

Dr. Osite Onyekw

P_q + II

Support Vector Machine

Vary Hyperplane

Optimization

04/19/2024

Original Optimization Problem:

$$\min_w \underbrace{\frac{1}{2} \|w\|^2}_{\text{Loss function } f(w)} \quad \text{s.t. } \underbrace{y^i (\vec{x}^{iT} \vec{w} + b) - 1}_{\substack{\text{estimate} \\ \nexists i \in \{1, \dots, n\}}} \geq 0$$

constraints

- We want the function $f_i(w) \geq 0$

We want to correctly classify each point i , in other words, we want to put it on the right side of the hyperplane.

$$\Rightarrow f_i(\vec{w}, b) \geq 0$$

\uparrow
w-Vec
 \nwarrow scalar b

$$\Rightarrow \vec{w} = \begin{bmatrix} b \\ \vec{w} \end{bmatrix}$$

- We have many constraints. We have a constraint for each data point " i ".

- This problem can be formulated in this way:

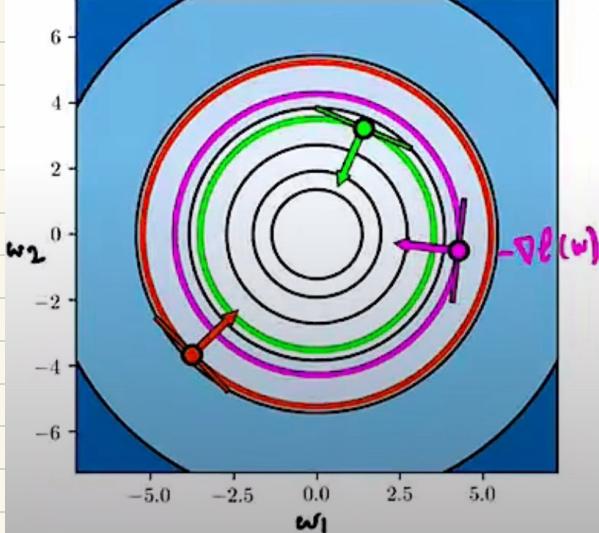
\downarrow are the space of parameters w

$$\underset{w}{\text{minimize}} \underbrace{f(w)}_{\text{loss function}}, \quad w \in \mathbb{R}^n$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

s.t. $f(w) \geq 0 \leftarrow$ This condition is fulfilled

(General formulation of a constrained optimization problem)

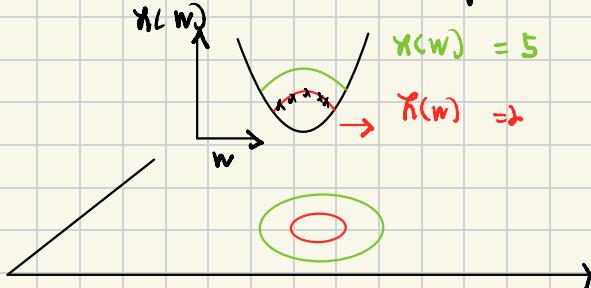


- Very Challenging Problem to solve mathematically.
- Gradient Descent : few steps/iterations about the gradient direction.
- When we plot the contours of our Loss function $\lambda(w)$

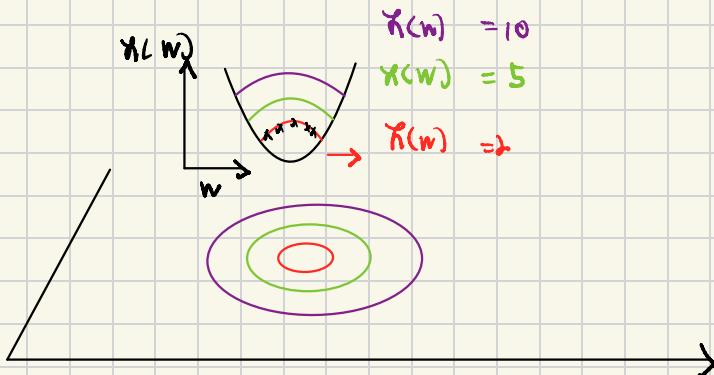


We defined the contour or the level contour of the loss function, where the loss function at each point, \Rightarrow when we get the w we know that all those points, in the loss function = constant, \Rightarrow

$\lambda(w) = \lambda$, and then we projected it. I imagine this is a 3-D space, and this we project it.

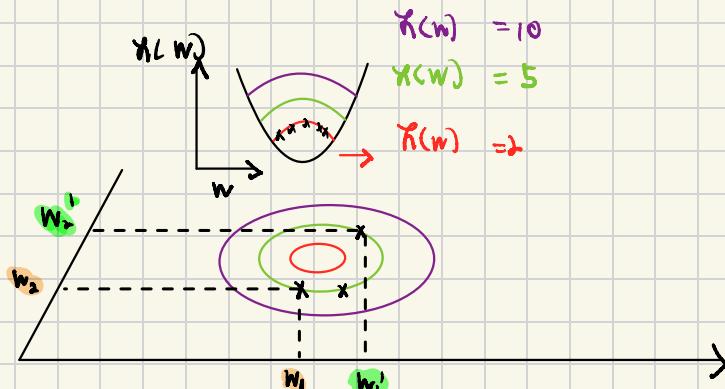


- We find the second contour, and this is when $\lambda(w)$ takes a higher value. and the λ is actually constant along the contour.
-



Then, we have another one, $f(w) = 10$

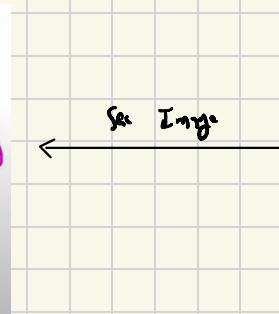
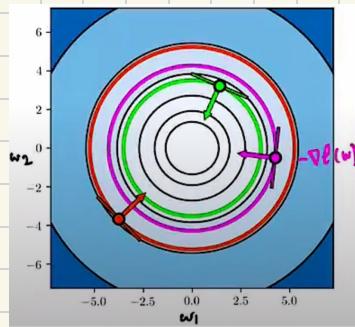
∴ If we select a random point:



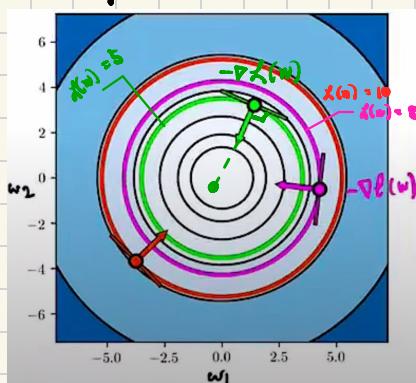
We know that all points have the same value, i.e., all equal to 5 because they all lie on the green contour.

①

In gradient Descent, the gradient at each point that lies on the contours, the gradient (regardless of it) is perpendicular to the contours.



① We also want to move along the negative direction of the gradient towards our optimum. i.e., because it is an minimum. because this is local for the values in say the training centers we go closer to the center.

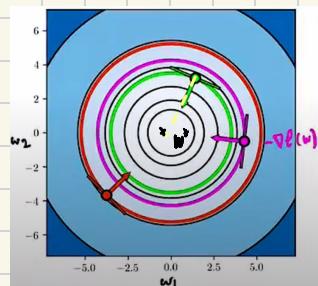


Remember, our goal is to minimize oil loss thickness of walls slowly jumping from one contour to another very gradually descent

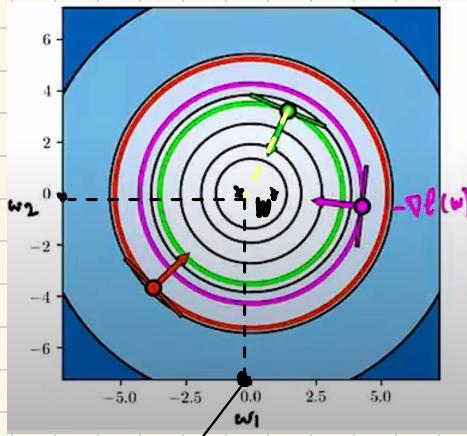
What we want to do is to minimize this.

$$\begin{array}{l} \text{minimize}_w \quad f(w), \\ \text{s.t.} \quad g_i(w) \leq 0, \quad i = 1, 2, \dots, m \end{array}$$

↳ If we remove the constraint, and we have a general loss function that we want to minimize, we will find a optimum.



original
contours. g
 $A(w)$



Item, it is located at $(0,0)$ ($w_1 = 0, w_2 = 0$) \Rightarrow
 w is a \mathbb{R}^2 vector $\rightarrow w \in \mathbb{R}^2$, however, we have
 an additional condition,

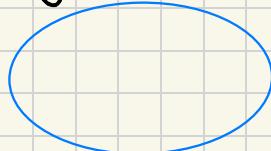
① If we solve this optimization problem it also boils down to
 solving this condition:

$$\begin{array}{ll} \text{minimize}_w & h(w), \quad w \in \mathbb{R}^2 \\ \text{s.t.} & \underbrace{f(w)}_{=} c \end{array}$$

$$\text{s.t. } f(w) = c \quad (\text{constant})$$

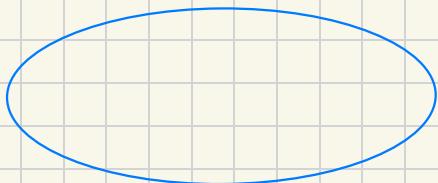
Value c .
 can be solved by setting the condition to a constant

→ This is exactly a contour line of the function.

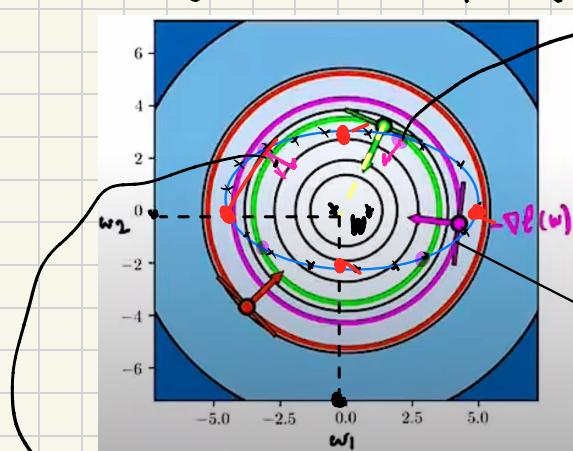


\Rightarrow The function is constant.

- We are trying to find all w 's, such that $f(w) = c$
- \therefore let's say that this is the contour, s.t $f(w) = c$



and let's say it lies exactly here.



Not drawn to scale
but assume that $f(w)$
intersects those contours of
 f at those points.

We want both curves/contours to be tangent.
We want the "w" to lie on both contours

$$\rightarrow f(w) = c$$

\Rightarrow not interested in other intersection points because of the directions of the gradients they make.

\therefore On this contour, we have that $f(w) = c$. \downarrow be want tangent curves.

\therefore We are looking for the vector \vec{w} (the set of parameters) s.t
we are minimizing L , and also our points should lie on the
contour.

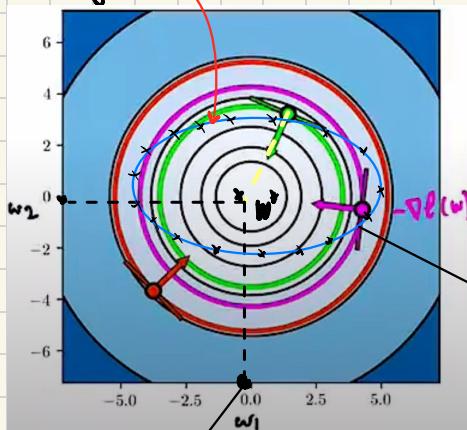
We want to find all possible solutions to this constraint optimization problem such that the contours of the loss function $l(w)$ are tangent to the contours of the loss $f(w)$.

- Imagine that the contour b takes on as many values.



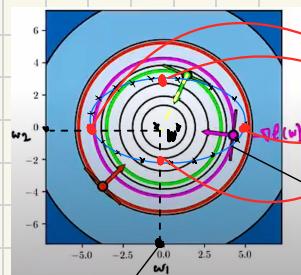
, but we can only intersect

in particular with the $\nabla l(w) \perp b$. because that is our condition, so we want to find an optimum W on this specific contour.

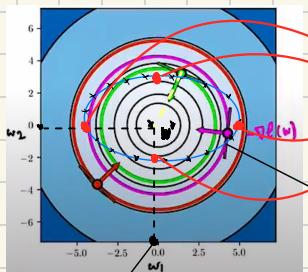


② We also want to lie on a minimal contour of h .

: we know that, as W should minimize the value of our loss function, it should lie on a large contour of h which by $h(w)$ takes on a minimum value.

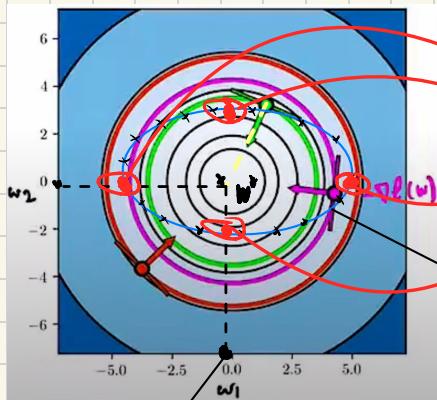


We want these key parts



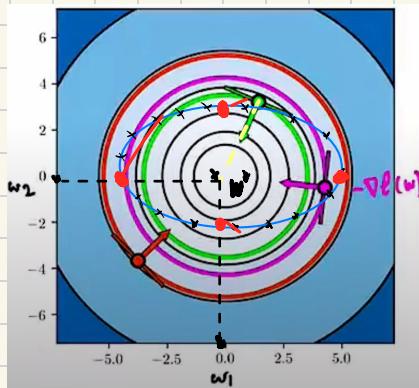
key parts.

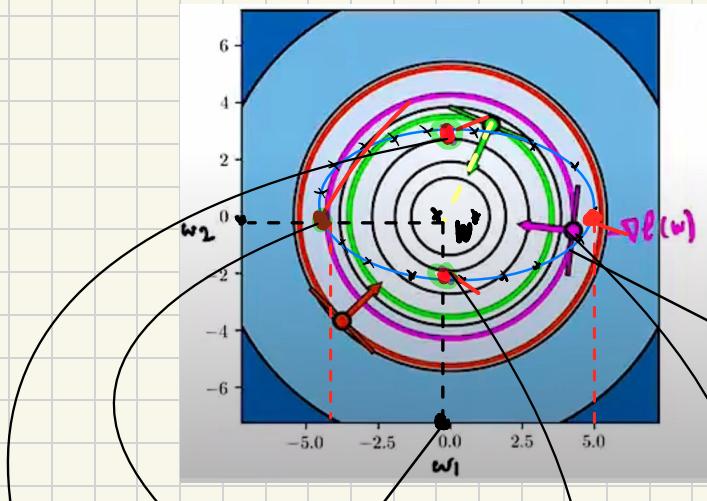
- Interest
- We want to find out where the contours for $f(w)$ intersect the contours of l .
 - We want last two contours to be tangent.



Here we have 4 candidate points.

find out
If we project them, we know their coordinates, but which one is then best.





$$f(w) = \underline{c}$$

For this one, I know that $f(w) = c$, however for $h(w)$, I know that this part lies on the purple constraint $\Rightarrow h(w) = 8$

$$h(w) = 8$$

$$h(w) = 2$$

$$h(w) = 2$$

Best candidate point here.

is looking at those 4 points, we know that the best one is the one that supports $h(w) = 2$.

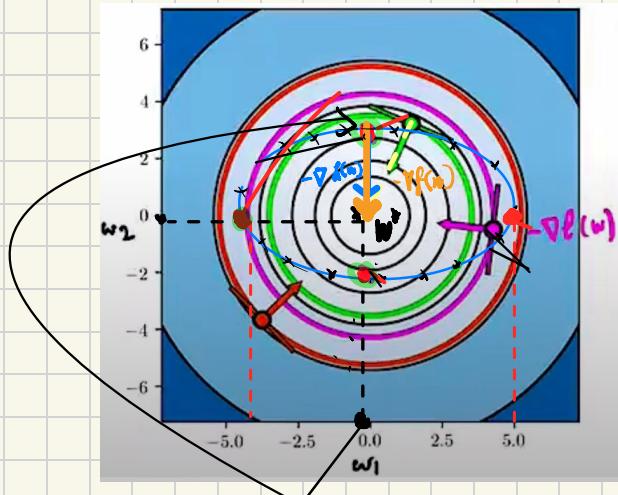
Those are our best candidates.

At the optimal point of w^* we should have the negative gradient of the loss function equal to the constraint.

- $\nabla f(w) = -\alpha \nabla f(w) - c$

colinear

Notice that the negative of the gradient of the loss function $-\nabla h(w)$ and if we calculate the negative of the gradient of the ℓ -function $-\nabla \ell(w)$, we



the gradient of the ℓ -function $-\nabla \ell(w)$, we notice that the vectors are collinear and they point towards the same direction, because the

curve and tangent is. Both of those vector gradients should be perpendicular to the contour of the loss function and the constraint.

We also want the 2 gradients to be colinear.

We know at optimal point, w^* !

$$-\nabla h(w) = -\alpha \underbrace{\nabla f(w) - c}_{\text{f-factor}} \quad \begin{matrix} c \text{ doesn't matter} \\ \text{because it is a} \\ \text{constant} \end{matrix}$$

(colinear)
sumulation
↓
The constraint itself.

We can remove the negative

$$\therefore \nabla h(w) = \alpha \nabla f(w)$$

∴ The Condition that the optimal points should satisfy is thus given by : $\nabla L(w) = \alpha \nabla f(w) \geq 0$

The most important thing to get out of this is that the gradients should be colinear and that $\alpha > 0$ because the 2 vectors should point to the same direction

Let's now define what we call the Lagrangian function $L(w)$ as :

We want $f(w) - c = 0$ ≠ because our "c" (in this case $c=0$)

$$\underbrace{y^i(x^i w + b) - 1}_{\text{constraints}} \geq 0,$$

$\rightarrow L(w) = L - \alpha [f(w) - c]$, we could get rid of the "c" because it is equal to 0, however this just leaves it like that for now.

∴ Writing it more formally we have that :

$$L(w) = L(w) - \alpha [f(w) - c]$$

Lagrangian of w Loss Lagrangian function

∴ If we take the gradient of the Lagrangian with respect to w_i , what do we find?

$$\nabla L = \nabla L(w) - \alpha \nabla f(w) - \nabla c$$

constant so ϕ

$$\therefore \nabla L = \nabla L(w) - \alpha \nabla f(w) \neq$$

some books write : $L(x, \lambda) = f(x) - \lambda g(x)$

$$\nabla \mathcal{L} = \nabla L(w) - \alpha \nabla f(w)$$

Remember that $\nabla L(w) = \begin{pmatrix} \frac{\partial \mathcal{L}}{\partial w_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w_d} \end{pmatrix}$
 We have 2 components & we have.

$\in \mathbb{R}^d$

, because our vector w belongs to \mathbb{R}^d ($w \in \mathbb{R}^d$)
 also, our samples belong to \mathbb{R}^d ($x \in \mathbb{R}^d$)

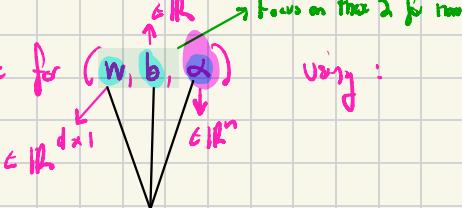
- We know that the gradient of the loss function = 0 [$\nabla \mathcal{L} = 0$]
- This is an important condition that will help us solve this problem.

① Original Optimization Problem: (Primal Optimization Problem):

$$\min_w \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad \begin{array}{l} y^T (X^T w + b) - 1 \geq 0 \\ \text{+ i.e. } \sum_i \text{ constraints} \end{array}$$

② Lagrangian: $L = \text{loss function to optimize} - \alpha x$ [constraint]

Lagrangian, $\alpha \geq 0$ (positive)
 multiplier
 per constraint

③ Solve for (w, b, α) using:


We need to solve for these parameters, and to do that we need to solve the following set of linear Equations using linear Algebraic tools.

$\left\{ \begin{array}{l} \nabla \mathcal{L} = 0 \quad (\text{ext. gradient} = 0 + \text{constant} \cdot \text{the gradient of the loss function}) \\ \text{+} \\ \text{all constraints} \end{array} \right.$
 We also need to use all the constraints!!
 We know our parameters $w + b$ satisfy particular constraints

$$\nabla L = \nabla R(w) - \alpha \nabla f(w) = 0$$

By adding the constraints we can find our optimal $\vec{w}^* \neq \vec{b}^*$

$$\min_w \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad \underbrace{y^i(x^{iT}w + b) - 1}_{\text{for } i \in \{1, \dots, n\}} \geq 0$$

constraints

L = loss function to optimize
 $-\alpha x$ [constraint]

Log/sgm, $\alpha \geq 0$
 multiplier
 per constraint

$$L = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^i(x^{iT}w + b) - 1] \quad \xrightarrow{\text{summing over all our constraints}}$$

Loss function
 $\alpha(x)$

① Compute gradient of L wrt w :

$$\frac{\partial L}{\partial w} = - \sum_{i=1}^n$$

$\in \mathbb{R}^n$

Hints:

$$\frac{\partial x^T w}{\partial w} = x$$

$$w^T w = \|w\|^2$$

$$\therefore \frac{\partial w^T w}{\partial w} = 2w$$

partial derivative
 of ReLU/margin
 with respect
 to w .

② Set gradient $\frac{\partial L}{\partial w} = 0$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow \underline{\underline{w}} \leftarrow \text{Optimal parameter}$$

Our goal is to \min_w : $\min_{\vec{w}, b} \max_{\vec{x} \geq 0} L(\vec{w}, \vec{x})$

Side Note : Dual SVM Definition : The Linearly Separable Case :

(Primal)

$$\min_{\vec{w}, b} \max_{\alpha \geq 0} \frac{1}{2} \|\vec{w}\|^2 - \sum_j \alpha_j [(\vec{w} \cdot \vec{x}_j + b)y_j - 1] \quad (1)$$

~~Swap min and max~~

(Dual)

$$\max_{\alpha \geq 0} \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 - \sum_j \alpha_j [(\vec{w} \cdot \vec{x}_j + b)y_j - 1] \quad (2)$$

Note that switching the $\vec{w} \cdot \vec{x}$ order would/could result in putting the terms on the $\vec{x} \rightarrow$
 $\rightarrow (\vec{x}^T \vec{w} + b)y_j$

Sufficient Condition from convex optimization guarantees that these two optimization problems are equivalent!

$$\therefore \frac{\partial L}{\partial \vec{w}} = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i y_i \vec{x}_i \in \mathbb{R}^d = 0 \star \leftarrow$$

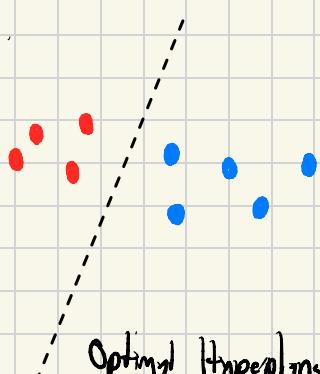
Set to 0

$$\Rightarrow \vec{w}^* = \sum_{i=1}^n \alpha_i y_i \vec{x}_i \quad \begin{matrix} \text{label} \\ \text{training sample} \end{matrix}$$

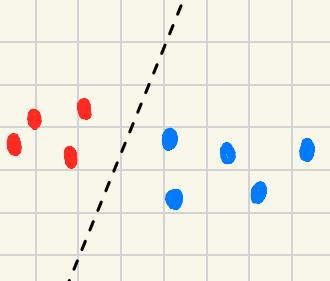
Optimal \vec{w} which is important because \vec{w} defines the margin and hence the hyperplane, so it carries all the information we need to carry on linear classification.

$\therefore \vec{w}$ depends on our training samples \vec{x}_i & their labels y_i .

E.g.:

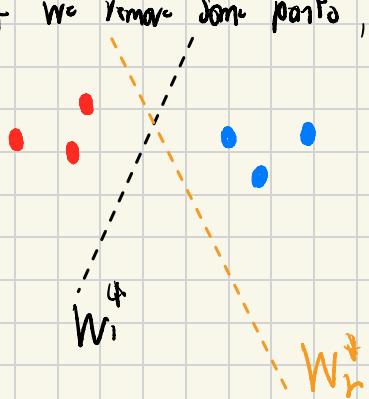


Optimal Hyperplane that Separates them (\vec{w}^*)

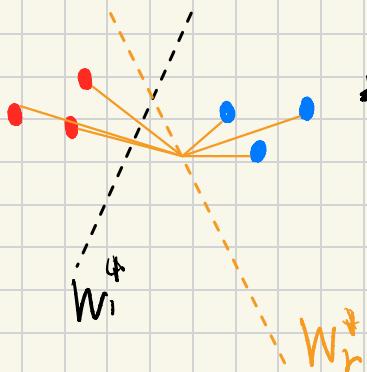


Optimal Hyperplane that Separates them (W^*)
 W^*

- If we remove some points, you can see that the optimal W is completely different.



This is very important because we know that W_r^* depends on the training samples.



* this is visible through this equation:

$$W^* = \sum_{i=1} \alpha_i y_i x^i$$

↑ training samples

So now, let's try to compute the "b"

$$\mathcal{L} = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^i (x^T w + b) - 1]$$

(3) Compute gradient of \mathcal{L} wrt to b :

$$\frac{\partial \mathcal{L}}{\partial b} = 0 - \sum_{i=1}^n \alpha_i y_i = 0$$

Note that $\frac{d \alpha x}{d x} = \alpha$

(4) Set gradient $\frac{\partial \mathcal{L}}{\partial b} = 0$

$$\frac{\partial \mathcal{L}}{\partial b} = 0$$

$$\Rightarrow b \Rightarrow$$

Labels

$$\sum_{i=1}^n \alpha_i y_i = 0$$

↑
sum over all the
training points

↑
This is a new
constraint that we
have found that

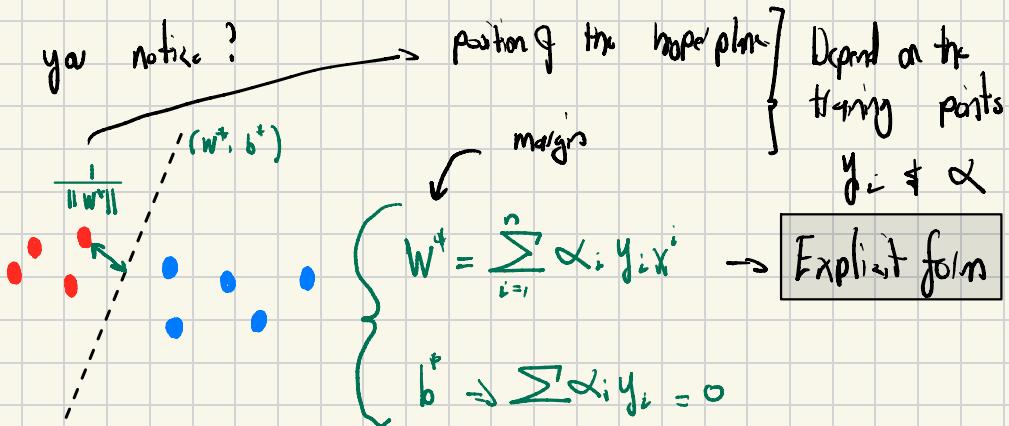
We can use to find the optimal $w^* \text{ and } b$

Remember that we need to add any new constraint that
pops up along the way e.g. $b \geq 0$.

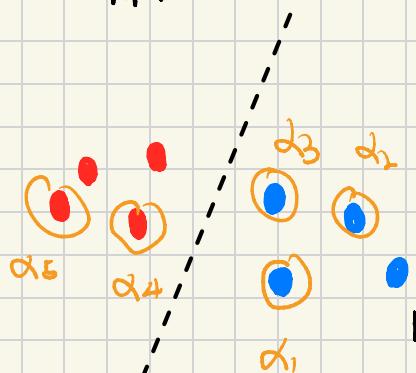
$$\left\{ \begin{array}{l} \nabla \mathcal{L} = 0 \\ \text{+} \\ \text{all constraints} \end{array} \right.$$

$$\longrightarrow \begin{array}{l} w^* \\ b^* \end{array}$$

What do you notice?



- w^* can be seen as a linear combination of all our training samples and the α_i is the weights. (Weighted sum of all training samples.)
- Not all of them will have the same weight (α s).



• We know that all the alphas are positive, but we know that they might take on different values.

• w^* depends on the samples but the samples do not have the same weight.

! finding the optimal parameters: \Rightarrow They contribute differently to finding w^* after finding the optimal parameters:

$$f(w^*, b^*) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i (w^T x_i + b) - 1]$$

$w^T w$

\Rightarrow plug w^*

$f(w, b)$ = w . Now need to compute the

Lagrangian: goal is to work with an PRIMAL

$$\min_w \mathcal{L}(w) \text{ s.t. } f(w) \geq c$$

→ Need to max the Dual Problem.

thru $\max_{\alpha} \mathcal{L}(w^*, \alpha)$ We want the value that Lagrangian takes at the optimal w to be max by

finding the best way

(Primal)

$$\min_{w, b} \max_{\alpha \geq 0} \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [(w \cdot \vec{x}_i + b)y_i - 1]$$

(1)

α^*

Swap min and max

(Dual)

$$\max_{\alpha \geq 0} \min_{w, b} \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [(w \cdot \vec{x}_i + b)y_i - 1]$$

(2)

↳ Finding the best set of Lagrangian Multipliers → The

first step will be to compute the Lagrangian at the optimum, & then we try to find the maximum of the Lagrangian with respect to the α .

$$\mathcal{L}(w^*, b) = \frac{1}{2} \left(\sum_{i=1}^n \alpha_i y_i \vec{x}^T \right) \left(\sum_{j=1}^n \alpha_j y_j \vec{x}^j \right) -$$

squared norm therefore we need to change the index:

$$\left(\sum_{i=1}^n \alpha_i y^i \vec{x}^{iT} \right) \sum_{j=1}^n \alpha_j y_j \vec{x}^j - \sum_{i=1}^n \alpha_i y_i b$$

$$+ \sum_{i=1}^n \alpha_i$$

so we know that the sum = 0 & "b" is a constant.

since it also does not depend on the index.

Result that: $\sum \alpha_i y_i = 0$ so that entire term goes away.

$$L(w^*, b) = \frac{1}{2} \left(\sum_{i=1}^n \alpha_i y_i x^i \right) \left(\sum_{j=1}^n \alpha_j y_j x^j \right) -$$

Squared norm therefore we need to change the index:

$$\left(\sum_{i=1}^n \alpha_i y_i x^i \right) \left(\sum_{j=1}^n \alpha_j y_j x^j \right) - \underbrace{\sum_{i=1}^n \alpha_i y_i b}_{+ \sum_{i=1}^n \alpha_i}$$

and also, then the terms are equal,

$$\begin{aligned} L(w^*, b^*) &= \frac{1}{2} \left(\sum_{i=1}^n \alpha_i y_i x^i \right) \left(\sum_{j=1}^n \alpha_j y_j x^j \right) - \\ &\quad \left(\sum_{i=1}^n \alpha_i y_i x^i \right) \left(\sum_{j=1}^n \alpha_j y_j x^j \right) - \underbrace{\sum_{i=1}^n \alpha_i y_i b}_{+ \sum_{i=1}^n \alpha_i} \\ \therefore L(w^*, b^*) &= -\frac{1}{2} \left(\underbrace{\sum_{i=1}^n \alpha_i y_i x^i}_{\text{Rebels}} \right) \left(\underbrace{\sum_{j=1}^n \alpha_j y_j x^j}_{\text{Dot Product}} \right) + \sum_{i=1}^n \alpha_i \end{aligned}$$

$$[\text{Ex: } Y_n x - x = -Y_n x]$$

$$\therefore L(w^*, b^*) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x^i x^j$$

\hookrightarrow [Explicit formula for regression evaluated at optimal parameters] between training product symbols

We now want

$$\max_{\alpha_i} L(w^*, b^*, \alpha_i) \Rightarrow w^* = \sum \alpha_i y_i x^i$$

↓
to find

We have not explicitly computed it.
Therefore dependencies

set α_i is We want to maximize it and find the best w . It depends on α_i , that will enable us to compute the bias as it depends on those α_i 's.

$$\therefore w^* = \sum \alpha_i y_i x^i$$

[Need to max the regions and then it will]
enable us to find α_i

∴ Looking at the regions at the optimal point:

$$L(w^*, b^*) = \sum_{i=1}^n \alpha_i - \gamma_2 \sum \sum_j \alpha_i \alpha_j y_i y_j x^{i\top} x^j$$

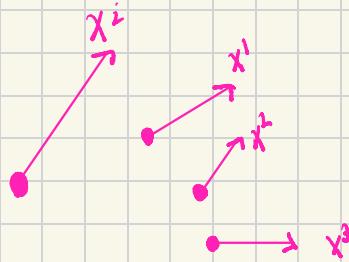
What does it depend on?

[It depends on the labels of the training points and dot product between two training samples.]

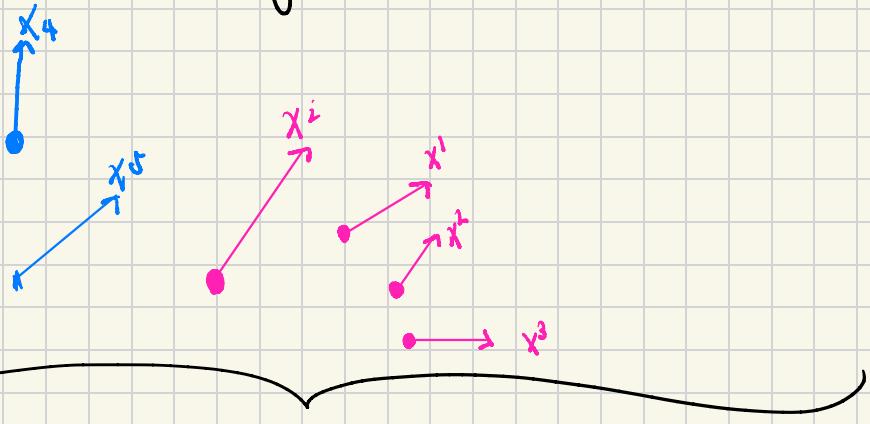
Dot product \Rightarrow multiplying one sample onto another sample.

Modelling the pairwise relationship between 2 samples is very important.

Eg:
This is our sample
 x^0 and other
samples.



And we have samples belonging to the other class.



We have a double summation, so this means that

$$L(w^*, b^*) = -\gamma_2 \left(\sum_{i=1}^n \alpha_i y_i x^{i \top} \right) \left(\sum_{j=1}^m \alpha_j y_j x^{j \top} \right) + \sum_{i=1}^n \alpha_i$$

$$L(w^*, b^*) = \sum_{i=1}^n \alpha_i - \gamma_2 \sum_{i=1}^n \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x^{i \top} x^j$$

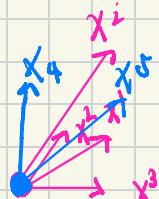
We are merging the double summations.

$$L(w^*, b^*) = -\gamma_2 \left(\sum_{i=1}^n \alpha_i y_i x_i^\top \right) \left(\sum_{i=1}^n \alpha_i y_i x_i^\top \right) + \sum_{i=1}^n \alpha_i$$

Merge them together.

$$L(w^*, b^*) = \sum_{i=1}^n \alpha_i - \gamma_2 \sum_i \sum_i \alpha_i \alpha_i' y_i y_i^\top x_i^\top x_i$$

∴ Considering the pairwise projection between all possible pairs of training samples.



∴ x^i , projecting all other samples onto sample x^i and summing all those projections for each point

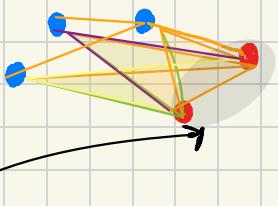
$$L(w^*, b^*) = \sum_{i=1}^n \alpha_i - \gamma_2 \sum_i \sum_i \underbrace{\alpha_i \alpha_i' y_i y_i^\top}_{\text{sum}} \underbrace{x_i^\top x_i}_{\text{Weight}}$$

The dot product tells us how two vectors are related. If they are close forming similar magnitude, opp/ same direction etc. The dot product between 2 vectors $x^i \cdot x^j$ models the relationship/captures the relationship between the features of 2 training points.

$$\stackrel{\circ}{\therefore} x^i \cdot x^j = x_{i1} x_{j1}, x_{i2} x_{j2}$$

↑ feature ↑ feature ↑ feature ↑ feature

$$\in \mathbb{R}^2$$

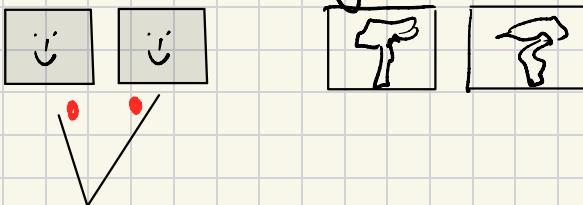


} comparing the pairwise relationship
between all the samples.

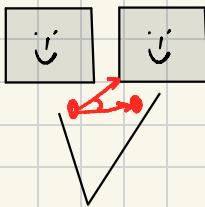
If they are too close, their dot product will tell us how close they are.

is 2 points that might be very similar,

for example classifying images of the same finger



These fingers will be similar, so I will expect the dot product between them to be quite large \rightarrow They are similar. is



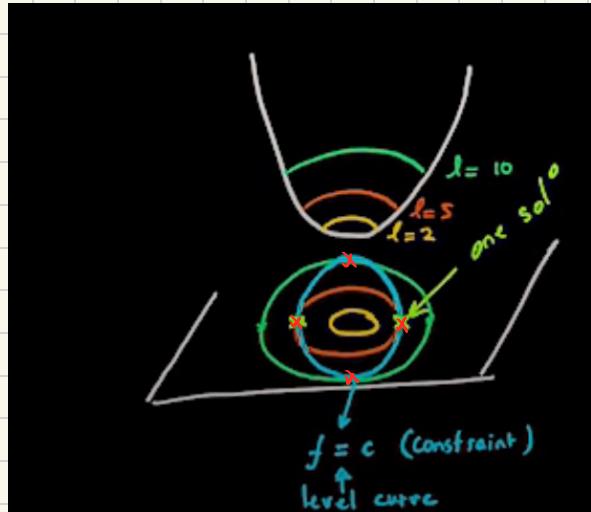
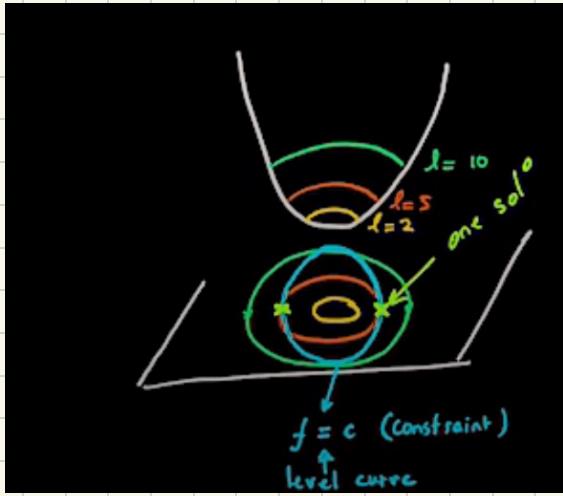
These 2 points shall also be similar.

What happens to L at the optimum (w^*, b^*) ?

Remember our goal is to find the α .

$$\Rightarrow L(w, b) = d(w, b) - \sum_{i=1}^n \alpha_i f_i(w, b)$$

$$L(w^*, b^*) = ?$$



We have different solutions and not all of them are right. So to find the optimal solution, we need to also get the α 's. That will help us find the optimal one because the primal problem has different solutions so by optimizing over the α , we will refine our first solution.

$$L(w^*, b^*) = ?$$

↑ ↑

We will find the right conditions to refine these solutions.

$$\max_{\alpha} L(w^*, b^*) = \underbrace{L(w^*, b^*)}_{\text{min}} \quad \text{solve the primal problem}$$

$$- \sum_{i=1}^n \alpha_i f_i(w^*, b^*)$$

↓
positive

strictly positive

∴ We know that (w^*, b^*) hasn't minimized our loss function

$$\max_{\alpha_i} L(w^*, b^*) = \underbrace{L(w^*, b^*)}_{\substack{\min \\ \text{solves the} \\ \text{primal problem}}} - \sum_{i=1}^n \alpha_i f_i(w^*, b^*)$$

↓
over all α_i

↑
strictly positive

What is the maximum that the Lagrangian at the optimum that we computed earlier can take?

$$\frac{1}{2} \|w\|^2$$

$$\max_{\alpha_i} L(w^*, b^*) = \underbrace{L(w^*, b^*)}_{\substack{\text{minimum } w \rightarrow 0}} - \sum_{i=1}^n \alpha_i f_i(w^*, b^*)$$

≥ 0

↑

$$\Rightarrow \underbrace{L(w^*, b^*) - \sum_{i=1}^n \alpha_i f_i(w^*, b^*)}_{\geq 0} \geq 0$$

(A) term

i. Subtracting something positive from the L quantity will always decrease it.

ii. To maximize the Lagrangian over the α , we ideally want to set the term to 0.

$$\sum_{i=1}^n \alpha_i f_i(w^*, b^*) \geq 0$$

] Set this term to zero.

(A) term

iii. The maximum value that it will take is going to be

W will

$$\{(w^*, b^*)\}$$

minimum value

≥ 0

is since we are multiplying two positive quantities α_i & f_i .
If we set them to zero, one of them should be equal to zero or both equal to zeros.

We can't have them both $= 0$. Not possible in our case.

$$\text{since } \begin{cases} \alpha_i = 0 \\ f_i = 0 \end{cases} \Rightarrow A = \begin{cases} \alpha_i = ? \\ f_i(w^*, b^*) = ? \end{cases} \quad \text{or}$$



We know that those conditions should always be satisfied (along positive constraint of the optimal point).

is for α_i strictly positive

$$\rightarrow \text{for } \alpha_i > 0 \text{ & } f_i(w^*, b^*) = 0$$

\Rightarrow

$$A = 0$$

$$\alpha_i = 0$$

or

$$f_i(w^*, b^*) = 0$$

$$\alpha_i > 0$$

$$A = 0$$

$$a_i = 0 \Rightarrow$$

$$f_i > 0$$

(2)

$$\underline{f_i(w^*, b^*) = 0}, \quad \alpha_i > 0$$



$$\text{for } \alpha_i > 0 \Rightarrow$$

$$f_i(w^*, b^*) = 0$$

(1)

$$w^* = \sum_{i=1}^n \alpha_i y_i x_i$$

\Rightarrow Result

$$y^*(x^{*T} w + b) - 1 \geq 0$$

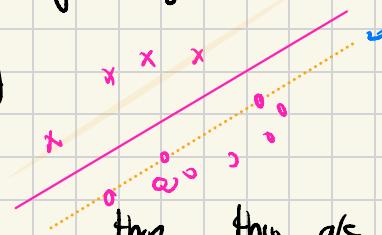
$\forall i \in \{1, \dots, m\}$

constraints

- $\downarrow = 0 \Rightarrow$ The points should lie on the hyperplane.

We know that point lie on the margin hyperplane

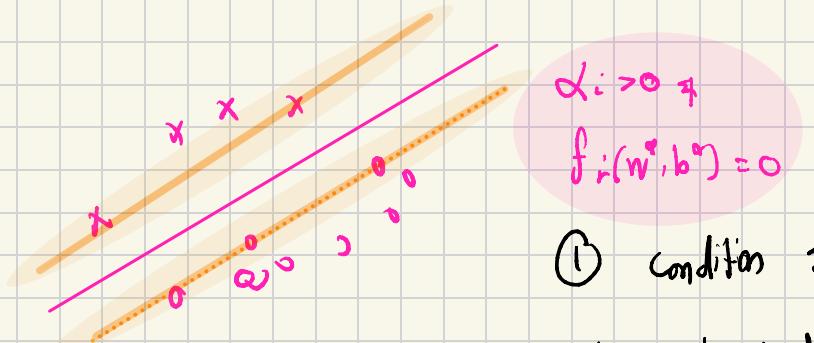
they satisfy that inequality by definition



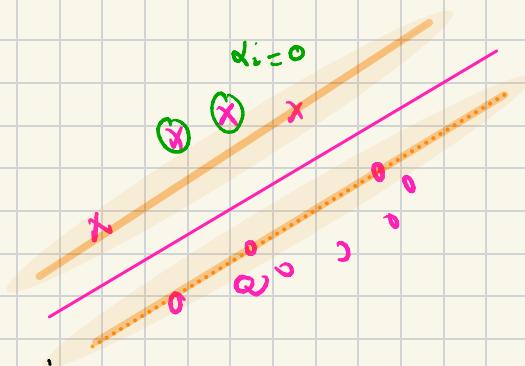
- If they are greater than zero, then they are correctly classified. By greater, we mean strictly greater than zero.

- This condition states that those training points x^* should lie on the margin plane.

- For points that lie on the margin hyperplane the α_i 's should be strictly positive



and for the others, when they are correctly classified,
we know that the $\alpha_i = 0$



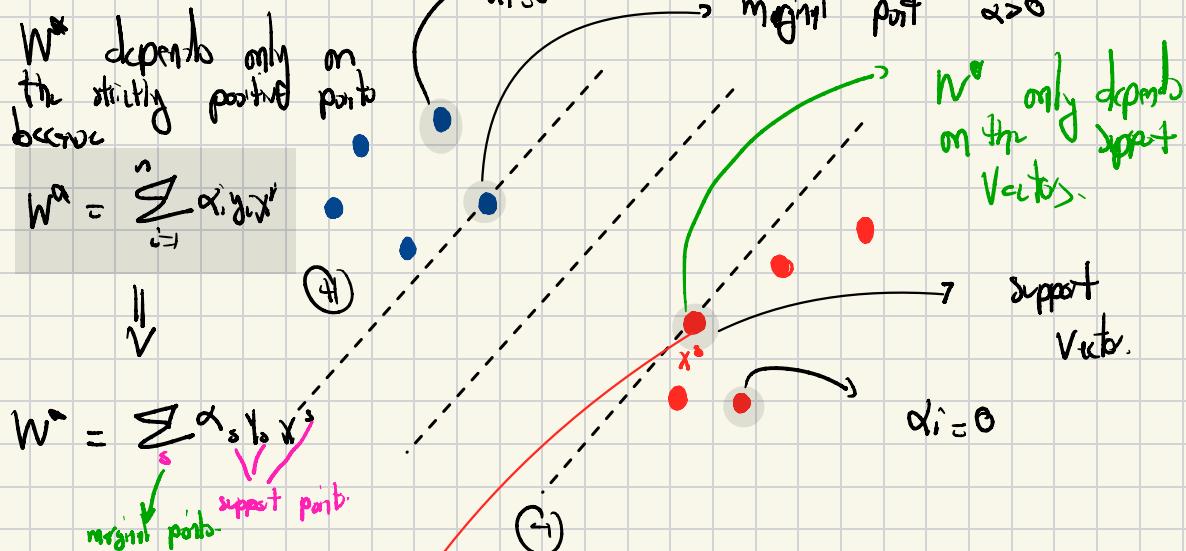
so all other points that do not lie on the margin planes
they have an $\alpha_i = 0$.

Again, looking at our optimal w^* which depends on the

α_i

$$w^* = \sum_{i=1}^n (\alpha_i y_i) e_i$$

know that if α_i is strictly positive ($\alpha_i > 0$), then the point x^i lies on the margin hyperplane and we know that if $\alpha_i = 0$, then x^i is not margin.



Now quality if we want this point to satisfy (support point)

Estimate b^* given $W^* = \sum_{i=1}^n \alpha_i y_i x_i$ using margin points.

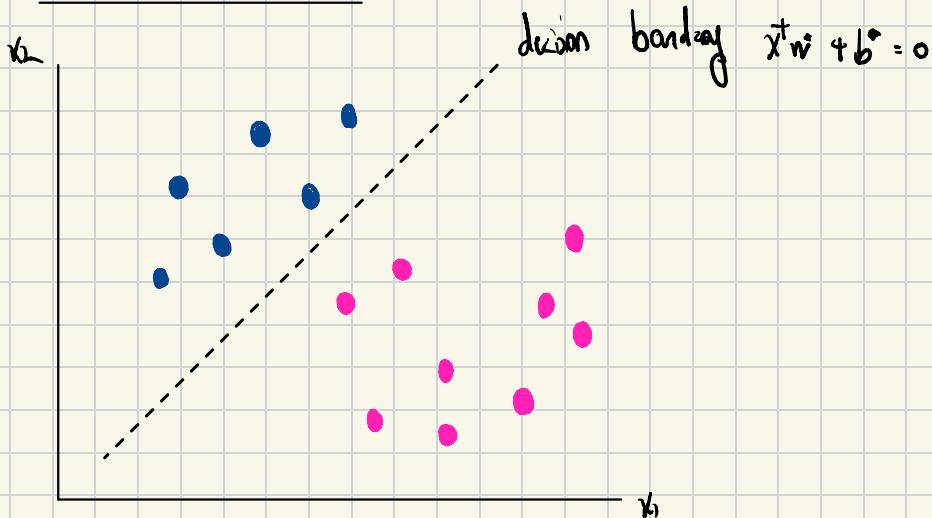
$$\Rightarrow x^T W^* + b^* = y \Rightarrow b^* = y - x^T W^*$$

A support vector example

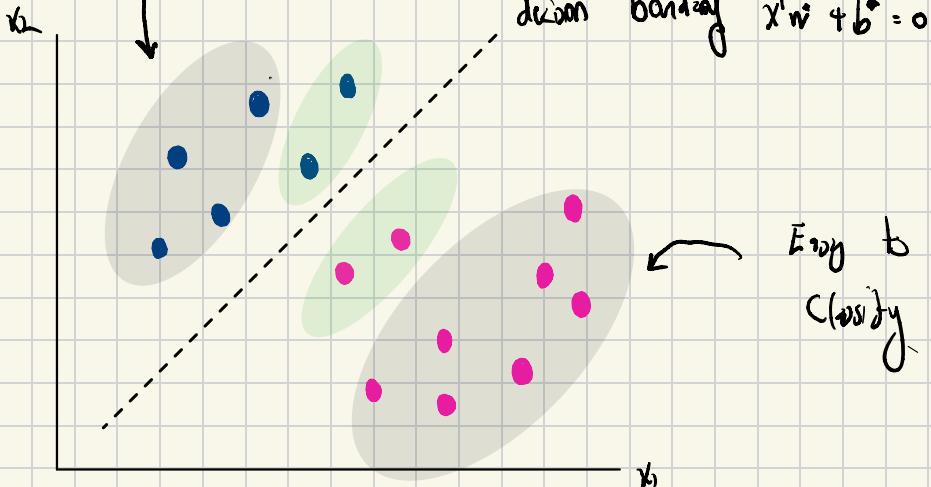
Dot product between the Support Vector / Sample and the optimal W^*

$$\Rightarrow W^* = \sum_s \alpha_s y_s x_s$$

This is intuitive



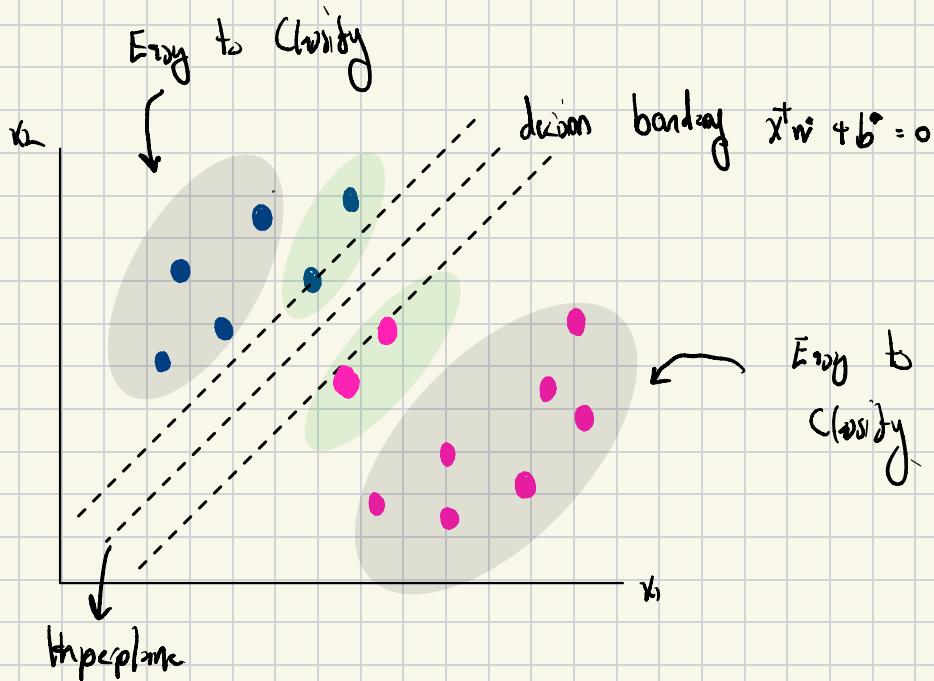
Easy to Classify



The most difficult ones are the points close to the boundary.

If you identify the marginal points (points close to the margin), then, you have solved the problem because you now know

that those vectors/points define the margin.



You are thus looking for the extreme point marginal points that lie on the hyperplane

- i. ① What are the samples that are most difficult to classify?
- support samples that lie on the margin hyperplane
- ② What directly affects the position of finding the optimum position of the decision boundary!
[w^* depends on?]
 - w^* only depends on the support vectors. If we can find them, then our problem is solved.

This is not always the case:

facto try that this is a training matrix: X^{tr}

$$X^{tr} = \begin{bmatrix} \text{class 1} \\ \vdots \\ \text{class 2} \end{bmatrix}$$

training sample

and we want to map them from the parameters W^*, b^*
which maps them into the fitted space.

$$X^{tr} = \begin{bmatrix} \text{class 1} \\ \vdots \\ \text{class 2} \end{bmatrix} \xrightarrow{W^*, b^*}$$

$$\begin{bmatrix} \text{fitted} \\ \vdots \\ \vdots \end{bmatrix}$$

but it's not
obvious... you can't

easily find the support points because when you look at them, there are
a plethora of things to keep in mind

Relationship

- ↳ feature vectors,
- ① samples, pairwise relationships

② We don't know which vectors lie on the margin ~ We are looking for those S.V.s.

so Only a very small subset of training samples (S.V.)
(n) fully specify the decision function.

$$W^* = \sum_s \alpha_s y_s x^s, \quad \delta = \delta \cdot n, \quad \alpha > 0$$

Testing stage:

• What is the decision rule for a new testing sample?

$$\text{Sign}(x^{\text{test}}^T w^\alpha + b^\alpha)$$

$$\text{Sign} \left(x^{\text{test}}^T \sum_s \alpha_s y_s x_s + b^\alpha \right)$$

↓ Compat. sign ↓

tells us if point

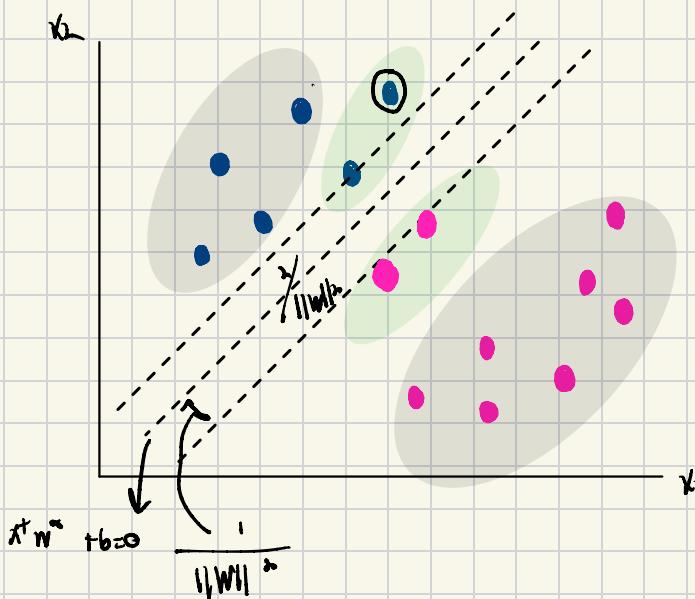
lies above

below H.P.

within the

marg

$$\frac{\sigma}{\|w\|_2}$$



- Key Take Aways :
- ① The dot product allows for the quantification of the similarity between two samples.
 - ② If 2 samples are very similar in features (dot product ~ 1), then they would have similar labels (ie lie on the same side with respect to the classifier hyperplane).
 - ③ Quantification of the similarity between samples remains a vibrant research area in machine learning.
 - ④ Kernels can be used to define a similarity between two samples.
 - ⑤ Once we find optimal w^* & b^* , you can determine if the support vector (sample) lies on the $+1$ or -1 margin hyperplane by checking the sign of $(x^T w^* + b^*)$. If it's positive then it lies on the $+1$ hyperplane, otherwise it lies on the -1 hyperplane.