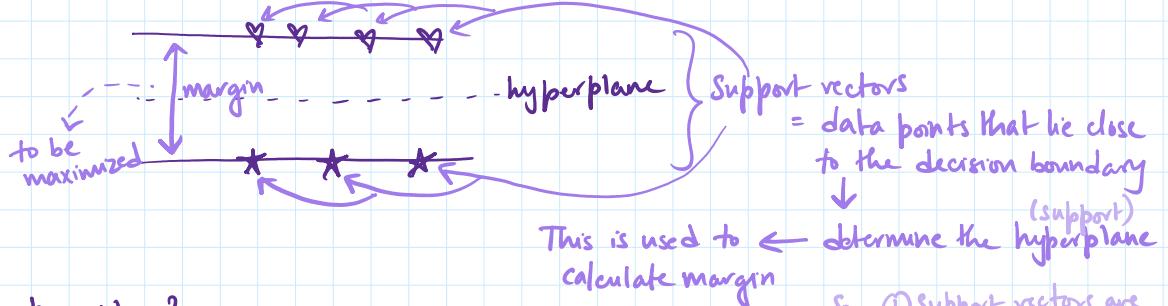


Margin = dist. b/w decision boundary + data point of each class.

SVM tries to MAXIMIZE MARGIN $\xrightarrow{\text{how}}$ FINDING OPTIMAL HYPERPLANE.



Remember the perceptron?

- It was guaranteed to converge
- gave us a LINEAR hyperplane
- acc. to the weights, all hypotheses treated equally.

Problem

The perceptron couldn't find an optimal hypothesis
Big issue.

This meant that each time you ran the perceptron,
you'd get a new hyperplane!

Say you're a two dimensional input pt (\bar{x}_1, \bar{x}_2)
 \Rightarrow One dimensional hyperplane (straight line)

But what does it mean for a random point (not) on hyperplane?

DISTANCE OF \bar{x} FROM HYPERPLANE

So say we've (\bar{x}, y)
+ hyperplane defined by $\bar{w} + b$.

$$\Rightarrow |\beta| = |\bar{w}\bar{x} + b| \quad \text{where } \beta \text{ is the distance of } \bar{x} \text{ from hyperplane}$$

$$\bar{w}_2 = \bar{w}_1 + b$$

$$\Rightarrow \bar{w}_1 - \bar{w}_2 + b = 0$$

$$\Rightarrow \begin{bmatrix} \bar{w}_1 \\ -1 \end{bmatrix} \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix} + b = 0$$

$$\Rightarrow \bar{w} \cdot \bar{x} + b = 0$$

Holds true for all points on the hyperplane.

Say you have a dataset $D = \{(\bar{x}_i, y_i) | \bar{x}_i \in \mathbb{R}^n, y_i \in \{-1, 1\}\}_{i=1}^m$

Let $B = \text{smallest } \beta$ over all m samples

Don't get lost:

β = distance of pt. from hyperplane.

B = min distance of all samples from given hyperplane

m = no. of samples

k = no. of hyperplanes.

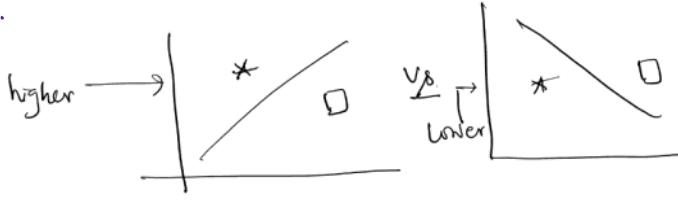
What are we doing?

for (hyperplane in hyperplanes):
for (sample in samples):

J. for (hyperplane in hyperplanes):
 for (sample in samples):
 calculate β .
 return min β as B . — ① get min distance B of all samples.
 return max β — ② get max β of all hyperplanes.

separation of the classes.

Now, look at this:



Both classify correctly.
 So what's up?

They're producing opposite signs!

Oki. Let's resolve.

Our training sample is of the form (\bar{x}, y)

What does f do?

- When we looked @ the margin, we were looking @ how well separated points were.

- f looks at that AND whether those classifications are correct.

More nuance:

- point is correctly classified + far from decision boundary
 ↳ large, +ve

- misclassified → -ve

- point is on or close to db. (correctly classif) → small, +ve

① Multiply β by y to get a new number f :

$$f = y \times \beta$$

$$f = y(\bar{w}\bar{x} + b)$$

} who is f ? THE FUNCTIONAL MARGIN OF THE SAMPLE.

② Now we compare:

y is +ve, $(\bar{w}\bar{x} + b)$ is +ve | f is +ve

y is -ve, $(\bar{w}\bar{x} + b)$ is -ve | f is +ve

Both cases imply correct classification

↓
 f is +ve for corr. classif.

y is -ve, $(\bar{w}\bar{x} + b)$ is +ve | f is -ve

y is +ve, $(\bar{w}\bar{x} + b)$ is -ve | f is -ve

Both cases imply wrong classification

↓
 f is -ve for wrong classif.

③ SO, for a given dataset D that has m samples:

$$F = \min_{i=1 \dots m} f_i$$

$$F = \min y_i (\bar{w}\bar{x}_i + b)$$

} who is F ? THE FUNCTIONAL MARGIN OF THE DATASET

④ But why take minimum?

First understand the functional margin.

- It's like a fence separating two groups of animals.
- It tells you how far each animal is from the fence + if they're on the correct side of it.
- The farther an animal is from the fence (+ on the correct side), the more confident you are that the fence is placed correctly.

Got that so far?

Now, for a classifier to be effective, it should:

- classify all training samples correctly.
- classify correctly WITH CONFIDENCE

Now, for a classifier to be effective, it should:

- classify all training samples correctly.
- classify correctly WITH CONFIDENCE

What does $\min_{i=1 \dots m} f_i$ imply?

f_i is essentially the 'confidence' of the SVM in its predictions.

- if $y_i = 1$:

large +ve value of $(\bar{w}\vec{x} + b)$ gives large, +ve f_i .
⇒ Strong confidence that point is +ve.

- if $y_i = -1$

large -ve value of $(\bar{w}\vec{x} + b)$ gives large -ve f_i .
⇒ Strong confidence that point is -ve.

Why $\min f_i$ then? Don't we want most confidence?

Why is the sample that's closest to hyperplane matter?

By finding $\min f_i$, we're finding the sample CLOSEST to the hyperplane.

a) It represents the hardest case for classification

I mean, it's cuz they're so close that the confidence is so low.

b) By focusing on this minimum margin, we ensure that challenging examples are classified correctly.

Now one more issue omg.

SCALE INVARIANCE & ISSUE.

The hyperplane is orthogonal to \bar{w} .

The hypothesis is $\bar{w}\vec{x} + b$.

If $|\bar{w}|$ or $|b|$ are of higher magnitude, the hypothesis doesn't change.
It just gets scaled ↑ or ↓.

\bar{w} still represents the same hyperplane.

To fix this scaling up & down business, we'll take the unit vector of \bar{w} .

From before, $f = y \cdot (\bar{w}\vec{x} + b)$

↓ get unit vector now

Magnitude wise,
 $GM = \frac{FM}{|\bar{w}|}$ functional margin
geometric margin

Geometric margin of sample

So per usual, for the dataset, it becomes:

Geometric margin of dataset

$$M = \min_{i=1 \dots m} \gamma_i$$

Smallest margin over dataset.

$$i.e. M = \min_{i=1 \dots m} y_i \left(\frac{\bar{w}}{|\bar{w}|} \cdot \vec{x} + \frac{b}{|\bar{w}|} \right)$$

So, net-net:

- a) functional margin of a data point wrt the hyperplane is
a measure of confidence in the classification of that point.

$$f_i = y_i (\bar{w} \bar{x}_i + b)$$

↑ weight vector ↑ bias term

confidence

fm tells us how far \bar{x}_i is from the decision boundary in terms of the raw output of the linear function.

- b) the geometric margin is the "normalized" version of PM for scale invariance
 (by dividing by $|\bar{w}|$)

This gives us the actual distance from the data point to decision boundary.

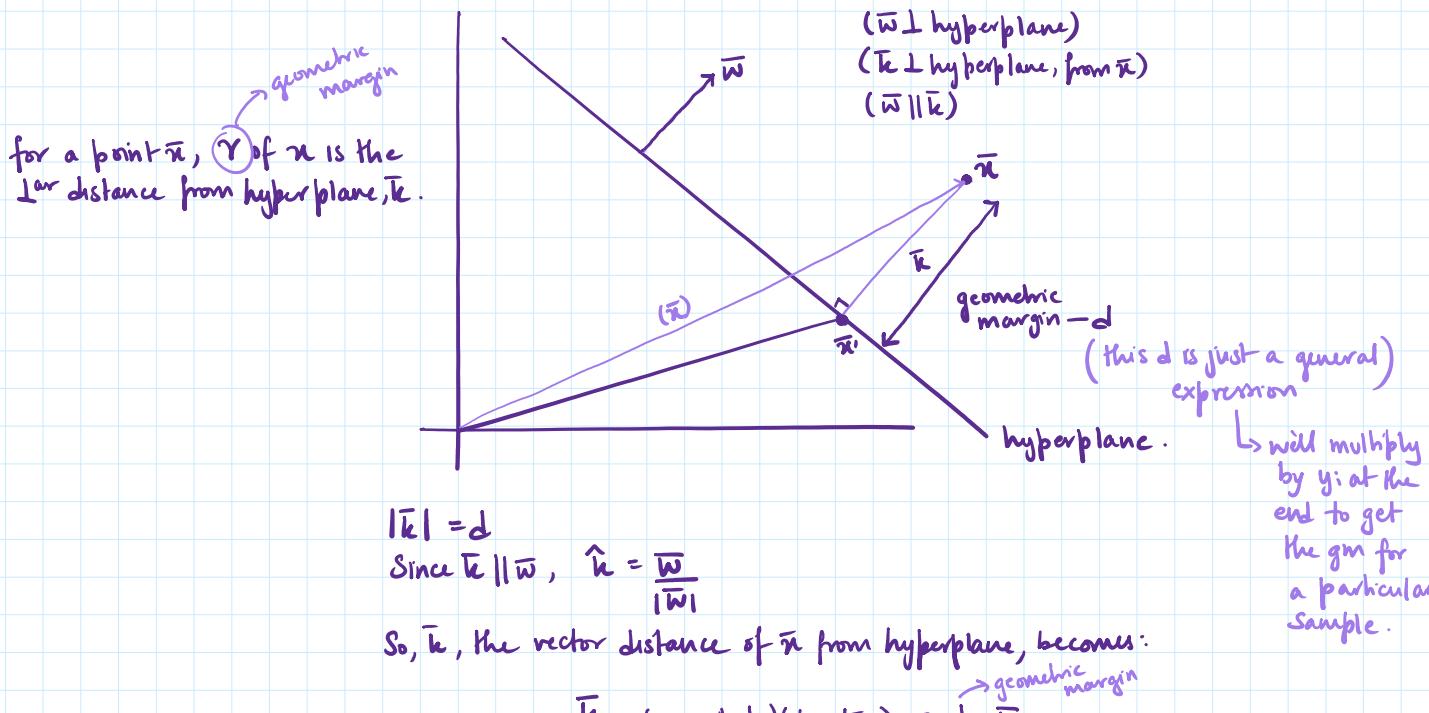
$$\gamma_i = y_i \left(\frac{\bar{w}}{|\bar{w}|} \bar{x}_i + \frac{b}{|\bar{w}|} \right)$$

distance

gives the 1st distance from the point to the hyperplane.

More direct measure of distance of pt. from hyperplane.

Now, let's look at it geometrically:



Using the Δ^k law of vector addition:

$$\bar{x} = \bar{x}' + \bar{k}$$

$$\Rightarrow \bar{x}' = \bar{x} - \bar{k}$$

$$= \bar{x} - d \cdot \frac{\bar{w}}{|\bar{w}|}$$

What \bar{x}' is on the hyperplane ($\bar{w} \bar{x}_i + b = 0$)
 plug in \bar{x}' for \bar{x} ,

$$\bar{w} \left(\bar{x} - d \cdot \frac{\bar{w}}{|\bar{w}|} \right) + b = 0$$

$$\bar{w}\bar{x} - d \frac{\bar{w} \cdot \bar{w}}{|\bar{w}|} + b = 0.$$

$$\bar{w}\bar{x} - d \cdot \frac{|\bar{w}|^2}{|\bar{w}|} + b = 0$$

$$\bar{w}\bar{x} - d |\bar{w}| + b = 0$$

$$d |\bar{w}| = \bar{w}\bar{x} + b$$

$$d = \left[\frac{\bar{w}\bar{x}}{|\bar{w}|} + \frac{b}{|\bar{w}|} \right] \quad \text{Now we multiply by } y \text{ so that we get a hyperplane that correctly classifies data.}$$

$$Y_i = y \left[\frac{\bar{w}\bar{x}}{|\bar{w}|} + \frac{b}{|\bar{w}|} \right]$$

And then when we find the g_m for another plane, we just change the values of \bar{w}, b .

Oof.

Now how do we find the \bar{w} that gives us the highest g_m ?

Let's look at optimizations :

Unconstrained — if we're to minimize $f(x)$, here, will search ALL possible values of x .

Constrained — let's work with the example, $\min f(x)$, such that $x=2$.
(st)

$$\begin{aligned} &\text{We write it as } \min f(x) \\ &\text{st } x=2 \end{aligned}$$

The feasible set will be the values of x that can optimize $f(x)$.

for $\min f(x)$ st $x=2$, the feasible set is the set of all points (x, y) such that $(x, y) \in (2, y)$

Now, we can have many constraints too,

look at this: $\min f(x, y, z)$ st $x=2, y=-3, z=-3$

now rewrite it as:

$$\min f(x, y, z)$$

$$\text{st } x=2, y=-3, z=-3$$

↓ or just represent it as a vector!

$$\min f(\bar{x})$$

$$\text{st } x_1=2$$

$$x_2=-3$$

$$x_3=-3$$

Now,

SVM Optimization Problem.

① What we have is a linearly separable dataset D .

$$D = \left[(x_i, y_i) \mid x_i \in \mathbb{R}^n, y_i \in \{+1, -1\} \right]$$

① What we have is a linearly separable dataset D .

$$D = \left\{ (\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in \mathbb{R}^n, y_i \in \{+1, -1\} \right\}$$

and also vector \bar{w} + bias b .

② Let the gm of the dataset $M = \min_{i=1 \dots m} Y_i \rightarrow$ dataset has m samples.

$$\text{and what's } Y_i? Y_i = y_i \left(\frac{\bar{w}}{\|\bar{w}\|} \cdot \mathbf{x}_i + \frac{b}{\|\bar{w}\|} \right)$$

gm of a sample
 i among the
 m samples in
the dataset.

③ We also know that the Optimal hyperplane is the hyperplane defined by $\bar{w} + b$ such that M is the LARGEST

where M is the gm over the dataset D .

Now, to find $\bar{w} + b$, we solve the optimization problem

↳ Constraint: margin of each sample Y_i should be greater than or equal to M .

$$\text{maximize } M \text{ such that } Y_i \geq M, i=1 \dots m$$

$$\text{From before, the gm is the fm divided by } \|\bar{w}\|. \rightarrow M = \frac{F}{\|\bar{w}\|}$$

Let's rewrite them.

The optimization problem becomes:

$$\text{Maximize } M \text{ such that } \frac{f_i}{\|\bar{w}\|} \geq \frac{F}{\|\bar{w}\|}, i=1 \dots m$$

We can simplify further.

$$\text{Maximize } M \text{ such that } f_i \geq F, i=1 \dots m.$$

Hurrah. What does this tell us, though?

In the hyperplane, if we scale up \bar{w}, b :

→ The hyperplane direction DOES NOT CHANGE.

→ Only FUNCTIONAL MARGIN gets bigger.

If we scale down \bar{w}, b such that $F=1$ — OPTIMAL GM REMAINS SAME.

Let's check it out:

$$\rightarrow \text{maximize } M \text{ st } f_i \geq F.$$

$$\rightarrow \text{maximize } \frac{F}{\|\bar{w}\|} \text{ st } f_i \geq F$$

! now, F is scaled to 1. !

$$\rightarrow \text{maximize } \frac{1}{\|\bar{w}\|} \text{ st } f_i \geq F$$

→ Maximizing $\frac{1}{\|\bar{w}\|} = \text{Minimizing } \|\bar{w}\|.$

$$\text{minimize } \|\bar{w}\| \text{ st } f_i \geq 1.$$

$$\text{Minimizing } \|\bar{w}\| = \text{Minimizing } \frac{1}{2} \|\bar{w}\|^2$$

minimize $\|\bar{w}\|$ st $f_i \geq 1$.

$$\text{Minimizing } \|\bar{w}\| = \text{Minimizing } \frac{1}{2} \|\bar{w}\|^2$$
$$\text{minimize } \frac{1}{2} \|\bar{w}\|^2 \text{ st } f_i \geq 1$$

$$f_i = y_i(\bar{w}x_i + b)$$

$$\Rightarrow \text{minimize } \frac{1}{2} \|\bar{w}\|^2 \text{ st } y_i(\bar{w}x_i + b) - 1 \geq 0$$

Let's look @ CONSTRAINED OPTIMIZATION now:
↓
Lagrange Multiplier.

$$\text{Ex 1: } f(x, y) = 3x + 4y \quad \rightarrow \text{Objective}$$
$$\text{st } g(x, y): x^2 + y^2 = 1 \quad \rightarrow \text{constraint}$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial x} \\ \frac{\partial g}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

A contour plot is a CONSTANT VALUE SURFACE.

- ↳ Since all points in the contour plot have the same value.
- ↳ rate of change in tangent to contour
- ↳ BUT, gradient \perp contour.

Since we're doing $f(x, y)$ and its constraint $g(x, y)$, $\nabla g \parallel \nabla f$

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

$$\nabla f(x, y) - \lambda \nabla g(x, y) = 0.$$

Now, let's assume $L(x, y, \lambda) = f(x, y) - \lambda g(x, y)$

$$\Rightarrow \nabla L(x, y, \lambda) = \nabla f(x, y) - \lambda \nabla g(x, y)$$

primal objective f^* — $f(x, y)$ ← We then solve for $\nabla L = 0$ since $\nabla f(x, y) - \lambda \nabla g(x, y) = 0$.
Lagrangian dual — $L(x, y, \lambda)$

Let's get back to our example now:

$$f(x, y) = 3x + 4y$$
$$g(x, y): x^2 + y^2 = 1$$

$$\nabla L = \nabla f - \lambda \nabla g = 0$$

$$[3] - \lambda [2x] = 0$$

$$\nabla L = \nabla f - \lambda \nabla g = 0$$

$$\begin{vmatrix} 3 \\ 4 \end{vmatrix} - \lambda \begin{bmatrix} 2x \\ 2y \end{bmatrix} = 0$$

$$\begin{vmatrix} 3 \\ 4 \end{vmatrix} = \lambda \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

$$\begin{aligned} 3 &= \lambda 2x \rightarrow x = \frac{3}{2\lambda} \\ 4 &= \lambda 2y \rightarrow y = \frac{4}{2\lambda} \end{aligned}$$

Now, what did the constraint say? $x^2 + y^2 = 1$
plug in the new values of x, y

$$\left(\frac{3}{2\lambda}\right)^2 + \left(\frac{4}{2\lambda}\right)^2 = 1$$

$$9 + 16 = 4\lambda^2$$

$$25 = 4\lambda^2$$

$$\lambda = \pm \frac{5}{2}$$

$$(x, y) = \left(\frac{3}{5}, \frac{4}{5}\right)$$

$$(x, y) = \left(-\frac{3}{5}, \frac{4}{5}\right)$$

Next: Try w/ 2 constraints:

$$\begin{aligned} f(x, y) &= x^2 + y^2 \\ g_1(x, y) &: x + y = 0 \\ g_2(x, y) &: y + 1 = 0 \end{aligned}$$

$$\text{Here, } L(x, y, \lambda) = f(x, y) - \lambda_1 g_1(x, y) - \lambda_2 g_2(x, y)$$

$$\Rightarrow \nabla L(x, y, \lambda) = \nabla f(x, y) - \lambda_1 \nabla g_1(x, y) - \lambda_2 \nabla g_2(x, y)$$

$$\underset{\nabla L \rightarrow 0}{\underset{\text{(use we set)}}{\nabla L \rightarrow 0}} = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} - \lambda_1 \begin{bmatrix} \frac{\partial g_1}{\partial x} \\ \frac{\partial g_1}{\partial y} \end{bmatrix} - \lambda_2 \begin{bmatrix} \frac{\partial g_2}{\partial x} \\ \frac{\partial g_2}{\partial y} \end{bmatrix}$$

$$\begin{bmatrix} 2x \\ 2y \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \lambda_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \lambda_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 2x \\ 2y \end{bmatrix} = \begin{bmatrix} \lambda_1 \\ \lambda_1 \end{bmatrix} + \begin{bmatrix} 0 \\ \lambda_2 \end{bmatrix}$$

$$\Rightarrow x = \lambda_1/2 \quad y = \lambda_1 + \lambda_2/2$$

Two constraints: $x + y = 0$ $y + 1 = 0$
substitute.

$$\frac{\lambda_1}{2} + \frac{\lambda_1 + \lambda_2}{2} = 0$$

$$2\lambda_1 + \lambda_2 = 0$$

$$\frac{\lambda_1 + \lambda_2}{2} = -1$$

$$\lambda_1 + \lambda_2 = -2$$

So we get a free variable extra of λ in each.

$$2\lambda_1 + \lambda_2 = 0$$

$$\lambda_1 + \lambda_2 = -2$$

so we got a two variable system of two eqns.

$$\begin{aligned} & \lambda_1 + \lambda_2 = -2 \\ & 2\lambda_1 + \lambda_2 = 0 \end{aligned}$$

$$-\lambda_1 = -2$$

$$\lambda_1 = 2$$

$$\lambda_2 = -4$$

next.

Wolfe Dual + Slater's Theorem

So from before, we're to minimize $f(\bar{w}) = \frac{1}{2}|\bar{w}|^2$ subject to m constraints of the form $g_i(\bar{w}, b)$

We'll call $y_i(\bar{w}\bar{x}_i + b)$ as $g_i(\bar{w}, b)$

$$g_i(\bar{w}, b) = y_i(\bar{w}\bar{x}_i + b) - 1 \quad i=1, 2, \dots, m$$

$$\begin{aligned} L(\bar{w}, b, \alpha) &= f(\bar{w}) - \sum \alpha_i g_i(\bar{w}, b) \\ &= \frac{1}{2}|\bar{w}|^2 - \sum \alpha_i [y_i(\bar{w}\bar{x}_i + b) - 1] \end{aligned}$$

Here, we have one lagrangian multiple α_i for every constraint g_i .

Let's solve using the dual.

The Duality Principle.

The Lagrangian problem has m inequality constraints (cuz m samples)

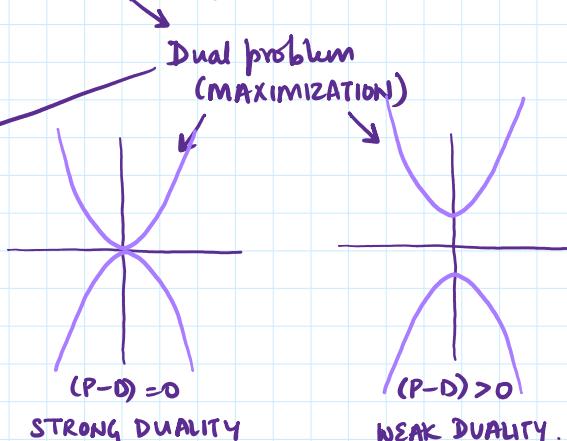
"An optimization can be solved in two ways."

Primal Problem
(MINIMIZATION)

Maximum of dual problem will be less than or equal to the minimal of the primal problem.

$$\max(\text{dual}) \leq \min(\text{primal})$$

dual provides lower bound to primal.



In convex optimization, an AFFINE CONSTRAINT is a constraint that can be expressed as

$$\underbrace{\bar{a}^T \bar{x}}_{\substack{\text{constant vector} \\ \downarrow \\ \text{decision variable}}} \leq \underbrace{b}_{\substack{\text{constant scalar} \\ \downarrow}} \quad \left. \right\} \text{An affine function is linear and constant.}$$

Slater's theorem says that strong duality holds for Slater's condition, which in turn holds for affine constraints.

Done?

Let's jump back real quick to the lagrangian function

$$L(\bar{w}, b, \alpha) = \frac{1}{2} |\bar{w}| \cdot |\bar{w}| - \sum \alpha_i [y_i (\bar{w} \bar{x}_i + b) - 1]$$

LAGRANGIAN PRIMAL PROBLEM

minimize wrt \bar{w}, b	maximize wrt α	At the same time!
---------------------------	-----------------------	-------------------

$$\min_{\bar{w}, b} \max_{\alpha} L(\bar{w}, b, \alpha) \quad \text{s.t. } \alpha_i \geq 0 \quad \forall i=1, 2, \dots, m$$

① To solve the minimization, we take PD wrt \bar{w}, b .

$$\nabla_{\bar{w}} L = \bar{w} - \sum \alpha_i y_i \bar{x}_i = 0 \quad \dots \quad (i)$$

$$\nabla_b L = -\sum \alpha_i y_i = 0$$

$$\Rightarrow \bar{w} = \sum \alpha_i y_i \bar{x}_i$$

Now plug this value into L

$$w(\alpha, b) = \frac{1}{2} \left(\sum \alpha_i y_i \bar{x}_i \right) \left(\sum \alpha_j y_j \bar{x}_j \right) - \sum \alpha_i [y_i (\sum \alpha_j y_j x_j) \cdot x_j + b] + \sum \alpha_i$$

↓ all the math later.

$$w = \sum \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \underbrace{\alpha_i \alpha_j y_i y_j}_{\text{their alphas.}} \underbrace{\bar{x}_i \cdot \bar{x}_j}_{\text{pair of samples}}$$

Optimization depends on the dot product of the pair of samples.

If $\frac{\partial L}{\partial b} = 0$, set $\sum \alpha_i y_i$ to zero!

$$w(\alpha) = \sum \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j$$

Wolfe Dual Lagrangian Function!

Use to solve primal problem.

Optimization becomes a WOLFE DUAL problem!

$$\underset{\alpha}{\text{maximize}} \sum \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j, \text{ s.t. } \alpha_i \geq 0 \text{ for } i=1 \dots m$$

and also $\sum \alpha_i y_i = 0$

$$\text{and also } \sum \alpha_i y_i = 0$$

Now, the main advantage of Wolfe-dual is that the objective function (\bar{w}) only depends on lagrangian multipliers!

ie $\bar{w} \leftrightarrow \text{lagrangian multipliers ONLY.}$

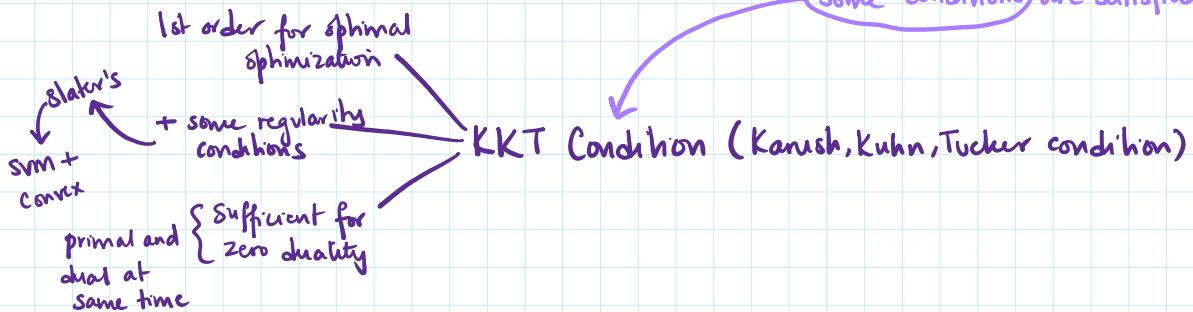
Important stuff to infer:

- ① Most of the data points really do not matter. only support vectors have non-zero α_i
- ② All datapoints \bar{x}_i are vectors. — few of them are support-vectors.
- ③ All $\alpha_i = 0$ belong to those datapoints that are FAR from decision boundary.

More about the lagrange multiplier:

- ① This method is used for problems w/ equality constraints So wassup?
BUT,
- ② We're using it with inequality constraints

This method still works for inequality provided some conditions are satisfied.



① STATIONARITY CONDITION

$$\nabla_w \alpha = w - \sum \alpha_i y_i \bar{x}_i = 0$$

- if no constraint, $\nabla_{\text{objective}} f^n$ is 0.
- w/ constraint, use lagrangian

$$\frac{\partial L}{\partial b} = - \sum \alpha_i y_i$$

② PRIMAL FEASIBILITY CONDITION

$$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 \geq 0 \quad \forall i=1 \dots m$$

- essentially, constraint of primal problem

③ DUAL FEASIBILITY CONDITION

$$\alpha_i \geq 0 \quad \forall i=1 \dots m$$

④ COMPLEMENTARY SLEEKNESS

$$\alpha_i [y_i(\bar{w} \cdot \bar{x}_i + b) - 1] = 0 \quad \forall i=1 \dots m$$

- either $\alpha_i = 0$ or $y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$
- $\alpha_i > 0$
- so

Constraint is active when $y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$.

Constraint is active when
 $y_i(\bar{w}\bar{x}_i + b) - 1 = 0$.

Now

After we solve Wolfe dual & get vector α w/ all α_i , we find $\bar{w} + b$.

Step 1: Calculate \bar{w}

$$\bar{w} = \sum_{i=1 \dots m} \alpha_i y_i \bar{x}_i$$

Step 2: calculate b

Here we use constraint of primal problem to find b .
 $y_i(\bar{w}\bar{x}_i + b) - 1 \geq 0$

For support vectors:

$$y_i(\bar{w}\bar{x}_i + b) = 1$$



$$(y_i) y_i (\bar{w}\bar{x}_i + b) = y_i$$

$$(y_i)^2 = 1$$

$$\bar{w}\bar{x}_i + b = y_i \rightarrow b = y_i - \bar{w}\bar{x}_i + b.$$

Ta daaaa

not done.

How do we find the class of an example x_j ?