pi^+ : & \omega^T x + b = +1 \\
\pi^- : & \omega^T x + b = -1
\end{bmatrix}
$$

Now the question here is why $[+1]$ & $[-1]$ on the "R.H.S" of __positive & negative hyperplanes__.

__Note__ :- Here we are not saying that "$\omega$" has to be unit vector, we are saying that "$\omega$" could be any vector (it could have any length).

All we are saying here is, "$\omega$" is $\perp$ to "$\pi$" but it could have any length. So the "$L_2$ Norm" or length of "$\omega$" need not be "1", it could be any value.

$$\Big\{ \| \omega \| \neq 1 \quad \text{(any length)} \text{ need not to be a unit vector} \Big\}$$

The whole margin: could be derived very simply as the distance of a point lying on ($\pi^+$ or $\pi^-$) from the alternate plane, hyperplane

The margin can be derived as :-

$$\left[ \text{Margin} :- \frac{2}{||\omega||} \right]$$

& Since we are saying that $||\omega||$ is not '1'. It could be any vector.
but it should be $\perp$ to '$\pi$' & '$\pi^+$' & '$\pi^-$' equivalently.

Our whole task is to maximize this margin

So, we want to find $[\omega^* \& b^*]$ in such a way that the
margin is maximized, with some constraints.

$$\left\{ \omega^*, b^* = \arg\max_{\omega, b} \frac{2}{||\omega||} \right\} \longrightarrow \begin{array}{c} \text{objective function} \\ \text{of SVM.} \end{array}$$

And the constraints are, all the +ve points should be on one side
& all the negative points should be on other side, but the object-
ive is to maximize the margin.

let's now look at the first approach to prove that "$+1$ & $-1$"
on R.H.S doesn't matter.

① .

Let $\pi^+$ : $\omega^T x + b = k$      | The only requirement is
& $\pi^-$ : $\omega^T x + b = -k$     | that "$k$" should be greater
                                                        | than "0" ($k > 0$) only then
                                                        $\pi^+$ & $\pi^-$ are away from $\pi$

& why are we taking "$+k$" & "$-k$" here. &
Why can't we take "$+k_1$ & $-k_2$" . The reason is we want
both of them +ve & -ve hyperplanes to be equally far away
from the margin hyperplane..

Note :→ '$k$' could be any number, as long as it is greater than "0"

PTO

Now when we have this, the margin will change to

margin :- $\dfrac{2k}{||w||}$

Here 'k' is a constant, It could have any value

If this is the margin, then from optimization's point of view, we know that the 'w + b' which maximizes $\left[\dfrac{2}{||w||}\right]$ is the same as the 'w + b' which maximizes $\left[\dfrac{2k}{||w||}\right]$

$$\left\{ \operatorname*{arg\,max}_{w,b} \dfrac{2}{||w||} = \operatorname*{arg\,max}_{w,b} \dfrac{2k}{||w||} \right\}$$

let $k = 4$. then we have

$$\operatorname*{arg\,max}_{w,b} \dfrac{2}{||w||} = \operatorname*{arg\,max}_{w,b} \dfrac{8}{||w||}$$

So, from optimization's point of view, if we take "1" or "k" it really does not matter.

So, we only took "1" here as a way of convenience, Nothing will change.

———— o ————