

## Linear Discriminant Models (Linear Models for Classification)

Definition: Linear Discriminant Model (LDM) is the following set:

$$\mathcal{F}_L = \left\{ f \mid \exists \omega = \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_n \end{bmatrix} \ni f(x) = \text{sign}(\omega^T \phi(x) - b) \neq x \in \mathcal{X} \right\}$$

→ functions in  $\mathcal{F}_L$  can be identified by

model parameters  $(\omega, b) \in \mathbb{R}^{n+1}$

→ Effective no. parameters =  $n$   
(One can be fixed arbitrarily)

→ Learning with this model can be posed as problem of obtaining the "best" parameters.

IDEAL OBJ.

$$P[\gamma \neq \text{sign}(\omega^T \phi(x) - b)]$$

$\min_{\omega \in \mathbb{R}^n, b \in \mathbb{R}}$   
s.t.

$$\|\omega\| = 1$$

$$\textcircled{I^*}$$

EMPIRICAL OBJ.

$$\approx \min_{\omega \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{y_i \neq \text{sign}(\omega^T \phi(x_i) - b)\}}$$

s.t.

$$\|\omega\| = 1$$

$$\textcircled{I}$$



→ Minimizer of Training error will go to that of true error, as  $m \rightarrow \infty$ .

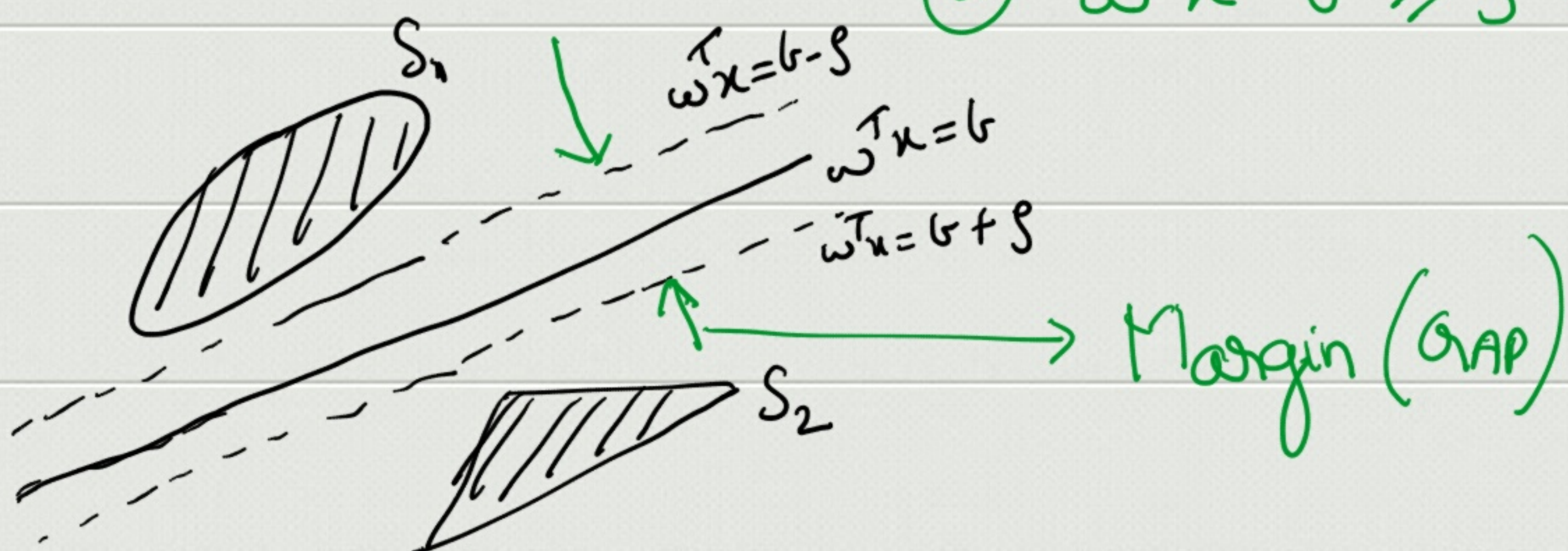
→ However, I is a non-convex problem with discontinuous objective!

→ In general, known to be Computationally hard

→ If, the two classes are strictly linearly separable, then polynomial time solver.

Strict Linear Separability: Two sets  $S_1, S_2$  are said to be strictly linearly separable iff:

$$\exists w \in \mathbb{R}^n, b \in \mathbb{R}, \rho > 0 \quad \text{such that} \quad \begin{aligned} \text{(i)} \quad & w^T x - b \leq -\rho \quad \forall x \in S_1 \\ \text{(ii)} \quad & w^T x - b \geq \rho \quad \forall x \in S_2 \end{aligned}$$





With the assumption,  $\textcircled{I}$  can be written as:

min  
 $w \in \mathbb{R}^n$ ,

$b \in \mathbb{R}$ ,

$\delta > 0$

s.t.

$$y_i (\omega^T \phi(x_i) - b) \geq \delta, \quad \forall i=1 \text{ to } m;$$

$$\|\omega\| = 1.$$

The constraints  $\|\omega\|=1$  &  $\delta > 0$  are spoilers:

→ Option(i): Maximize margin.

→ Option(ii): Different normalization.

$$\text{Margin} = \frac{2\delta}{\|\omega\|}$$

Option(i):

min

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

$b \in \mathbb{R}$ ,

$\delta \geq 0$

s.t.

$-\delta$

$$y_i (\omega^T \phi(x_i) - b) \geq \delta, \quad \forall i=1 \text{ to } m$$

$$\|\omega\| \leq 1.$$

$\textcircled{\text{II}_a}$



Option 2:

Select  $\beta = 1$ , then Margins =  $2/\|\omega\|$ .



min  
 $\omega \in \mathbb{R}^n$ ,  
 $b \in \mathbb{R}$   
s.t.

$$y_i (\omega^T \phi(x_i) - b) \geq 1 \quad \forall i=1 \text{ to } m$$

Option 1 meets Option 2:



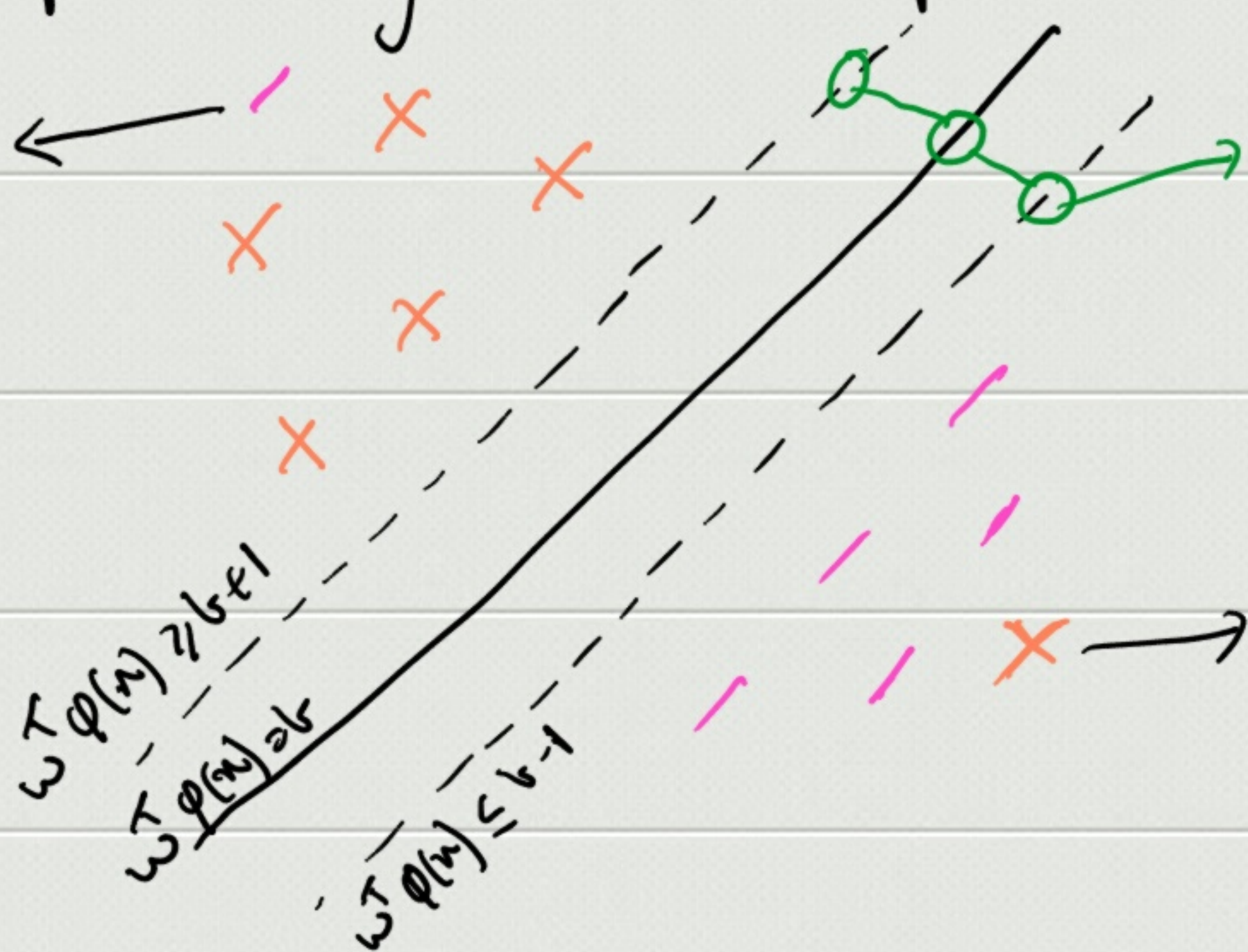
min  
 $\omega \in \mathbb{R}^n$ ,  
 $b \in \mathbb{R}$   
s.t.

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

$$y_i (\omega^T \phi(x_i) - b) \geq 1 \quad \forall i=1 \text{ to } m$$

If assumption of strict separation fails, then:

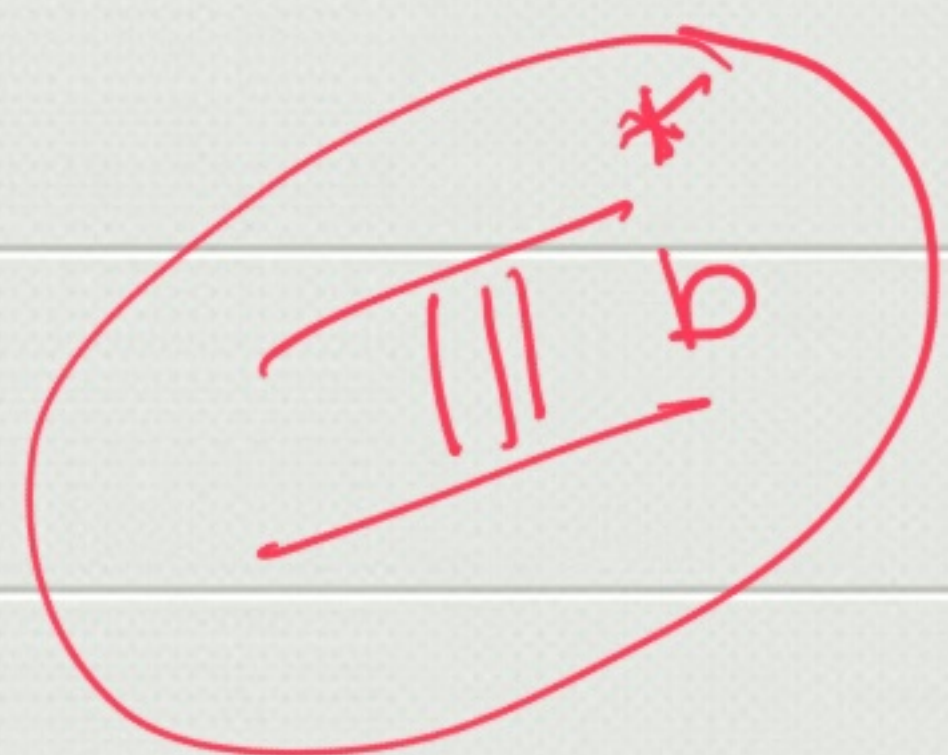
$\xi_i > 0$



Canonical  
hyperplane  
representation

$\xi_i > 0$





$$\min_{\omega \in \mathbb{R}^n, b \in \mathbb{R}, \xi \in \mathbb{R}^m}$$

$$\xi \in \mathbb{R}^m$$

s.t.

$$\sum_{i=1}^m 1_{\{\xi_i > 0\}}$$

$$y_i(\omega^T \phi(x_i) - b) \geq 1 - \xi_i, \xi_i \geq 0,$$

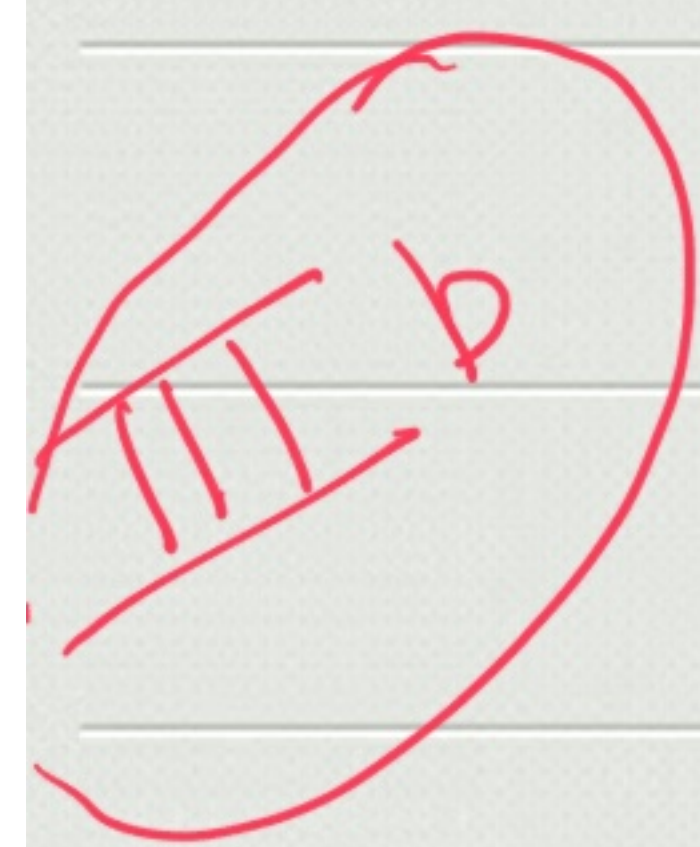
$\forall i=1 \text{ to } m.$

Combinatorial objective is reason for hardness.

→ Alternative perf. measure?

↓ usual trick

replace no. by sum. (tightest Convex Relaxation)



$$\min_{\omega \in \mathbb{R}^n, b \in \mathbb{R}, \xi \in \mathbb{R}^m}$$

$$\sum_{i=1}^m \xi_i$$

$$\underline{\underline{\text{s.t.}}} \quad y_i(\omega^T \phi(x_i) - b) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i=1 \text{ to } m$$

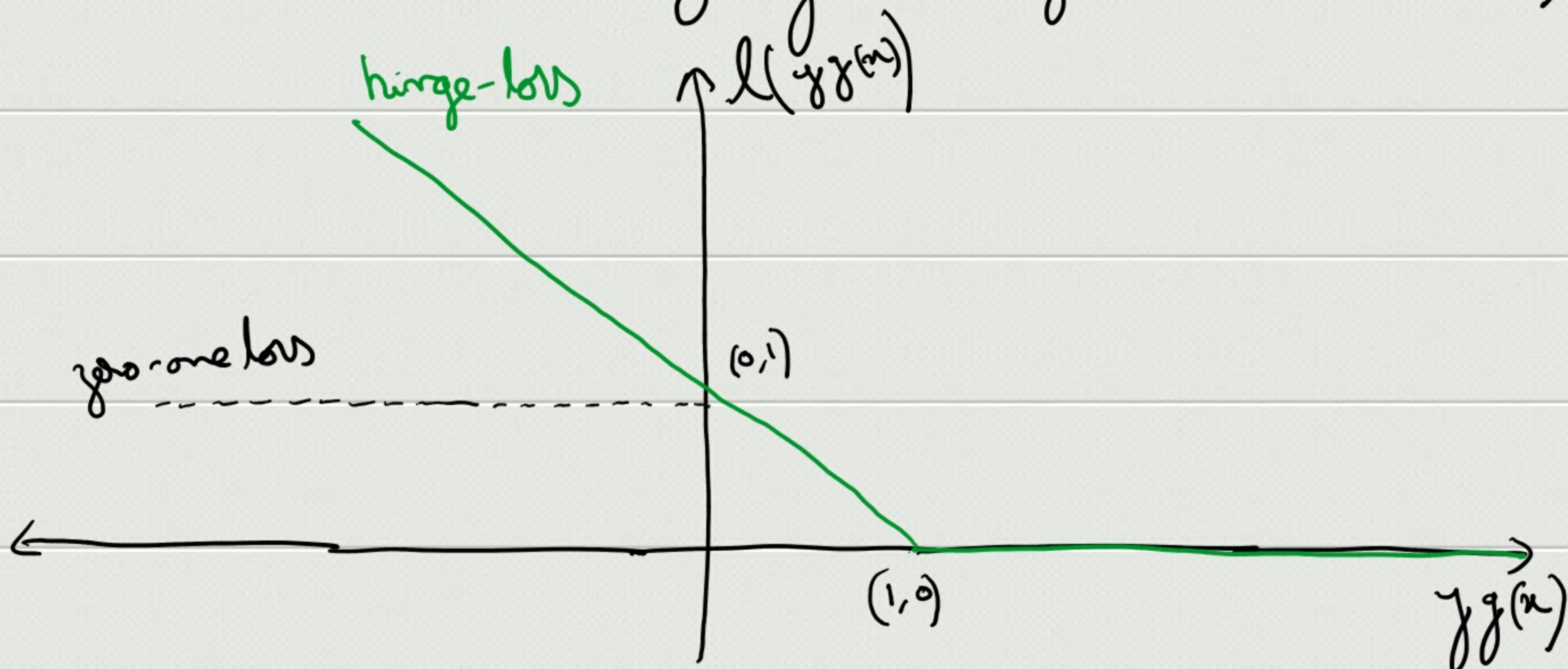


(III b) Can be re-written as:

$$\min_{\omega \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i(\omega^T \phi(x_i) - b))$$

$\downarrow$   
 $\ell(y_i g(x_i))$

Equivalent to minimizing average loss  $\ell(y g(x))$



→ Average 0-1 loss gives accuracy

→ Average hinge-loss gives obj. in (III b).

$\downarrow m \rightarrow \infty$

$$\mathbb{E}[\max(0, 1 - \gamma(\omega^T \phi(x) - b))]$$

expected hinge-loss.

(upper bound on minmax. error).



## Surrogate Loss Functions for Binary Classification

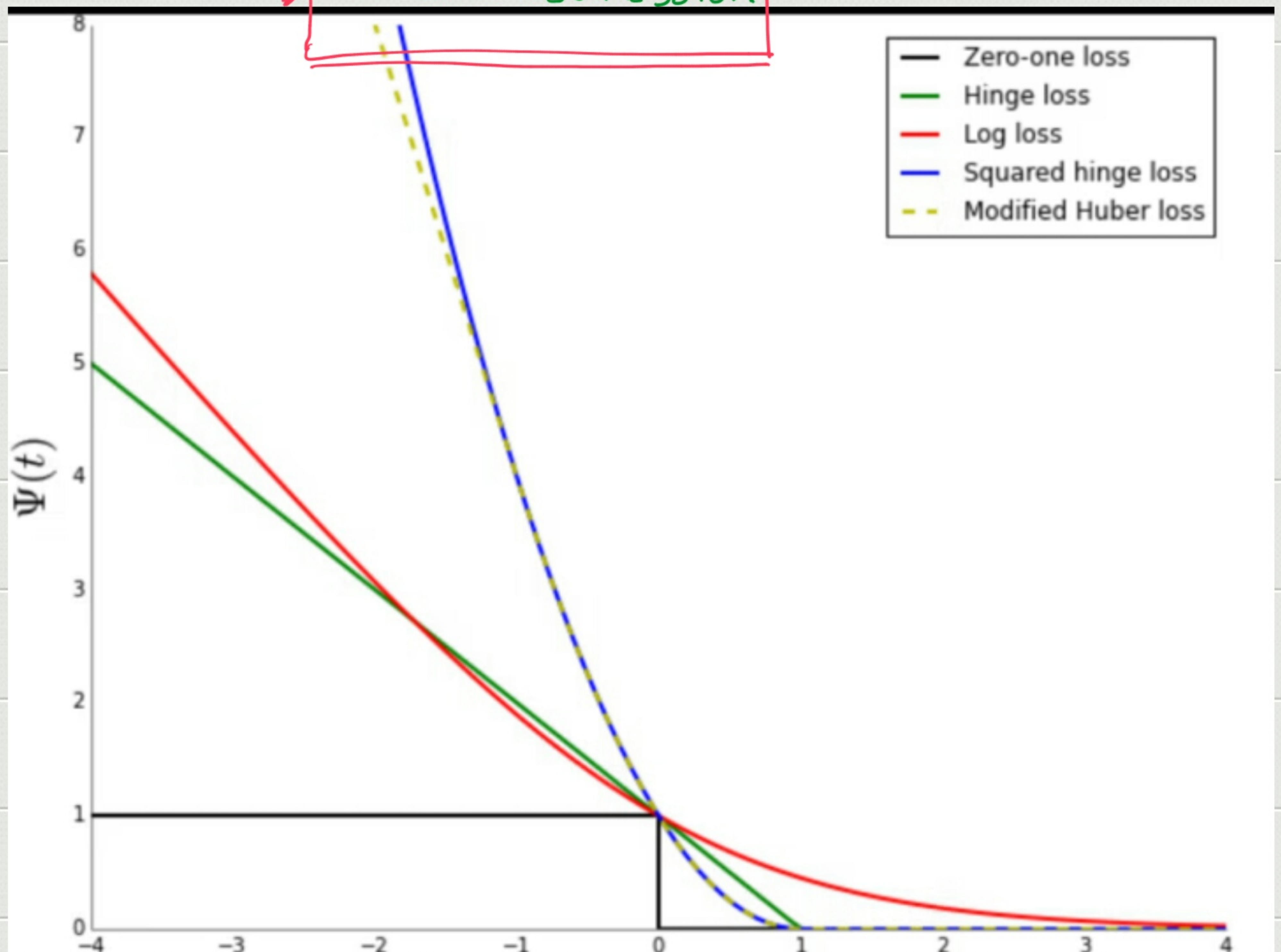
→ 0-1 loss:  $l(z) = \begin{cases} 0 & \text{if } z \geq 0 \\ 1 & \text{if } z < 0 \end{cases}$

→ hinge-loss:  $l(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ 1-z & \text{if } z < 1 \end{cases}$  L1-Loss

→ squared-hinge-loss:  $l(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ (1-z)^2 & \text{if } z < 1 \end{cases}$

→ logistic-loss:  $l(z) = \log(1 + e^{-z})$  L2-LOSS

→ LOGISTIC REGRESSION





Introducing surrogate loss functions in II IV will need a "regularization" parameter:

III b

min  
 $\omega \in \mathbb{R}^n$ ,  
 $b \in \mathbb{R}$ ,  
 $\xi \in \mathbb{R}^m$   
s.t.

$$\frac{1}{2} \|\omega\|^2 + \frac{C}{m} \sum_{i=1}^m \xi_i$$

$$y_i (\omega^T \phi(x_i) - b) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i=1 \text{ to } m$$

SVM (Support Vector Machine)

II a is called hard-margin SVM.

$C$  is reg. parameter

II b

min  
 $\omega \in \mathbb{R}^n, b \in \mathbb{R}$ ,  
 $\xi \in \mathbb{R}^m, \xi \geq 0$

$$\sum_{i=1}^m \xi_i - \gamma f$$

$\gamma$  is the reg. parameter.

s.t.

$$y_i (\omega^T \phi(x_i) - b) \geq \gamma - \xi_i, \xi_i \geq 0, \quad \forall i=1 \text{ to } m,$$

$$\|\omega\| \leq 1.$$

$\gamma$ -SVM ( $\gamma$ -Support Vector Machine)

II a is called hard-margin  $\gamma$ -SVM.



General form for maximum-margin lin. disc.:

$$\min_{w \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{C}{m} \sum_{i=1}^m \ell(\gamma_i(w^T \phi(x_i) - b))$$

'C' is called "hyperparameter"

and for linear discriminators is:

$$\min_{w \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m \ell(\gamma_i(w^T \phi(x_i) - b))$$

max-margin lin.  
+ sparse loss

= SUPPORT VECTOR

METHOD/MODEL

## Analysis of Linear Discriminant Models

① They are essentially nearest neighbour models

↑ This is more popularly known as Representer Theorem.

② Above + "sparse" loss = "smart" nearest neighbour.

That is where name Support Vector comes in



Representer Theorem: Both lin. disc. and

max-margin versions admit optimal solution  $\omega^*$  of the form:

$$\omega^* = \sum_{i=1}^m \alpha_i y_i \phi(x_i) \quad \text{for some } \alpha_i \in \mathbb{R}.$$

Inference:  $g(x) = \text{sign}(\omega^{*\top} \phi(x) - b)$

(nearest neighbour)  $= \begin{cases} 1 & \text{if } \sum_i \alpha_i y_i k(x, x_i) \geq b \\ 0 & \text{if } \sum_i \alpha_i y_i k(x, x_i) < b \end{cases}$

where,  $k(x, z) \equiv \phi(x)^\top \phi(z) \forall x, z \in \mathcal{X}$

With losses like hinge-loss

$\Rightarrow \omega^*$  does not depend on the "correct" examples

Further if training error is low,

$\Rightarrow$  many  $\alpha_i$  are zero. Non-zero  $\alpha_i$  examples are called Support Vectors

$\Rightarrow$  "Smart" nearest neighbours.



Representer result is useful in at least two more ways.

(i) Training and inferring with arbitrary  $\mathcal{X}$  as long as  $\phi(x)^T \phi(z)$  is known for  $x, z \in \mathcal{X}$ .

$$\rightarrow \omega^T \phi(x) = \sum_{i=1}^m \alpha_i y_i k(x_i, x)$$

(ii) EXPLOIT 'SPARSITY' in  $\alpha$ 's for better algorithms.

## On efficiently training linear models

$\rightarrow$  note that there may be multiple choices of  $\alpha$ 's such that

$$\omega = \sum_i \alpha_i y_i \phi(x_i)$$

$\rightarrow$  not all will be sparse.

HENCE, substituting in the optimization need not be the "best" way for obtaining  $\alpha$ 's



However, it is indeed a baseline:

$$\min_{\alpha \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{2} \alpha^T K \alpha + \frac{C}{m} \sum_{i=1}^m l_i \left( (K\alpha)_i - b \right),$$

$$\text{where, } (K)_{ij} = \varphi(x_i)^T \varphi(x_j)$$

Here,  $(K)_{ij}$  denotes the  $i, j^{\text{th}}$  entry in the matrix  $K$ .

$(K\alpha)_i$  denotes the  $i^{\text{th}}$  entry in the vector  $K\alpha$

One may solve this using cvx or gradient descent. However, since  $\alpha$  has no "special" property, there were observed to not scale well.

Whereas, we intuitively know that  $\alpha$  could be chosen to be sparse at least with losses like hinge-loss. Hence we aim at analyzing the optimal solution, more carefully.



→ So, analyze optimal soln. more carefully

→ Simple gradient cond. is not insightful:

$$\omega + \frac{c}{m} \sum_{i=1}^m \underbrace{l'(y_i \omega^T \phi(x_i))}_{\uparrow} \phi(x_i) = 0$$

Solve following instead

if decoupled then better  
insightful.

$$\omega = \sum_i \alpha_i \phi(x_i),$$

$$\frac{c}{m} l'(y_i \omega^T \phi(x_i)) = -\alpha_i$$

insightful, but how to solve?

Following is a generic technique  
that exposes both there:

Rewrite optimization prob. as:

$$\min_{\substack{\omega \in \mathbb{R}^n, \zeta \in \mathbb{R} \\ \zeta \in \mathbb{R}^m \\ \text{s.t.}}} \frac{1}{2} \|\omega\|^2 + \frac{c}{m} \sum_{i=1}^m l_i(\zeta_i)$$

$$\zeta_i = \omega^T \phi(x_i) - b_i \quad \forall i=1 \text{ to } m$$

$$l_i(\zeta_i) \equiv l(y_i \zeta_i)$$



$$= \min_{\substack{\omega \in \mathbb{R}^n, b \in \mathbb{R}, \\ \gamma \in \mathbb{R}^m}} \frac{1}{2} \|\omega\|^2 + \frac{c}{m} \sum_{i=1}^m l_i(\gamma_i) + \max_{\alpha \in \mathbb{R}^m} \sum_i \alpha_i (\gamma_i - \omega^T \phi(x_i) - b)$$

$$= \max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \alpha^T G \alpha - \frac{c}{m} \sum_{i=1}^m l_i^* \left( \frac{m \alpha_i \gamma_i}{c} \right),$$

s.t.

$$\sum_i \alpha_i \gamma_i = 0,$$

where,  $G_{ij} = \gamma_i \gamma_j k(x_i, x_j)$ ,

$l^*$  is known as  
 (i) Conjugate  
 (ii) Fenchel dual  
 (iii) Legendre transform of  $l$

$$l_i^*(\lambda) = \max_{\gamma \in \mathbb{R}} \lambda \gamma - l_i(\gamma)$$

and more importantly, at optimality, given for hinge loss as below

$$\rightarrow \omega = \sum_{i=1}^m \alpha_i \gamma_i \phi(x_i),$$

$$\rightarrow \alpha_i = 0 \Rightarrow \gamma_i (\omega^T \phi(x_i) - b) \geq 1$$

$$\rightarrow \alpha_i \in (0, \frac{c}{m}) \Rightarrow \gamma_i (\omega^T \phi(x_i) - b) = 1$$

$$\rightarrow \alpha_i = \frac{c}{m} \Rightarrow \gamma_i (\omega^T \phi(x_i) - b) \leq 1$$

Support Vectors



Exercise: Compute conjugate of all the surrogate losses you know.

Whenever, optimal soln. is known to be sparse, two algorithms come to mind:

(i) Co-ordinate descent

→ Optimize one or few variables iteratively.

(ii) Active-Set methods

→ Guess the right "active" variables.

→ In each case, the sub-problem at every iterate is small and no. iterations may not be high as solution is known to be sparse



→ Minimum no. variables to optimize is 2,  $\because \sum_i \alpha_i y_i = 0$

→ Optimization with such two variables  
→ has closed form solution

→ Selection of pairs by heuristic:  
→ Worst optimality violator.

→ If 'b' is set to zero (Say by  $\phi(x) \rightarrow \begin{bmatrix} \phi(x) \\ 1 \end{bmatrix}$ )  
→ Round robin variable by variable is OK.

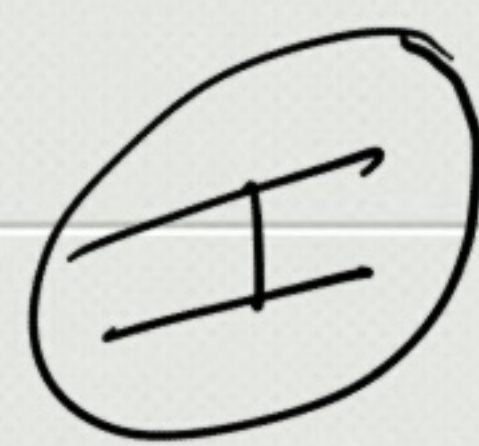
→ For non-sparse losses, above will be slow.  
→ Logistic regression will be covered later.



# LAGRANGE DUALITY

Given,

$$\begin{array}{ll} \min_{x \in X} & f(x) \\ \text{s.t.} & g_i(x) \leq 0 \quad \forall i=1 \text{ to } m \end{array}$$



PRIMAL

$x$  is primal variable.

Assumptions:

- (i) All  $f, g_i$  are convex
- (ii) Feasibility, Boundedness, Solvability
- (iii) Slater's cond. if  $g_i$ 's are non-linear.

Define: 
$$L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) \quad \forall x \in X, \lambda \geq 0$$

LAGRANGIAN FUNCTION

$$\overline{L}(x) \equiv \max_{\lambda \geq 0} L(x, \lambda) \quad \Bigg| \quad \underline{L}(\lambda) \equiv \min_{x \in X} L(x, \lambda)$$



# LAGRANGE DUALITY

$$\textcircled{I} = \min_{x \in X} \bar{L}(x) = \min_{x \in X} \max_{\lambda \geq 0} L(x, \lambda)$$

$$= \max_{\lambda \geq 0} \min_{x \in X} L(x, \lambda) = \max_{\lambda \geq 0} \underline{L}(\lambda) = \textcircled{II}$$

$\downarrow$  DUAL  $\downarrow$

$$\left. \begin{array}{l} x^* \text{ is an optimal soln. of } \textcircled{I} \\ \lambda^* \text{ is an optimal soln. of } \textcircled{II} \end{array} \right\} \bar{L}(x^*) = \underline{L}(\lambda^*)$$

TST:  $x^*$  is an optimal soln. of  $\min_{x \in X} L(x, \lambda^*)$ ;  
 $\lambda^*$  is an optimal soln. of  $\max_{\lambda \geq 0} L(x^*, \lambda)$

Proof:

$$\bar{L}(x^*) = \max_{\lambda \geq 0} L(x^*, \lambda) \geq L(x^*, \lambda^*) \geq \min_{x \in X} L(x, \lambda^*) = \underline{L}(\lambda^*)$$

$$f(x^*) = f(x^*) + \sum_{i=1}^m \lambda_i^* g_i(x^*) \Leftrightarrow \lambda_i^* g_i(x^*) = 0 \quad \forall i$$



# Optimality Conditions

Given  $\lambda^*$ , necessary conds. on  $x^*$  are:

(i)  $x^* \in \mathcal{X}, g_i(x^*) \leq 0 \quad \forall i$  (Feasibility)

(ii)  $x^* = \underset{x \in \mathcal{X}}{\operatorname{argmin}} L(x, \lambda^*) \xleftrightarrow[\text{L is diff.}]{\text{X is open}} \left. \nabla_x L(x, \lambda^*) \right|_{x=x^*} = 0$  (gradient)

(iii)  $\lambda_i^* g_i(x^*) = 0 \quad \forall i = 1 \text{ to } m$ . (Complementary Slackness)

Interestingly, above are sufficient for optimality of  $x^*$ .

→ in fact, sufficient if  $\lambda$  exists!

→ So (i), (ii), (iii) are nec. & suff. conds. for optimality of both  $x^*$  and  $\lambda^*$ .



# Dual of SVM + hinge-loss

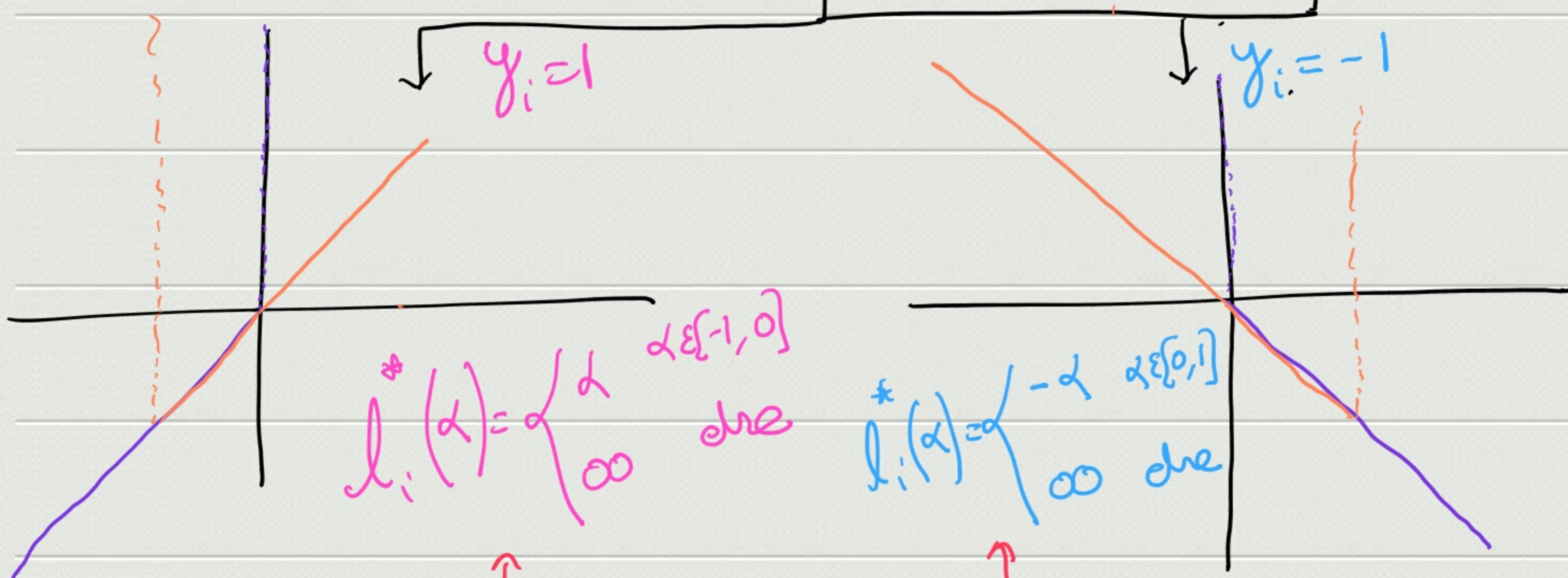
$$l_i(z) \equiv \max(0, 1 - y_i z)$$

$$l_i^*(\alpha) = \max_{z \in \mathbb{R}} \alpha z - l_i(z)$$

$$= \max_{z \in \mathbb{R}} \alpha z - \max(0, 1 - y_i z)$$

$$= \max \left( \begin{array}{l} \max_{z: 1 - y_i z \leq 0} \alpha z \\ \max_{z: 1 - y_i z > 0} (\alpha + y_i) z - 1 \end{array} \right)$$

$\begin{cases} -|\alpha| & \alpha y_i \leq 0 \\ \infty & \alpha y_i > 0 \end{cases}$ 
 $\begin{cases} \alpha y_i & (\alpha + y_i) y_i \geq 0 \\ \infty & (\alpha + y_i) y_i < 0 \end{cases}$



$$l_i^*(\alpha) = \begin{cases} -|\alpha| & \text{if } -\alpha y_i \in [0, 1] \\ \infty & \text{else} \end{cases}$$



From above (or otherwise), we can get  $z^*$  such that:

$$\alpha z^* - \max(0, 1 - y_i z^*) = -|\alpha|$$

$$(-\alpha y_i \in [0, 1])$$

$$\rightarrow \alpha = 0, \text{ then } 1 - y_i z^* \leq 0$$

$$\rightarrow -\alpha y_i = 1, \text{ then } -y_i z^* - \max(0, 1 - y_i z^*) = -1$$

$$\Leftrightarrow 1 - y_i z^* \geq 0$$

$$\rightarrow \alpha \in (-1, 0), \text{ then } \alpha(z^* - 1) = \max(0, 1 - z^*)$$

$$\Leftrightarrow z^* = 1$$

$$\rightarrow \alpha \in (0, 1), \text{ then } \alpha(z^* + 1) = \max(0, 1 + z^*)$$

$$\Leftrightarrow z^* = -1$$

$$\rightarrow -\alpha y_i \in (0, 1), \text{ then } 1 - y_i z^* = 0$$



Our dual:  $\max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \alpha^T K \alpha - \frac{c}{m} \sum_{i=1}^m \ell_i^* \left( -\frac{m \alpha_i}{c} \right)$

s.t.  $\sum_{i=1}^m \alpha_i = 0$

$(K)_{ij} = \phi(x_i)^T \phi(x_j)$

Let's substitute  $\alpha_i y_i \rightarrow \alpha_i$

$= \max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \alpha^T G \alpha - \frac{c}{m} \sum_{i=1}^m \ell_i^* \left( -\frac{m y_i \alpha_i}{c} \right)$

s.t.  $\sum_{i=1}^m \alpha_i y_i = 0$

$\begin{cases} \frac{m \alpha_i}{c} & \alpha_i \in [0, \frac{c}{m}] \\ \infty & \text{else} \end{cases}$

$(G)_{ij} = y_i y_j \phi(x_i)^T \phi(x_j)$

for hinge-loss

$= \max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \alpha^T G \alpha - \mathbf{1}^T \alpha$

s.t.  $0 \leq \alpha \leq \frac{c}{m} \mathbf{1}, \quad y^T \alpha = 0$



## Optimality cond. for hinge-loss:

(i)  $\omega \in \mathbb{R}^n, b \in \mathbb{R}, \gamma_i \in \mathbb{R}, \gamma_i = \omega^\top \phi(x_i) - b \quad \forall i$  (Feasibility)

(ii)  $\omega = \sum_{i=1}^m \alpha_i \gamma_i \phi(x_i), \sum_{i=1}^m \alpha_i \gamma_i = 0$ , (gradient)

$$\alpha_i = 0 \Rightarrow \gamma_i (\omega^\top \phi(x_i) - b) \geq 1$$

$$\alpha_i = \frac{c}{m} \Rightarrow \gamma_i (\omega^\top \phi(x_i) - b) \leq 1 \rightarrow \text{bounded Support Vectors}$$

$$\alpha_i \in (0, \frac{c}{m}) \Rightarrow \gamma_i (\omega^\top \phi(x_i) - b) = 1 \rightarrow \text{non-bound Support Vectors}$$

## Co-ordinate descent to exploit sparsity

→ assume  $b=0$ , so that  $\sum \alpha_i \gamma_i = 0$  does not appear!

→ Update sequentially one  $\alpha_i$  at a time.

→ repeat till convergence.



→ 1-d convex QP  
⇒ closed form solution

Use duality gap.

→ Get  $\omega^*, b^*$  from  $\alpha$

→  $\textcircled{\text{I}} - \textcircled{\text{II}} = \text{duality gap} \downarrow 0$