

Lecture 6

Kernel Machines – Part 1

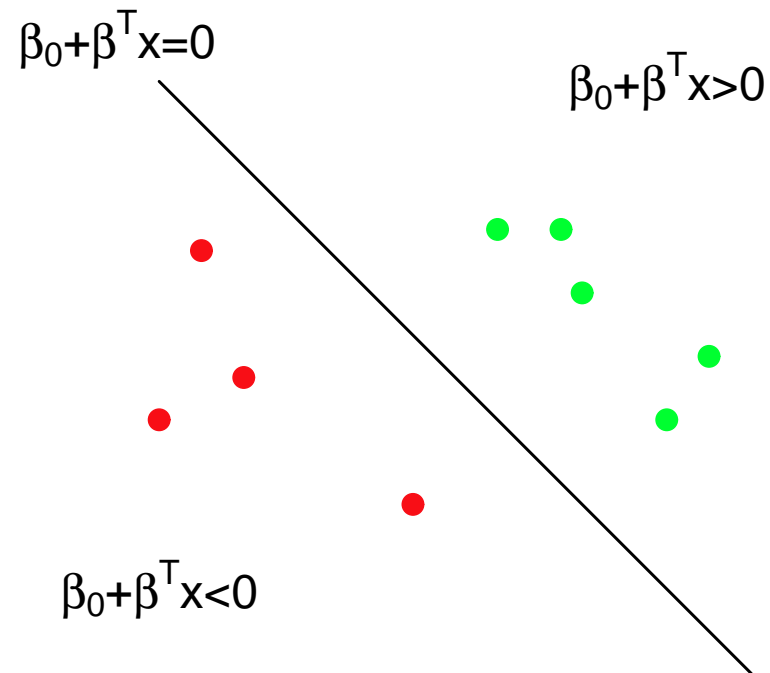
(February 13, 2015)

Mu Zhu
University of Waterloo

Outline

- Part 1
 - support vector machine (SVM)
- Part 2 (after short break)
 - kernel machines

Separating Hyperplane



$$y_i \in \{-1, +1\}$$

a **separating hyperplane** if $\exists c > 0$ s.t. $y_i(\beta_0 + \beta^T \mathbf{x}_i) \geq c \forall i$

Scaling

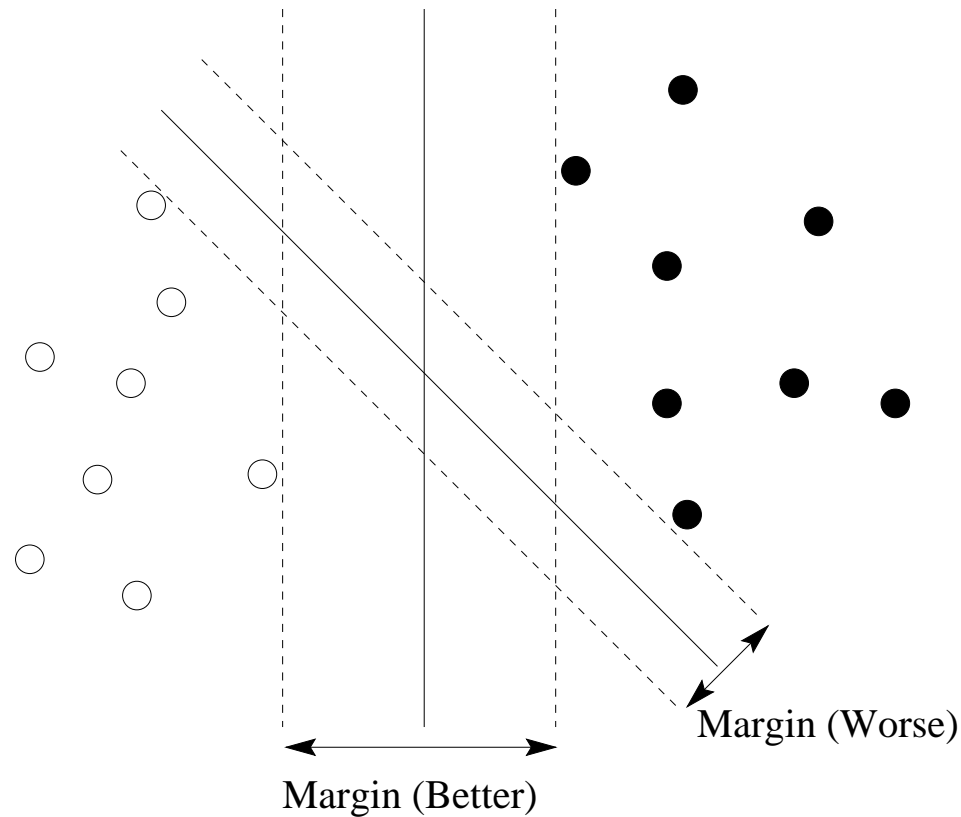
clearly, for any $s > 0$, $\beta_0 + \boldsymbol{\beta}^T \mathbf{x} = 0 \Leftrightarrow s(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}) = 0$

always choose s so that $c = 1$

call it a “canonical” separating hyperplane if

$$y_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) \geq 1 \quad \forall \quad i$$

Margin



if data are separable, \exists infinitely many separating hyperplanes
argument: the one with the largest **margin** is the best

Computing the Margin

(1) for any $\mathbf{x}_1, \mathbf{x}_2 \in \text{HP}$,

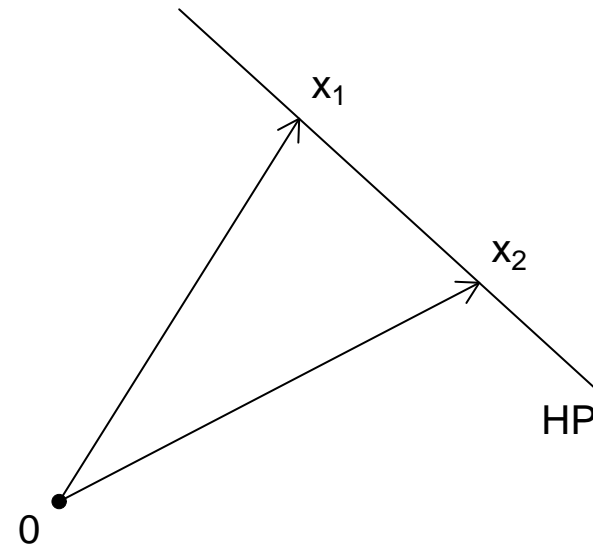
$$\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_1 = \beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_2 = 0,$$

so

$$\boldsymbol{\beta}^\top (\mathbf{x}_1 - \mathbf{x}_2) = 0 \quad \Rightarrow \quad \boldsymbol{\beta} \perp \text{HP}$$

(2) for any $\mathbf{x}_0 \in \text{HP}$,

$$\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_0 = 0 \quad \Rightarrow \quad \beta_0 = -\boldsymbol{\beta}^\top \mathbf{x}_0$$



Computing the Margin

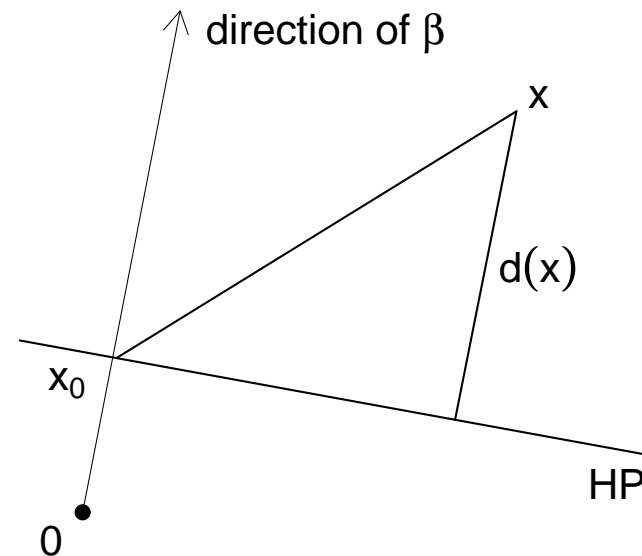
(3) for any \mathbf{x} , to compute

$d(\mathbf{x}) = \text{signed distance}$ to HP,

simply grab any $\mathbf{x}_0 \in \text{HP}$, and project $\mathbf{x} - \mathbf{x}_0$ onto direction of $\boldsymbol{\beta}$ (see figure), so

$$\begin{aligned}d(\mathbf{x}) &= \frac{\boldsymbol{\beta}^T}{\|\boldsymbol{\beta}\|} (\mathbf{x} - \mathbf{x}_0) \\ &= \frac{1}{\|\boldsymbol{\beta}\|} (\boldsymbol{\beta}^T \mathbf{x} + \beta_0)\end{aligned}$$

by (2) above



Computing the Margin

$$\begin{aligned}\text{margin} &= 2 \times \min_i [y_i d(\mathbf{x}_i)] \\ &= 2 \times \min_i \left[\frac{y_i (\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0)}{\|\boldsymbol{\beta}\|} \right] \\ &= \frac{2}{\|\boldsymbol{\beta}\|}\end{aligned}$$

since $\underbrace{y_i (\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0)}_{\text{separating HP}} \geq 1 \quad \forall \quad i$

Optimization Problem

$$\max_{\beta, \beta_0} \frac{2}{\|\beta\|} \Leftrightarrow \min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2$$

$$\text{s.t. } y_i(\beta^T \mathbf{x}_i + \beta_0) \geq 1 \quad \forall i$$

assumes classes are separable ... for now

extend to non-separable case later

Remark Going after maximum margin is not only intuitive, but also backed by rich **VC-theory** due to Vapnik & Chervonenkis.

Flavor of VC-Theory

- suppose [using hyperplanes (HP)]
 - all data $\mathbf{x} \in \mathbb{R}^d$ are confined within a ball of radius r
 - training data perfectly separated (ideal case) w/ margin Δ

- then, w/ probability $1 - \delta$, error (ε) on test data bounded by

$$\varepsilon \leq 4 \times \frac{h[1 + \log(2n/h)] - \log(\delta/4)}{n}$$

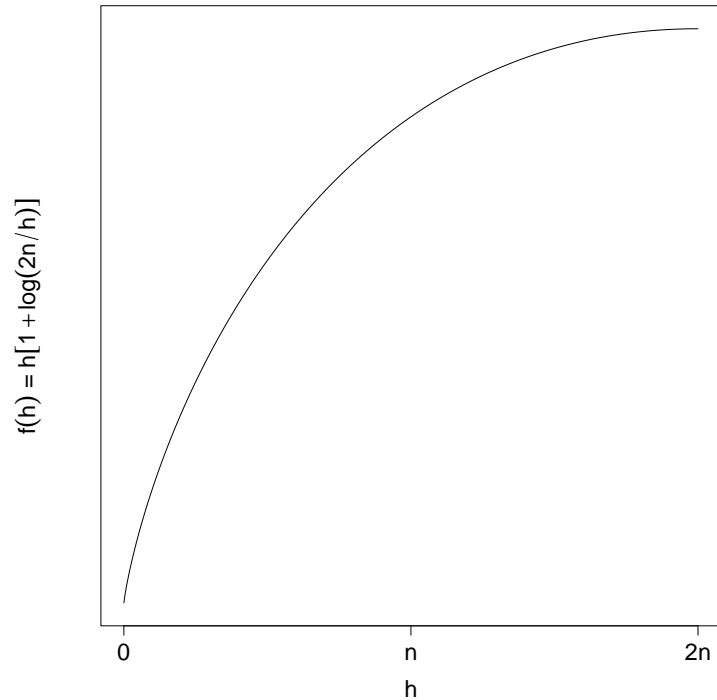
where

$$h \leq \min \left(\frac{r^2}{\Delta^2}, d \right) + 1$$

is the VC dimension of the separating HP

- basically, $\Delta \uparrow \Rightarrow h \downarrow \Rightarrow \varepsilon \downarrow \dots \exists$ many variations

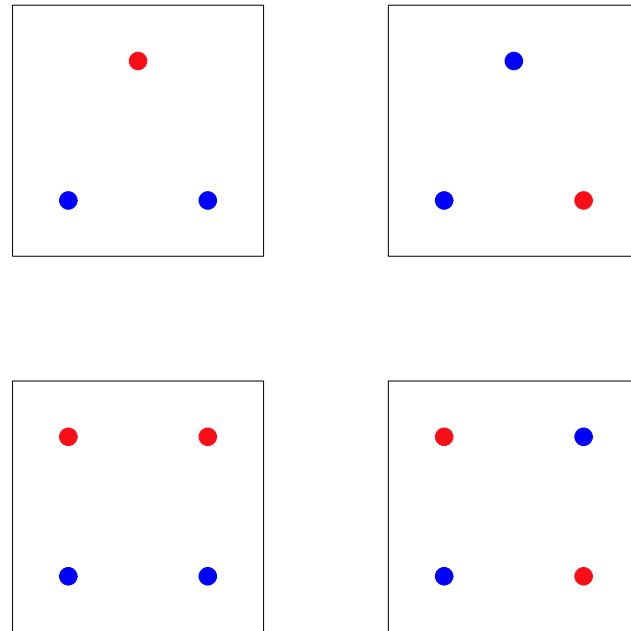
Some Details



Above: $f(h) = h[1 + \log(2n/h)]$ is an increasing function of h on $(0, 2n)$.

Below: VC dimension = maximum number of points a function can “shatter”, a measure of complexity.

In 2D, HPs have VC dim = 3.



Optimization

- Fermat (~ 1650): unconstrained optimization

$$\min_{\mathbf{x}} f(\mathbf{x})$$

- Lagrange (~ 1800): equality constraints

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad g_i(\mathbf{x}) = 0$$

- Kuhn-Tucker (~ 1950): inequality constraints

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \underbrace{g_i(\mathbf{x}) \geq 0}_{\text{main difficulty}} \quad \text{and} \quad \underbrace{h_j(\mathbf{x}) = 0}_{\text{treat "as usual" we'll ignore}}$$

The Primal Problem

- minimizing the Lagrangian function,

$$L(\mathbf{x}, \boldsymbol{\alpha}) \equiv f(\mathbf{x}) - \sum_i \alpha_i g_i(\mathbf{x}),$$

over \mathbf{x} , using $\alpha_i \geq 0$, gives lower bound on $f(\mathbf{x})$ for feasible \mathbf{x} :

$$\begin{aligned} L^*(\boldsymbol{\alpha}) &\equiv \inf_{\mathbf{x}} \left\{ f(\mathbf{x}) - \sum_i \alpha_i g_i(\mathbf{x}) \right\} \\ &\leq f(\mathbf{x}) - \sum_i \alpha_i g_i(\mathbf{x}) \quad [\text{having minimized over } \mathbf{x}] \\ &\leq f(\mathbf{x}) \quad [\alpha_i \geq 0 \text{ and } g_i(\mathbf{x}) \geq 0 \text{ for feasible } \mathbf{x}] \end{aligned}$$

- natural to ask for the best lower bound

The Dual Problem

- maximizing $L^*(\boldsymbol{\alpha})$ over $\alpha_i \geq 0$ gives the greatest lower bound
- under “good” conditions, can get

$$\max_{\alpha_i \geq 0 \ \forall i} L^*(\boldsymbol{\alpha}) \quad \Leftrightarrow \quad \min_{g_i(\boldsymbol{x}) \geq 0 \ \forall i} f(\boldsymbol{x}),$$

known as **strong duality**

- can see that, when this happens, must have

$$\alpha_i g(\boldsymbol{x}_i) = 0 \quad \forall \quad i,$$

known as **complementary slackness**

Exercise What’s the intuition behind complementary slackness?

Karush-Kuhn-Tucker Conditions

Let $f(\mathbf{x})$ be a continuously differentiable function for $\mathbf{x} \in \mathbb{R}^d$. Let $S = \{\mathbf{x} \in \mathbb{R}^d : g_i(\mathbf{x}) \geq 0\}$ be a set of **affine** constraints. Define

$$L \equiv f(\mathbf{x}) - \sum_i \alpha_i g_i(\mathbf{x})$$

to be the **Lagrangian**. If f has local minimum on S at \mathbf{x}^* , then

- (a) $\frac{\partial L}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^*} = 0 \quad \Rightarrow \quad f'(\mathbf{x}^*) = \sum_i \alpha_i g'_i(\mathbf{x}^*);$
- (b) $\alpha_i \geq 0$ for all i ;
- (c) $\alpha_i g_i(\mathbf{x}^*) = 0$ for all i .

If f is also (strictly) convex and (a)–(c) are satisfied, then \mathbf{x}^* is a (unique) global minimum.

The KKT Conditions

$$\min_{\boldsymbol{\beta}, \beta_0} \frac{1}{2} \|\boldsymbol{\beta}\|^2 \quad \text{s.t.} \quad y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) \geq 1 \quad \forall \quad i$$

$$L \equiv \frac{1}{2} \|\boldsymbol{\beta}\|^2 - \sum_i \alpha_i [y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1]$$

$$(a) \quad \frac{\partial L}{\partial \boldsymbol{\beta}} = 0 \Rightarrow \boldsymbol{\beta} = \sum_i \alpha_i y_i \mathbf{x}_i; \quad \frac{\partial L}{\partial \beta_0} = 0 \Rightarrow \sum_i \alpha_i y_i = 0$$

$$(b) \quad \alpha_i \geq 0 \quad \text{for all } i$$

$$(c) \quad \alpha_i [y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1] = 0 \quad \text{for all } i$$

plug (a) into L , and obtain $L^*(\boldsymbol{\alpha})$

The Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & L^*(\alpha) \\ &= \frac{1}{2} \|\beta\|^2 - \sum_i \alpha_i [y_i(\beta^T \mathbf{x}_i + \beta_0) - 1] \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \underbrace{\sum_i \alpha_i y_i = 0}_{\text{KKT (a)}} \quad \text{and} \quad \underbrace{\alpha_i \geq 0}_{\text{KKT (b)}} \end{aligned}$$

Quadratic Programming

- dual problem can be written as

$$\max_{\alpha} \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T D_y X X^T D_y \alpha \quad \text{s.t.} \quad \mathbf{y}^T \alpha = 0, \quad I \alpha \geq \mathbf{0}$$

$$\text{where } D_y = \begin{bmatrix} y_1 & - & - \\ - & \ddots & - \\ - & - & y_n \end{bmatrix}_{n \times n} \quad \text{and} \quad X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}_{n \times d}$$

- a standard [quadratic programming](#) problem of the form

$$\max_{\alpha} \mathbf{p}^T \alpha - \frac{1}{2} \alpha^T Q \alpha \quad \text{s.t.} \quad A \alpha = b, \quad C \alpha \geq d$$

“What’s in a name?”

“That which we call a **support vector machine** by any other name would sound as esoteric.”

- KKT(b)+(c) imply

$$\alpha_i > 0 \quad \Rightarrow \quad y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1 = 0$$

and

$$y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1 > 0 \quad \Rightarrow \quad \alpha_i = 0$$

- training observation \mathbf{x}_i called a **support vector** if $\alpha_i > 0$;

$$\boldsymbol{\beta} = \sum_i \alpha_i y_i \mathbf{x}_i$$

only affected by terms with $\alpha_i > 0$ — SVs are **on the margin**

SVM Steps

1. Solve dual problem (quadratic programming) for α_i .

2. Let

$$\boldsymbol{\beta} = \sum_i \alpha_i y_i \mathbf{x}_i.$$

3. Solve for β_0 . [How?]

4. Predict with

$$f(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{x}_i + \beta_0 = \beta_0 + \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}.$$

SVM Steps

1. Solve dual problem (quadratic programming) for α_i .

2. Let

$$\boldsymbol{\beta} = \sum_i \alpha_i y_i \mathbf{x}_i.$$

3. Solve for β_0 using any support vector [those with $\alpha_i > 0$, since all SVs satisfy $y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1 = 0$].

4. Predict with

$$f(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{x}_i + \beta_0 = \beta_0 + \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}.$$

Non-separable Case

$$\min_{\boldsymbol{\beta}, \beta_0} \frac{1}{2} \|\boldsymbol{\beta}\|^2 + \gamma \sum_i \xi_i$$

$$\text{s.t.} \quad y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) \geq 1 - \xi_i \quad \forall \quad i$$

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

[different formulations possible]

The KKT Conditions

$$L \equiv \frac{1}{2} \|\boldsymbol{\beta}\|^2 + \gamma \sum_i \xi_i - \sum_i \alpha_i [y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1 + \xi_i] - \sum_i \lambda_i \xi_i$$

$$(a) \quad \frac{\partial L}{\partial \boldsymbol{\beta}} = 0 \Rightarrow \boldsymbol{\beta} = \sum_i \alpha_i y_i \mathbf{x}_i; \quad \frac{\partial L}{\partial \beta_0} = 0 \Rightarrow \sum_i \alpha_i y_i = 0;$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \gamma - \alpha_i - \lambda_i = 0 \quad \text{or} \quad \alpha_i = \gamma - \lambda_i$$

$$(b) \quad \alpha_i \geq 0 \quad \text{for all } i; \quad \lambda_i \geq 0 \quad \text{for all } i$$

$$(c) \quad \alpha_i [y_i(\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0) - 1] = 0 \quad \text{for all } i; \\ \lambda_i \xi_i = 0 \quad \text{for all } i$$

Two Typs of Support Vectors

- extra KKT (a)+(b)+(c):

$$\alpha_i = \gamma - \lambda_i; \quad \lambda_i \geq 0; \quad \lambda_i \xi_i = 0$$

- together, they imply

$$\alpha_i < \gamma \quad \Rightarrow \quad \lambda_i > 0 \quad \Rightarrow \quad \xi_i = 0$$

and

$$\xi_i > 0 \quad \Rightarrow \quad \lambda_i = 0 \quad \Rightarrow \quad \alpha_i = \gamma$$

- support vectors

either on the margin $(\xi_i = 0, \alpha_i < \gamma)$

or over the margin $(\xi_i > 0, \alpha_i = \gamma)$

Non-separable Case

Exercises

- (a) Show that, in this particular formulation of the non-separable case, the **dual** remains the same as in the separable case, except for an additional set of constraints that $\alpha_i \leq \gamma$ for all i .
- (b) Specify the corresponding **dual** if the non-separable case is formulated as

$$\min_{\beta, \beta_0} \|\beta\|^2 + \gamma \sum_i \xi_i^2$$

$$\text{s.t.} \quad y_i(\beta^T \mathbf{x}_i + \beta_0) \geq 1 - \xi_i \quad \forall \quad i$$

instead. How to solve the corresponding **dual**?

- (c) Solve for β_0 in the non-separable case. [both formulations]

A Different Point of View

$$\begin{aligned} & \min_{\boldsymbol{\beta}, \beta_0} \quad \frac{1}{2} \|\boldsymbol{\beta}\|^2 + \gamma \sum_i \xi_i \\ \text{s.t.} \quad & y_i(\boldsymbol{\beta}^\top \mathbf{x}_i + \beta_0) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0 \quad \forall \quad i \end{aligned}$$

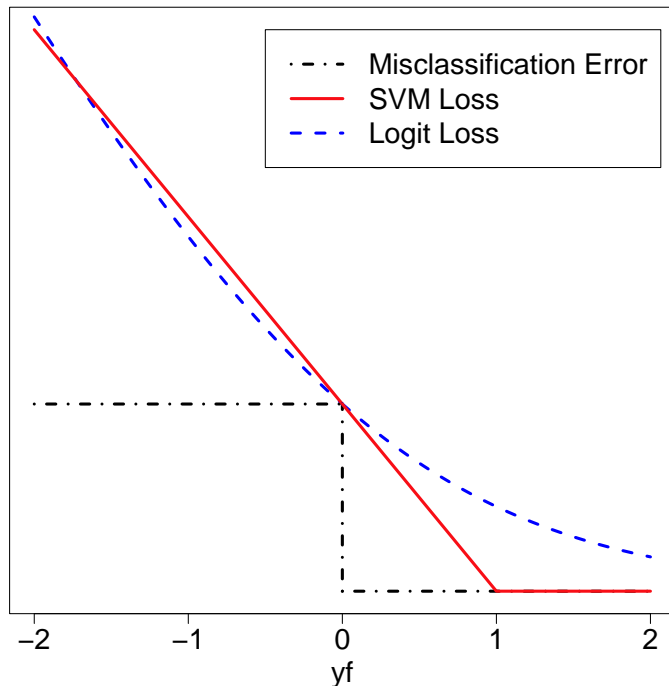
\Updownarrow

$$\min_{\boldsymbol{\beta}, \beta_0} \quad \underbrace{\sum_{i=1}^n [1 - y_i(\boldsymbol{\beta}^\top \mathbf{x}_i + \beta_0)]_+}_{\text{loss}} + \underbrace{\lambda \|\boldsymbol{\beta}\|^2}_{\text{penalty}}$$

$$\text{where } [z]_+ = \begin{cases} z, & z > 0; \\ 0, & z \leq 0 \end{cases}$$

The “Hinge” Loss

$L(y, f)$ as a function of yf



$$L[y, f(\mathbf{x})]$$

$$\text{Error : } I[yf(\mathbf{x}) < 1]$$

$$\text{Logit : } \log_2[1 + e^{-yf(\mathbf{x})}]$$

$$\text{SVM : } [1 - yf(\mathbf{x})]_+$$

in a way, SVM not too different from [ridge logistic regression](#)

Summary

- Part 1
 - separating hyperplanes; margin; VC-theory
 - optimization; primal; dual; KKT conditions
 - quadratic programming
 - support vector machine (SVM); “hinge” loss
- Part 2 (next, after short break)
 - kernel SVM; reproducing kernel Hilbert space
 - kernel ridge regression
 - kernel canonical correlation analysis