Consider an example (with equality constraints)

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

s.t. $x_1^2 + x_2^2 = a^2$

(The points are on a circle)

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

$\underline{x} \in \mathbb{R}^2 \qquad \nabla f = \left( \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right)$

$$C = x_1^2 + x_2^2 - a^2 \in \mathbb{E}$$

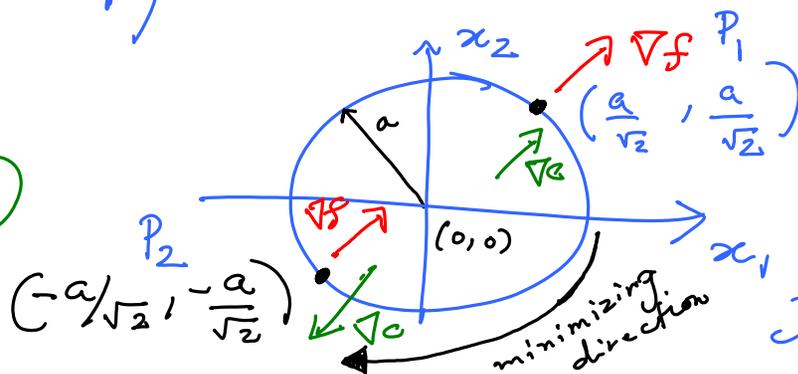$$\nabla f = \begin{pmatrix} 1 & 1 \end{pmatrix}$$

$\nabla C = \left( \frac{\partial C}{\partial x_1} \quad \frac{\partial C}{\partial x_2} \right) = 0$

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

$\nabla f_{(-\frac{a}{\sqrt{2}}, -\frac{a}{\sqrt{2}})} = (1, 1)$

$\nabla C = (-\sqrt{2}\,a, -\sqrt{2}\,a)$

$(-a/\sqrt{2}, -\frac{a}{\sqrt{2}})$



$P_2$

$(-a/\sqrt{2}, -\frac{a}{\sqrt{2}})$

$P_1$: $\left( \frac{a}{\sqrt{2}}, \frac{a}{\sqrt{2}} \right)$

$(0,0)$

minimizing direction

$\overline{III}^{rd}$ quadrant has both $x_1$ and $x_2$ $-ve$ $\Rightarrow$ Soln lies there

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

From the figure,

$$\nabla f(\underline{x}^*) = \lambda_1^* \nabla c(\underline{x}^*)$$

$$\lambda_1^* = \frac{-1}{a\sqrt{2}}$$

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

Let us __analyze__ this issue through a __Taylor__ Series expansion around the Constraint.

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

$$\left( \begin{array}{c} \because \text{ Equality Constraint} \end{array} \right)$$

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

$$\left( \begin{array}{c} \text{To maintain feasibility} \\ \text{w.r.t. } c(\underline{x}) = 0 \end{array} \right)$$

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

$$\left( \begin{array}{c} \text{With a first} \\ \text{order approx.} \end{array} \right)$$

<span style="color:red">$\underbrace{\hspace{3cm}}$ inner product</span>

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

$\nabla c$

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

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

$$\text{(A)}$$

Ill ly the <u>direction of optimization</u> <u>must produce</u>

a <u>decrease in f</u>

$$f(x+d) - f(x) < 0$$

$f(x) + \nabla^T f(x) \cdot d$

By doing a Taylor expansion around $\underline{x}$ using

a 1st order approx.

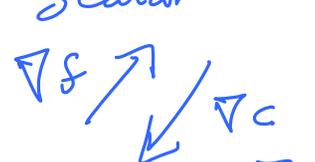$$\nabla^T f(\underline{x}) \underline{d} < 0 \qquad \text{———————} \qquad \text{Ⓑ}$$

If $\exists$ a $d$ satisfying Ⓐ and Ⓑ, an

<u>improvement</u> exists

There are 2 Cases to Consider here

(1) Such a direction does not exists

(2) Such a direction exists

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

$$\boxed{\nabla f = \lambda \nabla c}$$

Ponder why:     When $\nabla f = \lambda \nabla G$

Ⓐ and Ⓑ `do not` simultaneously hold

∴ forming a Lagrangian ⟵ Lagrange multiplier

$$L = f \pm \lambda G$$

$$\nabla L = 0 \implies \nabla f = \mp \lambda G$$

∴ 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

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

Let us verify if $\boxed{I}$ satisfies $\bigcirc\!\!\!A$ and $\bigcirc\!\!\!B$

$$d = -\nabla f + \frac{\overbrace{\nabla c \; \nabla c^T}^{\text{outer product}} \nabla f}{\underbrace{\nabla c^T \; \nabla c}_{\text{Inner product}}} \quad \text{---} \quad \bigcirc\!\!\!1$$

Let us consider $\bigcirc\!\!\!A$
Pre-multiply $\bigcirc\!\!\!1$ by $\nabla c^T$ ;

$$\nabla c^T \underline{d} = -\nabla c^T \nabla f + \frac{\overset{(Scalar)}{\nabla c^T \nabla c} \ \nabla c^T \nabla f}{\nabla c^T - \nabla c} = 0$$
$$\underset{(Scalar)}{}$$

Let us $\underline{\text{Consider}}$ $\boxed{B}$

$$\nabla f^T \underline{d}$$

Plug in $\underline{d}$

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

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

$$= - \|\nabla f\|^2 + \frac{\|\nabla f^T \nabla c\|^2}{\|\nabla c\|^2} \leq 0 \quad (\because \text{Cauchy Schwartz})$$
$$\text{inequality}$$

The equality is ruled out due to Case (A)

$(\because \nabla f \neq \lambda \nabla c)$

$\implies$ $\underline{d}$ is the direction satisfying the

constraints.

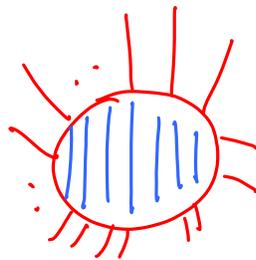# Single inequality constraint

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

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

Observe the '$\leq$' as against '$=$' in equality constraints

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

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

Considering the example from the circular constraints with inequality conditions, we are optimizing over all points lying on & inside the circle.

$$x_1^2 + x_2^2 \leq a^2$$

$$\underbrace{\phantom{xxxxxxx}}$$

$$-(x_1^2 + x_2^2) \geq -a^2$$

$$c(x) \geq 0 \quad \longrightarrow \quad a^2 - x_1^2 - x_2^2 \geq 0$$

We have 2 Cases

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

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

$$\boxed{II} \begin{cases} \nabla f^{T}(\underline{x}) \underline{d} < 0 & \underline{\hspace{2cm}} \quad (\because \textcircled{B}) \\ c(\underline{x}) + \nabla c^{T}(\underline{x}) \underline{d} \gtreqqless 0 & \underline{\hspace{2cm}} \quad (\because \textcircled{C}) \end{cases}$$

$A \quad \underline{d}$ that satisfies the constraints is

$$\underline{d} = - c(\underline{x}) \frac{\nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(\underline{x})\|} \underline{\hspace{2cm}} \textcircled{D}$$

We can verify that ⓓ Satisfies both the
Constraints in Ⓘ

(i) $\nabla^T f(x)\, \underline{d} = -\nabla^T f(\underline{x}) \cdot c(\underline{x}) \underbrace{\dfrac{\nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(x)\|}}_{\text{Scalar}}$

$= -c(\underline{x}) \dfrac{\nabla^T f(\underline{x})\, \nabla f(\underline{x})}{\|\nabla f(\underline{x})\| \|\nabla c(x)\|}$

Evaluates to $\|\nabla f\|$

$= -c(x) \dfrac{\|\nabla f\|}{\|\nabla c\|}$

$> 0$

$\Longrightarrow$ First Constraint in Ⓘ is
Satisfied          i.e., $< 0$

(ii)     Consider

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

$$= \quad c(\underline{x}) + \nabla^T c(\underline{x}) \left[ - \frac{c(\underline{x}) \, \nabla f(\underline{x})}{\| \nabla f(\underline{x}) \| \, \| \nabla c(\underline{x}) \|} \right]$$

$$= \quad c(\underline{x}) - c(\underline{x}) \underbrace{\frac{\nabla^T c(\underline{x}) \, \nabla f(\underline{x})}{\| \nabla f(\underline{x}) \| \, \| \nabla c(\underline{x}) \|}}_{\ge 0} < 1$$

$$| \cdot | < 1$$

Unless     $\nabla f(\underline{x}) \not\equiv \lambda \nabla c(\underline{x}),$

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

We have $\quad c(\underline{x}) + \nabla c(\underline{x})^T \underline{d}$

$$\implies \qquad c(\underline{x}) - c(\underline{x}) \, \alpha$$

$\alpha$ can be +ve or −ve

$$c(\underline{x})(1 - \alpha) \geq 0$$

$\left( \begin{array}{c} \text{The equality is} \\ \text{only over the case} \\ \text{when } \alpha = 1 \end{array} \right)$
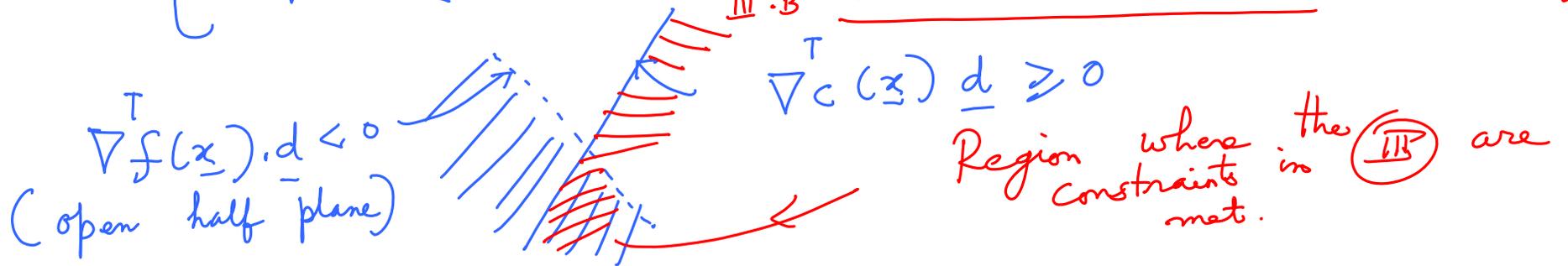
Case B : When $\underline{x}$ is on the boundary of the constraint eqn i.e., $c(\underline{x}) = 0$

We have

$\boxed{III}$
$\begin{cases} \nabla^T f(\underline{x})\, \underline{d} < 0 \quad \underline{\hspace{2cm}} \quad \underline{III}.A \\[2em] \nabla^T c(\underline{x})\, \underline{d} \geq 0 \quad \underline{\hspace{2cm}} \\ \overline{III \cdot B} \end{cases}$

$\textcircled{B}$ boundary case

$\left( \because c(\underline{x}) = 0 \atop \text{Plug into } \textcircled{G} \right)$

FEASIBLE SOLN REGION (GEOMETRY)

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



$\nabla^T f(\underline{x}) \cdot \underline{d} < 0$
(open half plane)

Region where the $\textcircled{III}$ are constraints in met.

When $\nabla f = \lambda \nabla c$ and

$\nabla f$ and $\nabla c$ point in the same direction

the regions from (III) do not intersect. !

$\nabla f \nearrow$ $\nearrow \nabla c$

When $\nabla f = -\lambda \nabla c$ where $\lambda > 0$

the constrained regions satisfying (IV) overlap into

an entire half space ! ( Fully intersect ! )

$\nabla f \nearrow$ $\searrow \nabla c$

Forming the _Lagrangian_ for $\lambda > 0$ $\left( \begin{array}{l} x = 0 \\ \Rightarrow \text{No} \\ \text{Constraint} \end{array} \right)$

If $L = f - \lambda c$

When $\lambda > 0$

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

$\nabla f = + \lambda \nabla G \Rightarrow$ The search stops since constraints are not met

$\Rightarrow$

With $c(x) \geq 0$ While forming the "Lagrangian" with inequality constraint, have a $= \prime$ sign before the constraint scaled by $\prime \lambda \gtrless 0 \prime$ !

If the inequality was $c(x) \leq 0$,

We can form a $g(x) \geq 0$ such that

$$g(x) = -c(x) \geq 0$$

# Support Vector Machines

**SVMs :** Another class of algorithms for pattern classification and non linear regression.
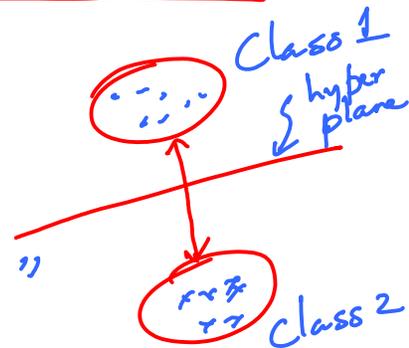
It is a <u>linear machine</u>

$$\omega^T x + b$$

**Roots to SVMs :** Vladimir Vapnik

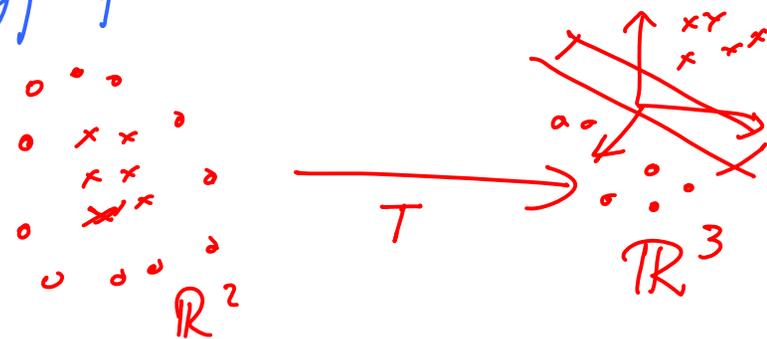Very elegant theory with firm roots in Convex optimization

**Idea:** Construct a hyperplane as the decision surface in such a way that the margin of separation between the 2 classes is <u>maximized</u>

Idea of deriving the hyperplane stems from 'structural risk minimization"

Class 1

hyperplane

Class 2

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

In the case of <u>non-linearly separable patterns,</u> 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 $\underline{x}$ 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 $\{\underline{x}_i, d_i\}_{i=1}^{N}$

← target

$i/p$ pattern
for the $i^{th}$ example

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

The eqn of the decision surface is

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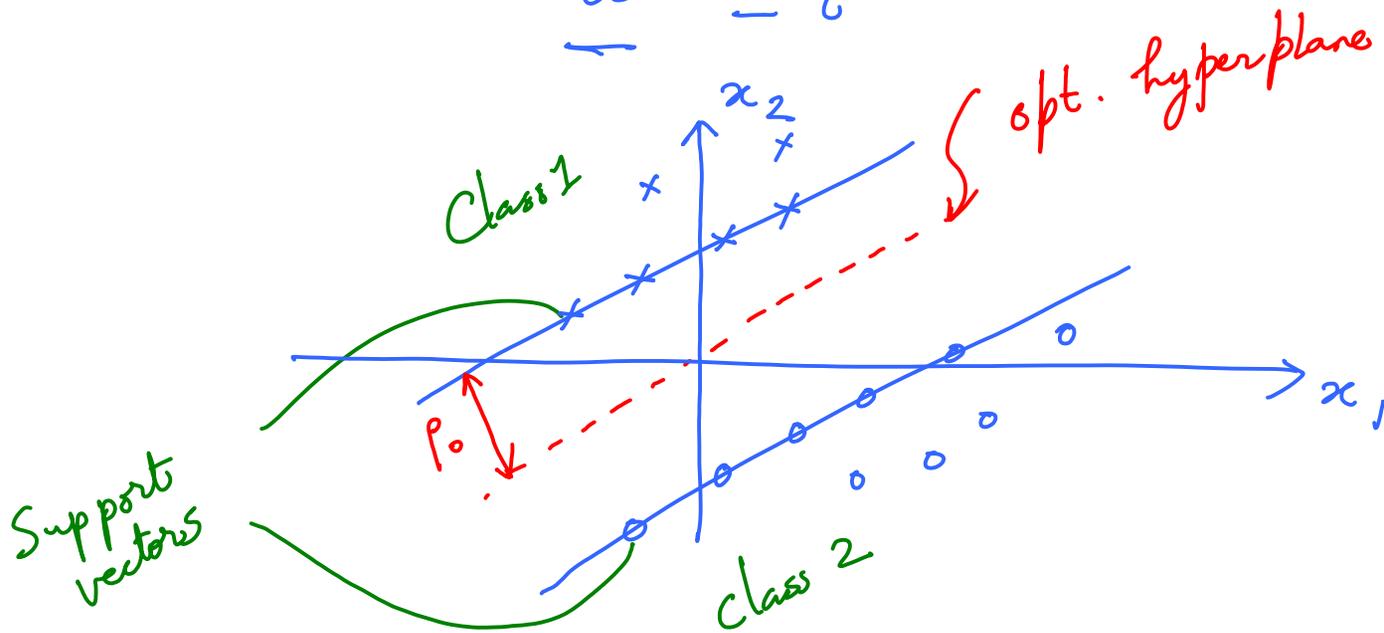
Now,

$$\omega^T x_i + b \geq 0 \qquad \text{for } d_i = +1$$

$$\omega^T x_i + b < 0 \qquad \text{for } d_i = -1$$



Class 1

Class 2

Support vectors

opt. hyperplane

$\rho_0$

$x_2$

$x_1$

Let $\underline{w}_0$ and $b_0$ be the opt. values of the weight vector and the bias
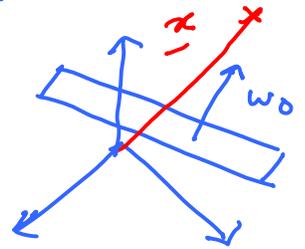
$$\underline{w}_0^T \underline{x} + b_0 = 0 \quad \longleftarrow \quad \text{Eqn of the decision boundary}$$
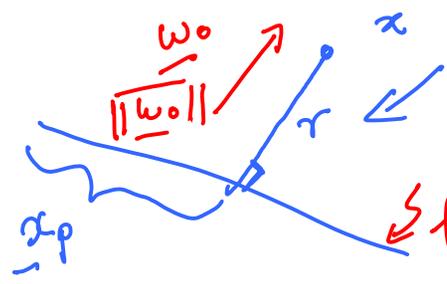
Let us write the discriminant function as

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

From our notion of the normal to a plane

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

$w_0$

$\dfrac{w_0}{\|w_0\|}$

$x$

$r$

$x_p$

hyperplane

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

$$\underline{x} = \underline{x}_p + r \dfrac{\underline{w_0}}{\|\underline{w_0}\|}$$

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

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

$g(\cdot)$ is an $\underline{\text{affine map}}$ $g(x) = \left( \underline{w}_o^T \underline{x} + \underline{b_o} \right)$ $\underset{\text{(Linear map)}}{b_o = 0}$

$$g(\underline{x}) = g\left( \underline{x}_p + r \frac{\underline{w}_o}{\|\underline{w}_o\|} \right) = \underline{w}_o^T \left( \underline{x}_p + r \frac{\underline{w}_o}{\|\underline{w}_o\|} \right) + b_o$$

$$g(\underline{x}) = \underbrace{\underline{w}_o^T \underline{x}_p + b_o}_{g(\underline{x}_p) = 0} + r \cdot \frac{\overbrace{\underline{w}_o^T \underline{w}_o}}{\|\underline{w}_o\|} \quad \Leftarrow \quad \|\underline{w}_o\|^2 = r \|\underline{w}_o\|$$

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

Relationship between the distance, $g(\underline{x})$, $\underline{w}_o$

Now, the distance from the origin to
the hyperplane

$$\frac{b_0}{\|\underline{w_0}\|}$$

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

$b_0 < 0$; the origin —"— —ve side
—"—

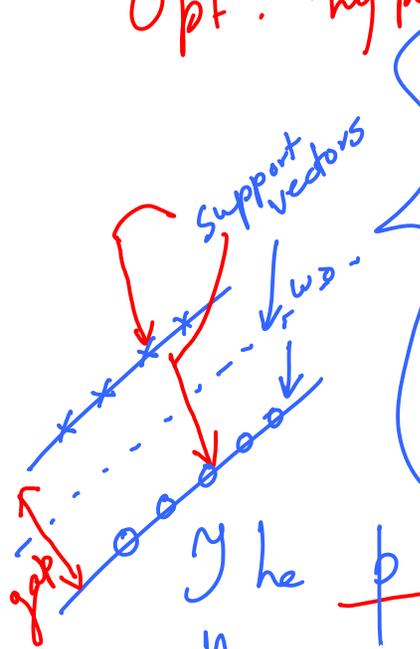If $b_0 = 0$, the opt. hyperplane passes through the
origin !

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

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

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

$$\underline{w}_0^T \underline{x}_i + b_0 \leqslant 1 \qquad d_i = -1 \qquad \text{(A)}$$

support vectors

The particular data points $(\underline{x}_i, d_i)$ for which the eqns in (A) are satisfied with equality are the "Support vectors"!

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

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

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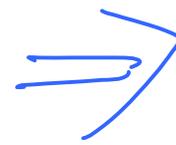
$$\underline{d}^{(s)} = \mp 1$$

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

$$r = \frac{g\left(\underline{x}^{(s)}\right)}{\|\underline{\omega}_0\|} = \begin{cases} \dfrac{1}{\|\underline{\omega}_0\|} & \text{if } d^{(s)} = +1 \\[3mm] -\dfrac{1}{\|\underline{\omega}_0\|} & \text{if } d^{(s)} = -1 \end{cases}$$

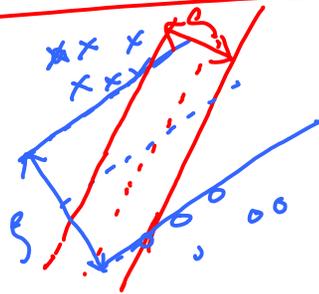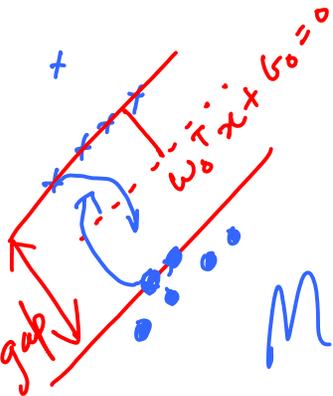Let $\rho$ be the opt. value of

margin of Separation

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

Max. margin of Separation "$\rho$" $\Longrightarrow$ Min the Euclidean norm of $w_o$

# Quadratic Optimization for finding opt. hyperplane

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

subject to

$$d_i(\underline{w}^T \underline{x}_i + b) \geq 1 \qquad \text{for } i = 1, 2, \ldots, N$$

and the $\underline{\text{weight vector}}$ that minimizes the cost function

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

NOTE:

a) Cost function is $\underline{\text{convex}}$

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

Set up the Lagrangian function

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

wt. vector · bias vector

of all
Lag. multipliers
for each constraint

Observe the
sign flip required
for inequality constraints

Lagrange
multiplier
for each
constraint
'i'

## Conditions

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

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

gives us the constraints

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

Evaluating the partial derivatives

Condition 1 gives us,

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

Condition 2 gives us,

$$\sum_{i=1}^{N} \alpha_i d_i = 0 \qquad \text{————————②}$$

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 <u>Lagrange multiplier</u> $\alpha_i$, the product of that multiplier with the constraint vanishes

i.e.,
$$\alpha_i \left[ d_i \left( \underline{w}^T \underline{x}_i + b \right) - 1 \right] = 0 \qquad \forall \ i = 1, \ldots, N$$

(<u>Home Work</u>)

$$\alpha_i \neq 0$$

$$\Rightarrow \boxed{d_i \left( \underline{w}^T \underline{x}_i + b \right) - 1 = 0}$$

# 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 $\underline{w}_{opt}$ for the primal problem, we may need to find an alternative variable that optimizes the dual problem

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

$$- b \underbrace{\sum_{i=1}^{N} \alpha_i d_i}_{③} + \underbrace{\sum_{i=1}^{N} \alpha_i}_{④}$$

$\left(\underline{\text{Expanding}} \text{ from the} \atop \underline{\text{primal} \quad \text{problem}}\right)$

From the optimality conditions,

$$\sum_{i=1}^{N} \alpha_i d_i = 0 \qquad \left(\frac{\partial J(\cdot)}{\partial b} = 0\right)$$

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

$$\left( \because \text{ Condition } 1 \quad \frac{\partial J(\cdot)}{\partial \underline{\omega}} = 0 \right)$$

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

$\alpha_i$'s are non-negative

Our dual objective function is $Q(\alpha)$ given by

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

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) $\alpha_i \geqslant 0$ $\qquad \forall\ i = 1, \ldots, N$

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 $i = 1, ..., N$ 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 $b_o = 1 - \underline{w}_o^T \, \underline{x}^{(s)}$ for $d^{(s)} = 1$