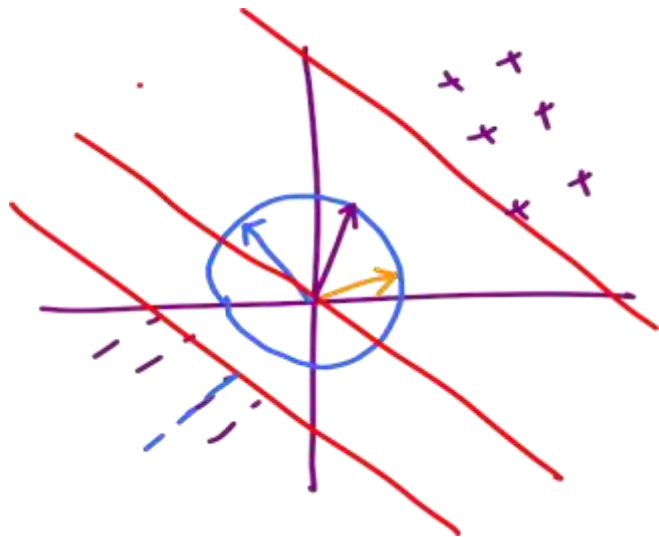**Goal:** Given a dataset $D = \{ (x_1, y_1) \cdots \cdots , (x_n, y_n) \}$

find a $w$ with maximum width s.t all points in dataset $D$ are classified correctly

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

s.t

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

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

Possible fix

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

$$(w^T x_i) y_i \geq \gamma$$

$$\|w\|^2 = 1$$

**Goal:** Given a dataset $D = \{ (x_1, y_1) \cdots , (x_n, y_n) \}$

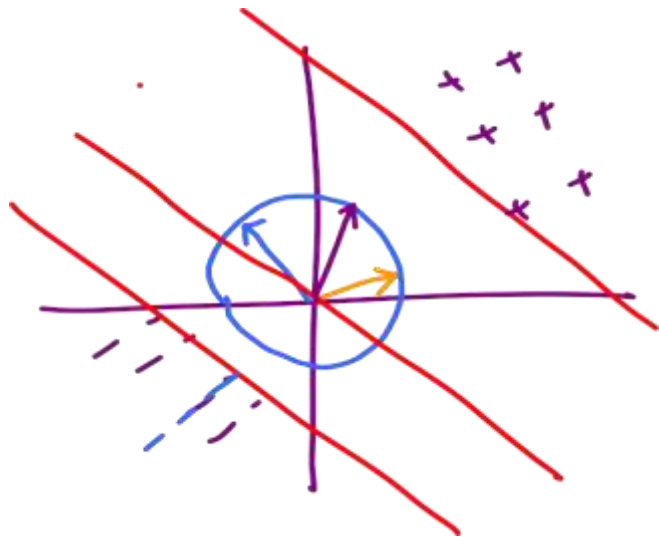find a $w$ with maximum width s.t all points in dataset $D$ are classified correctly

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

s.t

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

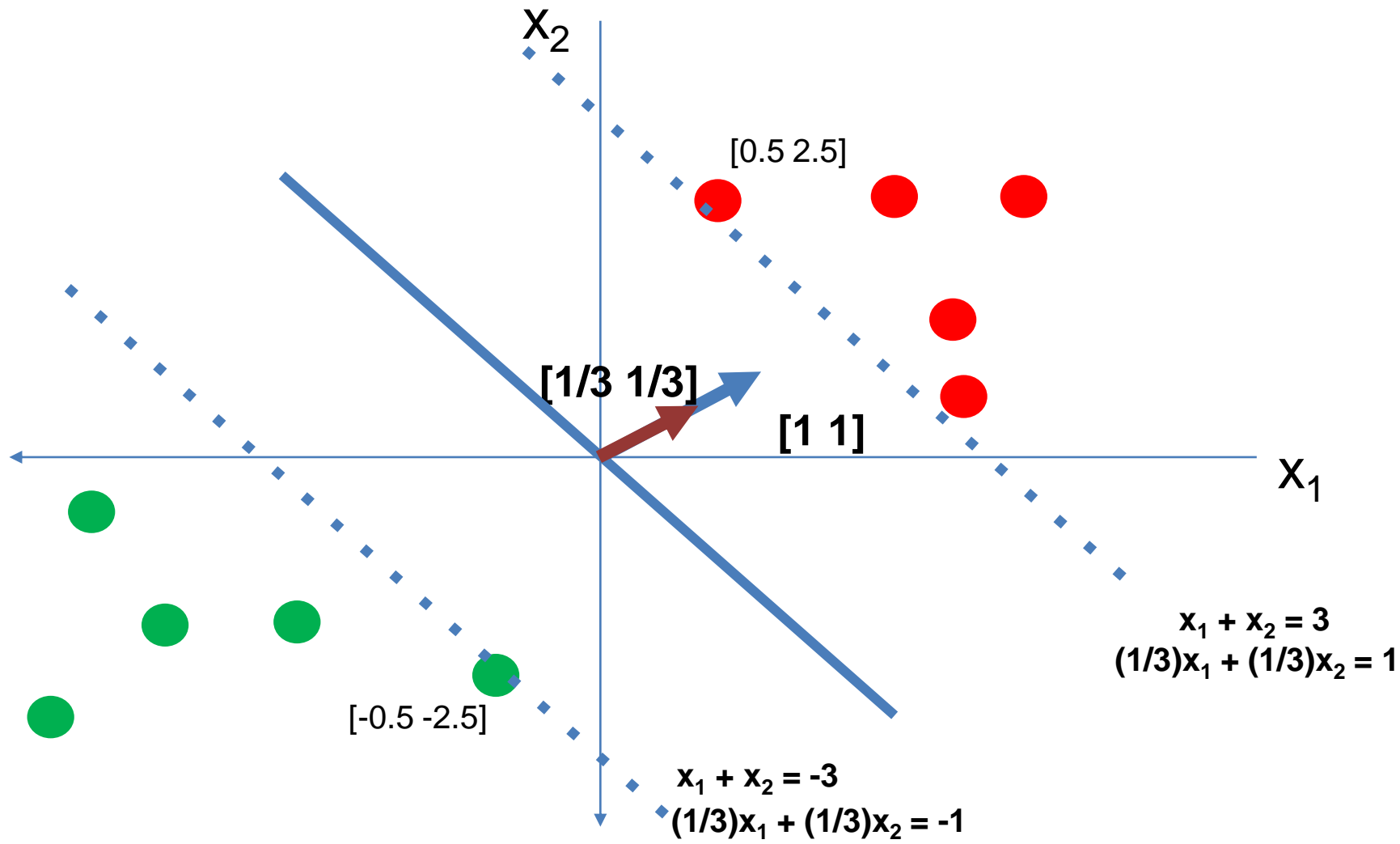**Issue:** Can scale $w$ arbitrarily.

Possible fix

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

$$(w^T x_i) y_i \geq \gamma$$

$$\|w\|^2 = 1$$

$X_2$

[0.5 2.5]

[1/3 1/3]

[1 1]

$X_1$

$x_1 + x_2 = 3$
$(1/3)x_1 + (1/3)x_2 = 1$

[-0.5 -2.5]

$x_1 + x_2 = -3$
$(1/3)x_1 + (1/3)x_2 = -1$

$x_2$

[0.5 2.5]

[1 1]

$x_1$

[-0.5 -2.5]

Notice what happens to the lengths of the w as we adjust it to have margin 1

# OBSERVATIONS

-> Once a direction is fixed,
the width between the margin lines
is fixed

-> If the width is large, then the w that
achieves margin 1 in that direction
has smaller length

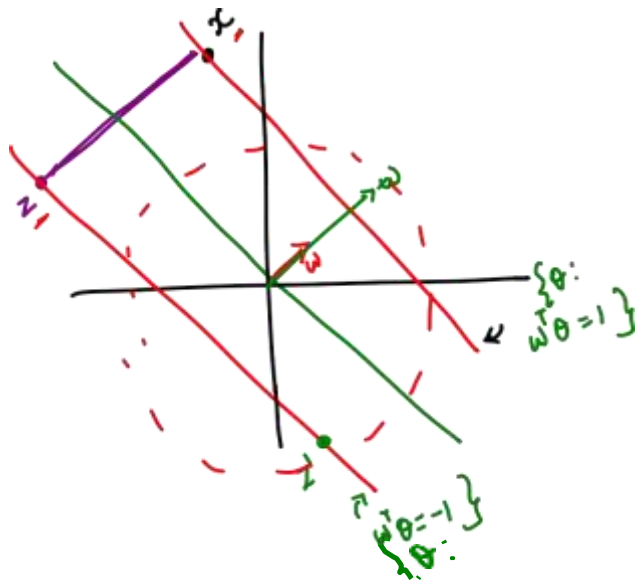-> If the width is small, then the w that
achieves margin 1 in that direction
has larger length

-> In general, width(w) seems to be
inversely proportional to length(w)

$$\max_{w} \quad \boxed{width(w)}$$

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

What is width(w) ?

$$\min_{z:} \quad \frac{1}{2} \| x - z \|^2 \leftarrow$$

$$s.t \quad \begin{matrix} w^{\mathsf{T}} x = 1 \\ w^{\mathsf{T}} z = -1 \end{matrix} \Bigg\}$$

[ Exercise ]

$$\text{width} (\omega) = \frac{2}{\|w\|^2}$$

$$\max_{w} \quad \boxed{\frac{2}{||w||^2}}$$

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



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

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

## Issues

- L·S is a strong assumption

- Non-linear structure?

$$\min_{\omega} \; f(\omega)$$

$$g(\omega) \quad \leq 0$$

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

Fix some $\omega$.

Consider

$$\boxed{\max_{\alpha \geq 0} \ L(w, \alpha)} \leftarrow$$

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

$$\begin{cases} \infty & \text{if} & g(w) > 0 \\ \\ f(w) & \text{if} & g(w) \leq 0 \end{cases}$$

$L(w_2, \alpha)$

$\cdot w_2$

$f(w)$

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

$$\min_{\omega} \left[ \max_{\alpha \geq 0} \frac{B(\omega)}{L(\omega, \alpha)} \right]$$

equivalent

$$\equiv \quad \min_{\omega} \quad F(\omega)$$

$$\text{s.t} \quad g(\omega) \leq 0.$$

- Can we swap min and max?

mi
W

# Multiple Constraints

$\rightarrow$ Same idea

$$\min_{\omega} \quad f(\omega)$$

$$\text{s.t} \quad g_i(\omega) \leq 0 \quad \substack{\forall i \\ = 1 \cdots k}$$

$$\equiv \quad \min_{\omega} \left[ \max_{\substack{\{\alpha_1, \cdots \\ \alpha_k \geq 0\}}} \left[ f(\omega) + \sum_{i=1}^{k} \alpha_i g_i(\omega) \right] \right]$$

$$|||\quad \text{Strong duality} \underset{\text{convex } f, g_i}{\text{for}}$$

$$\max_{\alpha_1, \cdots, \alpha_k \geq 0} \quad \min_{\omega} \quad f(\omega) + \sum_{i=1}^{k} \alpha_i g_i(\omega)$$

$$\min_{w} \quad \frac{1}{2} \|w\|^2 \quad \leftarrow f(w)$$

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

$$1 - (w^T x_i) y_i \leq 0$$

$$g_i(w) = 1 - (w^T x_i) y_i$$

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

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

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

$$\hookrightarrow \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} \geq 0$$

$$|||$$

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

Fix $\quad \alpha \geq 0$

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

Grad w.r.t w

$$w^* + \sum_{i=1}^{n} - \alpha_i x_i y_i = 0$$

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

$\in R^d$

$\{+1, -1\}$

Fixed Choice

In matrix notation

$$\boxed{w^* = X Y \alpha}$$

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

Substituting $\underset{=}{\text{soln}}$ back in the objective.

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

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

$$\frac{1}{2}(XY\alpha)^T(XY\alpha) \; + \; \alpha^T 1 \; - \; \sum_{i=1}^{n}(XY\alpha)^T x_i \, y_i \, \alpha_i$$

$$\begin{bmatrix} \alpha_i \\ \vdots \\ \alpha_n \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

On   Simplification   $\begin{bmatrix} \text{please-} \\ \text{do This} \end{bmatrix}$

$$\alpha^T 1 \; - \; \frac{1}{2}(XY\alpha)^T(XY\alpha)$$

Can be KERNELIZED!

__DUAL   PROBLEM__

Solving in n instead of d-

$$\max \quad \alpha^T 1 \; - \; \frac{1}{2}\,\alpha^T Y^T X^T X Y \alpha$$

$\boxed{\alpha \geq 0}$

↳ easy constraints

→ n×n.

# Revisiting The Lagrangian

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

PRIMAL                                    DUAL

$\omega^*$                                $\alpha^*$

$$\max_{\alpha \geq 0} \ f(\omega^*) + \alpha \, g(\omega^*) \quad = \quad \min_{w} \ f(\omega) + \alpha^* g(\omega)$$

$$f(\omega^*) \quad = \quad f(\dot\omega) + \alpha^* g(\omega')$$

$$f(\omega^*) \quad \leq \quad f(\omega^*) + \alpha^* g(\omega^*)$$

$$\Rightarrow \alpha^* g(\omega^*) \geq 0$$

But we know $\alpha^* g(\omega^*) \leq 0$

$$\Rightarrow \quad \alpha^* g(\omega^*) = 0 \quad \longrightarrow \quad \text{COMPLEMENTARY SLACKNESS}$$

# For multiple constraints

$$\alpha_i^* \, g_i(\omega^*) = 0 \quad \forall i$$

For our problem

$$\boxed{\alpha_i^* \left( 1 - (w^T x_i) y_i \right) = 0 \qquad \forall i}$$



Supporting hyperplanes

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

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



## Issues

- L·S is a strong assumption

- Non-linear structure ?

**So far**

Support Vector Machines
Primal Problem – Margin Maximization
Dual Problem
- Kernel Version

**Now**

- What if there are **outliers** in the problem?

**Idea** (to deal with outliers):

Fix any $w$. $w$ classifies some points correct and some incorrectly. Let the incorrect points pay "bribe" to get to the correct side.

# Modified formulation

$c \geqslant 0$ [hyper parameter]

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

$\rightarrow (\vec{w}^T x_i)\, y_i + \xi_i \geqslant 1 \quad \leftarrow \quad \forall i$

$\Rightarrow \quad \xi_i \geqslant 0 \quad \leftarrow \quad \forall i$

If $\quad c = 0 \quad \Rightarrow \quad$ Bribes don't cost $\Rightarrow \quad$ $w = 0$ is solution

$c \rightarrow \infty \quad \Rightarrow \quad$ Bribes are too costly $\Rightarrow \quad$ Linear separable case.

$$L\left(\omega, \xi, \alpha, \beta\right) = \frac{1}{2}\|w\|^2 + c\left(\underbrace{\sum_{i=1}^{n} \xi_i}_{\uparrow}\right) + \underline{\sum_{i=1}^{n}\alpha_i\left(1 - \underset{-\xi_i}{(\omega^T x_i)\, y_i}\right)}$$

$$+ \sum_{i=1}^{n} \beta_i \underset{\uparrow}{(-\xi_i)}$$

Dual:

$$\max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \quad \min_{\omega} \quad L\left(\omega, \xi, \alpha, \beta\right)$$

$$\frac{\partial L}{\partial \omega} = 0 \quad \Rightarrow \quad \omega^* = \sum_{i=1}^{n} \alpha_i\, x_i\, y_i \qquad \boxed{\omega^* = x\, y\, \alpha}$$

$$\frac{\partial L}{\partial \xi_i} = 0 \implies \boxed{C - \alpha_i - \beta_i = 0}$$

$$\boxed{\alpha_i + \beta_i = C} \quad \forall i$$

Substitute $\vec{v} = XY\alpha$ in the original objective

$$\frac{1}{2}(XY\alpha)^T(XY\alpha) + \sum_{i=1}^{n}(C - \alpha_i - \beta_i)\xi_i + \alpha^T 1$$
$$- (XY\alpha)^T(XY\alpha)$$

SOFT - MARGIN

SUPPORT

VECTOR

MACHINE

$$\begin{array}{l} \max \\ \alpha \geq 0 \\ \beta \geq 0 \\ \alpha_i + \beta_i = C \end{array} \qquad \alpha^T 1 \; - \; \frac{1}{2} (x y \alpha)^T (x y \alpha)$$

$\equiv$

$$\begin{array}{l} \max \\ 0 \leq \alpha \leq C \end{array} \qquad \alpha^T 1 \; - \frac{1}{2} \; \underbrace{\alpha^T y^T (x^T x) y \alpha}$$

BOX
CONSTRAINT.

## PRIMAL

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

$$\text{st} \quad \underbrace{(W^T x_i) y_i}_{1 - W^T x_i y_i \leq 0} \geq 1 \quad \ast i$$

## DUAL

$$\max_{\alpha \geq 0} \quad \alpha^T 1 - \alpha^T y^T x^T x y \alpha$$

$$\alpha_i^* (1 - W^{*T} x_i y_i) = 0 \qquad \ast i$$

$$\boxed{W^* = \sum_{i \geq 1}^{n} \alpha_i^* x_i y_i}$$

## PRIMAL ✓

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

$$\text{st} \quad \begin{cases} (W^T x_i) y_i + \xi_i \geq 1 \quad \ast i = 1, \dots n \\ \xi_i \geq 0 \quad \ast i = 1, \dots n \end{cases}$$

$\alpha$

$\beta$

## DUAL ✓

$$\max \quad \alpha^T 1 - \alpha^T y^T x^T x y \alpha$$

$$\boxed{\alpha + \beta = C}$$

$$\alpha \geq 0$$

$$\beta \geq 0$$

$$0 \leq \alpha \leq C$$

- Let $(\underline{w}^*, \underline{\xi}^*)$ be the primal optimal solution

- Let $(\underline{\alpha}^*, \underline{\beta}^*)$ be the dual optimal solution
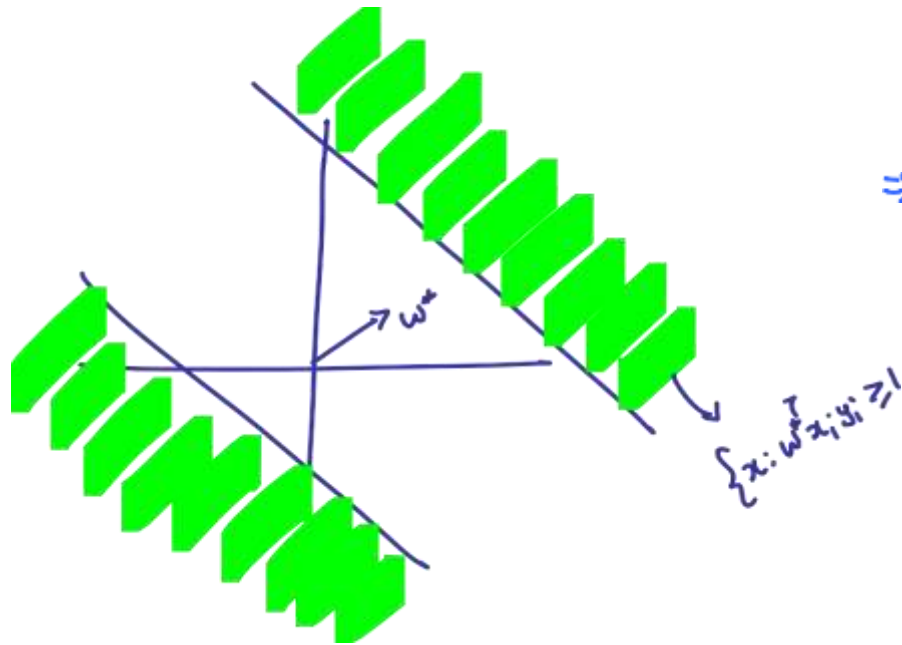
## COMPLEMENTARY SLACKNESS

$$\underline{\alpha_i^*} \left( 1 - (w^{*T} x_i) y_i - \xi_i^* \right) = 0 \qquad \forall i$$

$$\beta_i^* \xi_i^* = 0 \qquad \forall i$$

$\alpha_i^* + \beta_i^* = c$

$\forall i$

$\hookrightarrow$ Ⓐ

Various cases possible

**Case 1:** $\qquad \alpha_i^* = 0 \qquad \overset{\text{(A)}}{\Rightarrow} \qquad \beta_i^* = C \qquad \overset{\boxed{\text{CS}}}{\Rightarrow} \quad \xi_i^* = 0$

$$1 - (w^{*^T} x_i) y_i - \underset{=0}{\underline{\xi_i^*}} \leq 0 \quad [\text{Primal feasibility}]$$

$$\Rightarrow \quad 1 - (w^{*^T} x) y_i \leq 0$$

$$\Rightarrow \qquad w^{*^T} x_i \, y_i \geq 1$$

$$\Rightarrow \quad w^* \text{ classifies } (x_i, y_i) \text{ correctly.}$$

$\{ x : w^T x_i y_i \geq 1 \}$

$\to w^*$

Case 2:  $\quad 0 < \alpha_i^* < C \quad \overset{\textcircled{A}}{\Longrightarrow} \quad 0 < \beta_i^* < C \quad \overset{cs}{\Longrightarrow} \quad \xi_i^* = 0$

$\Big\Vert$ $\boxed{cs}$

$1 - (w^{*T} x_i) y_i - \xi_i^* = 0$

$\Downarrow$

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

$\Longrightarrow \quad$ $(x_i, y_i)$ lies on the supporting hyperplane.

Case 3:

$$\alpha_i^+ = C \implies \beta_i^+ = 0 \implies \xi_i^+ \geq 0$$

$$\Downarrow \boxed{CS}$$

$$1 - w^{*T} x_i y_i - \xi_i^+ = 0$$

$$\xi_i^+ = 1 - w^{*T} x_i y_i \geq 0$$

$$\implies \boxed{w^{*T} x_i y_i \leq 1}$$

Let's see this from P.o.v of data

CASE 1

$$\boxed{w^{*T} x_i y_i < 1}$$

$$1 - w^T x_i y_i - \xi_i^* \leq 0$$

$$w^{*T} x_i y_i \geq 1 - \xi_i^*$$

$$\boxed{\xi_i^* \geq 1 - w^{*T} x_i y_i}$$

$$\Rightarrow \xi_i^* > 0 \Rightarrow \beta_i^* = 0 \Rightarrow \alpha_i^* = C$$

$$\alpha_i^* \left( 1 - w^T x_i y_i - \xi_i^* \right) = 0$$

$$\beta_i^* \xi_i^* = 0.$$

CASE 2 : $\quad w^{*T} x_i y_i = 1$

$$\xi_i^* \geq 1 - w^{*T} x_i y_i$$

$$\Rightarrow \quad \xi_i^* \geq 0 \quad \Rightarrow \quad \alpha_i^* \in [0, C]$$

CASE 3 $\quad w^{*T} x_i y_i > 1$

$$1 - w^{*T} x_i y_i - \xi_i^* \leq 0 \quad [\text{Primal feasibility}]$$

$$\Rightarrow \quad 1 - w^{*T} x_i y_i - \xi_i^* < 0 \quad \overset{\text{c.s}}{\Rightarrow} \quad \alpha_i^* = 0$$

# SUMMARY

$\alpha_i^* = 0 \qquad \Rightarrow \qquad w^{*T} x_i y_i \geq 1$

$0 < \alpha_i^* < C \qquad \Rightarrow \qquad w^{*T} x_i y_i = 1$

$\alpha_i^* = \underline{C} \qquad \Rightarrow \qquad w^{*T} x_i y_i \leq 1$

---

✓ $w^T x_i y_i < 1 \qquad \Rightarrow \qquad \underline{\alpha_i^* = C}$

$\rightarrow w^T x_i y_i = 1 \qquad \Rightarrow \qquad \underline{\alpha_i^* \in [0, C]}$

$\rightarrow w^T x_i y_i > 1 \qquad \Rightarrow \qquad \underline{\alpha_i^* = 0.}$

Binary classification

GENERATIVE

Naïve Bayes

G·D·A

DISCRIMINATIVE

- k-NN
- Decision trees
- Perceptron
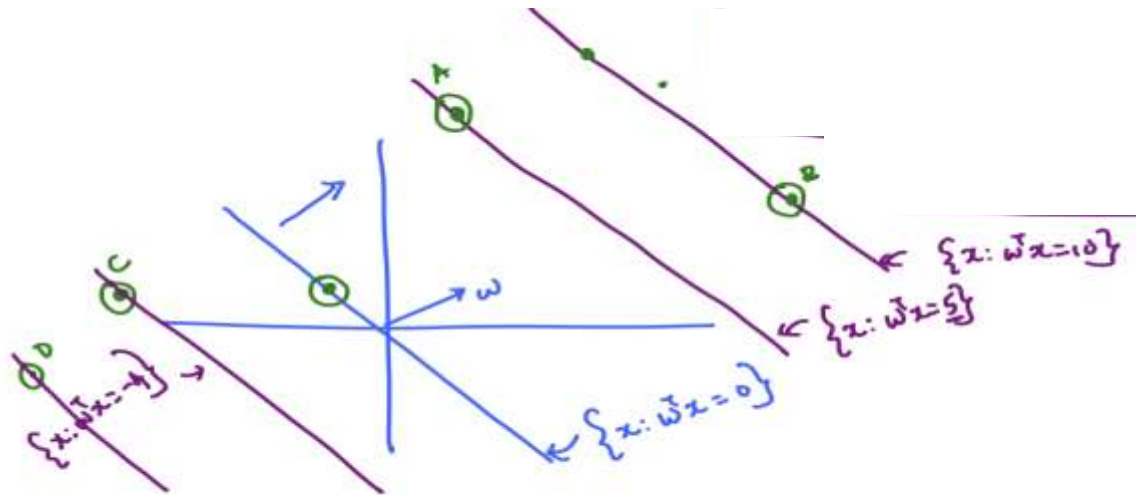- Support-vector-machines

Dosen't "really" model $\boxed{P(y/x)}$

$\rightarrow$ Just finds $f: \mathbb{R}^d \rightarrow \{\pm 1\}$

- Can we model $P\left(y = +1/x\right)$ differently?

Start with a simple model

Given $x \in \mathbb{R}^d$   $z = \vec{w}^T x$   $w \in \mathbb{R}^d$.



$\{x: \vec{w}^T x = 10\}$

$\{x: \vec{w}^T x = 5\}$

$\{x: \vec{w}^T x = 0\}$
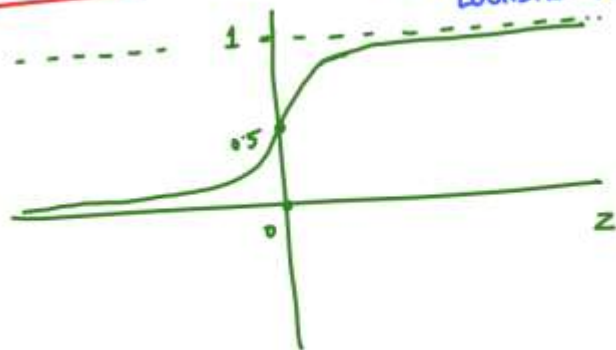
$\{x: \vec{w}^T x = -5\}$

$$P\left(y = +1 / x\right) = g(w^T x)$$

- $g(z) \in [0,1]$

- $g(z) \rightarrow \boxed{1}$ as $z \rightarrow \boxed{\infty}$

- $g(z) \rightarrow 0$ as $z \rightarrow -\infty$

LINK FUNCTION

- $g(z) = 0.5$ if $z = 0.$

ONE   POPOLAR  CHOICE

SIGMOID  FUNCTION
LOGISTIC  FUNCTION.

$$g(z) = \dfrac{1}{1 + e^{-z}}$$



MODEL :   LOGISTIC  REGRESSION

Data: $\{ (x_1, y_1) \cdots (x_n, y_n) \}$      $x_i \in \mathbb{R}^d$
                                              $y_i \in \{0, 1\}$

Max. Likelihood

$$L(w, \text{Data}) = \prod_{i=1}^{n} \left( g(w^T x_i) \right)^{y_i} \left( 1 - g(w^T x_i) \right)^{(1-y_i)}$$

$$\log L(w, \text{Data}) = \sum_{i=1}^{n} y_i \log\left( g(w^T x_i) \right) + (1-y_i) \log\left( 1 - g(w^T x_i) \right)$$

$$= \sum_{i=1}^{n} \left[ y_i \log \left( \frac{1}{1 + e^{-w^T x_i}} \right) + (1 - y_i) \log \left( \frac{e^{-w^T x_i} \times 1}{1 + e^{-w^T x_i}} \right) \right]$$

$$= \sum_{i=1}^{n} \left[ \log \left( \frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}} \right) - y_i (-w^T x_i) \right]$$

$$= \sum_{i=1}^{n} \left[ (1 - y_i)(-w^T x_i) - \log \left( 1 + e^{-w^T x_i} \right) \right]$$

- No closed form solution

- ## Gradient ascent

$$\nabla \log L(w) = \sum_{i=1}^{n} (1-y_i)(-x_i) - \frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}} (-x_i)$$

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

$$= \sum_{i=1}^{n} \overset{\underset{\mathbb{R}^d}{\downarrow}}{x_i} \left( \overset{\underset{\{0,1\}}{\downarrow}}{y_i} - \overset{\underset{g(w^T x_i)}{\downarrow}}{\frac{1}{1 + e^{-w^T x_i}}} \right)$$

$$W_{t+1} = W_t + \eta_t \nabla \log L(W_t)$$

## REGULARIZED VERSION

$$\min_{w} \quad \sum_{i=1}^{n} (1-y_i)\, w^T x_i \; + \; \log\left(1 + e^{-w^T x_i}\right) \quad + \quad \frac{\lambda}{2} \|w\|^2$$

---

## KERNEL VERSION

- Can argue $\quad w = \sum_{i=1}^{n} \alpha_i x_i$ 

  Formal Theorem

  Representer Theorem

Exercise: Derive the kernel version of logistic regression

# META CLASSIFIERS (or)

## ENSEMBLE CLASSIFIERS.

WEAK
CLASSIFIERS
[better than random]
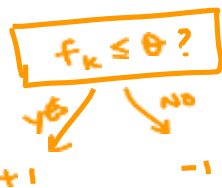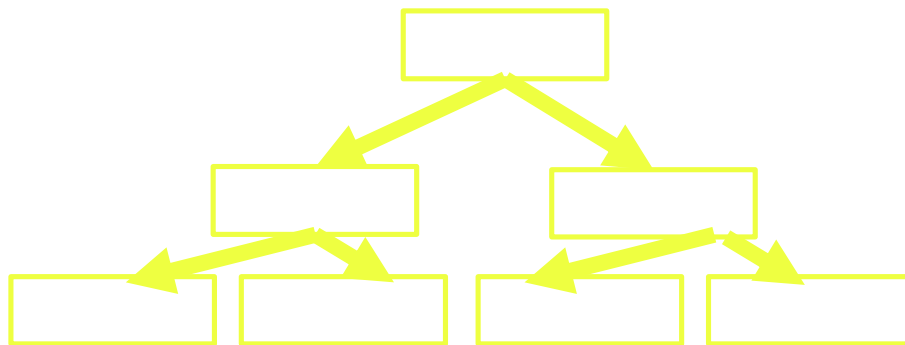
→

STRONG
CLASSIFIERS

# Weak classifiers

## Overfit decision tree

DECISION STUMP

$f_k \leq \theta$ ?

YES        NO

+1          −1

**high bias, low variance**

…

……………

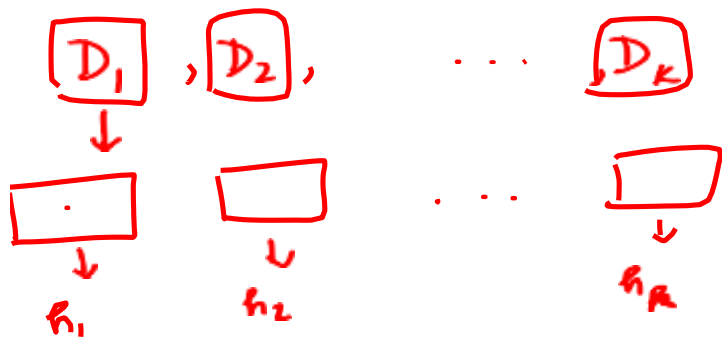**low bias, high variance**

$$X_1, X_2, \ldots X_n \sim N(M, 1)$$

$$\hat{M}_1 = X_1 \qquad \hat{M}_2 = X_2 , \ldots \hat{M}_N = X_n \qquad \hat{M}_{ML} = \frac{1}{n} \Sigma x_i$$

Overfit decision trees



$$h^*(x) = \text{majority}\left( h_1(x), \ldots , h_k(x) \right)$$

$$h_i : \mathbb{R}^d \to \{\pm 1\}$$

BAGGING — Bootstrap Aggregation.

$$D = \{ (x_1, y_1) \cdots (x_n, y_n) \}$$

Chance that a point appears in a dataset

$$1 - \left(1 - \frac{1}{n}\right)\left(1 - \frac{1}{n}\right)\cdots\left(1 - \frac{1}{n}\right)$$

$$1 - \left(1 - \frac{1}{n}\right)^n$$

$$1 - \frac{1}{e} \quad (\text{as } n \to \infty)$$

$$\approx 66\%$$

— Create datasets $D_1, \cdots, D_k$ from $D$ by

"Sampling with replacement"

— Run weak classifier on $D_1, \cdots, D_k$ to get $h_1, \cdots, h_k$

— Aggregate $h_1, \cdots, h_k$ using majority.

FEATURE BAGGING → Bag the features in addition to data points

Feature bagged decision trees -> RANDOM FOREST

BOOTSTRAP - Sampling with Replacement ?

?

AGGREGATION - Majority.

# BOOSTING

[ Freund & Schappire.
1995
Gödel Prize ]

ADA-BOOST

↑

---

Distribution $D$ over $(\mathcal{X} \times \mathcal{Y})$

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{+1, -1\}$

$\longrightarrow$ unknown but fixed.

$x_1, \dots, x_n$ are iid from $D$.

$$h : \underset{x}{\mathbb{R}^d} \longrightarrow \underset{y}{\{\pm 1\}}$$

Measure performance using

Misclassification probability.

$$P_{(x,y) \sim D} \left( h(x) \neq y \right)$$

A weak learner is one which outputs a Classifier

Strong

$h$ for which

$$P_{x,y \sim D} \left( h(x) = y \right) \geq \frac{1-\epsilon}{\frac{1}{2} + \gamma} \qquad \boxed{\gamma > 0}$$

for any unknown but fixed distribution $D$.

$x_2$

[0.5 2.5]

[1/3 1/3]

[1 1]

$x_1$

$x_1 + x_2 = 3$
$(1/3)x_1 + (1/3)x_2 = 1$

[-0.5 -2.5]

$x_1 + x_2 = -3$
$(1/3)x_1 + (1/3)x_2 = -1$

Notice what happens
to the lengths
of the w as we
adjust it to have
margin 1

$X_2$

[0.5 2.5]

[1 1]

$X_1$

[-0.5 -2.5]

## OBSERVATIONS

-> Once a direction is fixed,
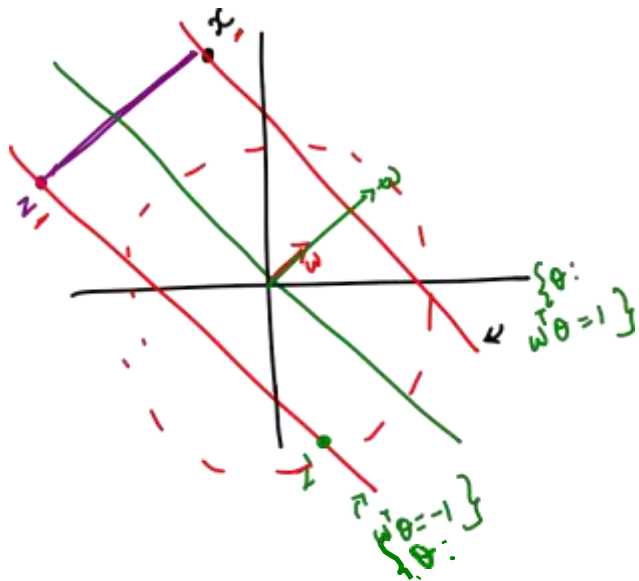the width between the margin lines
is fixed

-> If the width is large, then the w that
achieves margin 1 in that direction
has smaller length

-> If the width is small, then the w that
achieves margin 1 in that direction
has larger length

-> In general, width(w) seems to be
inversely proportional to length(w)

$$\max_{w} \boxed{\text{width}(w)}$$

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

What is width(w) ?

$$\min_{z:} \quad \frac{1}{2} \|x - z\|^2 \leftarrow$$

$$\text{s.t} \quad \left. \begin{array}{l} w^T x = 1 \\ w^T z = -1 \end{array} \right\}$$

[ Exercise ]

$$\text{width}(\omega) = \frac{2}{\|w\|^2}$$

$$\max_{\omega} \boxed{\frac{2}{||\omega||^2}}$$

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

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

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

<u>Issues</u>

- L.S is a strong assumption

- Non-linear structure?

DETOUR

$$\min_{\omega} \; f(\omega)$$

$$g(\omega) \quad \leq 0 \quad \leftarrow$$

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

Fix some $\omega$.

Consider

$$\max_{\alpha \geq 0} \; L(w, \alpha)$$

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



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

$$\begin{cases} \infty & \text{if} \quad g(w) > 0 \\ f(w) & \text{if} \quad g(w) \leq 0 \end{cases}$$

$$\min_{\omega} \left[ \max_{\alpha \geq 0} \frac{B(\omega)}{\mathcal{L}(\omega, \alpha)} \right] \underset{\equiv}{\text{equivalent}} \quad \min_{\omega} \; F(\omega)$$

$$\text{s.t} \quad g(\omega) \leq 0.$$

- Can we swap min and max?

mi
w

# Multiple Constraints

$\rightarrow$ Same idea

$$\min_{\omega} \quad f(\omega)$$

$$s.t \quad g_i(\omega) \leq 0 \quad \substack{\forall i \\ =1 \cdots k}$$

$\equiv$

$$\min_{\omega} \left[ \max_{\substack{\{\alpha_1, \cdots \\ \alpha_k \geq 0\}}} \left[ f(\omega) + \sum_{i=1}^{k} \alpha_i \, g_i(\omega) \right] \right]$$

$|||$ Strong duality $\substack{\text{for} \\ \text{convex } f, g_i}$

$$\max_{\alpha_1, \cdots, \alpha_k \geq 0} \quad \min_{\omega} \quad f(\omega) + \sum_{i=1}^{k} \alpha_i \, g_i(\omega)$$

$$\min_{w} \quad \frac{1}{2} \|w\|^2 \quad \leftarrow f(w)$$

$$g_i(w) = 1 - (w^T x_i) y_i$$

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

$$1 - (w^T x_i) y_i \leq 0$$

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

$$\uparrow$$

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

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

$$\hookrightarrow \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} \geq 0$$

$$|||$$

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

Fix $\quad \alpha \geq 0$

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

Grad  w.r.t  w

$$w^* + \sum_{i=1}^{n} -\alpha_i x_i y_i = 0$$

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

∈ R^d

{+1, -1}

Fixed
Choice

In  matrix  notation

$$\boxed{w^* = X Y \alpha}$$

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

Substituting $\frac{soln}{=}$ back in the objective.

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

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

$$\frac{1}{2}(XY\alpha)^T(XY\alpha) + \alpha^T 1 - \sum_{i=1}^{n}(XY\alpha)^T x_i \, y_i \, \alpha_i$$

$$\begin{bmatrix} \alpha_i \\ \vdots \\ \alpha_n \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

On Simplification $\begin{bmatrix} \text{please-} \\ \text{do This} \end{bmatrix}$

$$\alpha^T 1 - \frac{1}{2}(XY\alpha)^T(XY\alpha)$$

Can be KERNELIZED!

**DUAL PROBLEM**

Solving in n instead of d-

$$\max \quad \boxed{\alpha \geq 0} \quad \alpha^T 1 - \frac{1}{2}\alpha^T Y^T \overset{T}{\boxed{X^T X}} Y \alpha$$

↓ easy constraints

→ nxn.

# Revisiting The Lagrangian

$$\underset{\substack{w}}{min} \left[ \underset{\alpha \geq 0}{max} \quad f(w) + \alpha g(w) \right] \equiv \underset{\alpha \geq 0}{max} \left[ \underset{w}{min} \quad f(w) + \alpha g(w) \right]$$

PRIMAL                             DUAL

$w^*$                                          $\alpha^*$

$$\underset{\alpha \geq 0}{max} \quad f(w^*) + \alpha \, g(w^*) = \underset{w}{min} \quad f(w) + \alpha^* g(w)$$

$$f(\omega^*) \quad = \quad f(\dot{\omega}) + \alpha^* g(\omega')$$

$$f(\omega^*) \quad \leq \quad f(\omega^*) + \alpha^* g(\omega^*)$$

$$\Rightarrow \alpha^* g(\omega^*) \geq 0$$

But we know $\quad \alpha^* g(\omega^*) \leq 0$

$$\Rightarrow \quad \boxed{\alpha^* g(\omega^*) = 0} \quad \rightarrow \quad \text{COMPLEMENTARY SLACKNESS}$$

# For multiple constraints

$$\alpha_i^* \, g_i(\omega^*) = 0 \quad \forall i$$

For our problem

$$\boxed{\alpha_i^* \left( 1 - (\omega^T x_i) y_i \right) = 0 \qquad \forall i}$$



Supporting hyperplanes

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

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

Issues

- L·S is a strong assumption

- Non-linear structure?

**So far**

Support Vector Machines
Primal Problem – Margin Maximization
Dual Problem
- Kernel Version

**Now**

- What if there are **outliers** in the problem?

**Idea** (to deal with outliers):

Fix any $w$. $w$ classifies some points correct and some incorrectly. Let the incorrect points pay "bribe" to get to the correct side.

## Modified formulation

$C \geqslant 0$  [hyper parameter]

$$\min_{w, \varepsilon} \quad \frac{1}{2} \|w\|^2 \quad + C \sum_{i=1}^{n} \varepsilon_i$$

→ $(\vec{w}^T x_i) \, y_i + \varepsilon_i \geqslant 1$  ←  $\forall i$

⟹ $\varepsilon_i \geqslant 0$  ←  $\forall i$

If  $C = 0$  ⟹  Bribes don't cost ⟹  $w = 0$ is solution

$C \to \infty$  ⟹  Bribes are too costly  ⟹  Linear separable case.

$$L(\omega, \xi, \alpha, \beta) = \frac{1}{2}\|w\|^2 + c\left(\underbrace{\sum_{i=1}^{n}\xi_i}_{\uparrow}\right) + \sum_{i=1}^{n}\alpha_i\left(1 - (\omega^T x_i)\, y_i\right)$$

$$\underset{-\xi_i}{\phantom{x}}$$

$$+ \sum_{i=1}^{n}\beta_i\,(-\xi_i)$$

Dual: $\underset{\substack{\alpha \geq 0 \\ \beta \geq 0}}{max} \quad \underset{\omega}{min} \quad L(\omega, \xi, \alpha, \beta)$

$$\frac{\partial L}{\partial \omega} = 0 \implies \omega^* = \sum_{i=1}^{n}\alpha_i\, x_i\, y_i \qquad \boxed{\omega^* = x\, y\, \alpha}$$

$$\frac{\partial L}{\partial \xi_i} = 0 \quad \Rightarrow \quad \boxed{C - \alpha_i - \beta_i = 0}$$

$$\boxed{\alpha_i + \beta_i = C} \quad \forall i$$

Substitute $\vec{v} = XY\alpha$ in the original objective

$$\frac{1}{2}(XY\alpha)^T(XY\alpha) + \sum_{i=1}^{n}(C - \alpha_i - \beta_i)\xi_i + \alpha^T 1$$
$$- (XY\alpha)^T(XY\alpha)$$

SOFT - MARGIN

SUPPORT

VECTOR

MACHINE

$$\max_{\substack{\alpha \geq 0 \\ \beta \geq 0 \\ \alpha_i + \beta_i = C}} \quad \alpha^T 1 - \frac{1}{2}(xy\alpha)^T(xy\alpha)$$

$=$

$$\max_{0 \leq \alpha \leq C} \quad \alpha^T 1 - \frac{1}{2} \underbrace{\alpha^T y^T (x^T x) y \alpha}$$

BOX
CONSTRAINT.

PRIMAL

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

$$\text{st} \quad \underbrace{(w^T x_i) y_i}_{1 - w^T x_i y_i \le 0} \ge 1 \quad \forall i$$

DUAL

$$\max_{\alpha \ge 0} \quad \alpha^T 1 - \alpha^T y^T x^T x y \alpha$$

$$\alpha_i^* (1 - w^{*T} x_i y_i) = 0 \quad \forall i$$

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

PRIMAL ✓

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

$$\text{st} \quad \begin{cases} \alpha \to (w^T x_i) y_i + \xi_i \ge 1 \quad \forall i = 1, \dots n \\ \beta \to \xi_i \ge 0 \quad \forall i = 1, \dots n \end{cases}$$

DUAL ✓

$$\max \quad \alpha^T 1 - \alpha^T y^T x^T x y \alpha$$

$$\boxed{\alpha + \beta = C}$$

$$\alpha \ge 0$$

$$\beta \ge 0$$

$$0 \le \alpha \le C$$

- Let $(\underline{w}^*, \underline{\xi}^*)$ be the primal optimal solution

- Let $(\underline{\alpha}^*, \underline{\beta}^*)$ be the dual optimal solution
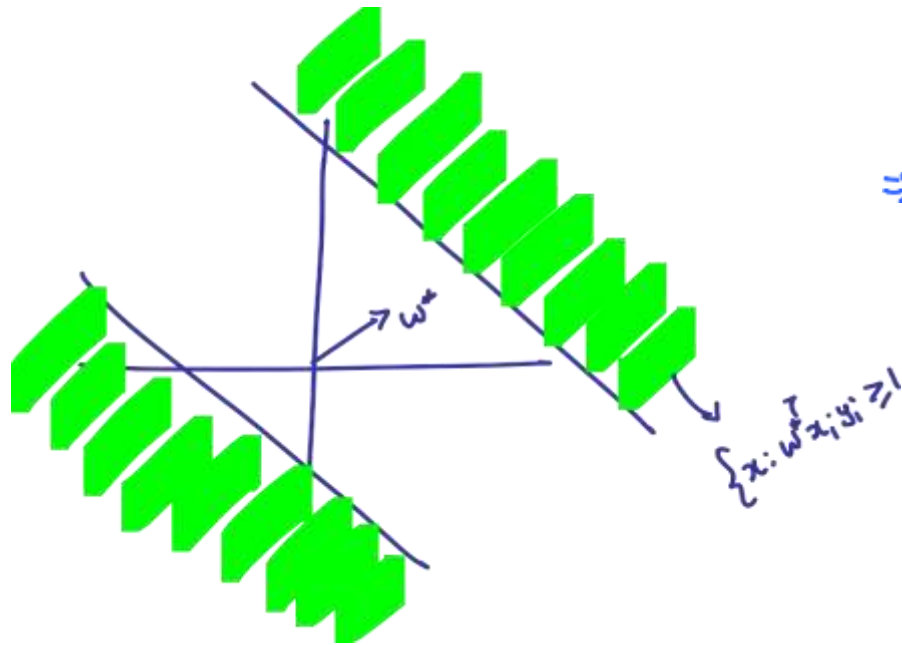
## COMPLEMENTARY SLACKNESS

$$\underline{\alpha_i^*} \left( 1 - (w^{*T} x_i) y_i - \xi_i^* \right) = 0 \qquad \forall i$$

$$\beta_i^* \xi_i^* = 0 \qquad \forall i$$

$$\alpha_i^* + \beta_i^* = c$$
$$\forall i$$
$$\hookrightarrow \text{(A)}$$

Various cases possible

**Case 1:** $\qquad$ $\alpha_i^* = 0$ $\qquad \overset{\textcircled{A}}{\Longrightarrow} \qquad \beta_i^* = C \qquad \overset{\boxed{CS}}{\Longrightarrow} \qquad \xi_i^* = 0$



$$1 - (w^{*T}x_i)y_i - \underset{=0}{\underline{\xi_i^*}} \leq 0 \quad \left[\text{Primal feasibility}\right]$$

$$\Rightarrow \qquad 1 - (w^{*T}x)y_i \leq 0$$

$$\Rightarrow \qquad w^{*T}x_i\, y_i \geq 1$$

$$\Rightarrow \quad w^* \text{ classifies } (x_i, y_i)$$
$$\text{correctly.}$$

$\{x_i : w^T x_i y_i \geq 1\}$

Case 2 :     $0 < \alpha_i^* < C$   $\overset{\text{\textcircled{A}}}{\Longrightarrow}$   $0 < \beta_i^* < C$   $\overset{CS}{\Rightarrow} \xi_i^* = 0$

$\Big\Vert \boxed{CS}$

$1 - (w^{*T} x_i) y_i - \xi_i^* = 0$

$\Downarrow$

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

$\Rightarrow$     $(x_i, y_i)$ lies on the

supporting hyperplane.

Case 3:
$$\alpha_i^+ = C \quad \Rightarrow \quad \beta_i^+ = 0 \quad \Rightarrow \quad \xi_i \geq 0$$

$$\Downarrow \boxed{CS}$$

$$1 - w^{*T} x_i y_i - \xi_i^+ = 0$$

$$\xi_i^+ = 1 - w^{*T} x_i y_i \geq 0$$

$$\Rightarrow \boxed{w^{*T} x_i y_i \leq 1}$$

Let's see this from P.o.v of data

CASE 1

$$\boxed{w^{*T} x_i y_i < 1}$$

$$1 - w^T x_i y_i - \xi_i^* \leq 0$$

$$w^{*T} x_i y_i \geq 1 - \xi_i^*$$

$$\boxed{\xi_i^* \geq 1 - w^{*T} x_i y_i}$$

$$\Rightarrow \xi_i^* > 0 \Rightarrow \beta_i^* = 0 \Rightarrow \alpha_i^* = C$$

$$\alpha_i^* \left( 1 - w^T x_i y_i - \xi_i^* \right) = 0$$

$$\beta_i^* \xi_i^* = 0.$$

CASE 2 : $\quad w^{*T} x_i y_i = 1$

$$\xi_i^* \geq 1 - w^{*T} x_i y_i$$

$$\Rightarrow \quad \xi_i^* \geq 0 \quad \Rightarrow \quad \alpha_i^* \in [0, c]$$

CASE 3 $\quad w^{*T} x_i y_i > 1$

$$1 - \underbrace{w^{*T} x_i y_i} - \xi_i^* \leq 0 \quad [\text{Primal feasibility}]$$

$$\Rightarrow \quad 1 - w^{*T} x_i y_i - \xi_i^* < 0 \quad \overset{\text{C.S}}{\Rightarrow} \quad \alpha_i^* = 0$$

# SUMMARY

$$\alpha_i^* = 0 \implies w^{*T} x_i y_i \geq 1$$

$$0 < \alpha_i^* < C \implies w^{*T} x_i y_i = 1$$

$$\alpha_i^* = \underline{C} \implies w^{*T} x_i y_i \leq 1$$

✓ $w^{*T} x_i y_i < 1 \implies \alpha_i^* = C$

$\rightarrow$ $w^{*T} x_i y_i = 1 \implies \alpha_i^* \in [0, C]$

$\rightarrow$ $w^{*T} x_i y_i > 1 \implies \alpha_i^* = 0.$

Binary classification

GENERATIVE

Naïve Bayes

G.D.A

DISCRIMINATIVE

k-NN
Decision trees
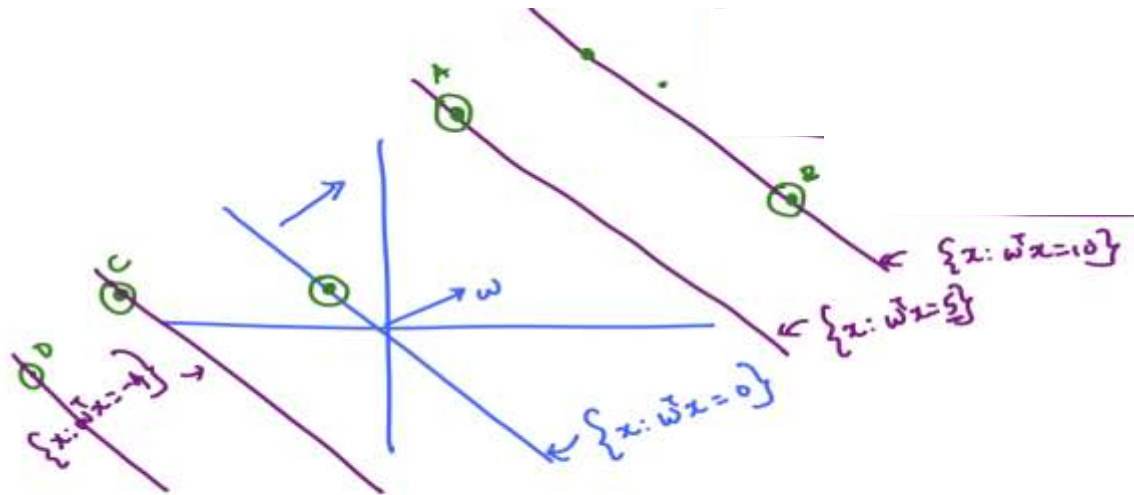Perceptron
Support-vector-machines

Dosen't "really" model $\boxed{P(y/x)}$

$\rightarrow$ Just finds $f: \mathbb{R}^d \rightarrow \{\pm 1\}$

- Can we model $P\left(y = +1/x\right)$ differently?

Start with a simple model

Given $x \in \mathbb{R}^d$     $z = \vec{w}^T x$     $w \in \mathbb{R}^d$.

$$\boxed{P\left(y = +1 \Big/ x\right) = g(w^T x)}$$
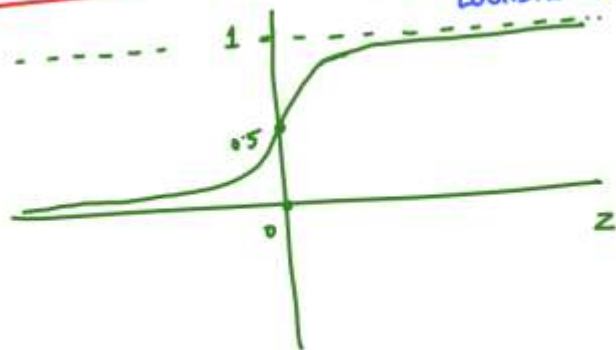
LINK
FUNCTION $\longrightarrow$

- $\underline{g(z)} \in \underline{[0,1]}$

- $g(z) \rightarrow \boxed{1}$ as $z \rightarrow \boxed{\infty}$

- $g(z) \rightarrow 0$ as $z \rightarrow -\infty$

- $g(z) = 0.5$ if $z = 0.$

ONE POPOLAR CHOICE

SIGMOID FUNCTION
LOGISTIC FUNCTION.

$$g(z) = \frac{1}{1 + e^{-z}}$$

MODEL : LOGISTIC REGRESSION

Data: $\{ (x_1, y_1) \ \cdots \ (x_n, y_n) \}$ 
$\qquad x_i \in \mathbb{R}^d$
$\qquad y_i \in \{0, 1\}$

Max. Likelihood

$$L(w, \text{Data}) = \prod_{i=1}^{n} \left( g(w^T x_i) \right)^{y_i} \left( 1 - g(w^T x_i) \right)^{(1-y_i)}$$

$$\log L(w, \text{Data}) = \sum_{i=1}^{n} y_i \ \log \left( g(w^T x_i) \right) + (1-y_i) \log \left( 1 - g(w^T x_i) \right)$$

$$= \sum_{i=1}^{n} \left[ y_i \; \log \left( \frac{1}{1 + e^{-w^T x_i}} \right) + (1 - y_i) \; \log \left( \frac{e^{-w^T x_i} \times 1}{1 + e^{-w^T x_i}} \right) \right]$$

$$= \sum_{i=1}^{n} \left[ \log \left( \frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}} \right) - y_i \left( -w^T x_i \right) \right]$$

$$= \sum_{i=1}^{n} \left[ (1 - y_i) \left( -w^T x_i \right) - \log \left( 1 + e^{-w^T x_i} \right) \right]$$

• No closed form solution

- ### Gradient ascent

$$\nabla \log L(w) = \sum_{i=1}^{n} (1-y_i)(-x_i) - \frac{e^{-w^{\mathsf{T}}x_i}}{1+e^{-w^{\mathsf{T}}x_i}}(-x_i)$$

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

$$= \sum_{i=1}^{n} \overset{\underset{\mathbb{R}^d}{\downarrow}}{x_i} \left( \overset{\underset{\{0,1\}}{\downarrow}}{y_i} - \overset{g(w^{\mathsf{T}}x_i)}{\frac{1}{1+e^{-w^{\mathsf{T}}x_i}}} \right)$$

$$w_{t+1} = w_t + \eta_t \, \nabla \log L(w_t)$$

## REGULARIZED VERSION

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

---

## KERNEL VERSION

- Can argue $\quad w = \sum_{i=1}^{n} \alpha_i x_i \quad$ $\begin{bmatrix} \text{Formal Theorem} \\ \text{Representer Theorem} \end{bmatrix}$

Exercise: Derive the kernel version of logistic regression

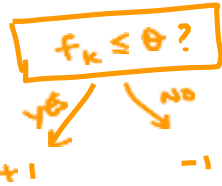# META CLASSIFIERS (or)

## ENSEMBLE CLASSIFIERS.

WEAK
CLASSIFIERS          $\longrightarrow$          STRONG
[better than                                    CLASSIFIERS
random]

# Weak classifiers

## Overfit decision tree

DECISION
STUMP

$f_k \leq \theta$ ?

YES        NO

+1        -1

**high bias, low variance**

…

……………

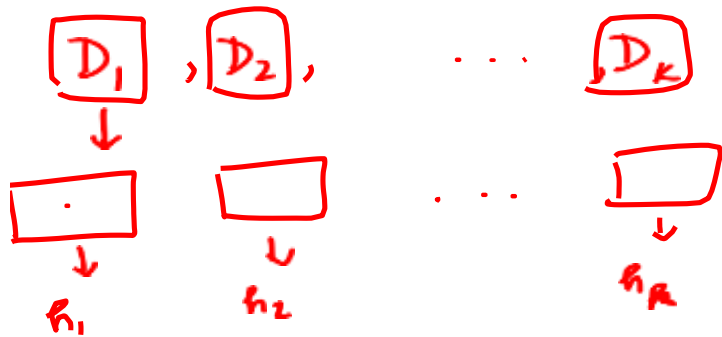**low bias, high variance**

$$X_1, X_2, \ldots X_n \sim N(M, 1)$$

$$\hat{M}_1 = X_1 \qquad \hat{M}_2 = X_2, \ldots \hat{M}_N = X_n \qquad \hat{M}_{ML} = \frac{1}{n}\Sigma x_i$$

$$D_1, D_2, \ldots D_k$$

Overfit decision trees

$$h_1 \qquad h_2 \qquad h_k$$

$$h_i : \mathbb{R}^d \to \{\pm 1\}$$

$$h^*(x) = \text{majority}(h_1(x), \ldots, h_k(x))$$

BAGGING — Bootstrap Aggregation.

$$D = \{ (x_1, y_1) \cdots (x_n, y_n) \}$$

Chance that a point appears in a dataset

$$1 - \left(1 - \frac{1}{n}\right)\left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{1}{n}\right)$$

$$1 - \left(1 - \frac{1}{n}\right)^n$$

$$1 - \frac{1}{e} \quad (\text{as } n \to \infty)$$

$$\simeq 66\%$$

— Create datasets $D_1, \cdots, D_k$ from $D$ by "Sampling with replacement".

— Run weak classifier on $D_1, \cdots, D_k$ to get $h_1, \cdots, h_k$

— Aggregate $h_1, \cdots, h_k$ using majority.

FEATURE BAGGING → Bag the features in addition to data points

Feature bagged decision trees -> RANDOM FOREST

BOOTSTRAP - Sampling with Replacement ?

?

AGGREGATION - Majority.

# BOOSTING

$\left[\begin{array}{l}\text{Freund \& Schappire.}\\ \qquad 1995 \\ \text{Godel Prize}\end{array}\right.$

↑
ADA-BOOST

---

Distribution $D$ over $(\underset{\mathbb{R}^d}{x} \times \underset{\{+1,-1\}}{y})$

$\longrightarrow$ unknown but fixed.

$x_1, \cdots, x_n$ are iid from $D$.

$h: \underset{x}{\mathbb{R}^d} \longrightarrow \underset{y}{\{\pm 1\}}$

Measure performance using

$$P_{(x,y) \sim D} \left( h(x) \neq y \right)$$

Misclassification probability.

A weak learner is one which outputs a Classifier
Strong
$h$ for which

$$P_{x,y \sim D} \left( h(x) = y \right) \geq \frac{1}{2} + \gamma$$

$1 - \epsilon$

$\gamma > 0$

for any unknown but fixed distribution $D$.