

## 1. Support Vector Machines

- a) Optimal separating hyperplane
- b) SVMs
- c) Interpretation as penalized ERM with hinge loss
- d) Dual Formulation
- e) Kernel trick
- f) Derivation of dual formulation

### Unsupervised Learning:

- 2. K-Means
- 3. EM for Gaussian Mixtures (studying 3 is optional)

Vapnik & Chervonenkis: foundational work on statistical learning starting in sixties and seventies, leading to SVMs ~1990.

Vapnik, 1998: "solve the problem directly and never solve a more general problem as an intermediate step"

## 1. SVMs

### a) Optimal separating hyperplane

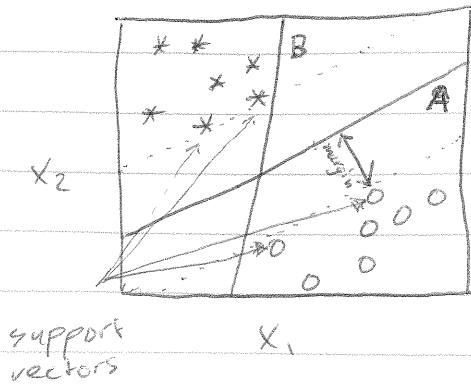
estimate  
 generative:  $P(x, y)$   
 discriminative:  $P(y|x)$   
 now: decision boundary directly

Assume: 2 classes, linearly separable

$$\begin{pmatrix} y_1 \\ x_1 \end{pmatrix}, \dots, \begin{pmatrix} y_N \\ x_N \end{pmatrix} \quad y \in \{-1, +1\}$$

Linear

Model: classifiers that compute  $x^T \beta + \beta_0$  ← writing intercept separately and return its sign



Which decision boundary do you like best? A or B?

A has larger "margin" = distance to nearest point

Fig. 12.1, left

"Optimal" separating hyperplane maximizes the margin.

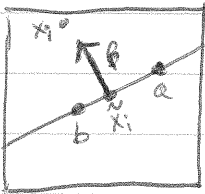
support vectors: points on the margin

decision boundary:  $\{x : x^T \beta + \beta_0 = 0\}$

Let  $f(x) = x^T \beta + \beta_0$ . Then distance of  $x_i$  to decision boundary is  $\frac{|f(x_i)|}{\|\beta\|}$

margin:  $\frac{y_i \cdot f(x_i)}{\|\beta\|}$  is distance if  $x_i$  classified correctly  
 is negative distance if classified incorrectly

Optional: for any  $a, b$  on decision boundary:



$(b-a)^T \beta = 0$ , so  $\frac{\beta}{\|\beta\|}$  is normal vector to the decision boundary of length 1.

suppose  $\hat{x}_i$  is projection of  $x_i$  onto decision boundary. Then

$$x_i = \hat{x}_i + \alpha \cdot \frac{\beta}{\|\beta\|} \quad |\alpha| \text{ is distance of } x_i \text{ from decision boundary}$$

$$\alpha \frac{\beta}{\|\beta\|} = x_i - \hat{x}_i$$

$$\alpha \cdot \|\beta\| = (x_i - \hat{x}_i)^T \beta \quad (\text{because } \beta^T \beta = \|\beta\|^2)$$

$$\alpha \cdot \|\beta\| = x_i^T \beta - \hat{x}_i^T \beta$$

$$= x_i^T \beta + \beta_0 \quad (\text{because } \hat{x}_i \text{ on decision boundary})$$

$$= f(x_i)$$

$$\alpha = \frac{f(x_i)}{\|\beta\|}$$

3

to maximize the margin  $M$ :

$$\max_{\beta, \beta_0, M} M$$

$$\text{subject to: } \frac{y_i (x_i^T \beta + \beta_0)}{\|\beta\|} \geq M \quad \text{for } i=1, \dots, N$$

$\uparrow$   
margin of  $(x_i)$

Same decision boundary and margin if we multiply  $\beta, \beta_0$  by a constant, so can always choose this constant such that  $\|\beta\| = \frac{1}{M}$ :

$$\max_{\beta, \beta_0, M} M$$

$$\|\beta\| = \frac{1}{M}$$

$$\text{s.t. } \frac{y_i (x_i^T \beta + \beta_0)}{\|\beta\|} \geq M \cdot \|\beta\| \quad \forall i$$

$$\max_{\beta, \beta_0} \frac{1}{\|\beta\|}$$

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

solution achieved by same parameters  $\beta, \beta_0$

convex function  $\rightarrow \min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2$

linear inequality constraints  $\rightarrow \text{s.t. } y_i (x_i^T \beta + \beta_0) \geq 1 \quad \forall i$

b) SVMs

What if classes are not linearly separable?

Greek letter "xi"

Fig. 12.1, right: introduce slack variable  $\xi_i \geq 0$  for each data point  $i = 1, \dots, N$

$$\begin{aligned} & \max_{\beta, \beta_0, M, \xi_i} M \\ & \text{subject to: } \xi_i \geq 0, \quad y_i (x_i^T \beta + \beta_0) \geq M(1 - \xi_i) \quad i = 1, \dots, N \\ & \quad \quad \quad \sum_{i=1}^N \xi_i \leq t \end{aligned}$$

parameter of alg. ← Assume we take large enough to have a solution

support vectors: all points inside the margin

$$\begin{aligned} & \min_{\beta, \beta_0, \xi_i} \frac{1}{2} \|\beta\|^2 \\ & \text{s.t. } \xi_i \geq 0, \quad y_i (x_i^T \beta + \beta_0) \geq 1 - \xi_i \quad i = 1, \dots, N \\ & \quad \quad \quad \sum_{i=1}^N \xi_i \leq t \end{aligned}$$

equivalent

$$\min_{\beta, \beta_0, \xi_i} \frac{1}{2} \|\beta\|^2 + C \cdot \sum_{i=1}^N \xi_i \quad (*)$$

Parameter

$$\text{s.t. } \xi_i \geq 0, \quad y_i (x_i^T \beta + \beta_0) \geq 1 - \xi_i \quad i = 1, \dots, N$$

c) Interpretation as penalized ERM

see slides 4

$$L(y_i, f(x_i)) = \max \{0, 1 - y_i \cdot f(x_i)\} \text{ is "hinge loss"}$$

$$\min_{\beta, \beta_0, \xi_i} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{s.t. } \xi_i \geq 0, \quad \xi_i \geq 1 - y_i (x_i^T \beta + \beta_0) \quad i = 1, \dots, N$$

$$\xi_i \geq L(y_i, x_i^T \beta + \beta_0)$$

$$\min_{\beta, \beta_0} \sum_{i=1}^N L(y_i, x_i^T \beta + \beta_0) + \frac{1}{2-C} \|\beta\|^2$$

ERM for hinge loss with  $L_2$ -penalty,  $\lambda = \frac{1}{2C}$ .

d) Dual Formulation

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k \langle x_i, x_k \rangle$$

$\langle x_i, x_k \rangle = x_i^T x_k$

$$\text{subject to } 0 \leq \alpha_i \leq C \quad \text{and} \quad \sum_{i=1}^N \alpha_i y_i = 0$$

Then solution to (\*) is:

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i \quad \leftarrow \begin{array}{l} \text{very different, but} \\ \text{reminiscent of nearest neighbor, because} \\ \text{also defined in terms of training data} \end{array}$$

$$\hat{\alpha}_i = 0 \quad \text{for } x_i \text{ outside margin}$$

$$0 < \hat{\alpha}_i < C \quad \text{for } x_i \text{ on margin}$$

$$\hat{\alpha}_i = C \quad \text{for } x_i \text{ inside margin}$$

$$\text{solve } \hat{\beta}_0 \text{ from } \hat{\alpha}_i [y_i (x_i^T \hat{\beta} + \hat{\beta}_0) - 1] = 0$$

for any  $i$  s.t.  $0 < \hat{\alpha}_i < C$

e) Kernel trick

Fig. 2.5: how to learn something like this with linear classifier?

map features  $x$  to a larger set of features  $h(x)$ !  
e.g.

$$h \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ x_1 \cdot x_2 \\ x_1^2 \\ x_2^2 \end{pmatrix}$$

Then  $\langle x_i, x_k \rangle$  in dual formulation becomes  $\langle h(x_i), h(x_k) \rangle$

kernel trick: don't need to specify  $h$ ,  
only need to know the kernel function

really nice if this were a simple function... so turn things around and start by defining  $K(x_i, x_k)$   $\rightarrow K(x_i, x_k) = \langle h(x_i), h(x_k) \rangle$  ← measure of similarity, larger if  $x_i, x_k$  more similar.  
 $h(x)$  may even be infinite-dimensional!

Examples:  $K(x, x') = (1 + \langle x, x' \rangle)^d$  :  $d$ -th degree polynomial

↳ e.g.  $d=2, x \in \mathbb{R}^2$ :  
$$K(x, x') = (1 + x_1 x'_1 + x_2 x'_2)^2 = \langle h(x), h(x') \rangle$$

for 
$$h(x) = \begin{pmatrix} 1 \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \end{pmatrix}$$

radial basis:  $K(x, x') = e^{-\gamma \|x - x'\|^2}$

neural network:  $K(x, x') = \tanh(a \langle x, x' \rangle + b)$

If  $K$  satisfies certain technical conditions (symmetric, positive definite) then there always exists some mapping  $h$  s.t.  $K(x, x') = \langle h(x), h(x') \rangle$ .  
 $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

Classifying a new  $x$ :

$$\begin{aligned} \hat{f}(x) &= h(x)^T \hat{\beta} + \hat{\beta}_0 = h(x)^T \sum_{i=1}^N \alpha_i y_i h(x_i) + \hat{\beta}_0 \\ &= \sum_{i=1}^N \alpha_i y_i \langle h(x), h(x_i) \rangle + \hat{\beta}_0 \\ &= \sum_{i=1}^N \alpha_i y_i K(x, x_i) + \hat{\beta}_0 \end{aligned}$$

Again no need to specify  $h$ ; only need to know kernel.

Fig. 12.3

f.) Derivation of Dual Formulation ← optional!

$$\begin{aligned} \min_{\beta, \beta_0, \xi_i} \quad & \frac{1}{2} \|\beta\|^2 + C \cdot \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \xi_i \geq 0, \quad y_i (x_i^T \beta + \beta_0) - (1 - \xi_i) \geq 0 \quad i=1, \dots, N \end{aligned}$$

$$\min_{\beta, \beta_0, \xi_i} \quad \sup_{\alpha_i, p_i \geq 0} \quad \underbrace{\frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (x_i^T \beta + \beta_0) - (1 - \xi_i)) - \sum_{i=1}^N p_i \xi_i}_A$$

(If  $\beta, \beta_0, \xi_i$  violate constraints, then  $\alpha_i$  or  $p_i$  becomes  $\infty$  and  $A = \infty$ , so  $\beta, \beta_0, \xi_i$  only achieve minimum while satisfying constraints. And then  $\alpha_i, p_i$  become 0, so the constraints ~~with~~ drop away and we are minimizing the previous objective.)

~~$$\nabla_{\beta} A = 0 \Rightarrow \beta = \sum_{i=1}^N \alpha_i y_i x_i \quad \nabla_{\beta_0} A = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$~~

~~$$\nabla_{\xi_i} A = 0 \Rightarrow p_i = C \alpha_i$$~~

plugging

⑧

$$\min_{\beta, \beta_0, \xi_i} \sup_{\alpha_i, p_i \geq 0} A = \sup_{\alpha_i, p_i \geq 0} \underbrace{\min_{\beta, \beta_0, \xi_i} A}_{\text{can solve this}} \quad \text{by convex optimization theory (Slater's condition)}$$

$$\nabla_{\beta} A = 0 \Rightarrow \beta = \sum_{i=1}^N \alpha_i y_i x_i \quad \nabla_{\beta_0} A = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$

$$\nabla_{\xi_i} A = 0 \Rightarrow p_i = C - \alpha_i$$

Plugging these in gives dual formulation.



# Unsupervised Learning

## 4.3 Intro

Supervised:  $T = \left( \begin{array}{c} y_1 \\ x_1 \end{array} \right), \dots, \left( \begin{array}{c} y_N \\ x_N \end{array} \right)$  ← can evaluate every method by EPE

Unsupervised:  $T = \{x_1, \dots, x_N\}$

What can we do?

Clustering: split data into groups of points (clusters) that are similar.

- ① What do you mean by "similar"?
- ② What kind of clusters are you looking for?

## K-means

① Squared Euclidean distance between  $x_i, x_{i'}$ :

$$d(x_i, x_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

② Find  $K$  clusters s.t. distance to mean within each cluster is small.

$c(i) \in \{1, \dots, K\}$  is cluster assigned to  $x_i$ .

$$\text{minimize } \sum_{k=1}^K \sum_{i: c(i)=k} d(x_i, \mu_k) \quad (*)$$

where  $\mu_k$  = mean of points in cluster  $k$

NB. Mistake in book! (multiplies by  $N_k$  in each cluster)

### Algorithm:

1. Initialize cluster assignment  $C$
2. Set  $\mu_k$  to current mean in cluster  $k$
3. Given means  $\mu_1, \dots, \mu_K$ , assign each  $x_i$  to nearest mean to get new  $C$
4. repeat from 2 until no changes in  $C$ .

- Converges to local minimum

- ~~often~~ <sup>may</sup> choose  $K$  s.t. increasing  $K$  does not reduce  $(*)$  very much.

### Gaussian Mixtures and EM

- like k-means, but with "soft" cluster assignments

- for simplicity: explain for two clusters

each cluster is a Gaussian modelled by

cluster 1:  $X \sim \mathcal{N}(\mu_1, \sigma_1^2)$  w.p.  $1-\pi$

cluster 2:  $X \sim \mathcal{N}(\mu_2, \sigma_2^2)$  w.p.  $\pi$

probability density:  $(1-\pi)\phi_{\mu_1, \sigma_1^2}(x) + \pi\phi_{\mu_2, \sigma_2^2}(x)$

To find parameters:

ML is difficult numerically

+ gives bad solution with  $\hat{\sigma}_1 = 0$

$\hat{\mu}_1 = x_i$  for some  $i$ .

Solution: algorithm very similar to K-means

Hidden variables  $\Delta_i = \begin{cases} 0 & \text{if } x_i \text{ in cluster 1} \\ 1 & \text{if } x_i \text{ in cluster 2} \end{cases}$

1. Initializes parameters  $\hat{\pi}, \hat{\mu}_1, \hat{\sigma}_1, \hat{\mu}_2, \hat{\sigma}_2$  ( $\hat{\sigma}_1 > 0.5, \hat{\sigma}_2 > 0.5$ )

2. Soft cluster assignment:

$$\hat{y}_i = E[\Delta_i | \hat{\theta}, T] = \frac{\hat{\pi} \phi_{\hat{\mu}_2, \hat{\sigma}_2}(x_i)}{\hat{\pi} \phi_{\hat{\mu}_2, \hat{\sigma}_2}(x_i) + (1 - \hat{\pi}) \phi_{\hat{\mu}_1, \hat{\sigma}_1}(x_i)}$$

3. Update parameter estimates for clusters:

$$\hat{\mu}_1 = \frac{\sum_{i=1}^N (1 - \hat{y}_i) x_i}{\sum_i (1 - \hat{y}_i)} \quad \hat{\mu}_2 = \frac{\sum_{i=1}^N \hat{y}_i x_i}{\sum_i \hat{y}_i}$$

$$\hat{\sigma}_1^2 = \frac{\sum_{i=1}^N (1 - \hat{y}_i) (x_i - \hat{\mu}_1)^2}{\sum_i (1 - \hat{y}_i)} \quad \hat{\sigma}_2^2 = \frac{\sum_{i=1}^N \hat{y}_i (x_i - \hat{\mu}_2)^2}{\sum_i \hat{y}_i}$$

$$\hat{\pi} = \frac{\sum_{i=1}^N \hat{y}_i}{\sum_{i=1}^N \hat{y}_i + \sum_{i=1}^N (1 - \hat{y}_i)} = \frac{\sum_{i=1}^N \hat{y}_i}{N}$$

4. repeat from 2 until convergence

- converges to local minimum.