

Consider an example (with equality constraints)

$$\min_{(x_1, x_2)} x_1 + x_2$$

s.t.  $x_1^2 + x_2^2 = a^2$   
 (The points are on a circle)

$$f(x) = x_1 + x_2$$

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix}$$

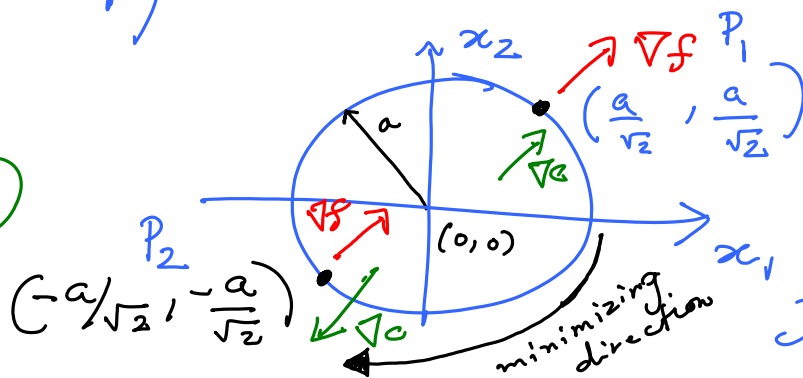
$$C = x_1^2 + x_2^2 - a^2 = 0$$

$$\nabla C = \begin{pmatrix} \frac{\partial C}{\partial x_1} \\ \frac{\partial C}{\partial x_2} \end{pmatrix} = \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix}$$

$$\nabla f = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\nabla f \begin{pmatrix} -a/\sqrt{2} \\ -a/\sqrt{2} \end{pmatrix} = (1, 1)$$

$$\nabla C \begin{pmatrix} -a/\sqrt{2} \\ -a/\sqrt{2} \end{pmatrix} = (-\sqrt{2}a, -\sqrt{2}a)$$



$\Rightarrow$  III<sup>rd</sup> quadrant has both  $x_1$  and  $x_2$  -ve  
 $\Rightarrow$  Soln lies there

Just  $\nabla f$  does not suffice for minima!

From the figure,

$$\nabla f(x^*) = \lambda_1^* \nabla c(x^*)$$
$$\lambda_1^* = \frac{-1}{a\sqrt{2}}$$

Note that:  $\nabla f$  is a scalar multiple of  $\nabla c$  @ the point of maxima as well i.e.  $(\frac{a}{\sqrt{2}}, \frac{a}{\sqrt{2}})$

Let us analyze this issue through a Jaylon Series expansion around the constraint.

---

$$c(\underline{x}) = 0 \quad ( \because \text{Equality constraint} )$$

$$c(\underline{x} + \underline{d}) = 0 \quad ( \text{To maintain feasibility w.r.t. } c(\underline{x}) = 0 )$$

$$c(\underline{x} + \underline{d}) \approx c(\underline{x}) + \underbrace{\nabla c^T(\underline{x}) \underline{d}}_{\text{inner product}} \quad ( \text{With a first order approx.} )$$

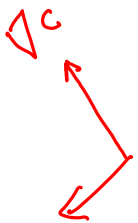
$$c(\underline{x}) + \nabla c^T(\underline{x}) \underline{d} = 0$$



$$\nabla c^T(\underline{x}) \underline{d} = 0$$

$$( \because c(\underline{x}) = 0 )$$

(A)



It by the direction of optimization must produce  
a decrease in  $f$

$f(\underline{x}) + \nabla^T f(\underline{x}) \cdot d$   
By doing a Taylor expansion around  $\underline{x}$  using  
a 1st order approx.

$$\nabla^T f(\underline{x}) \underline{d} < 0$$

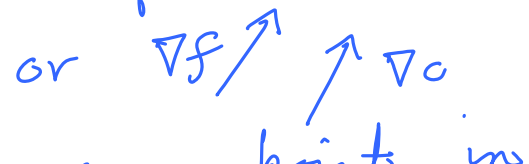

If  $\exists$  a  $\underline{d}$  satisfying (A) and (B), an  
improvement exists

There are 2 cases to consider here

(1) Such a direction does not exist

(2) Such a direction exists

Case 1 : When such a direction  $\nabla f$  and  $\nabla c$  are scalar multiples of each other  
i.e.  $\nabla f$  and  $\nabla c$  can point in the same or opposite directions



$\nabla f = \lambda \nabla c$

Ponder why: When  $\nabla f = \lambda \nabla c$   
(A) and (B) 'do not' simultaneously hold

$\therefore$  forming a Lagrangian  $\swarrow$  Lagrange multiplier

$$L = f \pm \lambda c$$

$$\nabla L = 0 \Rightarrow \nabla f = \mp \lambda c$$

$\therefore$  Sign of the "constraint" in the Lagrangian  
"does not" matter!  
Sign does not matter

We can arrive at a saddle point here

We still need the sign of the Hessian  
to proceed & assess the validity.

---

Case 2 : When such a direction exists

$$\underline{d} = - \left( I - \frac{\nabla c \nabla c^T}{\|\nabla c\|^2} \right) \nabla f \quad \textcircled{I}$$

Let us verify if  $\textcircled{I}$  satisfies  $\textcircled{A}$  and  $\textcircled{B}$

$$\underline{d} = -\nabla f + \frac{\nabla c \nabla c^T \nabla f}{\nabla c^T \nabla c} \quad \textcircled{v}$$

outer product
inner product

Let us consider  $\textcircled{A}$   
 Pre-multiply  $\textcircled{I}$  by  $\nabla c^T$ ;



$$\nabla_c^T \underline{d} = -\nabla_c^T \nabla f + \frac{\overset{\text{(Scalar)}}{\cancel{\nabla_c^T} \cancel{\nabla_c}} \nabla_c^T \nabla f}{\cancel{\nabla_c^T} \cancel{\nabla_c} \text{(Scalar)}} = 0$$

Let us consider (B)

$$\nabla_f^T \underline{d}$$

Plug in  $\underline{d}$

$$= -\nabla_f^T \left( \nabla f - \frac{\nabla_c \nabla_c^T \nabla f}{\|\nabla_c\|^2} \right)$$

$$= -\underbrace{\nabla_f^T \nabla f}_{\text{1st term}} + \underbrace{\frac{\nabla_f^T \nabla_c \nabla_c^T \nabla f}{\|\nabla_c\|^2}}_{\text{2nd term}}$$

$$= - \overset{\downarrow}{\|\nabla f\|^2} + \frac{\|\nabla^T f \nabla c\|^2}{\|\nabla c\|^2} < 0$$

(∵ Cauchy Schwartz inequality)

The equality is ruled out due to Case (A)  
 (∵  $\nabla f \neq \lambda \nabla c$ )

⇒ d is the direction satisfying the constraints.

# Single inequality constraint

$$\underline{c}(\underline{x}) \geq 0$$

$$0 \leq c(\underline{x} + \underline{d}) \approx c(\underline{x}) + \nabla^T c(\underline{x}) \underline{d}$$

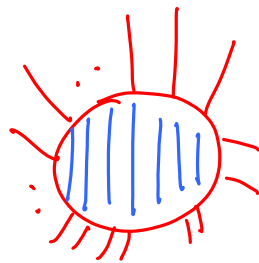
Feasibility of  $\underline{d}$  is retained while still improving the objective if

$$\underbrace{c(\underline{x})}_{\geq 0} + \nabla^T c(\underline{x}) \underline{d} \geq 0 \quad \textcircled{C}$$

Observe the  $\leq$  as against  $\geq$  in equality constraints

Considering the example from the conditions,  
circular constraints with inequality  
we are optimizing over all points lying on  $\mathcal{E}$   
inside the circle.

---



$$c(x) \geq 0$$



$$\begin{aligned} x_1^2 + x_2^2 &\leq a^2 \\ -(x_1^2 + x_2^2) &\geq -a^2 \\ a^2 - x_1^2 - x_2^2 &\geq 0 \end{aligned}$$

We have 2 cases

---

Case A : The strict inequality holds i.e.,  $c(\underline{x}) > 0$

Whenever  $\nabla f(\underline{x}) \neq 0$  i.e., when we have not yet reached optimum points

(II)

$$\left\{ \begin{array}{l} \nabla f(\underline{x})^T \underline{d} < 0 \quad \text{---} \quad (:\text{ (B)}) \\ c(\underline{x}) + \nabla c(\underline{x})^T \underline{d} \geq 0 \quad \text{---} \quad (:\text{ (C)}) \end{array} \right.$$

A  $\underline{d}$  that satisfies the constraints is  $\underline{d} = - \frac{c(\underline{x}) \nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|}$  (D)

We can verify that  $\textcircled{D}$  satisfies both the constraints in  $\textcircled{II}$

$$\begin{aligned}
 (i) \quad \nabla^T f(\underline{x}) \underline{d} &= - \nabla^T f(\underline{x}) \cdot c(\underline{x}) \frac{\nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|} \\
 &= - c(\underline{x}) \frac{\nabla^T f(\underline{x}) \nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|}
 \end{aligned}$$

Annotations:
 

- A red arrow points from the word "Scalar" to  $c(\underline{x})$ .
- A red arrow points from the text "Evaluates to  $\frac{\|\nabla f\|}{\|\nabla f\|}$ " to the fraction  $\frac{\nabla^T f(\underline{x}) \nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|}$ .
- A red arrow points from the text " $> 0$ " to  $c(\underline{x})$ .
- A red arrow points from the text " $-\frac{c(\underline{x}) \|\nabla f\|}{\|\nabla c\|}$ " to the final simplified expression.

$\Rightarrow$  First constraint in  $\textcircled{II}$  is satisfied i.e.,  $< 0$

(ii)

Consider

$$\begin{aligned} & c(\underline{x}) + \nabla^T c(\underline{x}) \underline{d} \\ \approx & c(\underline{x}) + \nabla^T c(\underline{x}) \left[ \frac{-c(\underline{x}) \nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|} \right] \end{aligned}$$

$$\approx c(\underline{x}) - c(\underline{x}) \frac{\nabla^T c(\underline{x}) \nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|} < 1$$

Unless

$$\nabla f(\underline{x}) \neq \lambda \nabla c(\underline{x}),$$

$$|\cdot| < 1$$

$$\nabla c^T(\underline{x}) \nabla f(\underline{x}) < \|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|$$

We have

$$c(\underline{x}) + \nabla c^T(\underline{x}) \underline{d}$$
$$\Rightarrow c(\underline{x}) - c(\underline{x}) \alpha$$
$$c(\underline{x})(1 - \alpha) \geq 0$$

$\alpha$  can be +ve or -ve

(The equality is only over the case when  $\alpha = 1$ )



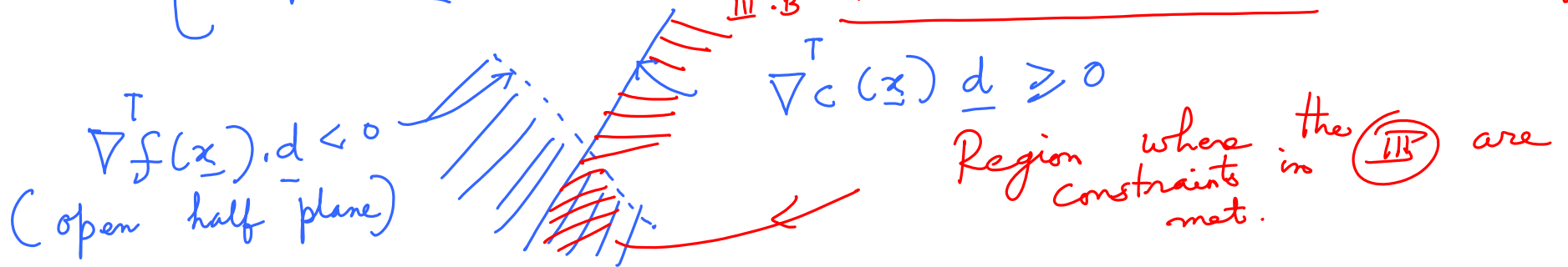
Case B : When  $\underline{x}$  is on the boundary  
of the constraints eqn i.e.,  $C(\underline{x}) = 0$

We have

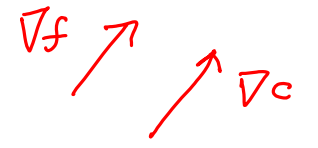
(III)  $\left\{ \begin{array}{l} \nabla^T f(\underline{x}) \underline{d} < 0 \quad \text{--- III.A} \\ \nabla^T C(\underline{x}) \underline{d} \geq 0 \quad \text{--- III.B} \end{array} \right.$

(B) boundary case  $C(\underline{x}) = 0$   
 $\left( \because \text{Plug into (C)} \right)$

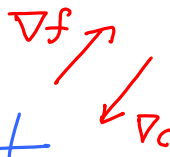
FEASIBLE SOLN REGION (GEOMETRY)



When  $\nabla f = \lambda \nabla c$  and  $\nabla f$  and  $\nabla c$  point in the same direction the regions from III do not intersect!



When  $\nabla f = -\lambda \nabla c$  where  $\lambda > 0$  the constrained regions satisfying III overlap into an entire half space! (Fully intersect!)



$\lambda = 0$   
 $\Rightarrow$  No  
Constraint

Forming the Lagrangian

for  $\lambda > 0$

If  $L = f - \lambda c$

When  $\lambda > 0$

$\nabla L = \nabla f - \lambda \nabla c = 0$   
 $\Rightarrow \nabla f = + \lambda \nabla c$

The search stops since constraints are not met

With  $c(x) \geq 0$

While forming the "Lagrangian" with inequality constraints, have a "-1" sign before the constraint scaled by  $\lambda > 0$ !

If the inequality was  $c(\underline{x}) \leq 0$ ,  
We can form a  $g(\underline{x}) \geq 0$  such that  
 $g(\underline{x}) = -c(\underline{x}) \geq 0$

# Support Vector Machines

SVMs : Another class of algorithms for pattern classification and non linear regression.

It is a linear machine

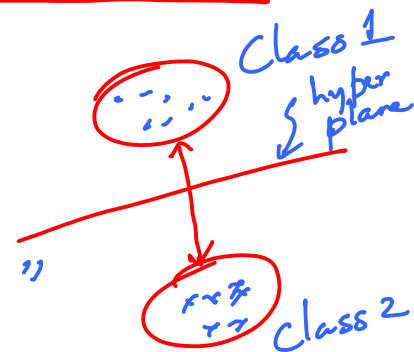
$$w^T x + b$$

Roots to SVMs : Vladimir Vapnik  
Very elegant theory with firm roots  
in Convex optimization

Idea:

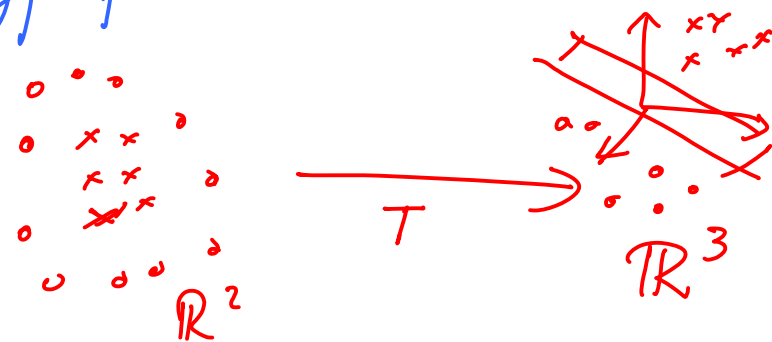
Construct a hyperplane as the decision surface  
in such a way that the margin of separation  
between the 2 classes is maximized

Idea of deriving the hyperplane  
Stems from "structural risk minimization"



In the case of linearly separable patterns, we need to derive a hyperplane that solves our objective.

In the case of non-linearly separable patterns, we need to lift the data points to a higher dimension so that we can still derive a hyperplane that solves our objective.



A notion central to the SVM is the "inner product kernel" between a support vector  $x_i^{(s)}$  and a vector  $x$  drawn from the input space.

The support vectors are a small subset of vectors extracted of the training set by the algo.



# Optimal hyperplane for linearly separable patterns

Consider the training samples  $\{ \underline{x}_i, d_i \}_{i=1}^N$   
i/p pattern for the  $i^{\text{th}}$  example target

Assume that the patterns represented by  $d_i = \{ +1, -1 \}$  is linearly separable

The eq<sup>n</sup> of the decision surface is

$$\underline{w}^T \underline{x} + b = 0$$

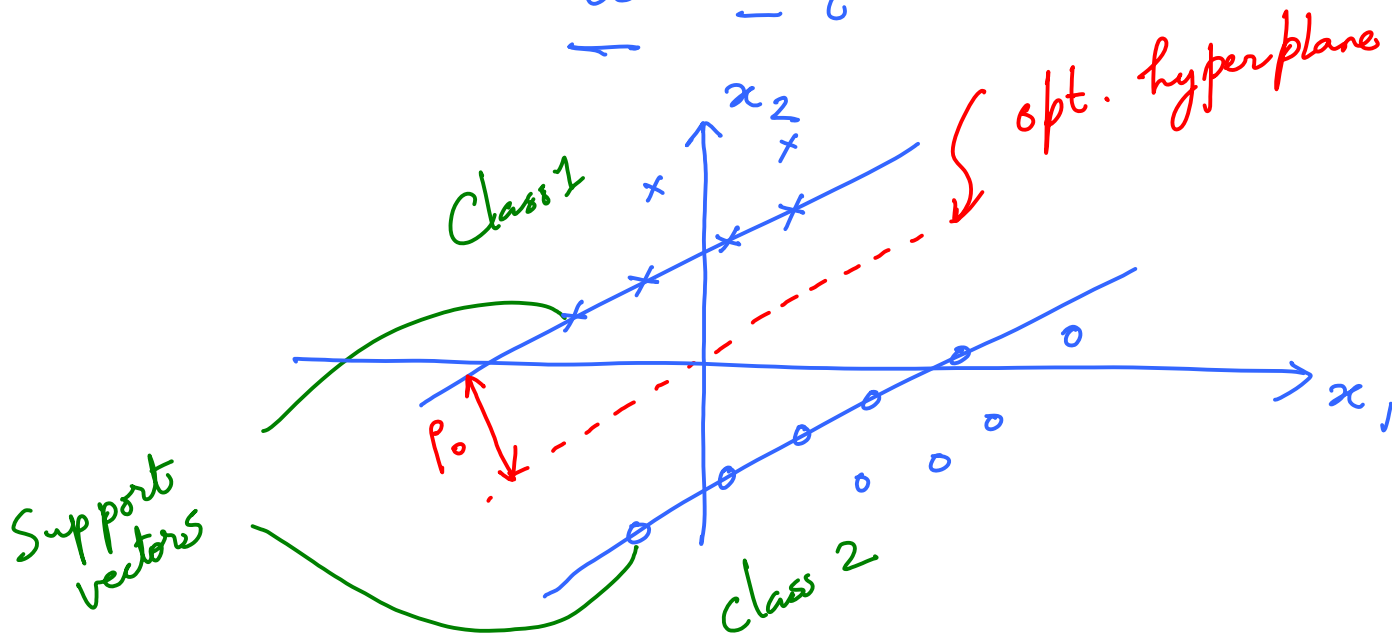
Now,

$$\underline{w}^T \underline{x}_i + b \geq 0$$

for  $d_i = +1$

$$\underline{w}^T \underline{x}_i + b < 0$$

for  $d_i = -1$



Let  $\underline{w}_0$  and  $b_0$  be the opt. values of the weight vector and the bias

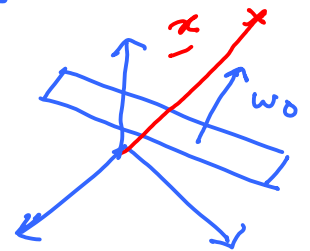
$$\underline{w}_0^T \underline{x} + b_0 = 0 \quad \leftarrow \text{Eqn of the decision boundary}$$

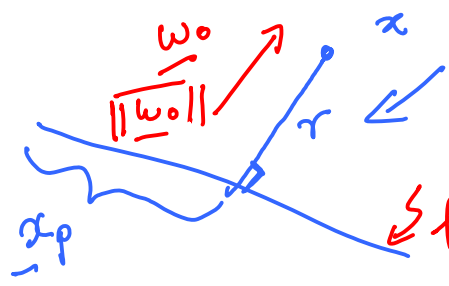
Let us write the discriminant function as

$$g(\underline{x}) = \underline{w}_0^T \underline{x} + b_0$$

From our notion of the normal to a plane

$$\underline{x} = \underline{x}_p + \frac{\underline{w}_0}{\|\underline{w}_0\|}$$





algebraic distance of the point  $\underline{x}$   
w.r.t plane

hyperplane

$$\underline{x} = \underline{x}_p + r \frac{\underline{w}_0}{\|\underline{w}_0\|}$$

normal projection of  $\underline{x}$   
on to the hyperplane

$r$  is +ve if  $\underline{x}$  is on the +ve side of the hyperplane  
 ———— || ———— -ve ———— || ———— -ve side of the hyperplane

$$g(\underline{x}_p) = 0 \quad \left( \because \underline{x}_p \text{ lies on the } \underline{\text{discriminant boundary}} \right)$$

$g(\cdot)$  is an affine map  $g(\underline{x}) = (\underline{w}_0^T \underline{x} + b_0)$   $b_0 = 0$   
(Linear map)

$$g(\underline{x}) = g\left(\underline{x}_p + r \frac{\underline{w}_0}{\|\underline{w}_0\|}\right) = \underbrace{\underline{w}_0^T \left(\underline{x}_p + r \frac{\underline{w}_0}{\|\underline{w}_0\|}\right) + b_0}_{\text{green bracket}}$$

$$g(\underline{x}) = \underbrace{\underline{w}_0^T \underline{x}_p + b_0}_{g(\underline{x}_p) = 0} + r \frac{\underline{w}_0^T \underline{w}_0}{\|\underline{w}_0\|} \quad \leftarrow \|\underline{w}_0\|^2 = r \|\underline{w}_0\|$$

$$\therefore r = \frac{g(\underline{x})}{\|\underline{w}_0\|}$$

Relationship  
between the alg.  
distance,  $g(\underline{x})$ ,  
 $\underline{w}_0$

Now, the distance from the origin to  
the hyperplane  $\frac{b_0}{\|w_0\|}$

If  $b_0 > 0$ ; the origin is on the +ve side  
of the hyperplane

||  $b_0 < 0$ ; the origin || -ve side

If  $b_0 = 0$ , the opt. hyperplane passes through the origin!

Our training set comprises of  $\mathcal{F} = \{ \underline{x}_i, d_i \}_{i=1}^N$

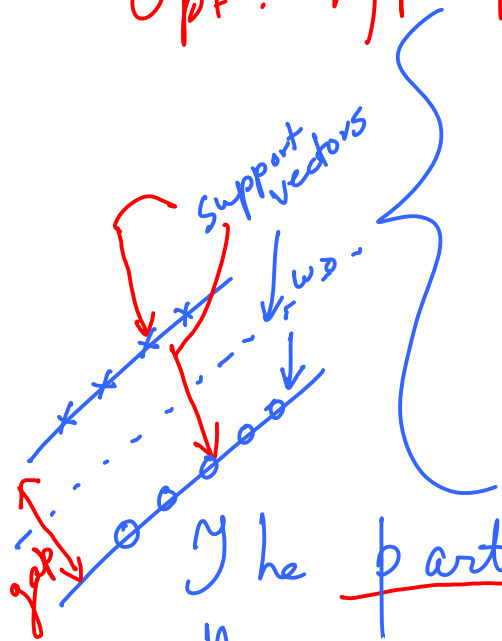
Opt. hyper plane  $\underline{w}_0^T \underline{x} + b_0 = 0$

$$\underline{w}_0^T \underline{x}_i + b_0 \geq 1 \quad d_i = +1$$

$$\underline{w}_0^T \underline{x}_i + b_0 \leq -1 \quad d_i = -1 \quad \textcircled{A}$$

$(\underline{x}_i, d_i)$  for which  
satisfied with equality are

$\underline{x}_i^{(s)}$  ← Support vect.



The particular data points  
The eqns in  $\textcircled{A}$  are  
"Support vectors"!

Consider a support vector  $\underline{x}^{(s)}$

$$g(\underline{x}^{(s)}) = \underline{w}_0^T \underline{x}^{(s)} + b_0 = \mp 1 \quad \text{for} \quad \underline{d}^{(s)} = \mp 1$$

The algebraic distance from the support vector  $\underline{x}^{(s)}$  to the opt. hyperplane is

$$r = \frac{g(\underline{x}^{(s)})}{\|\underline{w}_0\|} = \begin{cases} \frac{1}{\|\underline{w}_0\|} & \text{if } d^{(s)} = +1 \\ -\frac{1}{\|\underline{w}_0\|} & \text{if } d^{(s)} = -1 \end{cases}$$



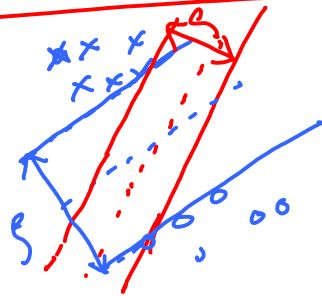
Let  $\rho$  be the opt. value of  
margin of separation

$$\rho = 2r \quad \text{where} \quad r = \frac{1}{\|w_0\|}$$

Max. margin of separation " $\rho$ "



Min the  
Euclidean norm  
of  $w_0$



## Quadratic Optimization for finding opt. hyperplane

---

Given  $\mathcal{J} = \{ \underline{x}_i, d_i \}_{i=1}^N$ , find the opt. hyperplane

subject to  $d_i (\underline{w}^T \underline{x}_i + b) \geq 1$  for  $i = 1, 2, \dots, N$

and the weight vector that minimizes the cost function

$$\phi(\underline{w}) = \frac{1}{2} \underline{w}^T \underline{w}$$

a) Cost function is convex

b) Constraints are linear in  $\underline{w}$

NOTE :

Set up the Lagrangian function

$$J(\underline{w}, b, \underline{\alpha}) = \frac{1}{2} \underline{w}^T \underline{w} - \sum_{i=1}^N \alpha_i [d_i (\underline{w}^T \underline{x}_i + b) - 1]$$

wt. vector  $\uparrow$   
 bias  $\uparrow$   
 vector of all Lag. multipliers for each constraint  $\uparrow$

Observe the sign flip required for inequality constraints  
 Lagrange multiplier for each constraint 'i'

Conditions

1)  $\frac{\partial J(\underline{w}, b, \underline{\alpha})}{\partial \underline{w}} = 0$

2)  $\frac{\partial J(\underline{w}, b, \underline{\alpha})}{\partial b} = 0$

3) Initially  $\frac{\partial J(\underline{w}, b, \underline{\alpha})}{\partial \alpha_i}$

gives us the constraints

## Evaluating the partial derivatives

Condition 1 gives  $w_s$ ,

$$\underline{w} = \sum_{i=1}^N \alpha_i d_i \underline{x}_i$$

\_\_\_\_\_ (1)

Condition 2 gives  $w_s$ ,

$$\sum_{i=1}^N \alpha_i d_i = 0$$

\_\_\_\_\_ (2)

Due to the nature of the convex opt. set up, soln is unique

NOTE :

1) It is important to note that, at the saddle point, for each Lagrange multiplier  $\alpha_i$ , the product of that multiplier with the constraint vanishes

i.e.,  $\alpha_i [d_i (\underline{\omega}^T \underline{x}_i + b) - 1] = 0 \quad \forall i = 1, \dots, N$

$$\alpha_i \neq 0$$

$\Rightarrow$

$$d_i (\underline{\omega}^T \underline{x}_i + b) - 1 = 0$$

(Home Work)

## Primal & dual problems

- 1) If the primal problem has an optimal solution, the dual too has, and the corresponding opt. values are equal. (For convex problems)
- 2) In order to find  $w_{opt}$  for the primal problem, we may need to find an alternative variable that optimizes the dual problem

$$J(\underline{w}, b, \underline{\alpha}) = \frac{1}{2} \underline{w}^T \underline{w} - \sum_{i=1}^N \alpha_i d_i \underline{w}^T \underline{x}_i$$

$\underline{w}$  (1)
 $\underline{x}_i$  (2)

$$- b \sum_{i=1}^N \alpha_i d_i + \sum_{i=1}^N \alpha_i$$

$\sum_{i=1}^N \alpha_i d_i$  (3)
 $\sum_{i=1}^N \alpha_i$  (4)

(Expanding from the primal problem)

conditions,

From the optimality

$$\sum_{i=1}^N \alpha_i d_i = 0$$

$$\left( \frac{\partial J(\cdot)}{\partial b} = 0 \right)$$

Also,

$$\underline{w}^T \underline{w} = \sum_{i=1}^N \alpha_i d_i \underline{w}^T \underline{x}_i$$

(∵ Condition 1)  
 $\frac{\partial J(\cdot)}{\partial \underline{w}} = 0$ )

$$\therefore \underline{w}^T \underline{w} = \sum_{i=1}^N \sum_{j=1}^N \alpha_i d_i \alpha_j d_j \underline{x}_i^T \underline{x}_j$$

Our dual objective function is  $Q(\alpha)$  given by  $\alpha_i$ 's are non-negative



$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \underline{x}_i^T \underline{x}_j$$

Statement of the dual problem

---

Given training samples  $\{ \underline{x}_i, d_i \}_{i=1}^N$ , find  
Lagrange multipliers  $\{ \alpha_i \}_{i=1}^N$  that maximize  
 $Q(\alpha)$

Subject to the conditions

$$1) \sum_{i=1}^N \alpha_i d_i = 0$$

$$2) \alpha_i \geq 0$$

$$\forall i = 1, \dots, N$$

Note that the dual problem is recast completely  
in terms of training data!

Having obtained the opt. Lagrange multipliers, denoted by  $\alpha_{opt, i}$ , each constraint  $i = 1, \dots, N$  we may compute the opt. weight  $\underline{w}_{opt}$  and write it as

$$\underline{w}_{opt} = \sum_{i=1}^N \alpha_{opt, i} d_i \underline{x}_i$$

Opt. bias  $b_0 = 1 - \underline{w}_0^T \underline{x}^{(s)}$  for  $d^{(s)} = 1$