**PERCEPTRON**

$$\text{\# mistakes} \leq \frac{R^2}{\gamma^2} \leftarrow$$
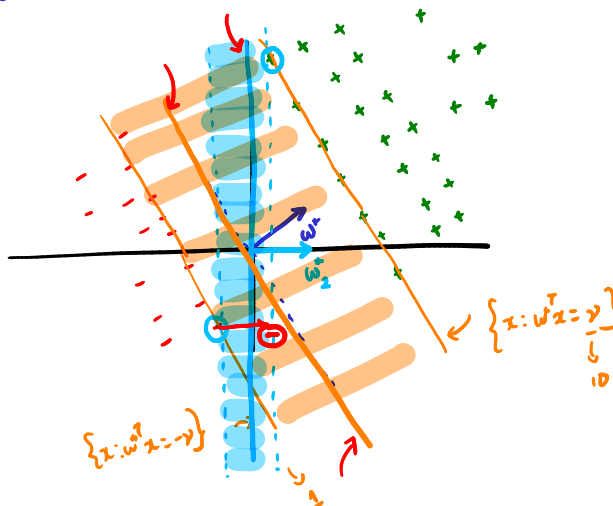
$$\|x_i\|^2 \leq R^2$$

Dataset - L·S with margin $\gamma$

$$(w^{*T} x_i) y_i \geq \gamma \quad \forall i$$
$$\gamma > 0$$

---

"QUALITY" OF FINAL SOLUTION

**Question**

Given that we prefer classifies with large margin, can we directly find them?



$\{x : w^T x = \gamma\}$
$\downarrow$
10

$\{x : w^2 x = -\gamma\}$

**Observation**
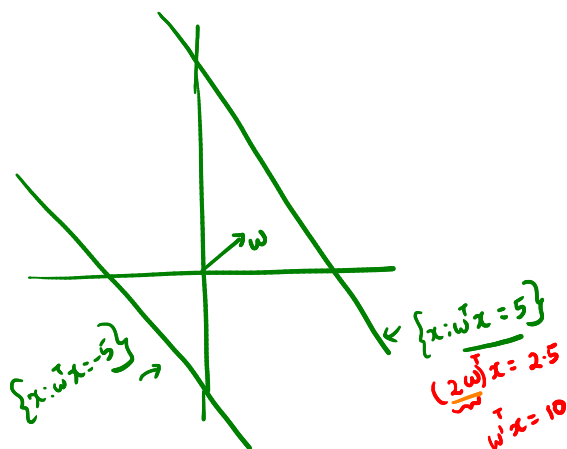
① # mistakes depends on the best possible $w^*$'s margin.

② $w_{perc}$ need not necessarily be $w^*$. It could be $w_2^*$ also (blue line)

---

**Goal**: To come up with a formulation that maximizes "margin"



$\{x : w^T x = 5\}$
$(2w)^T x = 2.5$
$w^T x = 10$

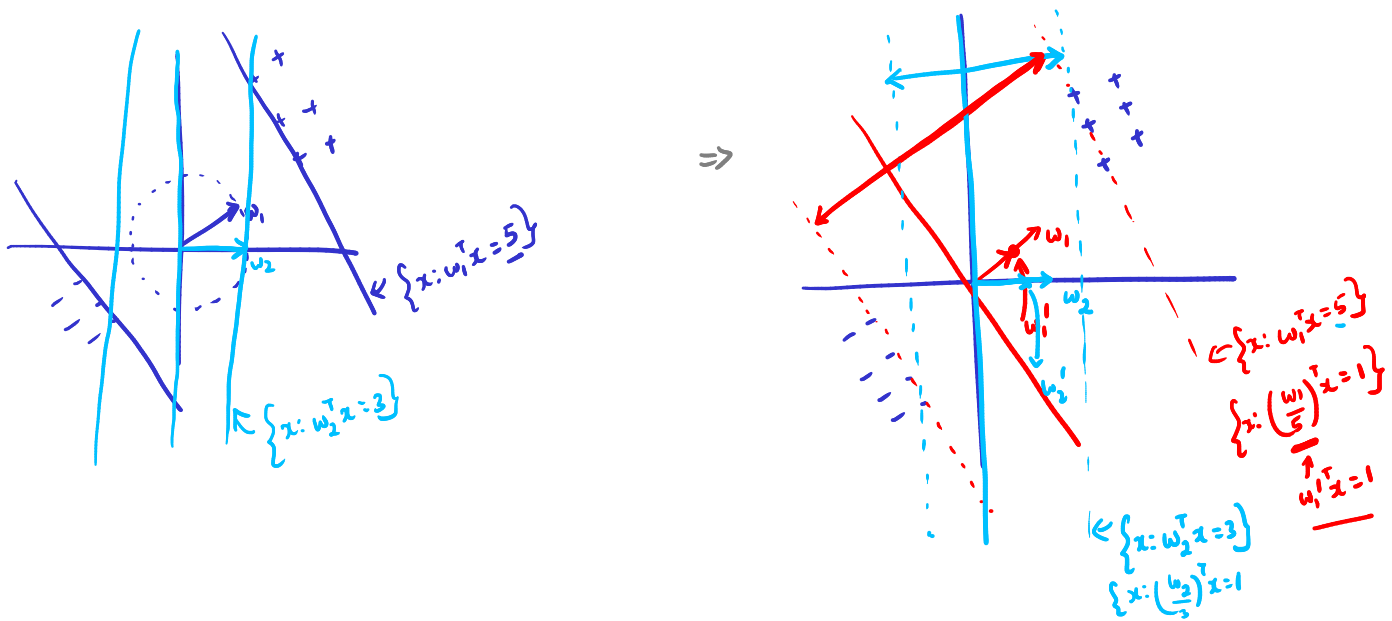$$\max_{w, \gamma} \quad \gamma$$

Such that

$$(w^T x_i) y_i \geq \gamma \quad \forall i$$

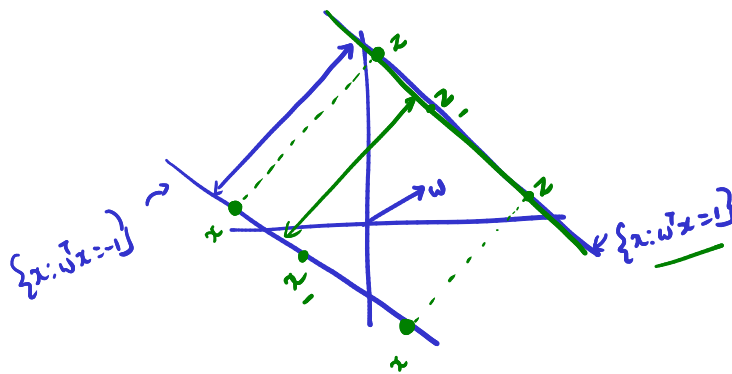**ISSUE**: Can scale $w$ arbitrarily.

$$\max_{w, \gamma} \quad \gamma$$
$$(w^T x_i) y_i \geq \gamma \quad \forall i$$
$$\|w\|^2 = 1$$

Top-left diagram labels: $w_1$, $w_2$, $\{x: w_1^T x = 5\}$, $\{x: w_2^T x = 3\}$

$\Rightarrow$

Top-right diagram labels: $w_1$, $w_2$, $\leftarrow \{x: w_1^T x = 5\}$, $\{x: \left(\frac{w_1}{5}\right)^T x = 1\}$, $w_1^T x = 1$, $\{x: w_2^T x = 3\}$, $\{x: \left(\frac{w_2}{3}\right)^T x = 1\}$

---

$$\max_{w} \quad \text{width}(\omega)$$
$$\text{s.t.} \quad (w^T x_i)\, y_i \geq 1 \quad \forall i$$

---

**What is width($\omega$)?**



Diagram labels: $z$, $z_1$, $\{x: w^T x = -1\}$, $\{x: w^T x = 1\}$, $w$

$$\boxed{w^T x = -1}$$

$$\min_{z} \quad \frac{1}{2}\|x - z\|^2$$
$$\text{s.t.} \quad w^T z = +1$$

Solution: $\quad \text{width}(w) = \boxed{\dfrac{2}{\|w\|^2}}$

---

$$\max_{w} \quad \frac{2}{\|w\|^2}$$
$$\text{s.t.} \quad (w^T x_i)\, y_i \geq 1$$

Bottom-right diagram labels: $w$, $w_2$, $\{x: w_2^T x = 1\}$, $\{x: w^T x = 1\}$

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

$$s.t \quad (w^T x_i) y_i \geq 1 \quad \forall i$$

$$\min_{w \in \mathbb{R}^d} \underbrace{\frac{1}{2} \|w\|^2}$$

$$\text{s.t } \underline{\pm i} \; (w^T x_i)\, y_i \geq 1$$

$- \text{(A)}$

DETOUR

$$\min_w \quad f(w)$$
$$\text{s.t} \quad g(w) \leq 0$$

$\leftarrow$

$$\mathcal{L}(w, \alpha) \;=\; f(w) + \alpha\, g(w)$$

Fix any $w$.

Consider $\quad \max_{\alpha \geq 0} \mathcal{L}(w, \alpha) \;=\; \max_{\alpha \geq 0} \underbrace{f(w)} + \underbrace{\alpha\, g(w)} \leftarrow$

| | $f(w)$ | $g(w)$ |
|---|---|---|
| $w$ $[1\,2\,3\,4]$ | $-100$ | $5$ |

$$\max_{\alpha \geq 0} \; -100 + \alpha 5 \uparrow$$

$\alpha = 1 \quad -95$
$\alpha = 10 \quad -50$
$\alpha = 100 \quad 400$

| | $\boxed{f(w)}$ | $g(w)$ |
|---|---|---|
| $w$ $[3\,4\,5]$ | $\boxed{100}$ | $-5$ |

$100 - 5\boxed{\alpha}$

$\alpha = 1 \Rightarrow 95$
$\alpha = 10 \Rightarrow 50 \downarrow$
$\alpha = 0 \Rightarrow \boxed{100}$

$= \begin{cases} \underline{\infty} & g(w) > 0 \\[1em] \underline{f(w)} & g(w) \leq 0 \end{cases}$



$\max_{\alpha \geq 0} \; f(w) + \alpha\, g(w)$

$f(w)$

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

$\{w : g(w) \leq 0\}$

$g(w_1) \leq 0.$

$w_1 \to \boxed{\phantom{x}} \to f(w)$

$w_2 \to \boxed{\phantom{x}} \to \infty$

$g(w_2) > 0$

$\min_w \; f(w)$

$\equiv$

$\min_w \left[ \max_{\alpha \geq 0} \; f(w) + \alpha\, g(w) \right]$

- Can we swap min and max in (B) ?

- In general, No! But if f and g are "nice" functions [convex functions], then yes!
  $\hookrightarrow$ [Quadratic / Linear]

$$\min_{w} \left[ \max_{\alpha \geq 0} f(w) + \alpha g(w) \right] \equiv \max_{\alpha \geq 0} \left[ \min_{w} f(w) + \alpha g(w) \right]$$

For convex f and g.

---

For multiple constraints

$$\begin{array}{|c|}\hline \min_{w} f(w) \\[2mm] s.t \quad g_i(w) \leq 0 \quad \forall i \\ {\scriptstyle = 1 \dots k} \\ \hline \end{array}$$

$\equiv$

$$\min_{w} \left[ \max_{\substack{\alpha_1, \dots, \alpha_k \\ \geq 0 \; \geq 0 \; \geq 0}} f(w) + \alpha_1 g_1(w) + \alpha_2 g_2(w) + \dots + \alpha_k g_k(w) \right]$$

$\equiv$

$$\max_{\alpha_1 \geq 0, \dots \alpha_k \geq 0} \left[ \min_{w} \left[ f(w) + \alpha_1 g_1(w) + \dots + \alpha_k g_k(w) \right] \right]$$

---

$$\min_{w} \frac{1}{2} \|w\|^2 \quad \leftarrow \text{Quadratic in } w$$

$$s.t \quad (w^T x_i) y_i \geq 1 \quad \forall i = 1, \dots, n \quad \leftarrow \text{Linear in } w$$

$$\equiv \quad \underbrace{1 - (w^T x_i) y_i \leq 0}_{g_i(w)} \quad \forall i = 1, \dots \underline{n}$$

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}$$

$$L(w, \overset{\downarrow}{\alpha}) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^{n} \alpha_i \left( 1 - (w^T x_i) y_i \right)$$

$$\min_{w} \max_{\alpha \geq 0} \left[ \frac{1}{2} \|w\|^2 + \sum_{i=1}^{n} \alpha_i \left( 1 - (w^T x_i) y_i \right) \right] \equiv \max_{\alpha \geq 0} \left[ \min_{w} \left[ \frac{1}{2} \|w\|^2 + \sum_{i=1}^{n} \alpha_i \left( 1 - (w^T x_i) y_i \right) \right] \right]$$

Fix Some $\alpha \geq 0$. $\qquad \alpha = \begin{bmatrix} 1 \\ 5 \\ 1 \end{bmatrix}$

$$\min_{w} \quad \frac{1}{2}\|w\|^2 + \underbrace{\sum_{i=1}^{n} \alpha_i \left(1 - (w^T x_i) y_i\right)}$$

$$w_{\alpha}^* + \sum_{i=1}^{n} \alpha_i (-x_i y_i) = 0$$

$$\boxed{w_{\alpha}^* = \sum_{i=1}^{n} \alpha_i x_i y_i} \quad \leftarrow \qquad \in \mathbb{R}^d$$

$\geq 0 \qquad \{\pm 1\}$

→ Substitute back value of $w_{\alpha}^*$ in the objective

$$w_{\alpha}^* = X y \alpha$$

$$X = \begin{bmatrix} 1 & 1 & & 1 \\ x_1 & x_2 & \cdots & x_n \\ 1 & 1 & & 1 \end{bmatrix}_{d \times n} \begin{bmatrix} y_1 & & 0 \\ & \ddots & \\ 0 & & y_n \end{bmatrix}_{n \times n} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

$$\frac{1}{2}\|w\|^2 + \sum_{i=1}^{n} \alpha_i \left(1 - (w^T x_i) y_i\right)$$

Substitute $w_{\alpha}^* = X y \alpha$ into

on Simplification

$$= \quad \alpha^T 1 - \frac{1}{2}(Xy\alpha)^T (Xy\alpha) \qquad \begin{bmatrix} \vdots \end{bmatrix}_{n \times 1}$$

DUAL PROBLEM

PRIMAL

$$\min_{w \in \mathbb{R}^d} \frac{1}{2}\|w\|^2$$

$$\max_{\alpha \geq 0} \quad \alpha^T 1 - \frac{1}{2} \alpha^T y^T \boxed{X^T X} y \alpha \qquad \mathbb{R}^{n \times n}$$

kernel k

## What have we gained?

- Dual variable dimension in $\mathbb{R}_+^{(n)}$ while primal problem dimension is $\mathbb{R}^{(d)}$

- Dual constraints are "easier"

- More importantly dual depends on $\bar{x}^T x$ and so can be "KERNELIZED"!

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

→ This says optimal $w^*$ is a linear combination of the data points where importance of a datapoint is given by $\alpha_i^*$ (for $i^{th}$ data point)

→ Question: Where are the "IMPORTANT" points? (i.e., points for which $\alpha_i^* > 0$)

---

### REVISITING THE LAGRANGIAN

Primal
$$\min_{w} \left[ \max_{\alpha \geq 0} \; f(w) + \alpha \, g(w) \right]$$

$\equiv$

Dual
$$\max_{\alpha \geq 0} \left[ \min_{w} \; f(w) + \alpha \, g(w) \right]$$

$w^*$ is the primal solution

$\alpha^*$ is the dual solution

$$\max f(w^*) + \alpha \, g(w^*) \quad = \quad \min_{w} f(w) + \alpha^* \, g(w)$$

$$f(w^*) = \min_w f(w) + \alpha^* g(w)$$
$$\leq f(\hat{w}^*) + \alpha^* g(\hat{w}^*)$$

$$\Rightarrow \quad f(\cancel{w}^*) \leq f(\cancel{w}^*) + \alpha^* g(\hat{w}^*)$$

$$\Rightarrow \quad \boxed{\alpha^* g(\hat{w}^*) \geq 0} \qquad - \text{①}$$

But we already know $\underline{\alpha^* \geq 0}$ & $\underline{g(w^*) \leq 0}$

$$\Rightarrow \quad \boxed{\alpha^* g(\hat{w}^*) \leq 0} \quad \text{②}$$

$$\text{①} \quad \& \quad \text{②} \quad \Rightarrow \quad \boxed{\alpha^* g(w^*) = 0} \rightarrow \text{COMPLEMENTARY SLACKNESS}$$

For multiple constraints,
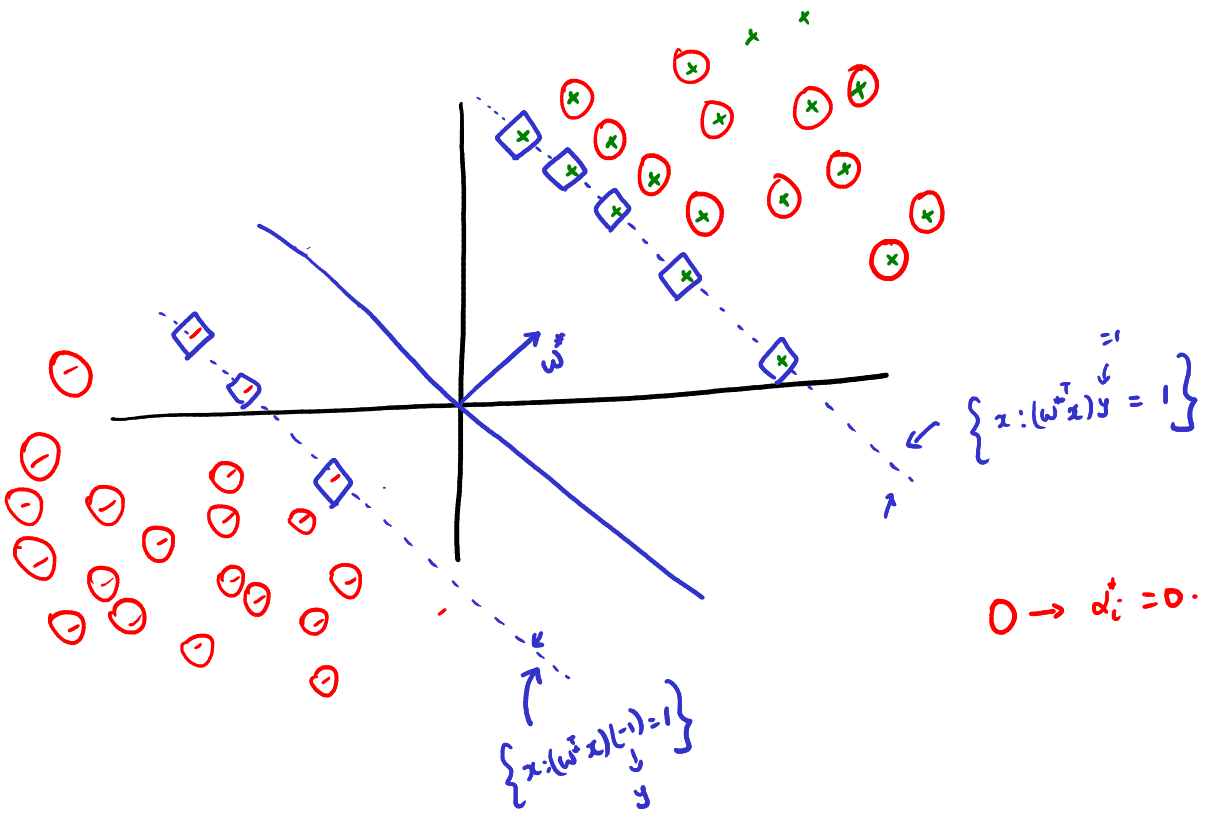
$$\boxed{\alpha_i^* g_i(w^*) = 0 \qquad \forall i}$$

In our problem

$$\alpha_i^* \underbrace{\left( 1 - (w^T x_i) y_i \right)}_{g_i(w^*)} = 0 \qquad \forall i \qquad \left[ \begin{array}{c} \text{by complementary} \\ \text{slackness} \end{array} \right]$$

$$\Rightarrow \quad \text{If} \quad \underline{\alpha_i^* > 0} \quad \overset{C.S}{\Rightarrow} \quad 1 - (w^T x_i) y_i = 0$$

$$\Downarrow$$

$$\boxed{(w^T x_i) y_i = 1}$$

$\{ z : (w^T z) \overset{=1}{y} = 1 \}$

$0 \rightarrow \alpha_i^* = 0.$

$\left\{ z : (w^T z)(\underset{y}{-1}) = 1 \right\}$

▶ Only the points that are on the "SUPPORTING" hyperplane can contribute to $w^*$

▶ These special points are called "SUPPORT VECTORS"

▶ ALGORITHM → SUPPORT VECTOR MACHINE [Vapnik et. al]
(SVM)

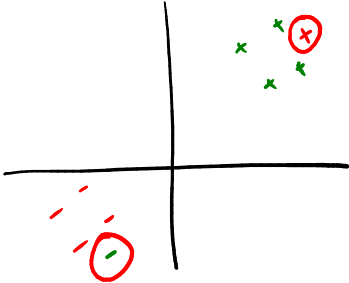▶ $w^*$ is a sparse linear combination of the data points.

Given

$x_{test}$ =

$$w^{*T} x_{test} = \left( \sum_{i=1}^{n} \alpha_i^* z_i y_i \right)^T x_{test}$$

$$= \sum_{i=1}^{n} \alpha_i^* y_i (z_i^T x_{test})$$

$x_{test}$    $w^{*T} \phi(x_{test}) = \sum_{i=1}^{n} \alpha_i^* y_i \, k(z_i, x_{test})$

# QUESTIONS

- How to adapt the SVM algorithm when data has Outliers.



- KERNELS Can help but is not the right way to solve this !

---

$$\min_{\omega} \quad \frac{1}{2}\|\omega\|^2$$

$$\text{s.t} \quad (\vec{\omega}x_i)y_i \geq 1 \quad \forall i$$

**Insight** : Make every $\omega$ feasible.

- Fix any $\omega$. $\omega$ classifies some points correctly and mis classifies some points

- The incorrectly classified points "pay bribe" to go to the "correct" side!

---

**MODIFIED FORMULATION**

$\geq 0 \rightarrow$ HYPER PARAMETER.

**SOFT MARGIN PRIMAL FORMULATION** $\leftarrow$

$$\min_{\omega,\xi} \quad \frac{1}{2}\|\omega\|^2 + C\sum_{i=1}^{n}\xi_i$$

$$\text{s.t} \quad (\omega^T x_i)y_i + \xi_i \geq 1 \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

$C = 0$
$\Rightarrow$ Bribes don't cost
$\Rightarrow \omega = 0 \in \mathbb{R}^d$ is the solution

$C = \infty$
$\Rightarrow$ Linear separable case.