

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

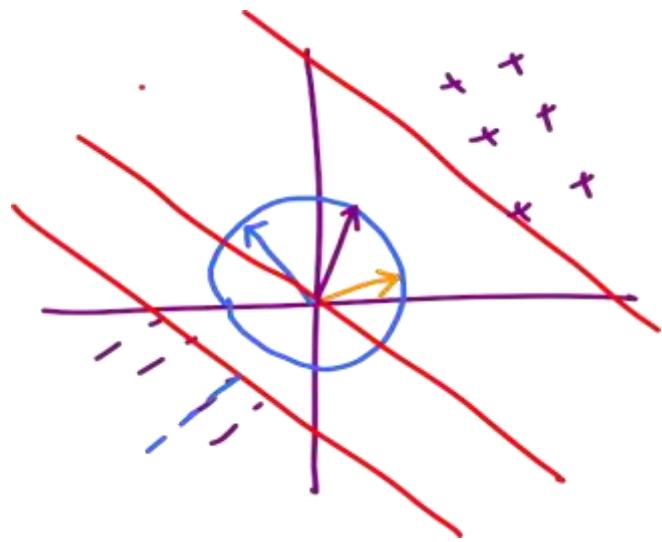
find a w with maximum width s.t all
points in dataset D are classified *wrong*

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

s.t.

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

Issue:
Can scale w
arbitrarily



Possible fix

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

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

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

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

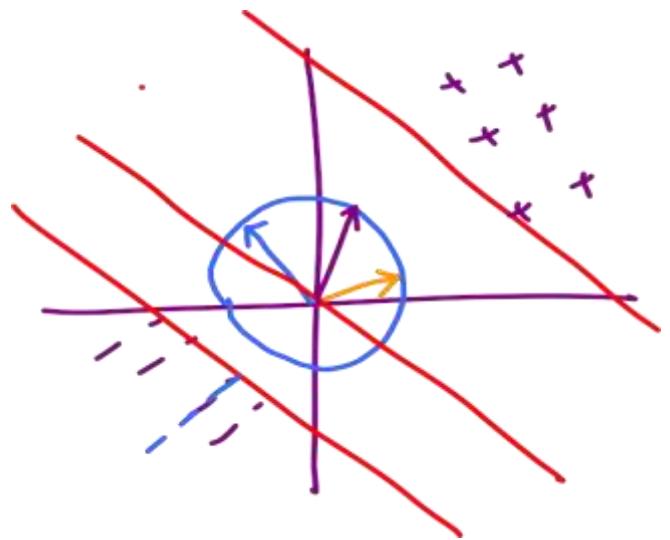
find a w with maximum width s.t all
points in dataset D are classified *wrong*

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

s.t.

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

Issue:
Can scale w
arbitrarily

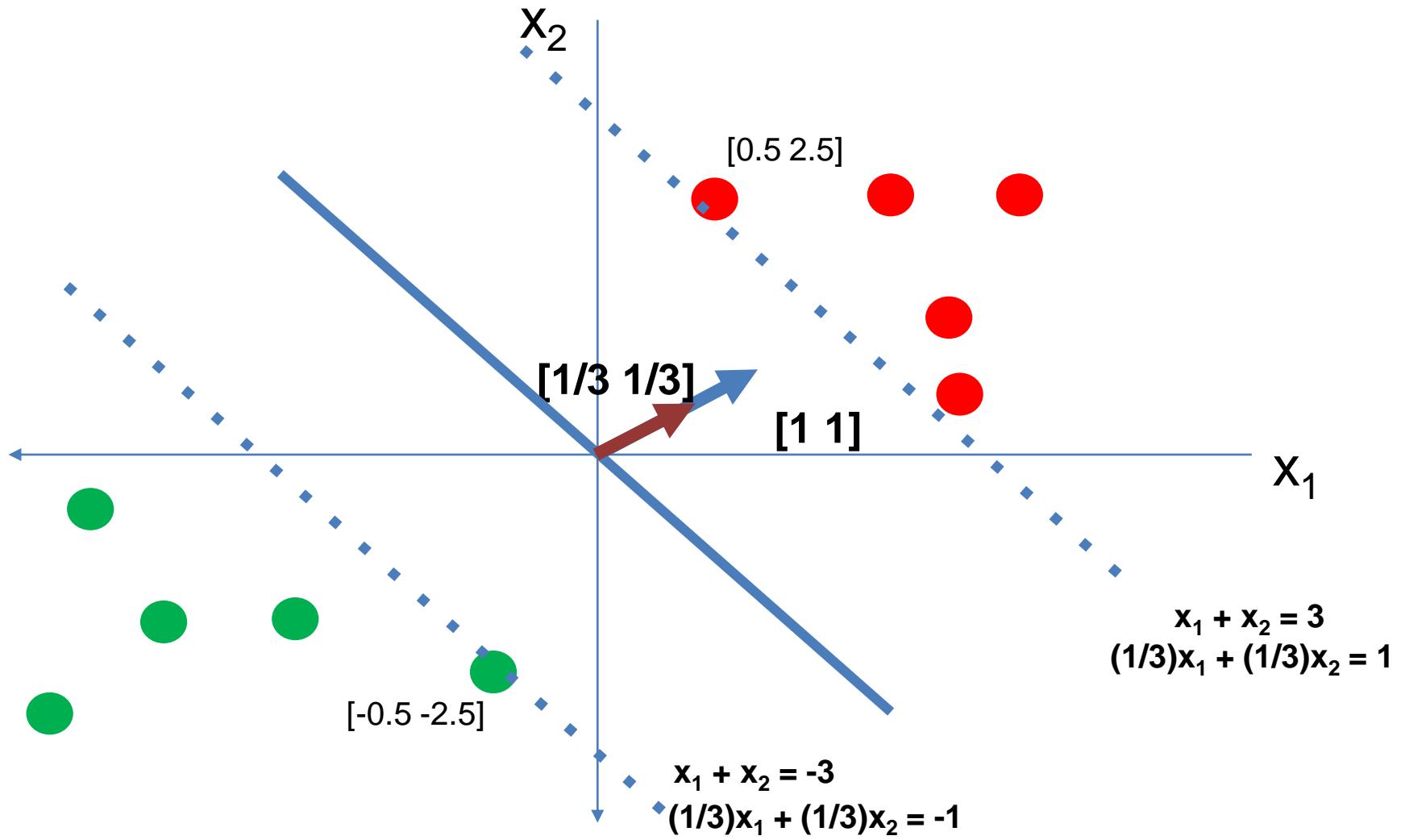


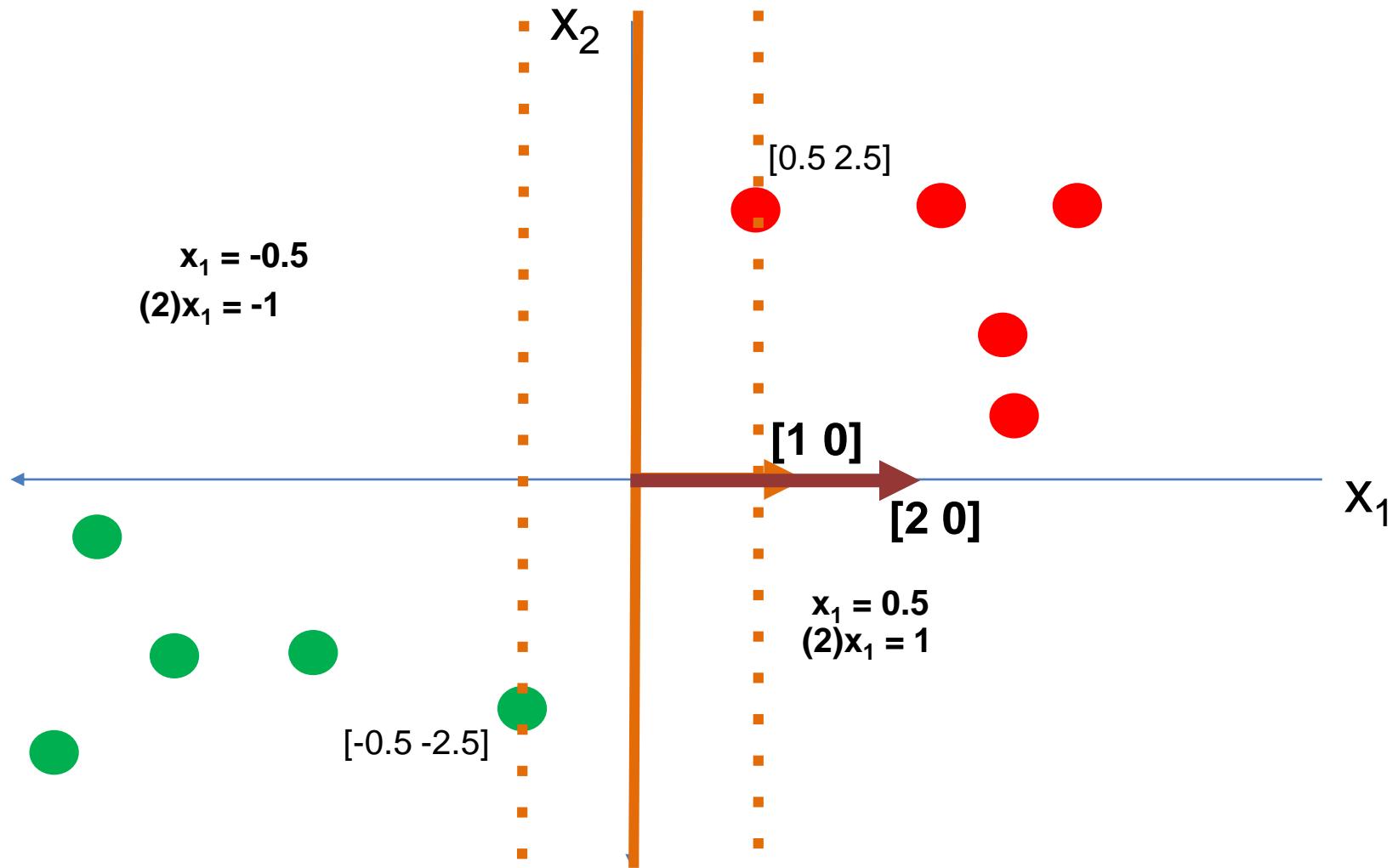
Possible fix

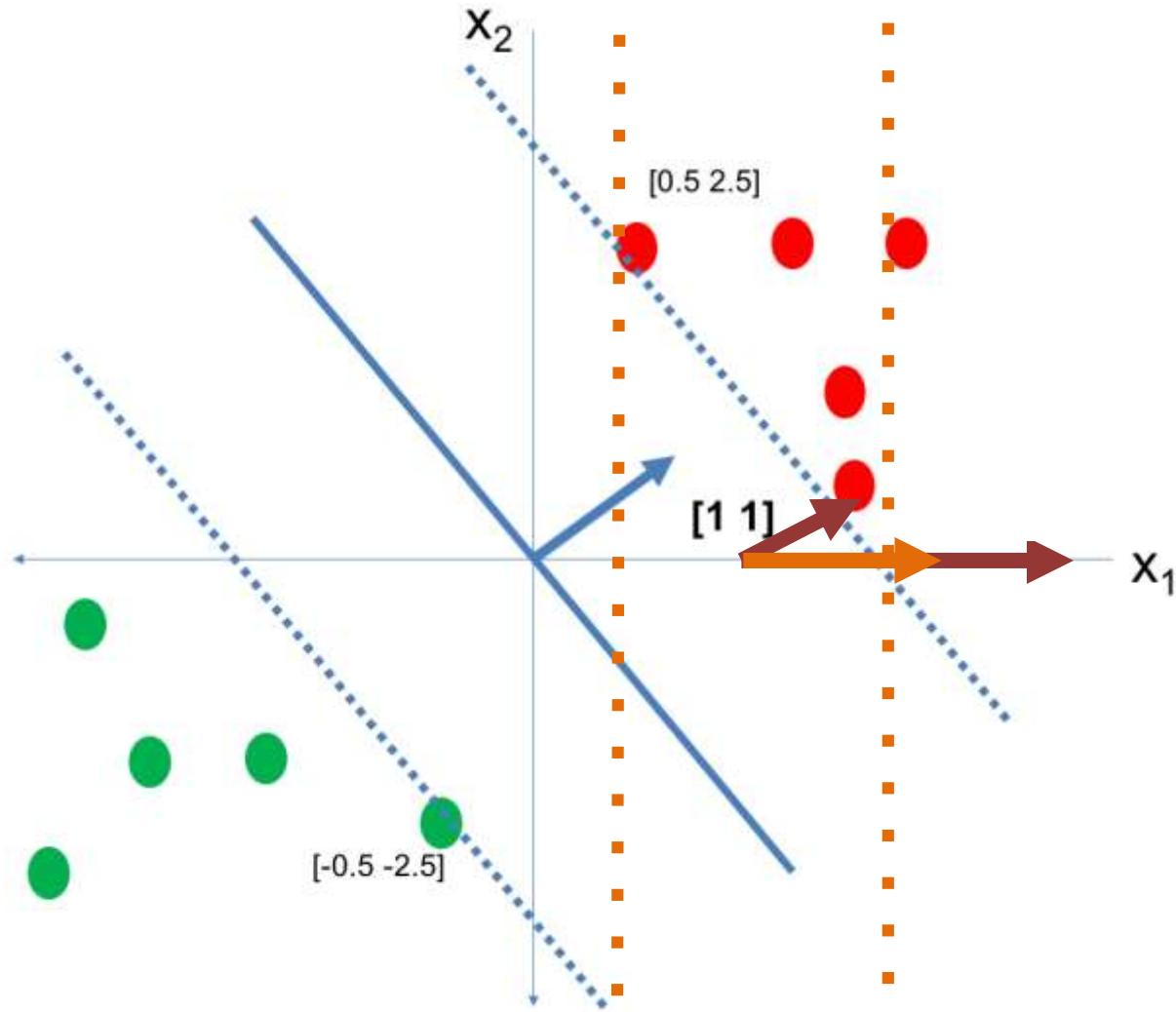
$$\max_{w, \gamma} \gamma$$

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

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







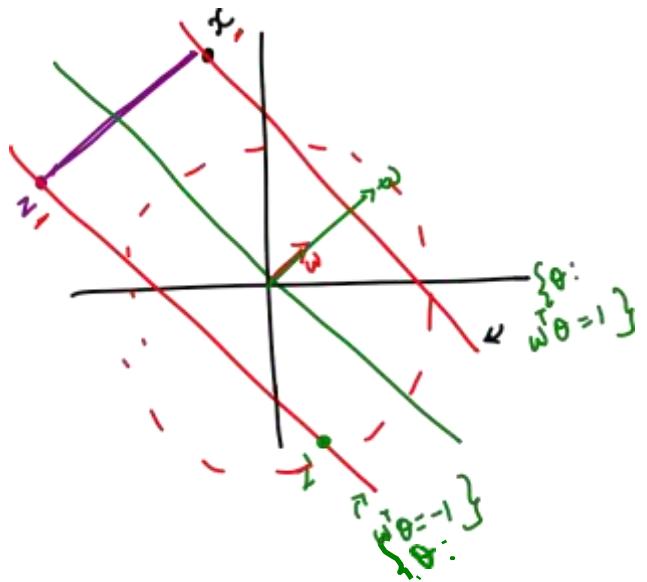
Notice what happens
to the lengths
of the w as we
adjust it to have
margin 1

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, $\text{width}(w)$ seems to be inversely proportional to $\|w\|$

$$\begin{aligned} & \max_w \\ & \text{width}(w) \\ \text{s.t. } & (\mathbf{w}^T \mathbf{x}_i) y_i \geq 1 \end{aligned}$$

What is $\text{width}(w)$?



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

[Exercise]

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

$$\max_w$$

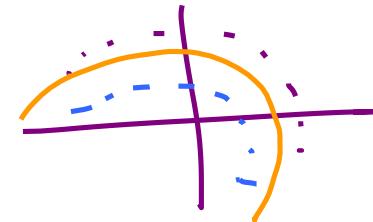
$$\frac{2}{\|\omega\|^2}$$

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

$$\min_w$$

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

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



Issues

- L.S is a strong assumption

- Non-linear structure?

DETOUR

$$\begin{array}{ll}\min_{\omega} & f(\omega) \\ \text{subject to} & g(\omega) \leq 0\end{array}$$



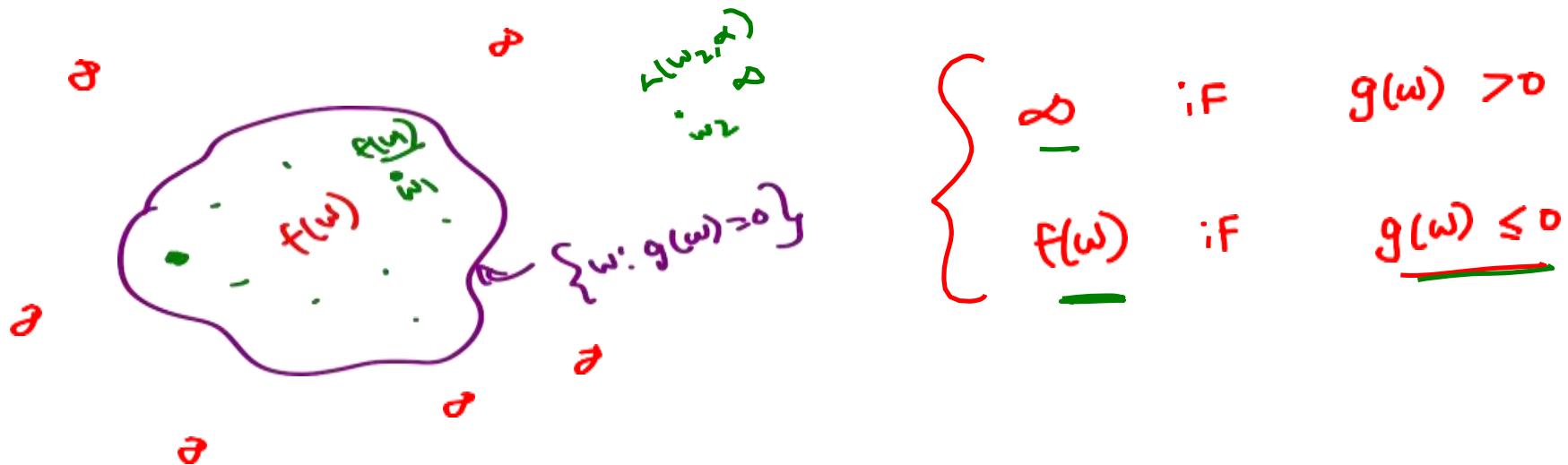
$$\underline{L}(\omega, \alpha) = f(\omega) + \alpha \cdot g(\omega)$$

Fix some ω .

Consider

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

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



$$\min_{\omega} \left[\max_{\alpha \geq 0} \frac{R(\omega)}{\lambda(\omega, \alpha)} \right] \stackrel{\text{equivalent}}{\equiv} \min_{\omega} f(\omega) \text{ s.t. } g(\omega) \leq 0.$$

- Can we swap min and max?

mi
n

Multiple Constraints

→ Same idea

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

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

=

$$\min_{\omega} \left[\max_{\substack{\{x_1, \dots \\ x_k \geq 0\}}} \left[f(\omega) + \underbrace{\sum_{i=1}^k \alpha_i g_i(\omega)}_{\text{---}} \right] \right]$$

III Strong duality for convex f, g_i

$$\max_{\substack{x_1, \dots, x_k \geq 0}} \min_{\omega} f(\omega) + \sum_{i=1}^k \alpha_i g_i(\omega)$$

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

s.t.

$$\underbrace{(\omega^T x_i) y_i}_{+i} \geq 1$$

$$1 - \underbrace{(\omega^T x_i) y_i}_{-i} \leq 0$$

$$g_i(\omega) = 1 - (\omega^T x_i) y_i$$

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

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

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

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

|||

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

Fix $\alpha \geq 0$

$$\min_w \left[\underbrace{\frac{1}{2} \|w\|^2}_{\text{Grad w.r.t } w} + \sum_{i=1}^n \alpha_i (1 - w^T x_i) y_i \right]$$

Grad w.r.t w

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

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

$\in \mathbb{R}^d$
 $\{\alpha_i\}$
 Fixed
Choice

In matrix notation

$$w^* = X Y \alpha$$

$$x = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}_{d \times n} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

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

$$\begin{aligned} & \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - w^T x_i) y_i \\ &= \frac{1}{2} \underline{w^T w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i (w^T x_i) y_i \end{aligned}$$

$$\underbrace{\frac{1}{2} (\mathbf{x}^T \boldsymbol{\alpha})^T (\mathbf{x}^T \boldsymbol{\alpha})}_{\text{matrix multiplication}} + \underbrace{\boldsymbol{\alpha}^T \mathbf{1}}_{\begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}} - \underbrace{\sum_{i=1}^n (\mathbf{x}^T \boldsymbol{\alpha})^T \mathbf{x}_i y_i \alpha_i}_{\text{sum}}$$

on Simplification [please do this]

$$\boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} (\mathbf{x}^T \boldsymbol{\alpha})^T (\mathbf{x}^T \boldsymbol{\alpha})$$

DUAL PROBLEM

Solving in $\boldsymbol{\alpha}$ instead of \mathbf{d}

$$\max_{\boldsymbol{\alpha} \geq 0} \quad \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i y_i \alpha_i$$

\mathbb{R} easy constraints

can be KERNELIZED!

Revisiting the Lagrangian

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

PRIMAL

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

DUAL

w^*

$$\boxed{\max_{\alpha \geq 0} f(w^*) + \alpha g(w^*)}$$

$$= \min_w f(w) + \alpha^* g(w)$$

α^*

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

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

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

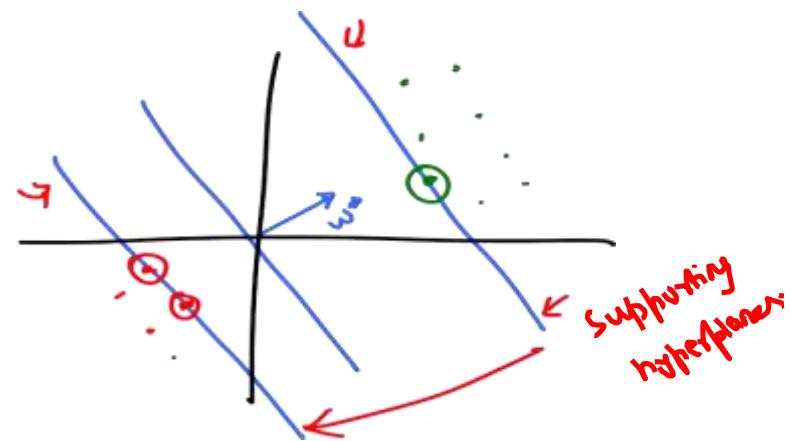
$$\Rightarrow \alpha^* g(\omega^*) = 0 \rightarrow \text{COMPLEMENTARY SLACKNESS}$$

For multiple constraints

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

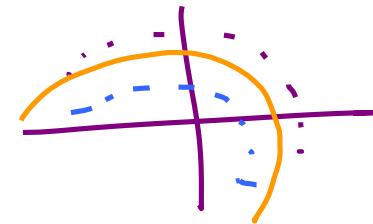
For our problem

$$(\alpha_i^*) (1 - (\omega^T x_i) y_i) = 0 \quad +i$$



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

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



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

$$\min_{\mathbf{w}, \boldsymbol{\varepsilon}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \boldsymbol{\varepsilon}_i$$

$C > 0$ [hyper parameter]

$$\rightarrow (\mathbf{w}^\top \mathbf{x}_i) y_i + \underline{\boldsymbol{\varepsilon}_i} \geq 1 \leftarrow +i$$

$$\rightarrow \underline{\boldsymbol{\varepsilon}_i} \geq 0 \leftarrow +i$$

if $C = 0 \Rightarrow$ Bribes don't cost $\Rightarrow \mathbf{w} = 0$ is solution

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

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

Dual :

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

$$\frac{\partial L}{\partial \omega} = 0 \Rightarrow \tilde{\omega} = \sum_{i=1}^n \alpha_i x_i y_i$$

$$\tilde{\omega}^* = X Y \alpha$$

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

$$\alpha_i + \beta_i = C \quad \forall i$$

Substitute $w = xy\alpha$ in the original objective

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

SOFT-MARGIN

SUPPORT
VECTOR
MACHING

$$\begin{aligned} \text{max } \\ \alpha > 0 \\ p > 0 \\ \underline{\alpha + p = C} \end{aligned}$$

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

=

$$\begin{aligned} \text{max } \\ 0 \leq \alpha \leq C \end{aligned}$$

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

BOX
CONSTRAINT

HARD-MARGIN SVM

PRIMAL

$$\min_w \frac{1}{2} \|w\|^2$$

st $(w^T x_i) y_i \geq 1 - \xi_i$

$\underbrace{(w^T x_i) y_i}_{1 - w^T x_i y_i \leq 0} \leq \xi_i$

DUAL

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

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

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

SOFT-MARGIN SVM

PRIMAL ✓

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

st $(w^T x_i) y_i + \xi_i \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

α, β

$\underbrace{(w^T x_i) y_i}_{\xi_i \geq 0} + \xi_i \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

DUAL ✓

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

$$\alpha + \beta = C$$

$$\alpha \geq 0$$

$$\beta \geq 0$$

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

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

COMPLEMENTARY SLACKNESS

$$\underline{\alpha}_i^* \left(1 - (\underline{w}^{*T} \underline{x}_i) \underline{y}_i - \underline{\xi}_i^* \right) = 0 \quad \forall i$$

$$\underline{\beta}_i^* \underline{\xi}_i^* = 0 \quad \forall i$$

$\alpha_i^* + \beta_i^* = c$
 $\forall i$
 $\hookrightarrow A$

Various cases possible

Case 1:

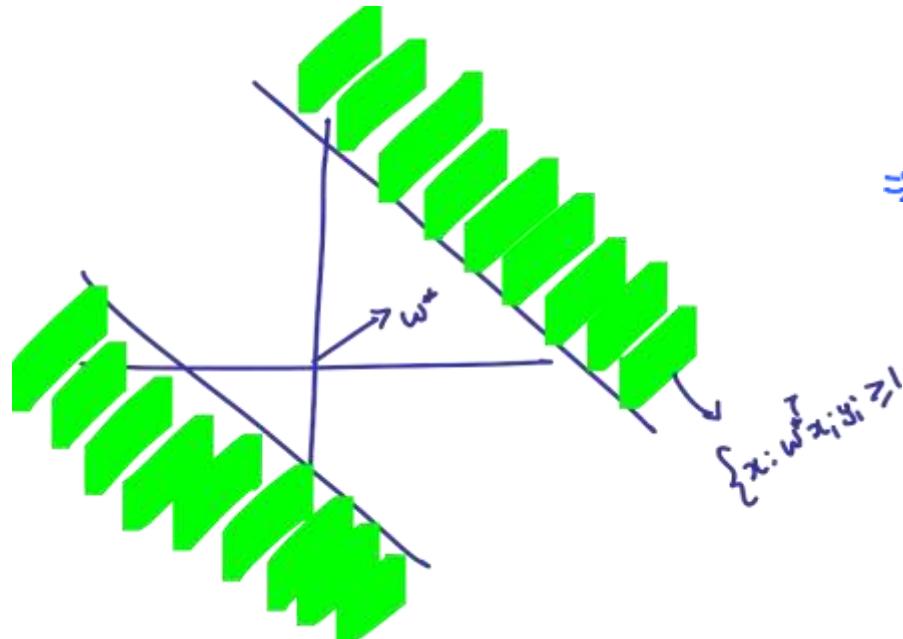
$$\alpha_i^* = 0$$

④

$$\beta_i^* = C$$

cs

$$\sum_i \underline{\xi}_i^* = 0$$



$$1 - (\mathbf{w}^\top \mathbf{x}_i) y_i - \underline{\xi}_i^* \leq 0 \quad [\text{Primal feasibility}]$$

$$\Rightarrow 1 - (\mathbf{w}^\top \mathbf{x}_i) y_i \leq 0$$

$$\Rightarrow \mathbf{w}^\top \mathbf{x}_i y_i \geq 1$$

$\Rightarrow \mathbf{w}^*$ classifies (\mathbf{x}_i, y_i) correctly.

Case 2: $0 < \alpha_i^* < C \Rightarrow 0 < \beta_i^* < C \Rightarrow \underline{\xi_i^*} = 0$

\Downarrow Case 2

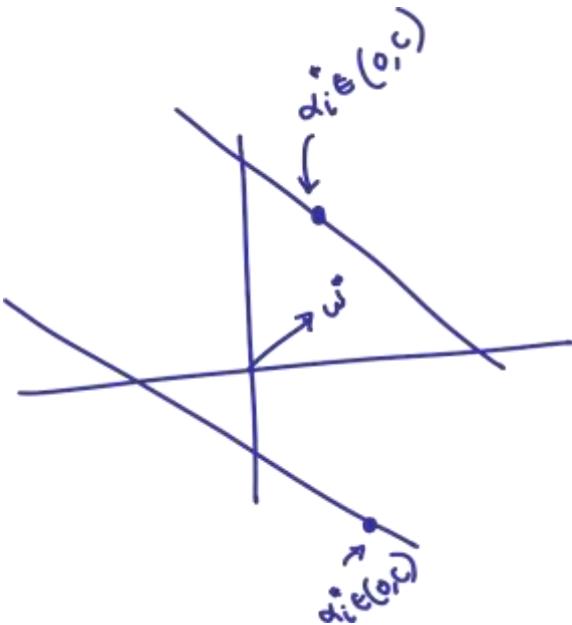
$$1 - (\omega^T x_i) y_i - \xi_i^* = 0$$

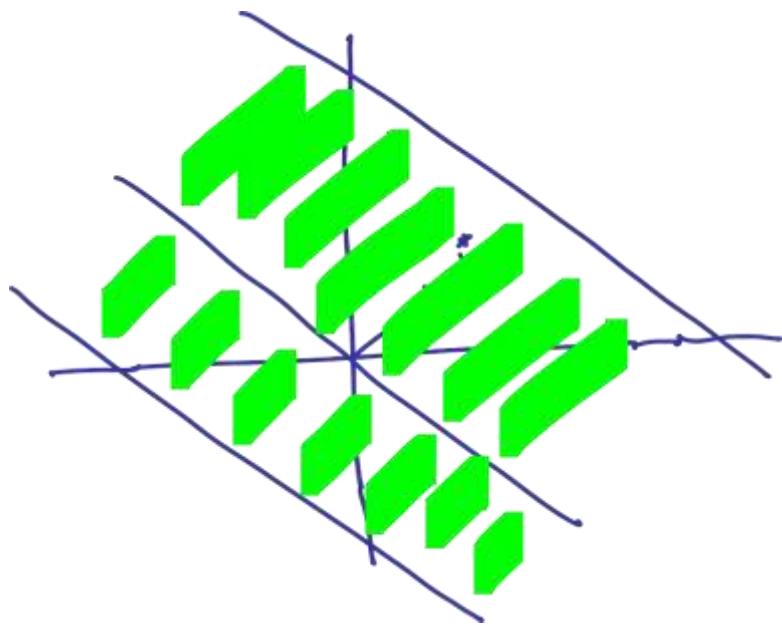
\Downarrow

$$(\omega^T x_i) y_i = 1$$

\Rightarrow

(x_i, y_i) lies on the
Supporting hyperplane.





Case 3:

$$\frac{\alpha_i^* = C}{\Downarrow \boxed{CS}} \Rightarrow \beta_i^* = 0 \Rightarrow \xi_i^* \geq 0$$

$$1 - \omega^T x_i y_i - \xi_i^* = 0$$

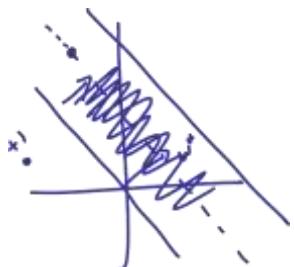
$$\xi_i^* = 1 - \omega^T x_i y_i \geq 0$$

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

Let's see this from P.O.V of data

CASE I

$$\boxed{\omega^T x_i y_i \leq 1}$$



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

$$\omega^T x_i y_i \geq 1 - \xi_i^*$$

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

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

$$\alpha_i^* \left(\frac{1 - \omega^T x_i y_i - \xi_i^*}{\beta_i^* \xi_i^*} \right) = 0$$

CASE 2:

$$\omega^T x_i y_i = 1$$

$$\xi_i^* \geq 1 - \frac{\omega^T x_i y_i}{\alpha_i}$$

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

CASE 3

$$\omega^T x_i y_i > 1$$

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

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

SUMMARY

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

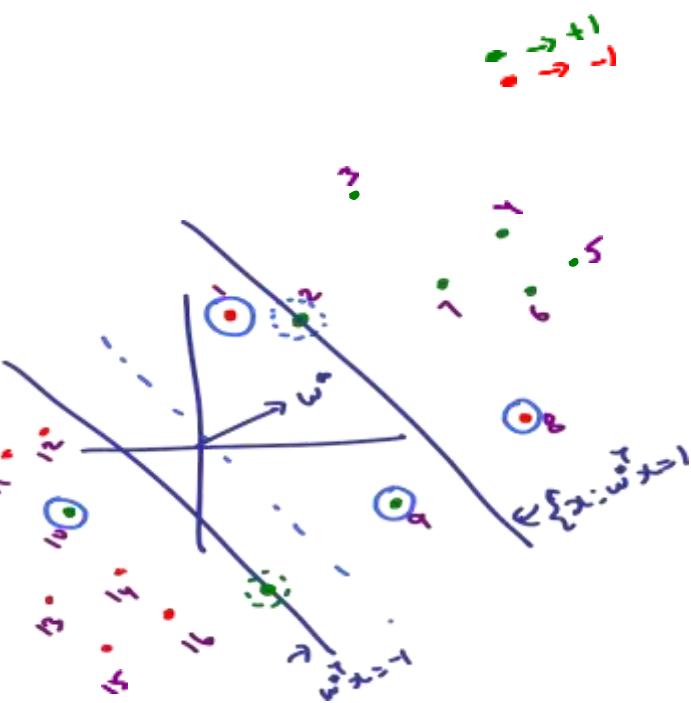
$$0 < \alpha_i^* < C \Rightarrow w^T x_i y_i = 1$$

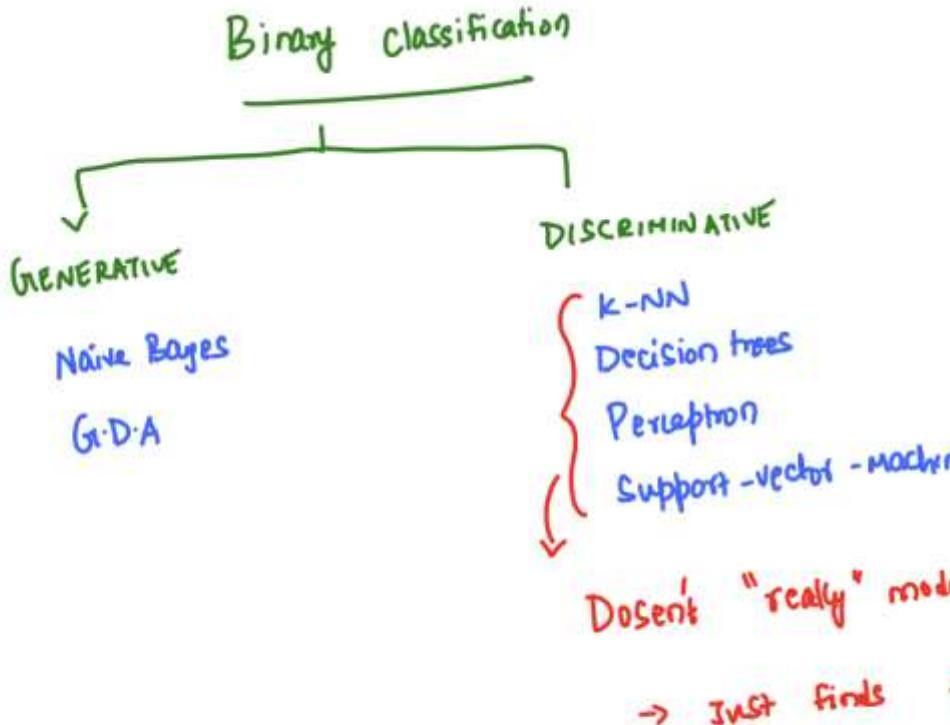
$$\alpha_i^* = C \Rightarrow w^T x_i y_i \leq 1$$

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

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

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





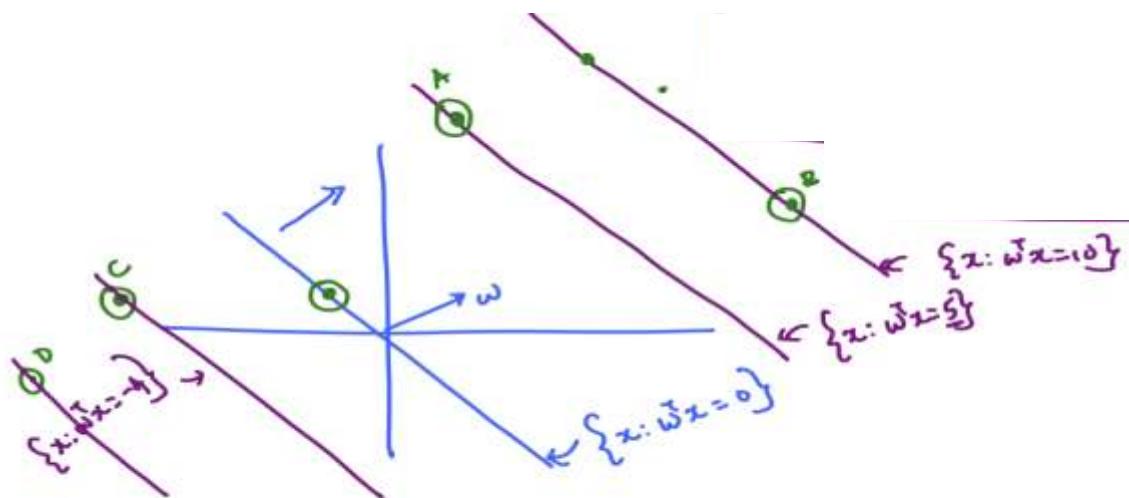
- Can we model $P(y=+1/x)$ differently?

| Start with a simple model

Given $x \in \mathbb{R}^d$

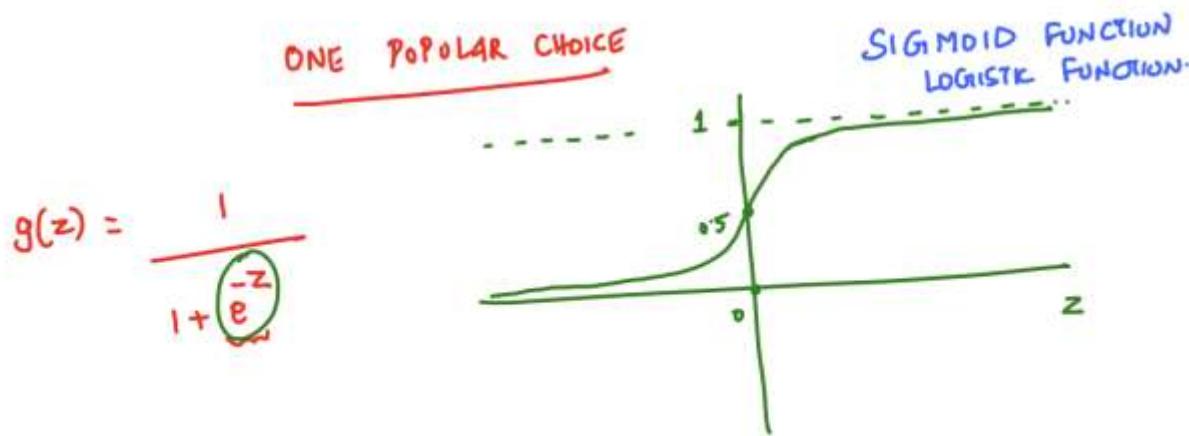
$$z = w^T x$$

$w \in \mathbb{R}^d$



$$P(y=+1|x) = g(w^T x)$$

- LINK FUNCTION
- $\underline{g(z)} \in [0, 1]$
 - $g(z) \rightarrow 1 \text{ as } z \rightarrow \infty$
 - $g(z) \rightarrow 0 \text{ as } z \rightarrow -\infty$
 - $g(z) = 0.5 \text{ if } z=0$



MODEL: LOGISTIC REGRESSION

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

Max. Likelihood

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

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

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

v

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

...

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

- No closed form solution

- Gradient ascent

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

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

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

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

REGULARIZED VERSION

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

KERNEL VERSION

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

[Formal Theorem
Representer Theorem]

Exercise: Derive the kernel version of logistic regression

META CLASSIFIERS (or)

ENSEMBLE CLASSIFIERS.

WEAK
CLASSIFIERS

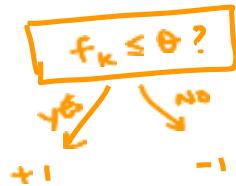
[better than
random]



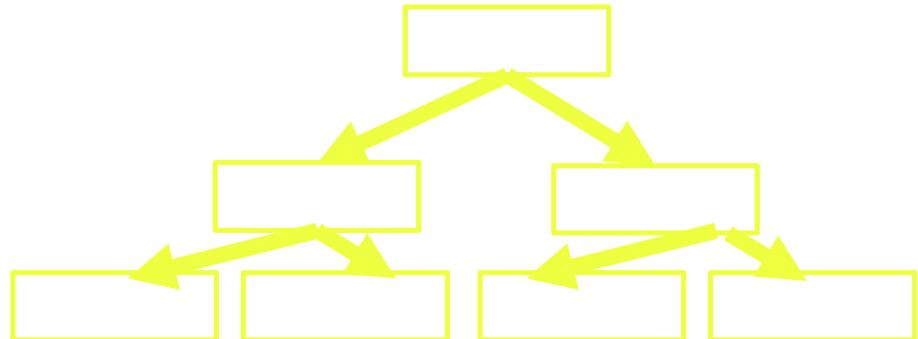
STRONG
CLASSIFIERS

Weak classifiers

DECISION STUMP



Overfit decision tree



high bias, low variance

...



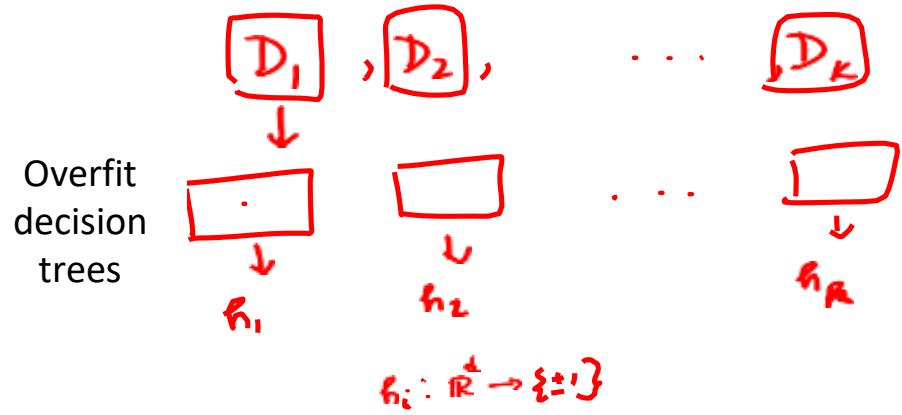
.....



low bias, high variance

$$x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$$

$$\hat{\mu}_1 = x_1 \quad \hat{\mu}_2 = x_2, \dots, \hat{\mu}_n = x_n \quad \hat{\mu}_{ML} = \frac{1}{n} \sum x_i$$



$$h^*(x) = \text{majority}(h_1(x), \dots, h_K(x))$$

BAGGING - Bootstrap Aggregation

Chance that a point appears in a dataset

$$1 - \underbrace{\left(1 - \frac{1}{n}\right)}_{\text{as } n \rightarrow \infty} \underbrace{\left(1 - \frac{1}{n}\right)}_{\text{as } n \rightarrow \infty} \cdots \underbrace{\left(1 - \frac{1}{n}\right)}_{\text{as } n \rightarrow \infty}$$

$$1 - \underbrace{\left(1 - \frac{1}{n}\right)}_{\text{as } n \rightarrow \infty}^n$$

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

$$\approx 63.2\%$$

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

- Create datasets D_1, \dots, D_k from D by
Sampling with replacement.

- Run weak classifier on D_1, \dots, D_k to get f_1, \dots, f_k

- Aggregate f_1, \dots, f_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

ADA-BOOST

[Freund & Schapire ·
1995
Gödel Prize]

Distribution D over $(\mathbb{R}^d \times \{+1, -1\})$
Unknown but fixed.

x_1, \dots, x_n are iid from D .

$$f_i: \mathbb{R}^d \rightarrow \{+1, -1\}$$

Measure performance using

$$P_{\substack{(x,y) \sim D}}(r(x) \neq y)$$

Misclassification probability.

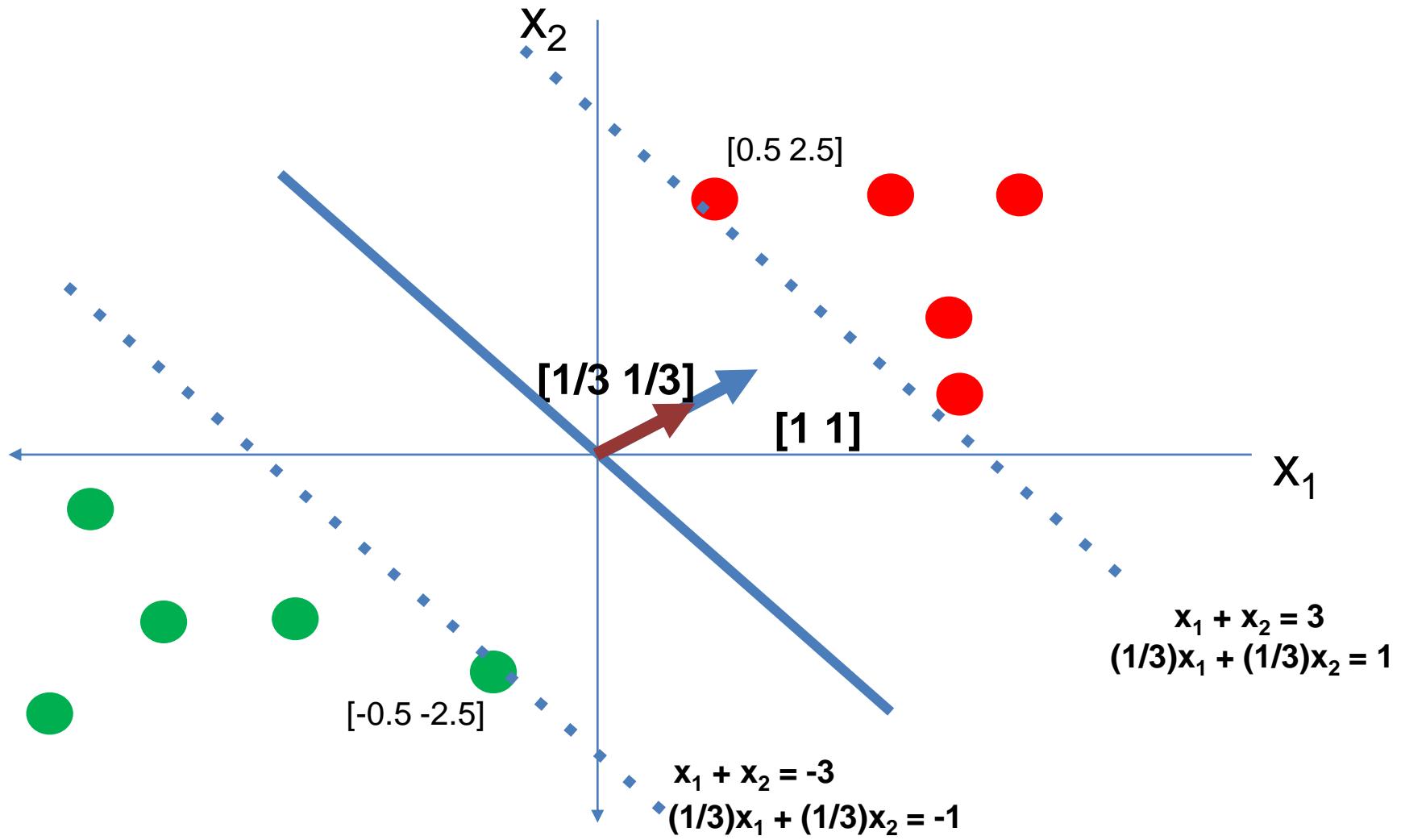
A weak learner is one which outputs a classifier
Strong

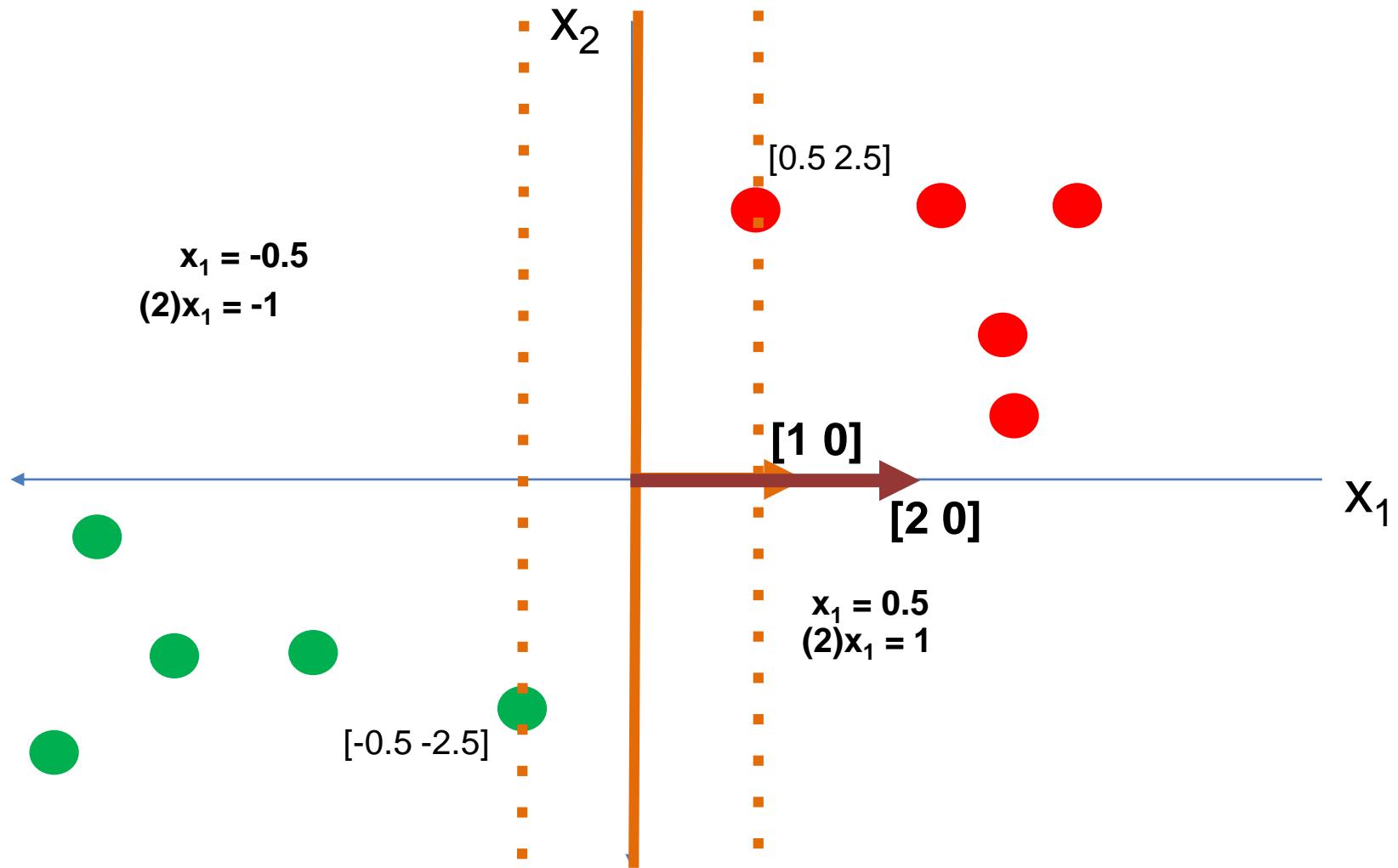
for which

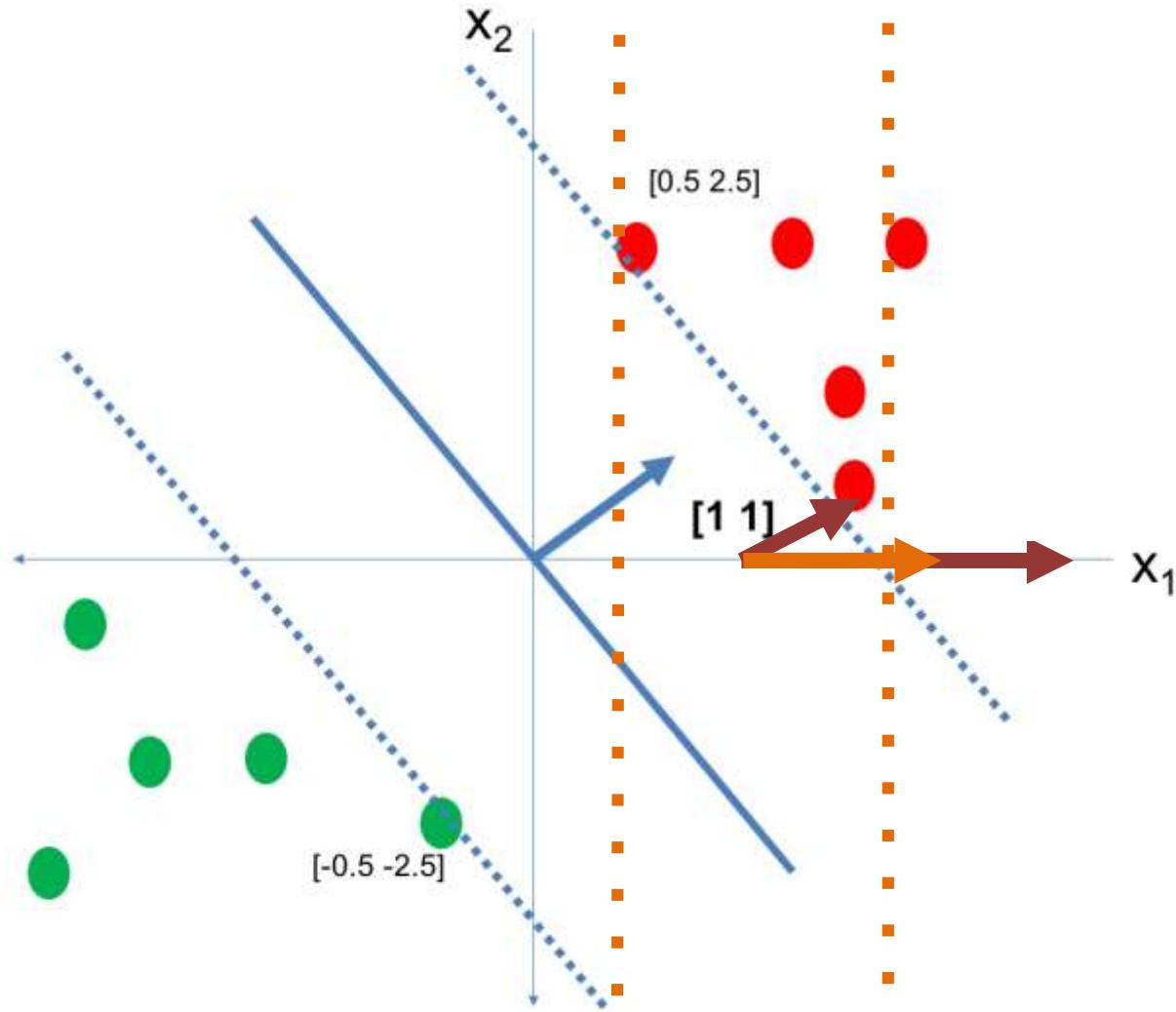
$$P_{\substack{x,y \sim D}}(r(x) = y) \geq \frac{1-\epsilon}{2} + \gamma$$

$$\gamma > 0$$

for any unknown but fixed distribution D .







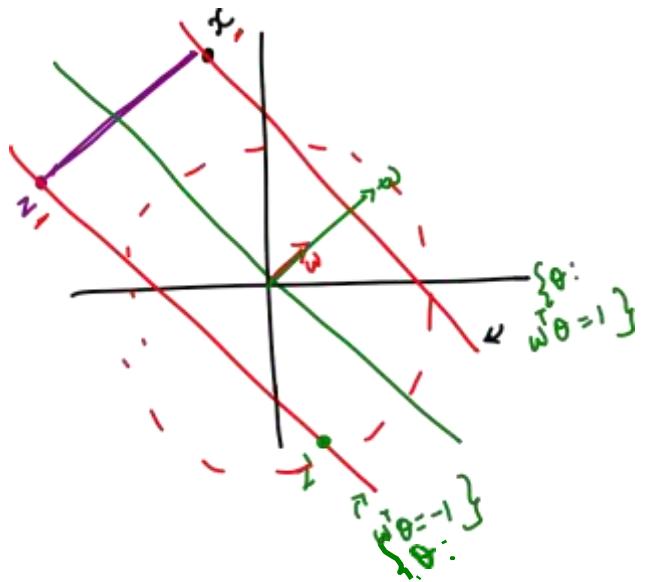
Notice what happens
to the lengths
of the w as we
adjust it to have
margin 1

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, $\text{width}(w)$ seems to be inversely proportional to $\|w\|$

$$\begin{aligned} & \max_w \\ & \text{width}(w) \\ \text{s.t. } & (\mathbf{w}^T \mathbf{x}_i) y_i \geq 1 \end{aligned}$$

What is $\text{width}(w)$?



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

[Exercise]

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

$$\max_w$$

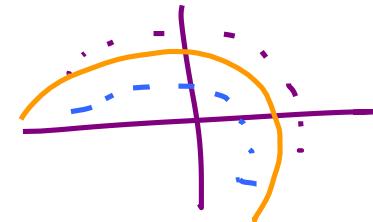
$$\frac{2}{\|\omega\|^2}$$

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

$$\min_w$$

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

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



Issues

- L.S is a strong assumption

- Non-linear structure?

DETOUR

$$\begin{array}{ll}\min_{\omega} & f(\omega) \\ \text{subject to} & g(\omega) \leq 0\end{array}$$



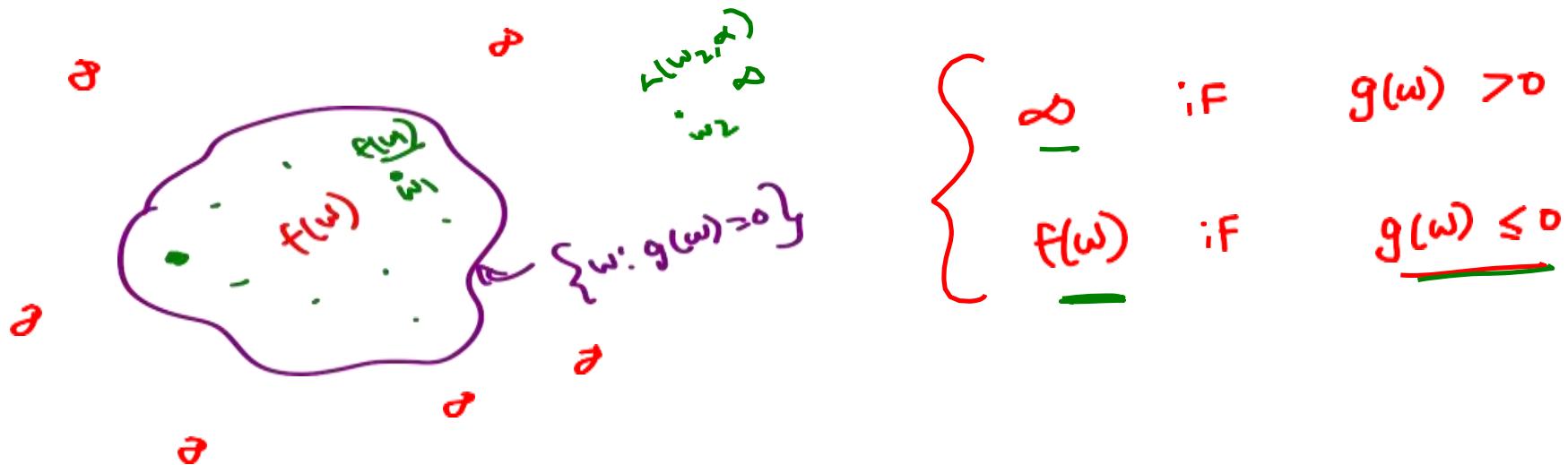
$$\underline{L}(\omega, \alpha) = f(\omega) + \alpha \cdot g(\omega)$$

Fix some ω .

Consider

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

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



$$\min_{\omega} \left[\max_{\alpha \geq 0} \frac{R(\omega)}{\lambda(\omega, \alpha)} \right] \stackrel{\text{equivalent}}{\equiv} \min_{\omega} f(\omega) \text{ s.t. } g(\omega) \leq 0.$$

- Can we swap min and max?

mi
n

Multiple Constraints

→ Same idea

$$\begin{aligned} \min_{\omega} \quad & f(\omega) \\ \text{s.t.} \quad & g_i(\omega) \leq 0 \quad i=1 \dots k \end{aligned}$$

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

III Strong duality for convex f, g_i

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

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

s.t.

$$\underbrace{(\omega^T x_i) y_i}_{+i} \geq 1$$

$$1 - \underbrace{(\omega^T x_i) y_i}_{-i} \leq 0$$

$$g_i(\omega) = 1 - (\omega^T x_i) y_i$$

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

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

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

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

|||

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

Fix $\alpha \geq 0$

$$\min_w \left[\underbrace{\frac{1}{2} \|w\|^2}_{\text{Grad w.r.t } w} + \sum_{i=1}^n \alpha_i (1 - w^T x_i) y_i \right]$$

Grad w.r.t w

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

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

$\in \mathbb{R}^d$

$\{\pm 1\}$

Fixed
choice

In matrix notation

$$w^* = X Y \alpha$$

$$x = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}_{d \times n} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

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

$$\begin{aligned} & \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - w^T x_i) y_i \\ &= \frac{1}{2} \underline{w^T w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i (w^T x_i) y_i \end{aligned}$$

$$\underbrace{\frac{1}{2} (\mathbf{x}^T \boldsymbol{\alpha})^T (\mathbf{x}^T \boldsymbol{\alpha})}_{\text{matrix multiplication}} + \underbrace{\boldsymbol{\alpha}^T \mathbf{1}}_{\begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}} - \underbrace{\sum_{i=1}^n (\mathbf{x}^T \boldsymbol{\alpha})^T \mathbf{x}_i y_i \alpha_i}_{\text{sum}}$$

On Simplification [please do this]

$$\boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} (\mathbf{x}^T \boldsymbol{\alpha})^T (\mathbf{x}^T \boldsymbol{\alpha})$$

DUAL PROBLEM

Solving in $\boldsymbol{\alpha}$ instead of \mathbf{x}

$$\max_{\boldsymbol{\alpha} \geq 0} \quad \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{y}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\alpha}$$

\mathbf{X} easy constraints

can be KERNELIZED!

Revisiting the Lagrangian

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

PRIMAL

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

DUAL

w^*

$$\boxed{\max_{\alpha \geq 0} f(w^*) + \alpha g(w^*)}$$

$$= \min_w f(w) + \alpha^* g(w)$$

α^*

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

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

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

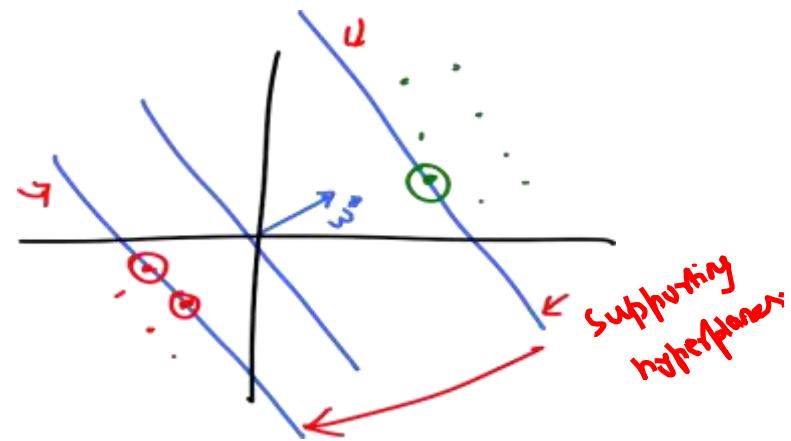
$$\Rightarrow \alpha^* g(\omega^*) = 0 \rightarrow \text{COMPLEMENTARY SLACKNESS}$$

For multiple constraints

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

For our problem

$$(\alpha_i^*) (1 - (\omega^T x_i) y_i) = 0 \quad +i$$

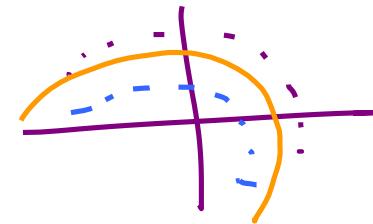


Supporting
hyperplanes

Margin

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

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



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

$$\min_{\mathbf{w}, \boldsymbol{\varepsilon}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \boldsymbol{\varepsilon}_i$$

$C > 0$ [hyper parameter]

$$\rightarrow (\mathbf{w}^\top \mathbf{x}_i) y_i + \underline{\boldsymbol{\varepsilon}_i} \geq 1 \leftarrow +i$$

$$\rightarrow \underline{\boldsymbol{\varepsilon}_i} \geq 0 \leftarrow +i$$

if $C = 0 \Rightarrow$ Bribes don't cost $\Rightarrow \mathbf{w} = 0$ is solution

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

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

Dual :

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

$$\frac{\partial L}{\partial \omega} = 0 \Rightarrow \tilde{\omega} = \sum_{i=1}^n \alpha_i x_i y_i$$

$\tilde{\omega}^* = X Y \alpha$

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

$$\alpha_i + \beta_i = C \quad \forall i$$

Substitute $w = xy\alpha$ in the original objective

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

SOFT-MARGIN

SUPPORT
VECTOR
MACHING

$$\begin{aligned} \text{max } & \\ \alpha > 0 & \\ \beta > 0 & \\ \underline{\alpha + \beta = C} & \end{aligned}$$

$$\alpha^T 1 - \frac{1}{2} (\mathbf{x}^T \mathbf{y} \alpha)^T (\mathbf{x}^T \mathbf{y} \alpha)$$

=

$$\begin{aligned} \text{max } & \\ 0 \leq \alpha \leq C & \end{aligned}$$

$$\alpha^T 1 - \frac{1}{2} \alpha^T \mathbf{y}^T (\mathbf{x}^T \mathbf{x}) \mathbf{y} \alpha$$

BOX
CONSTRAINT

HARD-MARGIN SVM

PRIMAL

$$\min_w \frac{1}{2} \|w\|^2$$

st $(w^T x_i) y_i \geq 1 - \xi_i$

$\underbrace{(w^T x_i) y_i}_{1 - w^T x_i y_i \leq 0} \leq \xi_i$

DUAL

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

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

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

SOFT-MARGIN SVM

PRIMAL ✓

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

st $(w^T x_i) y_i + \xi_i \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

α, β

$\underbrace{(w^T x_i) y_i}_{\xi_i \geq 0} + \xi_i \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

DUAL ✓

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

$$\alpha + \beta = C$$

$$\alpha \geq 0$$

$$\beta \geq 0$$

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

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

COMPLEMENTARY SLACKNESS

$$\underline{\alpha}_i^* \left(1 - (\underline{w}^{*T} \underline{x}_i) \underline{y}_i - \underline{\xi}_i^* \right) = 0 \quad \forall i$$

$$\underline{\beta}_i^* \underline{\xi}_i^* = 0 \quad \forall i$$

$\alpha_i^* + \beta_i^* = c$
 $\forall i$
 $\hookrightarrow A$

Various cases possible

Case 1:

$$\alpha_i^* = 0$$

$$\Leftrightarrow$$

$$\beta_i^* = C$$

$$\Rightarrow$$

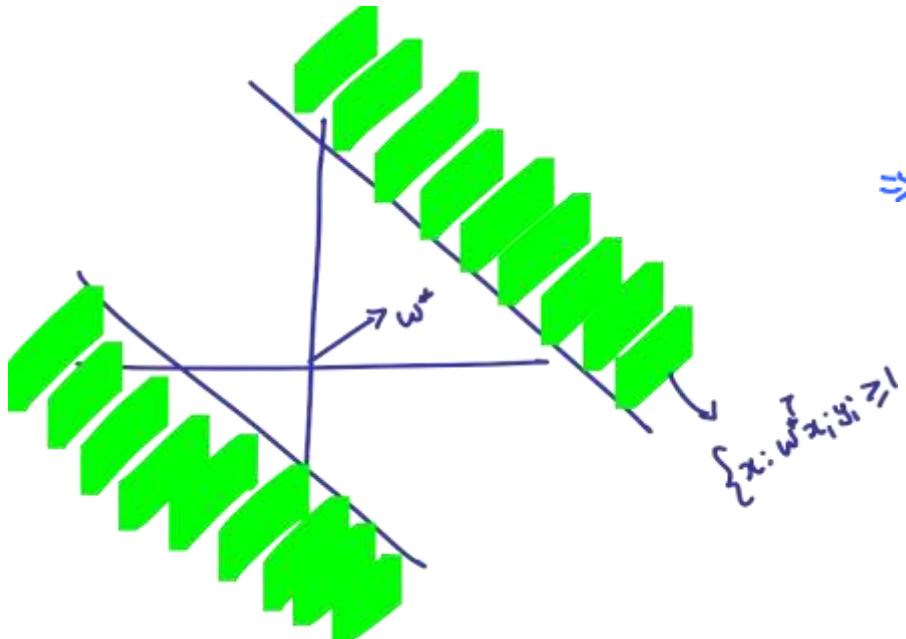
$$\sum_i \xi_i^* = 0$$

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

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

$$\Rightarrow \omega^T x_i y_i \geq 1$$

$\Rightarrow \omega^*$ classifies (x_i, y_i) correctly.



Case 2: $0 < \alpha_i^* < C \Rightarrow 0 < \beta_i^* < C \Rightarrow \underline{\xi_i^*} = 0$

\Downarrow Case 2

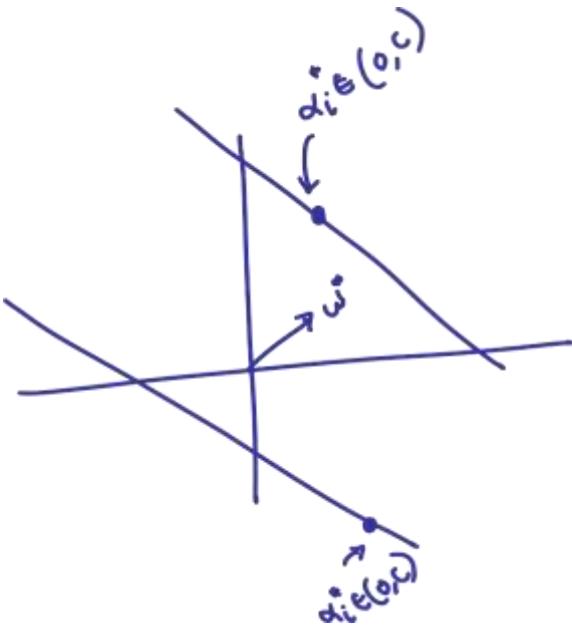
$$1 - (\omega^T x_i) y_i - \xi_i^* = 0$$

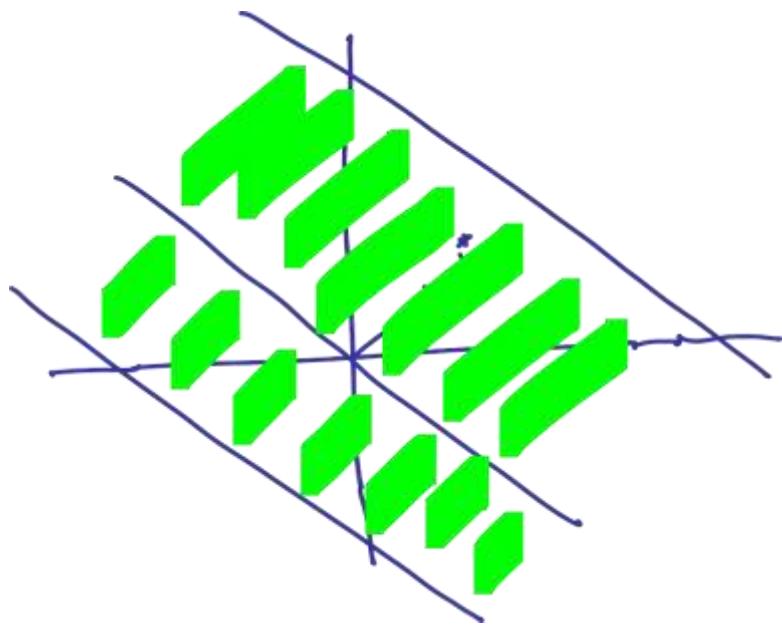
\Downarrow

$$(\omega^T x_i) y_i = 1$$

\Rightarrow

(x_i, y_i) lies on the
Supporting hyperplane.





Case 3:

$$\frac{\alpha_i^* = C}{\Downarrow \boxed{CS}} \Rightarrow \beta_i^* = 0 \Rightarrow \xi_i^* \geq 0$$

$$1 - \omega^T x_i y_i - \xi_i^* = 0$$

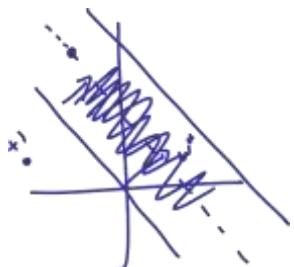
$$\xi_i^* = 1 - \omega^T x_i y_i \geq 0$$

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

Let's see this from P.O.V of data

CASE I

$$\boxed{\omega^T x_i y_i \leq 1}$$



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

$$\omega^T x_i y_i \geq 1 - \xi_i^*$$

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

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

$$\alpha_i^* \left(\frac{1 - \omega^T x_i y_i - \xi_i^*}{\beta_i^* \xi_i^*} \right) = 0$$

CASE 2:

$$\omega^T x_i y_i = 1$$

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

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

CASE 3

$$\omega^T x_i y_i > 1$$

$$1 - \underbrace{\omega^T x_i y_i}_{\text{C.S.}} - \xi_i^* \leq 0 \quad [\text{Primal feasibility}]$$

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

SUMMARY

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

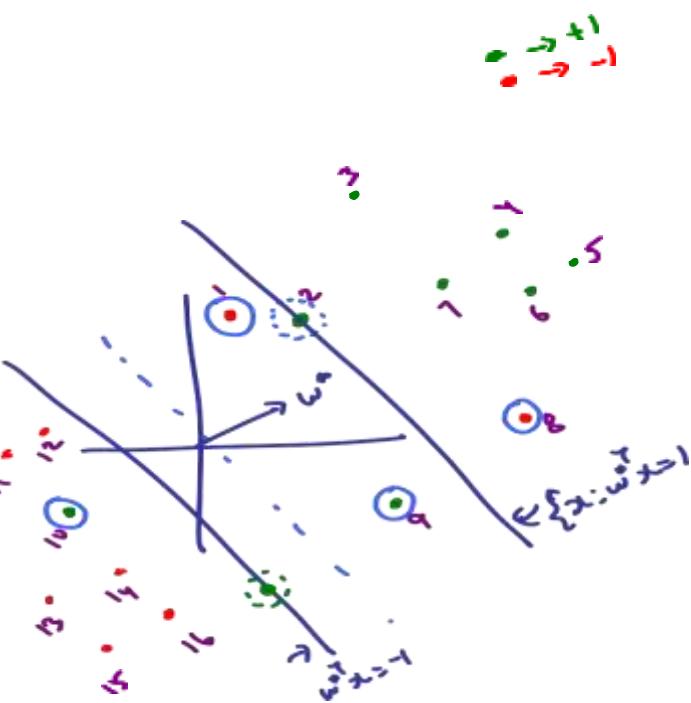
$$0 < \alpha_i^* < C \Rightarrow w^T x_i y_i = 1$$

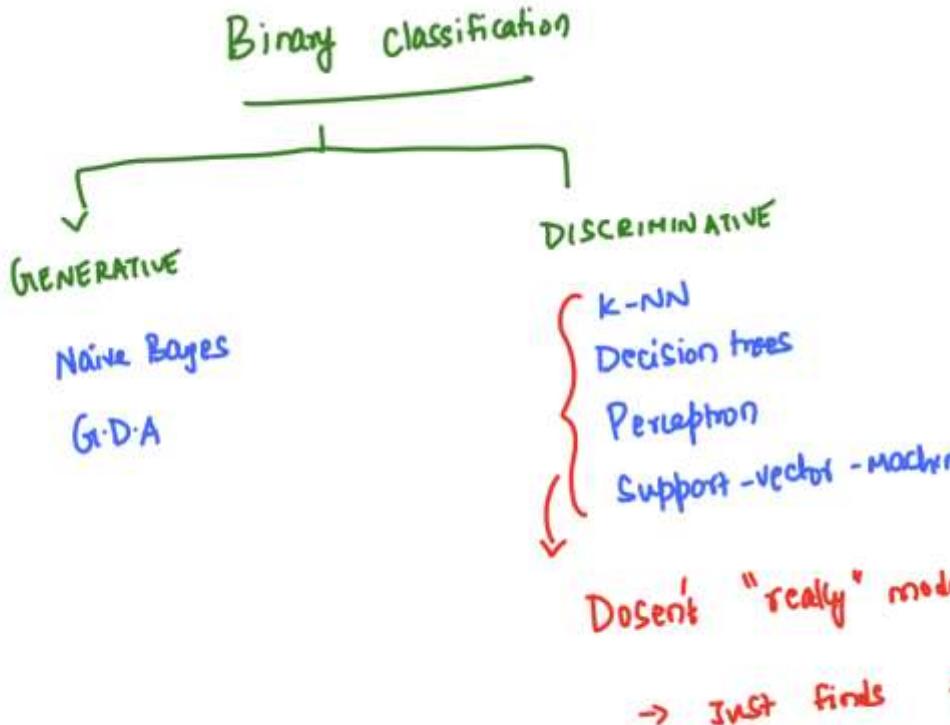
$$\alpha_i^* = C \Rightarrow w^T x_i y_i \leq 1$$

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

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

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





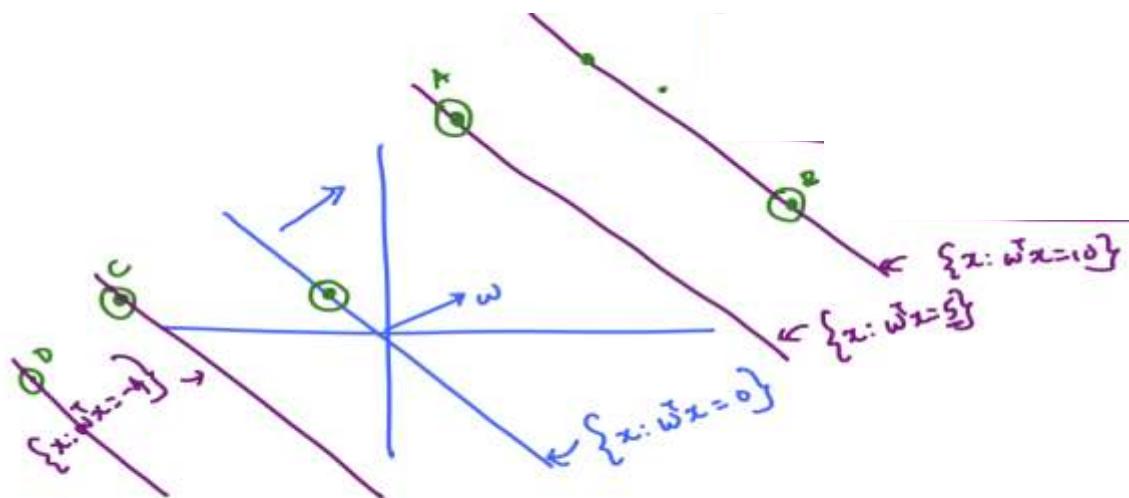
- Can we model $P(y = +1/x)$ differently?

| Start with a simple model

Given $x \in \mathbb{R}^d$

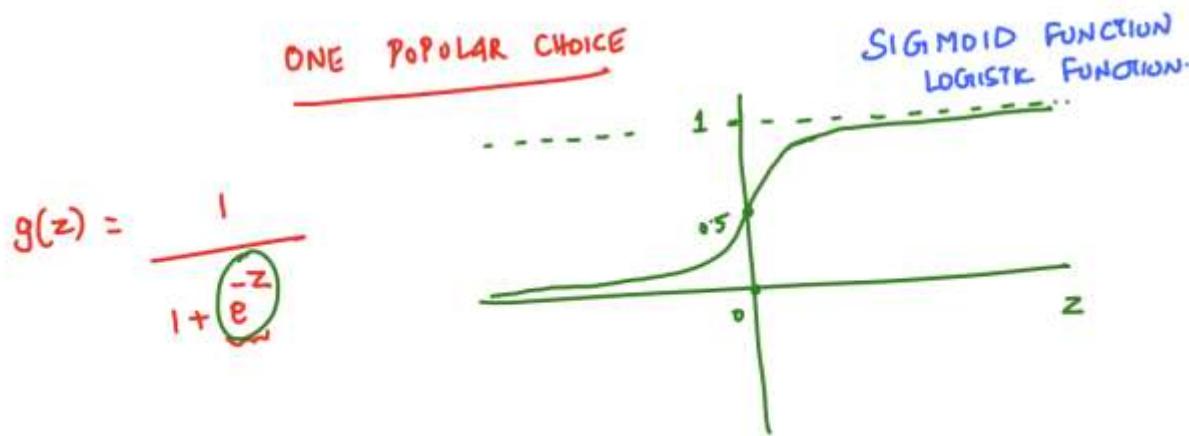
$$z = w^T x$$

$w \in \mathbb{R}^d$



$$P(y=+1|x) = g(w^T x)$$

- LINK FUNCTION
- $\underline{g(z)} \in [0, 1]$
 - $g(z) \rightarrow 1 \text{ as } z \rightarrow \infty$
 - $g(z) \rightarrow 0 \text{ as } z \rightarrow -\infty$
 - $g(z) = 0.5 \text{ if } z=0$



MODEL: LOGISTIC REGRESSION

Data: $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$x_i \in \mathbb{R}^d$$

$$y_i \in \{0, 1\}$$

Max. Likelihood

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

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

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

v

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

...

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

- No closed form solution

- Gradient ascent

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

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

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

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

REGULARIZED VERSION

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

KERNEL VERSION

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

Formal Theorem
Representer Theorem

Exercise: Derive the kernel version of logistic regression

META CLASSIFIERS (or)

ENSEMBLE CLASSIFIERS.

WEAK
CLASSIFIERS

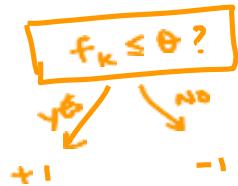
[better than
random]



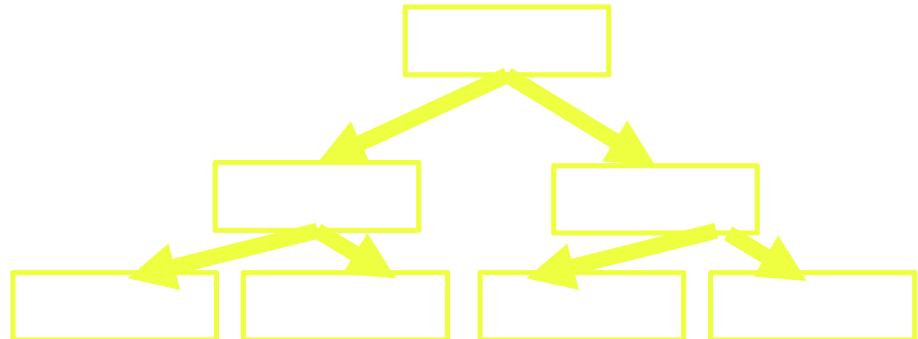
STRONG
CLASSIFIERS

Weak classifiers

DECISION STUMP



Overfit decision tree



high bias, low variance

...



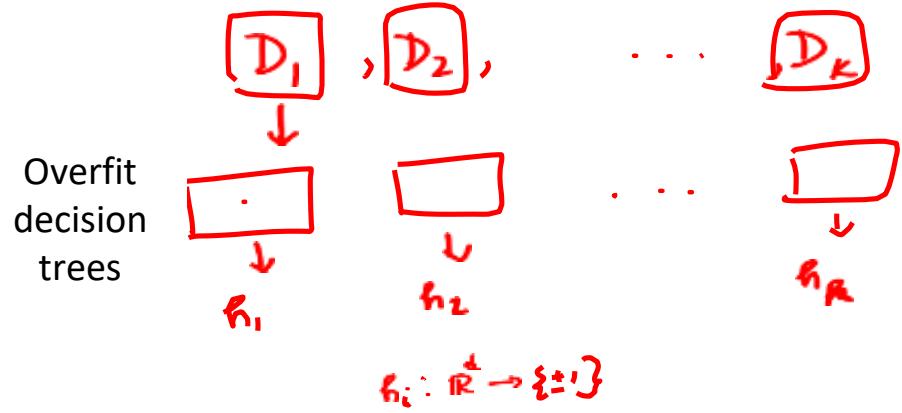
.....



low bias, high variance

$$x_1, x_2, \dots, x_n \sim N(\mu, \sigma^2)$$

$$\hat{\mu}_1 = x_1 \quad \hat{\mu}_2 = x_2, \dots, \hat{\mu}_n = x_n \quad \hat{\mu}_{ML} = \frac{1}{n} \sum x_i$$



$$h^*(x) = \text{majority}(h_1(x), \dots, h_K(x))$$

BAGGING - Bootstrap Aggregation

Chance that a point appears in a dataset

$$1 - \underbrace{\left(1 - \frac{1}{n}\right)}_{\text{as } n \rightarrow \infty} \underbrace{\left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{1}{n}\right)}_{n}$$

$$1 - \underbrace{\left(1 - \frac{1}{n}\right)}_{\text{as } n \rightarrow \infty}^n$$

$$1 - \frac{1}{e}$$

$$\approx 63.2\%$$

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

- Create datasets D_1, \dots, D_k from D by
Sampling with replacement.

- Run weak classifier on D_1, \dots, D_k to get f_1, \dots, f_k

- Aggregate f_1, \dots, f_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

ADA-BOOST

[Freund & Schapire ·
1995
Gödel Prize]

Distribution D over $(\mathbb{R}^d \times \{+1, -1\})$
Unknown but fixed.

x_1, \dots, x_n are iid from D .

$$f_i: \mathbb{R}^d \rightarrow \{+1, -1\}$$

Measure performance using

$$P_{\substack{(x,y) \sim D}}(r(x) \neq y)$$

Misclassification probability.

A weak learner is one which outputs a classifier
Strong

for which

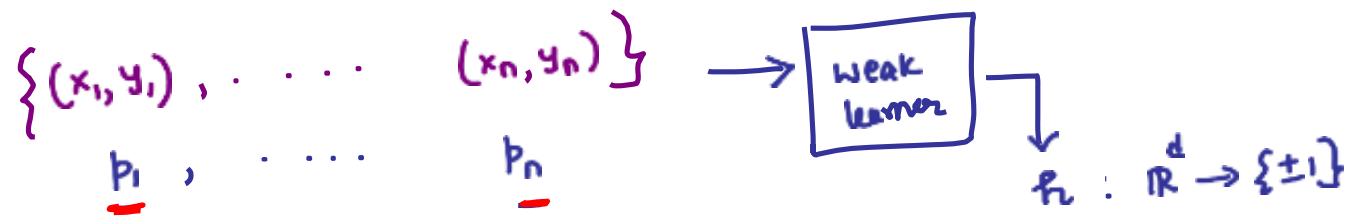
$$P_{\substack{x,y \sim D}}(r(x) = y) \geq \frac{1-\epsilon}{2} + \gamma$$

$$\gamma > 0$$

for any unknown but fixed distribution D .

BOOSTING

Weak learner \rightarrow Strong learner.



$$\boxed{\sum_i p_i = 1}$$

$$\boxed{\sum_{i=1}^n p_i \mathbf{1}(f(x_i) \neq y_i) \leq \frac{1}{2} - \gamma} \quad \boxed{\gamma > 0}.$$

Strong learner classifies all data points correctly.

BOOSTING

[ADA BOOST]

$$S = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$\begin{aligned} x_i &\in \mathbb{R}^d \\ y_i &\in \{-1\} \end{aligned}$$

Initialize

$$D_t(i) = \frac{1}{n} \quad \forall i$$

iteration

for $t = 1, \dots, T$

$$f_t = \left[\text{Input } \boxed{S, D_t} \text{ to get } f_t \right]$$

the weak learner to

$$D_{t+1}(i) = \frac{D_t(i) \cdot e^{-\alpha_t y_i R_t(x_i)}}{\sum_{j=1}^n D_t(j) e^{-\alpha_t y_j R_t(x_j)}} + i$$

$\alpha_t > 0$

end.

$$\underline{H(x)} = \sum_{t=1}^T \alpha_t \boxed{R_t(x)}$$

$C(x) = \text{sign}(H(x))$

ANALYSIS OF BOOSTING

$$D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i R_t(x_i)}}{z_t} \rightarrow \text{update rule}$$

$$z_t = \sum_{j=1}^n D_t(j) e^{-\alpha_t y_j R_t(x_j)}$$

$$D_{t+1}(i) = \left(\frac{D_{t-1}(i) e^{-\alpha_{t-1} y_i R_{t-1}(x_i)}}{z_{t-1}} \right) \cdot \frac{e^{-\alpha_t y_i R_t(x_i)}}{z_t}$$

$$D_{t+1}(i) = \frac{\frac{1}{n} e^{-\sum_{k=1}^t \alpha_k y_i R_k(x_i)}}{\prod_{k=1}^t z_k}$$

$$D_{T+1}(i) = \frac{\frac{1}{n} e^{-y_i H(x_i)}}{\prod_{k=1}^T z_k} \xrightarrow{\quad} \sum_{k=1}^T \alpha_k R_k(x_i)$$

$$\underbrace{\sum_{i=1}^n D_{t+1}(i)}_{=1} = \frac{\sum_{i=1}^n \frac{1}{n} e^{-y_i H(x_i)}}{\prod_{k=1}^T z_k}$$

$$\boxed{\prod_{k=1}^T z_k = \underbrace{\frac{1}{n} \sum_{i=1}^n e^{-y_i H(x_i)}}_{\text{---(A)}}}$$

Recall, $c(x) = \text{Sign}(H(x))$

$$\mathbb{1}(c(x_i) \neq y_i) = \frac{1}{\underbrace{e^{-H(x_i)y_i} + e^{H(x_i)y_i}}_{>0}} \leq e^{-H(x_i)y_i} \quad \text{---(B)}$$

$$\Rightarrow \prod_{k=1}^T z_k = \frac{1}{n} \sum_{i=1}^n e^{-y_i H(x_i)} \quad (\text{from ④})$$

$$\geq \frac{1}{n} \sum_{i=1}^n \mathbf{1}(H(x_i) y_i < 0) \quad (\text{from ⑤})$$

$$\geq \frac{1}{n} \sum_{i=1}^n \mathbf{1}(c(x) \neq y_i)$$

$\underbrace{\phantom{\sum_{i=1}^n \mathbf{1}(c(x) \neq y_i)}}_{\text{error}(c)}$

$\text{error}(c) \leq \prod_{k=1}^T z_k \leq \square$

$$\begin{aligned}
 \prod_{k=1}^T z_k &= \prod_{k=1}^T \left[\sum_{j=1}^n D_k(j) e^{-\alpha_k y_j R_k(x_j)} \right] \\
 &= \prod_{k=1}^T \left[\frac{\sum_{j=1}^n e^{-\alpha_k D_k(j)} \mathbb{1}(y_j = R_k(x_j))}{\sum_{j=1}^n e^{-\alpha_k D_k(j)} \mathbb{1}(y_j \neq R_k(x_j))} \right] \\
 &= \prod_{k=1}^T e^{-\alpha_k} \left(1 - \text{error}(R_k) \right) + e^{\alpha_k} \left(\text{error}(R_k) \right)
 \end{aligned}$$

Choose α_L & ϵ

$e^{-\alpha_L} (1 - \text{error}(h_L)) + \frac{\epsilon^2}{2} (\text{error}(h_L))^2$ is as small as possible.

$$\alpha_L = \ln \sqrt{\frac{1 - \text{error}(h_L)}{\text{error}(h_L)}}$$

$$= \prod_{L=1}^T 2 \sqrt{\text{error}(h_L) (1 - \text{error}(h_L))}$$

$$\leq \prod_{L=1}^T 2 \sqrt{\left(\frac{1}{2} - \gamma\right) \left(\frac{1}{2} + \gamma\right)}$$

[Show this]

$$e^{-\alpha_L} (1 - \gamma) + \frac{\epsilon^2}{2} \gamma$$

$$(1 - \gamma) e^{-\alpha_L} (-1) + \frac{\epsilon^2}{2} \gamma = 0$$
$$-1 + \gamma e^{-\alpha_L} + \frac{\epsilon^2}{2} \gamma = 0$$

$$= \prod_{L=1}^T 2 \sqrt{\frac{1 - 4\gamma^2}{4}}$$

$$\begin{aligned}
 &= \prod_{k=1}^T \sqrt{1 - e^{-2\gamma^2}} \\
 &\leq \prod_{k=1}^T (e^{-4\gamma^2})^{1/2} = \prod_{k=1}^T e^{-2\gamma^2} \\
 &= e^{-2\gamma^2 T} \quad \text{--- (2)}
 \end{aligned}$$

$\frac{1}{n} \sum_i z_i (c x_i + y_i)$

$$\begin{aligned}
 \underline{\text{error}(c)} &\leq \prod_{k=1}^T z_k \leq \underbrace{e^{-2\gamma^2 T}}_{\uparrow} \leq \frac{1}{2^n}
 \end{aligned}$$

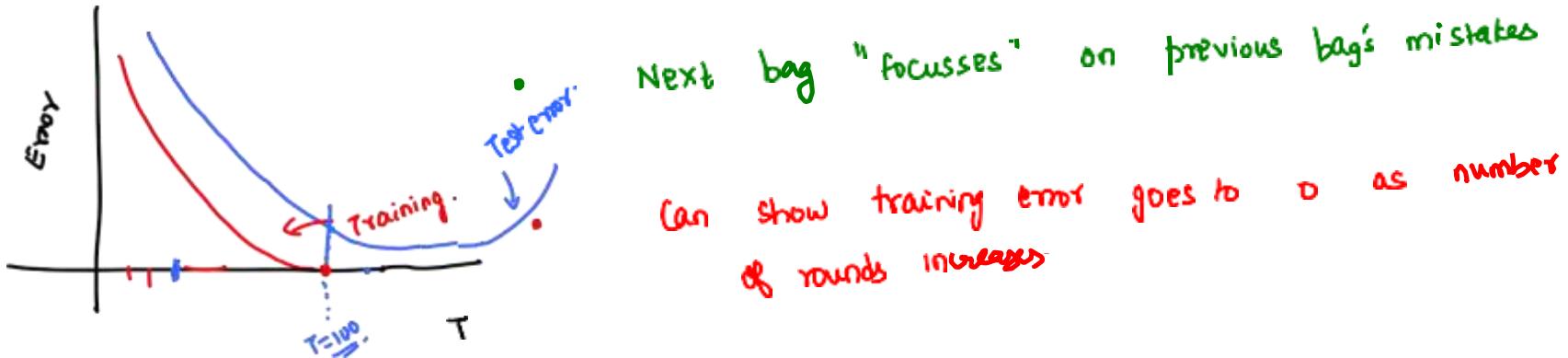
$\frac{1}{n} \geq \left(\frac{1}{2\gamma^2} \right) \ln(2n)$

$$\Rightarrow \underline{\text{error}(c) = 0}$$

Strong Learner!

Boosting-Summary

- Every bag depends on previous bag's performance



Cannot run in parallel – hence sometimes bagging is preferred in practice

Usually weak learners are decision stumps \rightarrow high bias, low variance

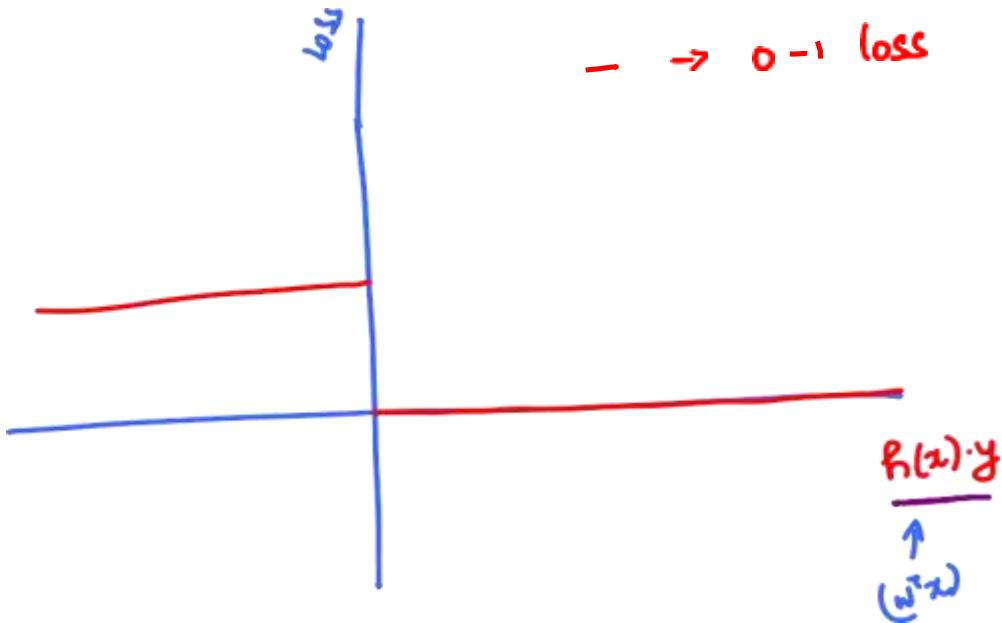
Boosting reduces bias without affecting variance a lot

Why are there so many binary classification algorithms?

$$\begin{matrix} x \in \mathbb{R}^d \\ y \in \{\pm 1\} \end{matrix}$$

$$\min_w \sum_{i=1}^n \mathbb{I}(w^T x_i \cdot y_i < 0)$$

NP-Hard problem



Linear Regression for classification

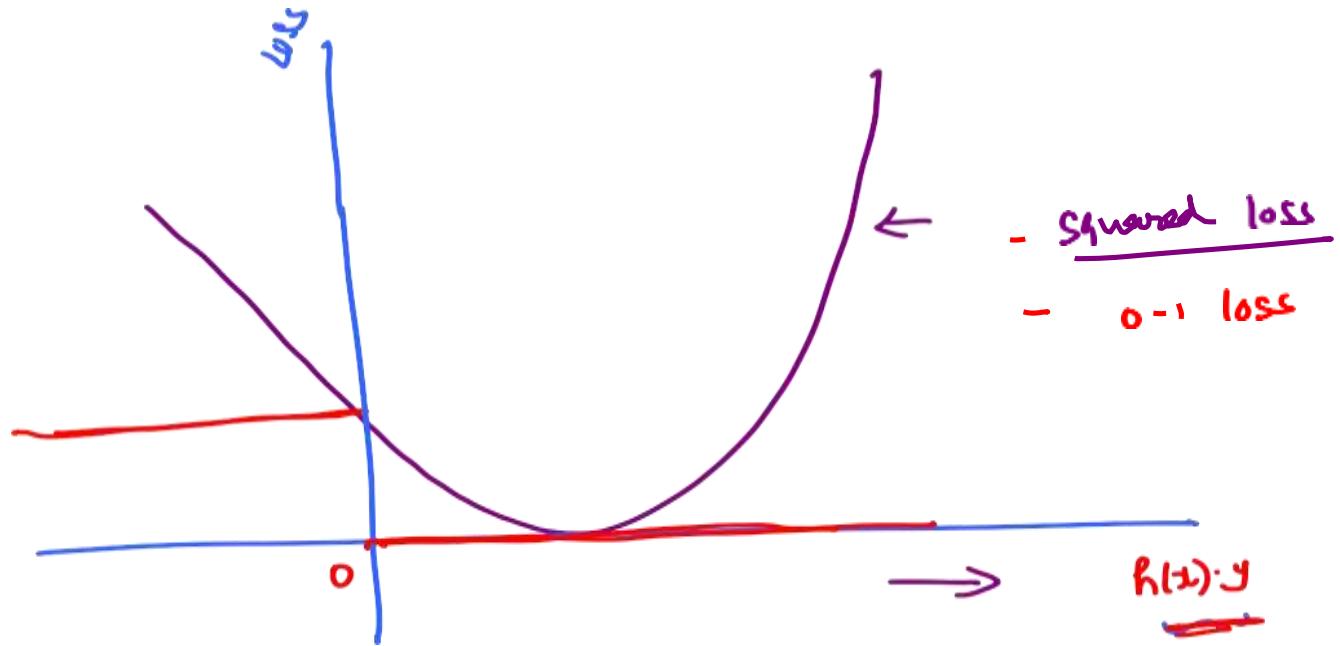
$$\min_{\mathbf{h}} \sum_{i=1}^n (h(x) - y_i)^2$$

Final classifier: $\text{sign}(h(x))$

Loss per point for h

$$(h(x) - y)^2 = \begin{cases} (h(x) - 1)^2 & \text{if } y = 1 \\ (h(x) + 1)^2 & \text{if } y = -1 \end{cases}$$

$$\begin{aligned} \text{if } y = 1 & \quad \frac{(h(x))^2 + 1 - 2h(x)}{2} = (h(x))^2 + 1 - \frac{2h(x)y}{2} \\ \text{if } y = -1 & \quad \frac{(h(x))^2 + 1 + 2h(x)}{2} = (h(x))^2 + 1 - \frac{2h(x)\cdot y}{2} \\ & = (h(x) \cdot y - 1)^2 \end{aligned}$$



SUPPORT VECTOR MACHINES

$$\min_{w, \xi}$$

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

$$\begin{aligned} w^T x_i y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

Equivalently

$$\begin{aligned} \xi_i &\geq 1 - w^T x_i y_i \\ \xi_i &\geq 0 \end{aligned}$$

Equivalently

$$\xi_i \geq \max(1 - w^T x_i y_i, 0)$$

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

$$\xi_i \geq \max(1 - w^\top x_i y_i, 0)$$

Equivalently,

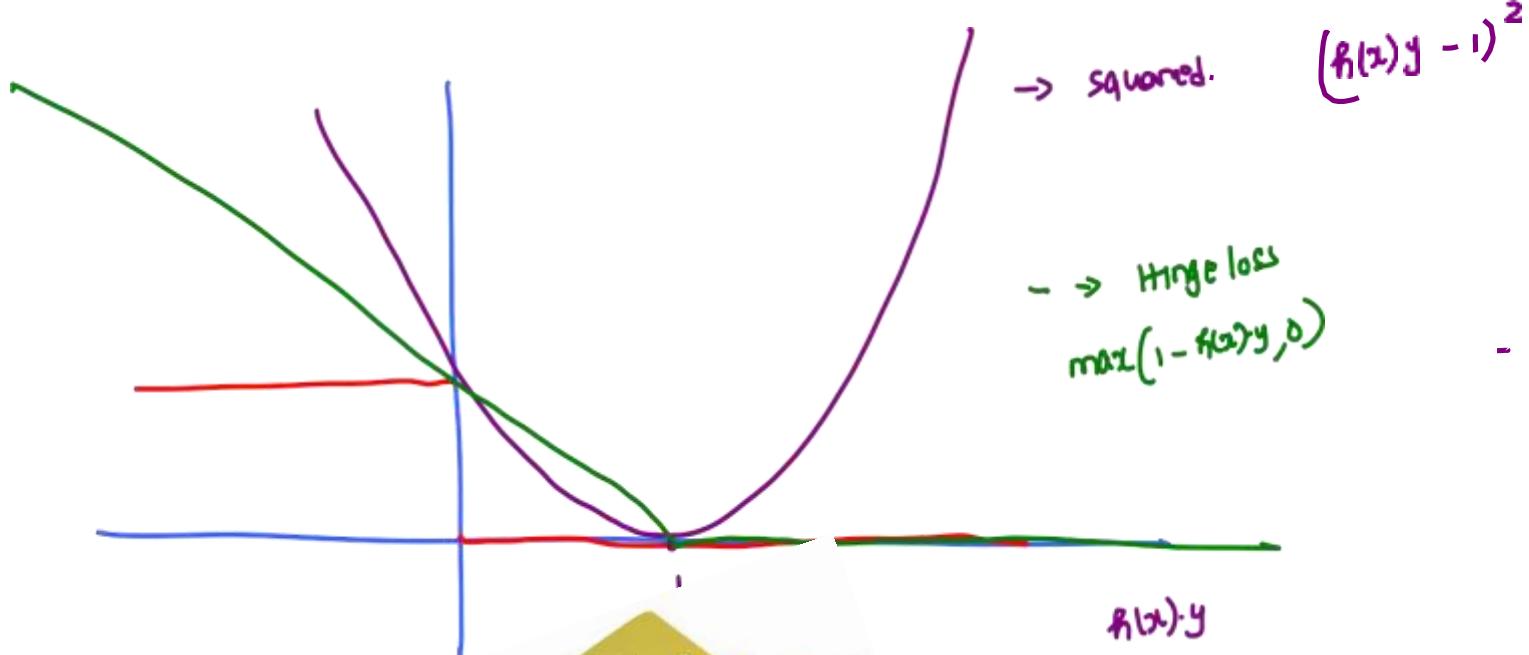
$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(1 - w^\top x_i y_i, 0)$$

Model Regularization term	Data Loss term
---------------------------------	----------------------

So what exactly is the loss?

$$l(w, (x, y)) = \max(1 - w^\top x y, 0)$$

HINGE LOSS



Hinge

Courtesy: Google images

Logistic regression

$$\max_{\omega} \prod_{i=1}^n \left(\frac{g(\omega^T x_i)}{1 + g(\omega^T x_i)} \right)^{z_i} \left(\cdot - \frac{g(\omega^T x_i)}{1 + g(\omega^T x_i)} \right)^{1-z_i}$$

$z_i = 1 \quad \text{if} \quad y_i = 1$
 $z_i = 0 \quad \text{if} \quad y_i = -1$

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

$$\max_{\omega} \prod_{i=1}^n \frac{g(\omega^T x_i)}{1 - g(\omega^T x_i)}^{z_i}$$

$$\max_{\omega} \sum_{i=1}^n z_i \log(g(\omega^T x_i)) + (1-z_i) \log(1-g(\omega^T x_i))$$

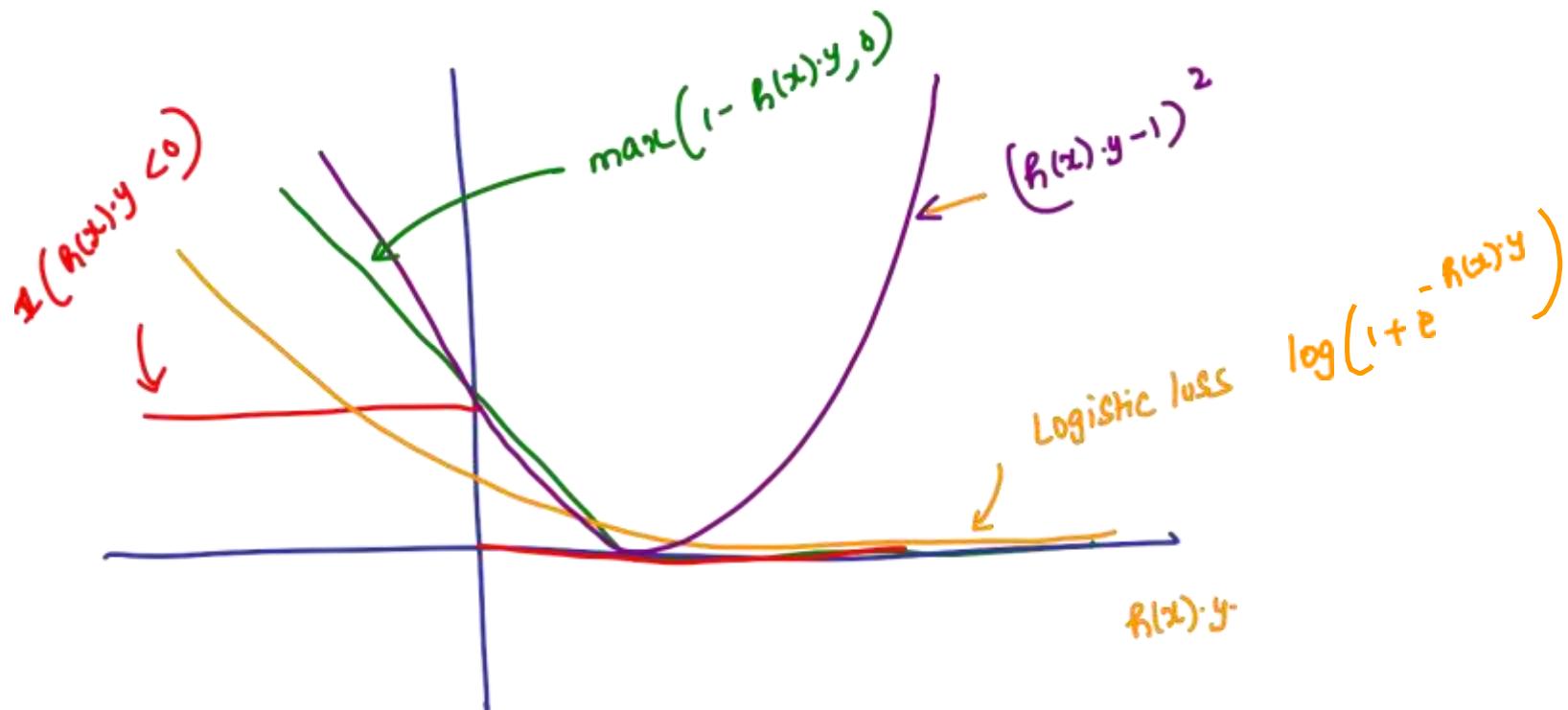
$$\equiv \min_{\omega} \sum_{i=1}^n \left[-z_i \log(g(\omega^T x_i)) + (1-z_i) \log(1-g(\omega^T x_i)) \right]$$

Loss for a single point when $z_i = 1$ ($y_i = 1$)

$$\begin{aligned}-\log(g(w^T x_i)) &= -\log\left(\frac{1}{1+e^{-w^T x_i}}\right) \\ &= \log\left(1 + e^{-w^T x_i}\right) = \boxed{\log\left(\frac{-y_i w^T x_i}{1+e^{-w^T x_i}}\right)}\end{aligned}$$

Loss for a single point when $z_i = 0$ ($y_i = -1$)

$$\begin{aligned}= -\log(1 - g(w^T x_i)) &= -\log\left(1 - \frac{1}{1+e^{-w^T x_i}}\right) \\ = -\log\left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}}\right) &= \log\left(1 + e^{-w^T x_i}\right) = \boxed{\log\left(\frac{-y_i w^T x_i}{1+e^{-w^T x_i}}\right)}\end{aligned}$$

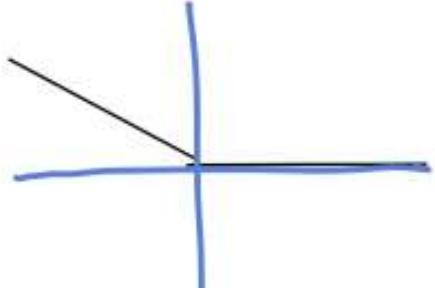


Perceptron

$$w_{t+1} = w_t + x_t y_t$$

Consider the hinge loss

$$\begin{aligned}\text{loss}(w, (x, y)) \\ = \max(0, 1 - w^T x y)\end{aligned}$$



$$\text{loss}(w, (x, y)) = \max(0, 1 - \vec{w}^T x y)$$

Sub-gradient

$$\frac{\partial \text{loss}}{\partial w} = \begin{cases} -x y & \text{if } (\vec{w}^T x) y < 0 \\ 0 & \text{if } (\vec{w}^T x) y > 0 \\ [-1, 0] x y & \text{if } (\vec{w}^T x) y = 0 \end{cases}$$

↳ choose $-x y$ if
 mistake.

when mistake

mistake

$$w_{t+1} = w_t - \eta_t (-x_t y_t)$$

$$\hookrightarrow = 1$$

• Perceptron can be viewed as

S.G.D with hinge loss with step size
= 1

BOOSTING

Recall

$$\prod_{t=1}^T z_t = \frac{1}{n} \sum_{i=1}^n e^{-\left(\sum_{t=1}^T \alpha_t h_t(x_i) y_i \right)}$$
$$:= \frac{1}{n} \sum_{i=1}^n e^{-H_T(x_i) y_i}$$

at round T , we so far have

$$\sum_{t=1}^{T-1} \alpha_t h_t(x)$$

We need to add one more
classifier in round T to get

$$\sum_{t=1}^{T-1} \alpha_t h_t(x) + \underline{\alpha_T h_T(x)}$$

↗ greedy choice
to minimize

$$\frac{1}{n} \sum_{i=1}^n e^{\left[\sum_{t=1}^{T-1} \alpha_t h_t(x_i) + \underline{\alpha_T h_T(x)} \right] y_i}$$

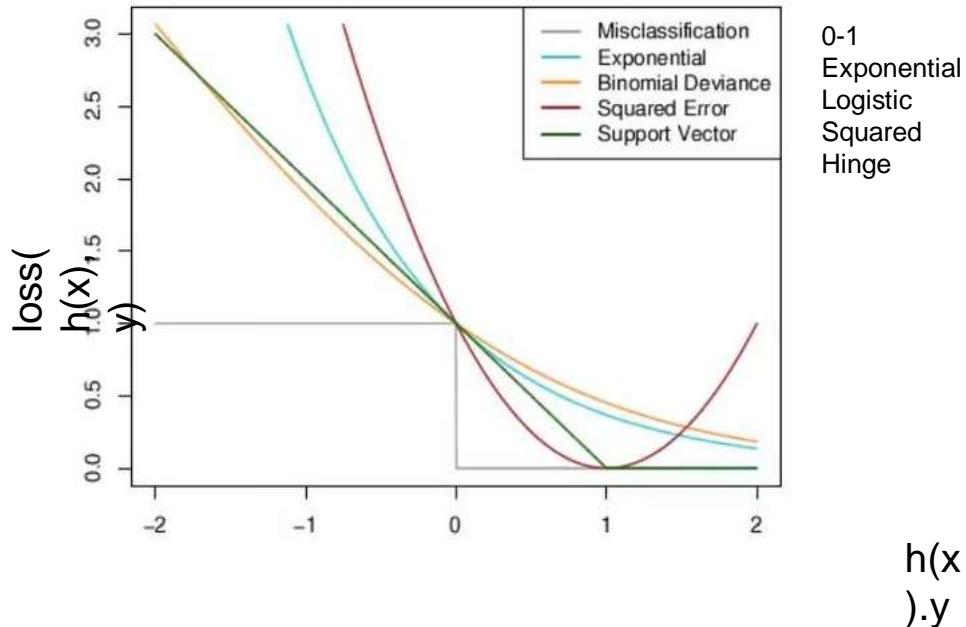
$$\text{Define } \text{loss}(h, y) = e^{-y h(x)}$$

- At every round, one chooses α_t to minimize Z_t .
- Thus Ada-boost can be viewed as "greedy / co-ordinate descent" on exponential loss

What about Perceptron and Boosting?

- one can argue perceptron update rule is equivalent to doing a SGD on hinge loss with stepsize = 1
↳ Stochastic sub-gradient descent.
- $w_{t+1} = w_t + \frac{x_i y_i}{B}$
$$-\frac{f(x) - y}{B}$$
- Boosting (AdaBoost) can be viewed as "greedy / coordinate wise" descent of exponential loss ↑

SUMMARY



-> The 0-1 loss is NP-hard to optimize even for linear classifiers

-> Different algorithms get around this by using a “surrogate” loss function

-> Surrogates are usually **“convex” surrogates** that are easy to optimize