

Course ID 80250993 Dates: 9/16-12/23/2021 @ I-205
Tencent/VooV Meeting 521 4678 6257 (passwd 1205)



Lecture 6

Support Vector Machines

Xuegong Zhang
Oct. 14, 2020



Xuegong Zhang



6.1

Large Margin and Optimal Hyperplane

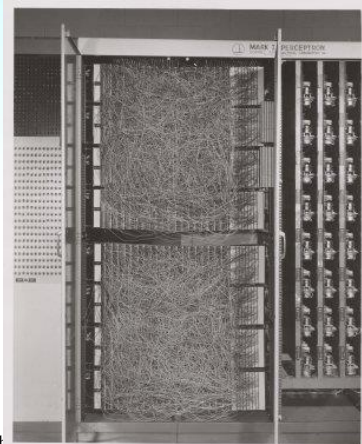
Xuegong Zhang

2

Perceptron

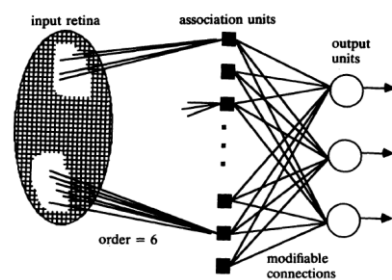
Frank Rosenblatt, *The Perceptron – a perceiving and recognizing automaton*,
Report 85-460-1, Cornell Aeronautical Laboratory, Jan. 1957

$$y = \text{sgn}\left(\sum_{i=1}^d w_i x_i + w_0\right)$$



Xuegong Zi

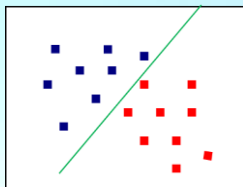
<https://en.wikipedia.org/wiki/Perceptron>



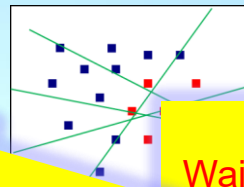
M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996

3

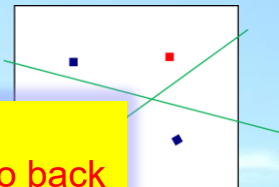
Linearly non-separable cases



Linearly separable



Not perfectly separable
but can achieve a
reasonable separation
with a linear classifier.



achievable linear
approximation.

Wait, let's go back
to consider the
easiest case.

Xuegong Zhang

4



Perceptron

- To train the machine $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$ with data $\{(x_1, y_1), \dots, (x_N, y_N)\}$.
- Goal: to optimize an objective function $J_P(\alpha) = \sum_{y_j \in \mathcal{Y}^k} (-\alpha^T y_j)$

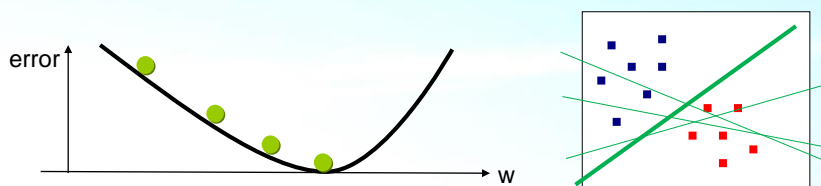


– Learning algorithm

- Basic Gradient Descent

$$\alpha(k+1) = \alpha(k) - \rho_k \nabla J$$

Any question here?



Xuegong Zhang

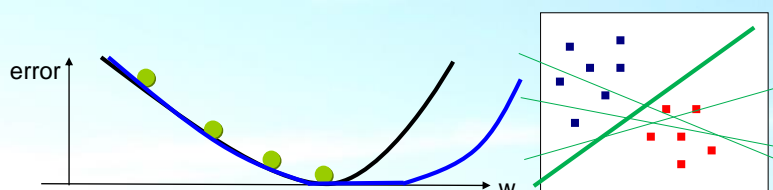
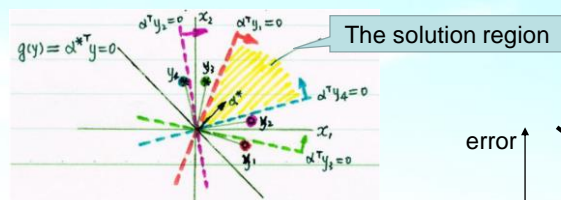
5

Question of Perceptron



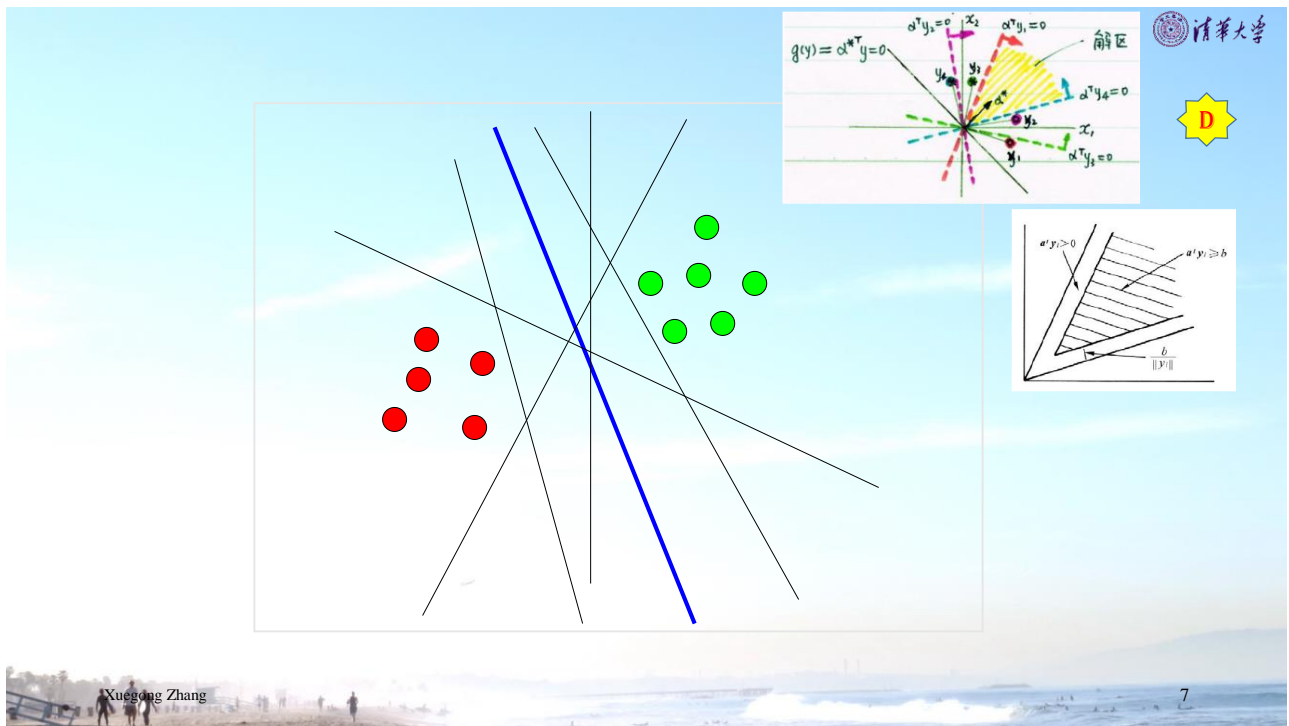
$$g(x) = w^T x + w_0 = \alpha^T y \quad \min J_P(\alpha) = \sum_{y_j \in \mathcal{Y}^k} (-\alpha^T y_j)$$

- There are multiple solutions, which one is the best?



Xuegong Zhang

6



- How can we get a "best" linear classifier?

- What is "the best" when no one is perfect?
 - A basic principle that differentiates different methods.
- What is "the best" when everyone looks perfect?

Top right corner: Tsinghua University logo and name in Chinese (清华大学).

Bottom left corner: Xuegong Zhang.

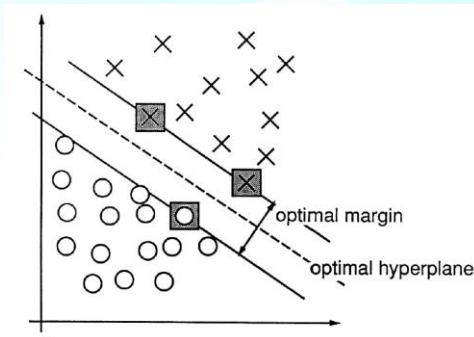
Bottom right corner: 8.

The Optimal Hyperplane



Vapnik

For sample set $(x_1, y_1), \dots, (x_l, y_l)$, $x \in R^d$, $y \in \{+1, -1\}$ that can be separated with a hyperplane $(w \cdot x) + b = 0$, **the optimal hyperplane** is defined as the linear decision function (hyperplane) $f(x) = \text{sgn}((w \cdot x) + b)$ with **maximal margin** between the vectors (samples) of the two classes among all hyperplanes that **separate the two classes without error**.



Separation:

$$(w \cdot x_i) + b > 0, \quad \text{if } y_i = 1$$

$$(w \cdot x_i) + b < 0, \quad \text{if } y_i = -1$$

$$y_i(w \cdot x_i + b) > 0, \quad i = 1, \dots, l$$

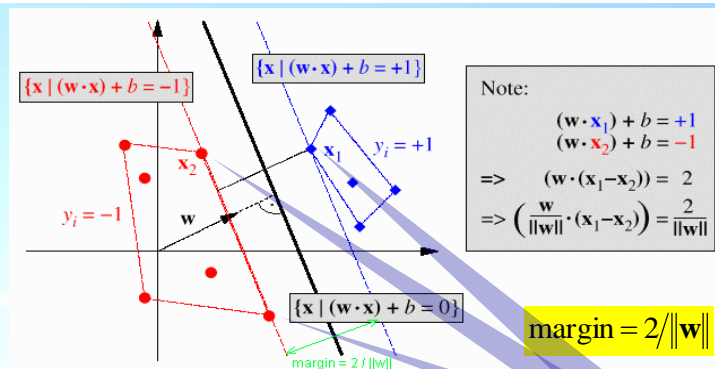


To fix the scale, we require

$$y_i(w \cdot x_i + b) \geq 1, \quad i = 1, \dots, l$$

Xuegong Zhang

9



• Normalization

$$(w \cdot x_i) + b \geq 1, \quad \text{if } y_i = 1$$

$$(w \cdot x_i) + b \leq -1, \quad \text{if } y_i = -1$$

$$\text{i.e., } y_i[(w \cdot x_i) + b] \geq 1, \quad i = 1, \dots, l$$

Support Vectors
(the ones with “=”)

-- The canonical form of the separating hyperplane

Xuegong Zhang

10

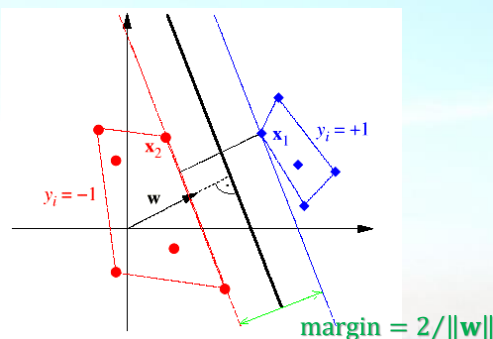


The Optimal Hyperplane

$$\min \Phi(\mathbf{w}) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) \text{ w.r.t. } \mathbf{w}$$

$$\text{s.t. } y_i[(\mathbf{x}_i \cdot \mathbf{w}) + b] \geq 1, \quad i = 1, 2, \dots, l$$

for training samples $(y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l), y \in \{-1, 1\}$



Xuegong Zhang

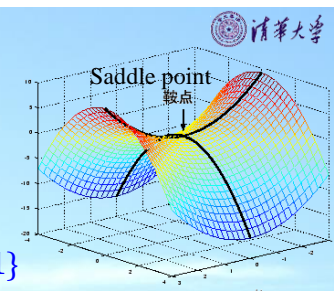
11

The solution: the **saddle point** of the Lagrangian

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) + b]y_i - 1\}$$

i.e.,

$$\min_{\mathbf{w}, b} \max_{\alpha} L(\mathbf{w}, b, \alpha) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) + b]y_i - 1\}$$



At the saddle point, we have

$$\frac{\partial L(\mathbf{w}_0, b_0, \alpha^0)}{\partial b} = 0 \quad \longrightarrow \quad \sum_{i=1}^l \alpha_i^0 y_i = 0$$

$$\frac{\partial L(\mathbf{w}_0, b_0, \alpha^0)}{\partial \mathbf{w}} = 0 \quad \longrightarrow \quad \mathbf{w}_0 = \sum_{i=1}^l y_i \alpha_i^0 \mathbf{x}_i$$

H

R

Xuegong Zhang

12

$$\min_{\mathbf{w}, b} \max_{\alpha} L(\mathbf{w}, b, \alpha) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) + b]y_i - 1\}$$



We have

(1) For the optimal hyperplane, α_i^0 must meet

$$\sum_{i=1}^l \alpha_i^0 y_i = 0, \quad \alpha_i^0 \geq 0, \quad i = 1, \dots, l$$

(2) and \mathbf{w}_0 must be the linear combination of training samples as:

$$\mathbf{w}_0 = \sum_{i=1}^l y_i \alpha_i^0 \mathbf{x}_i, \quad \alpha_i^0 \geq 0, \quad i = 1, \dots, l$$

Xuegong Zhang

13

$$\min_{\mathbf{w}, b} \max_{\alpha} L(\mathbf{w}, b, \alpha) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) + b]y_i - 1\}$$



(3) Only the support vectors have non-zero coefficients α_i^0 in \mathbf{w}_0 .

This is because, according to the Kuhn-Tucker theorem (Karush-Kuhn-Tucker (KKT) conditions), at the saddle point, the following equalities hold:

$$\alpha_i^0 \{[(\mathbf{x}_i \cdot \mathbf{w}_0) + b_0]y_i - 1\} = 0, \quad i = 1, \dots, l$$

Consider $y_i[(\mathbf{x}_i \cdot \mathbf{w}) + b] \geq 1$, $i = 1, 2, \dots, l$, $\alpha_i \neq 0$ only for cases where $y_i[(\mathbf{x}_i \cdot \mathbf{w}) + b] = 1$, which we call “support vectors” or SVs.

• The solution:

$$\mathbf{w}_0 = \sum_{SV_s} y_i \alpha_i^0 \mathbf{x}_i, \quad \alpha_i^0 \geq 0$$

Xuegong Zhang

14

Substituting $\mathbf{w}_0 = \sum_{i=1}^l y_i \alpha_i^0 \mathbf{x}_i$ into the Lagrangian, we get

$$W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

- **The dual problem** of the optimal hyperplane

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

s.t

$$\alpha_i \geq 0, \quad i = 1, \dots, l$$

and

$$\sum_{i=1}^l \alpha_i y_i = 0$$

- The original problem is called **the Primal Problem**.

$$\min \Phi(\mathbf{w}) = \frac{1}{2} (\mathbf{w} \cdot \mathbf{w}) \text{ w.r.t. } \mathbf{w}$$

$$\text{s.t.} \quad y_i [(\mathbf{x}_i \cdot \mathbf{w}) + b] \geq 1, \quad i = 1, 2, \dots, l$$

$$\text{for training samples } (y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l), \quad y \in \{-1, 1\}$$

Xuegong Zhang

15

At the solution $\alpha_0 = (\alpha_1^0, \dots, \alpha_l^0)$, we have

$$\|\mathbf{w}_0\|^2 = 2W(\alpha_0) = \sum_{\text{SVs}} \alpha_i^0 \alpha_j^0 (\mathbf{x}_i \cdot \mathbf{x}_j) y_i y_j = \sum_{i=1}^l \alpha_i^0$$

and the decision function is

$$f(\mathbf{x}) = \text{sign} \left(\sum_{\text{SVs}} y_i \alpha_i^0 (\mathbf{x}_i \cdot \mathbf{x}) + b_0 \right)$$

The threshold b_0 can be obtained from support vectors of the two classes:

$$b_0 = \frac{1}{2} [(\mathbf{w}_0 \cdot \mathbf{x}^*(1)) + (\mathbf{w}_0 \cdot \mathbf{x}^*(-1))]$$

Xuegong Zhang Tsinghua University

16

1-minute break

Any questions?



Xuegong Zhang

17

Questions

- Why can we call it “optimal”?
- What if the samples are not separable?
- Can we build optimal nonlinear machines like this?

Xuegong Zhang

18



6.2 Generalized Optimal Hyperplane

Xuegong Zhang

19

Non-separable Cases



For l samples of two classes: $(y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l)$, $y \in \{-1, 1\}$

$$\min \Phi(\mathbf{w}) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) \quad \text{w.r.t. } \mathbf{w}$$

$$\text{s.t.} \quad y_i[(\mathbf{x}_i \cdot \mathbf{w}) + b] \geq 1, \quad i = 1, 2, \dots, l$$

- Not linearly separable \leftrightarrow the inequalities cannot be met by all samples
- Introducing a slack variable $\xi_i \geq 0$, we define new constraints

$$y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, l$$

- Function $F_\sigma(\xi) = \sum_{i=1}^l \xi_i^\sigma$, $\sigma > 0$ can reflect the severity the original constraints are violated

Xuegong Zhang

20



Soft-margin Optimal Hyperplane (Generalized optimal hyperplane)

$$\begin{aligned} \min \quad & \Phi(\mathbf{w}, \xi) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \left(\sum_{i=1}^l \xi_i \right) \quad \text{w.r.t. } \mathbf{w} \\ \text{s.t.} \quad & y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, l \\ & \text{where parameter } C \text{ controls the penalty on errors.} \end{aligned}$$

D

The solution can be obtained with the same Lagrange optimization technology.

H

Xuegong Zhang

21

The dual problem of the generalized optimal hyperplane

Weights of the generalized optimal hyperplane are

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i$$

where parameters $\alpha_i, i = 1, \dots, l$ are the solution of the quadratic optimization problem

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

s.t.

$$\begin{aligned} & \sum_{i=1}^l y_i \alpha_i = 0, \\ & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l \end{aligned}$$

H

Xuegong Zhang

22



Similar to the linearly-separable cases, at the solution, we have

$$\alpha_i^0 \{[(\mathbf{x}_i \cdot \mathbf{w}_0) + b_0]y_i - 1 + \xi_i\} = 0, \quad i = 1, \dots, l$$

Only some α_i are non-zero. They corresponds to those that the equality holds in

$$y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, l$$

- The SVs also include the wrongly classified samples and their $\alpha_i = C$.

Xuegong Zhang

23

Summary on the generalized optimal hyperplane



The Primal Problem:

$$\begin{aligned} \min \quad & \psi(\mathbf{w}, \xi) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \left(\sum_{i=1}^n \xi_i \right) \\ \text{s.t.} \quad & y_i[(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 + \xi_i \geq 0, \quad \xi_i \geq 0, \quad i = 1, \dots, l \end{aligned}$$

The Dual Problem:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^l y_i \alpha_i = 0, \end{aligned}$$

and

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, l$$

The Solution:

$$f(\mathbf{x}) = \text{sgn} \left\{ \sum_{i=1}^n \alpha_i^* y_i (\mathbf{x}_i \cdot \mathbf{x}) + b^* \right\}$$

Xuegong Zhang

24

单选题 1分

⚙ 设置

5-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



Questions

- Why can we call it “optimal”?
- What if the samples are not separable?
- Can we build optimal nonlinear machines like this?





6.3

Nonlinear Classification via the Kernel Trick

Xuegong Zhang

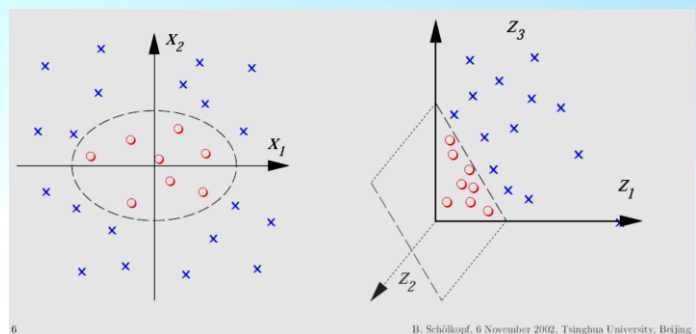
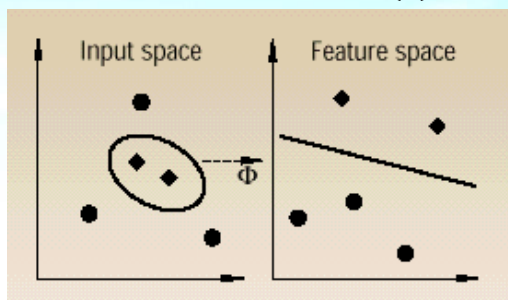
27

Generalization to nonlinear machines



Generalized linear discriminant

$$g(x) = c_0 + c_1 x + c_2 x^2 \quad \xrightarrow{z = [1 \quad x \quad x^2]^T} \quad g(z) = w \cdot z$$

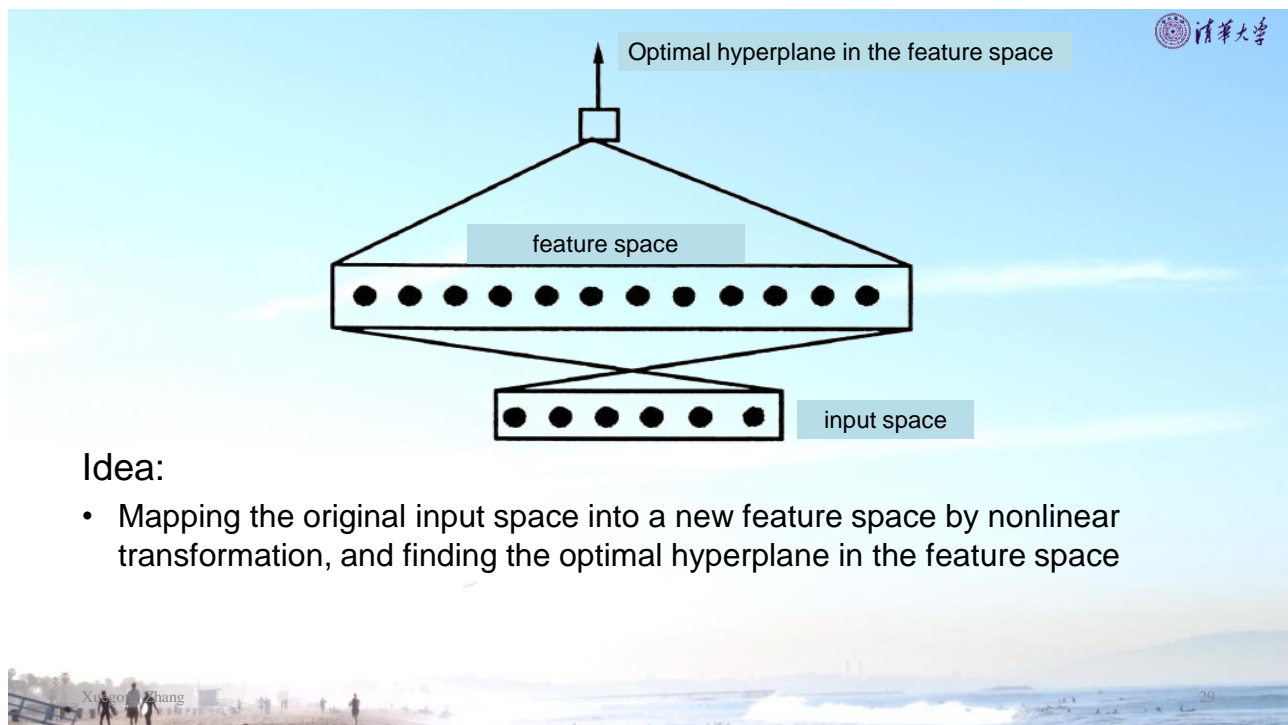
Nonlinear transform $z = \Phi(x)$ 

6

B. Schölkopf, 6 November 2002, Tsinghua University, Beijing

Xuegong Zhang

28



Idea:

- Mapping the original input space into a new feature space by nonlinear transformation, and finding the optimal hyperplane in the feature space

For example:

For d -dimensional samples, to construct a 2nd-order polynomial decision function, we can transform the original $X \subset R^d$ space to a feature space $Z \subset R^D$ with $N=d(d+3)/2$ coordinates:

$$\begin{aligned} z^1 = x^1, \dots, z^d = x^d & \quad d \text{ coordinates,} \\ z^{d+1} = (x^1)^2, \dots, z^{2d} = (x^d)^2 & \quad d \text{ coordinates,} \\ z^{2d+1} = x^1 x^2, \dots, z^D = x^d x^{d-1} & \quad d(d-1)/2 \text{ coordinates.} \end{aligned}$$

The linear discriminant function constructed in the feature space is a quadratic discriminant function in the original space.

Issues:

- Complicated transformation
- Dimensionality increasing



What problems we may meet with this approach?

(1) Can we handle the computation in feature space?

(technical issue)

- e.g, a 4th or 5th order polynomial function on 200d data (like 16x16 pixels) needs feature spaces of ~billion dimensions.

(curse of dimensionality)

(2) Will the “optimal hyperplane” still be optimal or good in the feature space?

(conceptual issue)

- Can it still work or can it generalize?

Let's leave it for future classes.

Xuegong Zhang

31

The Algorithm:

$$\begin{aligned} \min \quad & \psi(\mathbf{w}, \xi) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \left(\sum_{i=1}^n \xi_i \right) \\ \text{s.t.} \quad & y_i[(\mathbf{w} \cdot \mathbf{x}_i) + b] - 1 + \xi_i \geq 0, \quad \xi_i \geq 0, \quad i = 1, \dots, l \end{aligned}$$



$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^l y_i \alpha_i = 0, \end{aligned}$$

and

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, l$$

The Solution:

$$f(\mathbf{x}) = \text{sgn} \left\{ \sum_{i=1}^n \alpha_i^* y_i (\mathbf{x}_i \cdot \mathbf{x}) + b^* \right\}$$



$$\mathbf{z} = \Phi(\mathbf{x})$$



Look carefully



Xuegong Zhang

32



Only the inner-products of the vectors are involved!

$$\max_{\alpha} Q(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\mathbf{x} \rightarrow \Phi(\mathbf{x})$$

$$f(\mathbf{x}) = \text{sgn} \left\{ \sum_{i=1}^l \alpha_i^* y_i (\mathbf{x}_i \cdot \mathbf{x}) + b^* \right\}$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j))$$

Actually we can directly define the kernel without defining the transformation, like:

$$(\mathbf{x} \cdot \mathbf{y})^2 = \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \right)^2$$

$$\begin{aligned} \Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) &\longleftrightarrow \left(\begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix} \cdot \begin{pmatrix} y_1^2 \\ \sqrt{2}y_1y_2 \\ y_2^2 \end{pmatrix} \right) \\ &= (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})), \end{aligned}$$

Xuegong Zhang

33



Support Vector Machine (SVM)

The Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & Q(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s. t.} \quad & \sum_{i=1}^l y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l \end{aligned}$$

Decision Function:

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^l \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + b^* \right)$$

Xuegong Zhang

34

The Kernel Trick

(Boser, Guyon, and Vapnik, 1992)

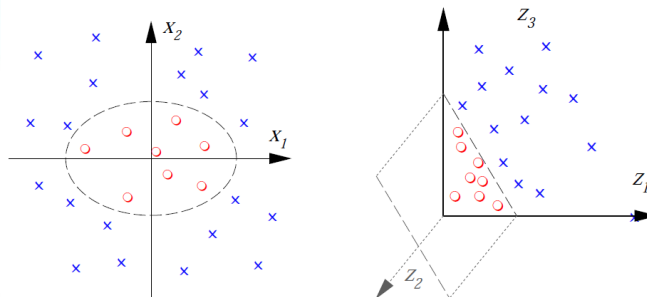


$$\begin{aligned}\langle \Phi(x), \Phi(y) \rangle &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)(y_1^2, \sqrt{2} y_1 y_2, y_2^2)^\top \\ &= \langle x, y \rangle^2 \\ &=: k(x, y)\end{aligned}$$

→ the dot product in F can be computed in \mathbb{R}^2

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$



Xu, Geon, Zhang

6

B. Schölkopf, 6 November 2002, Tsinghua University, Beijing

35

The Kernel Trick

(Boser, Guyon, and Vapnik, 1992)



$$\begin{aligned}\langle \Phi(x), \Phi(y) \rangle &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)(y_1^2, \sqrt{2} y_1 y_2, y_2^2)^\top \\ &= \langle x, y \rangle^2 \\ &=: k(x, y)\end{aligned}$$

→ the dot product in F can be computed in \mathbb{R}^2

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$

A proper transform, then optimal hyperplane in feature space

A proper kernel, then optimal hyperplane defined on kernels

"Convolution of the inner product"

What is a proper kernel?

—— It has to be the inner-product of a certain space.

Xu, Geon, Zhang

6

B. Schölkopf, 6 November 2002, Tsinghua University, Beijing

36

According to the Hilbert-Schmidt theory, the inner product in a Hilbert space has an equivalent representation



$$(z_1, z_2) = \sum_{r=1}^{\infty} a_r z_r(x_1) z_r(x_2) \Leftrightarrow K(x_1, x_2), \quad a_r \geq 0$$

Mercer Theorem

To guarantee that a continuous symmetric function $K(u, v)$ in $L_2(C)$ has an expansion

$$K(u, v) = \sum_{k=1}^{\infty} a_k \psi_k(u) \psi_k(v)$$

with positive coefficients $a_k > 0$ (i.e., $K(u, v)$ describes an inner product in some feature space), it is necessary and sufficient that the condition

$$\int K(u, v) g(u) g(v) du dv > 0$$

be valid for all $g \in L_2(C)$ (C being a compact subset of R^d).



Xu, Ge, Zhang

37

Mercer's Theorem



If k is a continuous kernel of a positive definite integral operator on $L_2(\mathcal{X})$ (where \mathcal{X} is some compact space),

$$\int_{\mathcal{X}} k(x, x') f(x) f(x') dx dx' \geq 0,$$

it can be expanded as

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(x')$$

using eigenfunctions ψ_i and eigenvalues $\lambda_i \geq 0$ [34].

In that case

$$\Phi(x) := \begin{pmatrix} \sqrt{\lambda_1} \psi_1(x) \\ \sqrt{\lambda_2} \psi_2(x) \\ \vdots \end{pmatrix}$$

satisfies $\langle \Phi(x), \Phi(x') \rangle = k(x, x')$.



B. Schölkopf, NIPS, 3 December 2001

Xu, Ge, Zhang

38

Positive Definite Kernels



Let \mathcal{X} be a nonempty set. The following two are equivalent:

- k is *positive definite (pd)*, i.e., k is symmetric, and for
 - any set of training points $x_1, \dots, x_m \in \mathcal{X}$ and
 - any $a_1, \dots, a_m \in \mathbb{R}$

we have

$$\sum_{i,j} a_i a_j K_{ij} \geq 0, \text{ where } K_{ij} := k(x_i, x_j)$$

- there exists a map Φ into a dot product space \mathcal{H} such that

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

\mathcal{H} is a so-called reproducing kernel Hilbert space.

Special case of positive definite kernels: “Mercer kernels”

39

H. Schölkopf, 6 November 2002, Technische Universität Berlin

Xuegong Zhang

39

- Got lost?
- It's ok. Let's skip it and just be assured that there is a theorem to safeguard the choosing of kernels.



Given training data $(x_1, y_1), (x_2, y_2), \dots, (x_l, y_l)$, choose some kernel $K(\cdot, \cdot)$ and solve the dual problem:

$$\max_{\alpha} Q(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{s.t.} \quad \sum_{i=1}^n y_i \alpha_i = 0$$

$$\text{and} \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l$$

to get solution $f(x) = \text{sgn}(\sum_{i=1}^l \alpha_i^* y_i K(x_i, x) + b^*)$

Xuegong Zhang

40



Most commonly used kernels

- Inner product: $K(\mathbf{x}, \mathbf{x}_i) = (\mathbf{x} \cdot \mathbf{x}_i)$ --- Linear SVM
- Polynomial: $K(\mathbf{x}, \mathbf{x}_i) = [(\mathbf{x} \cdot \mathbf{x}_i) + 1]^q$
--- SVM version of polynomial discriminant
- Radial Basis Function (RBF):

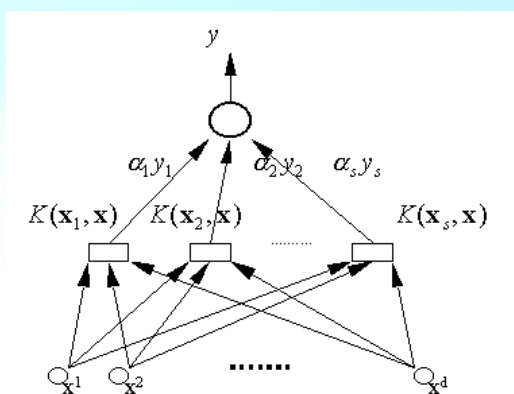
$$K(\mathbf{x}, \mathbf{x}_i) = \exp \left\{ -\frac{|\mathbf{x} - \mathbf{x}_i|^2}{\sigma^2} \right\}$$
 --- SVM version of RBF Network
- Sigmoid: $K(\mathbf{x}, \mathbf{x}_i) = \tanh(\nu(\mathbf{x} \cdot \mathbf{x}_i) + c)$
 --- SVM version of 3-layered MLP (one hidden layer)
 (with certain constraints on the choice of ν and c)

Xuefeng Zhang

41



The “Support Vector Network”



Ref. $g_k(\mathbf{x}) \equiv y_k = f \left(\sum_j w_{jk} f \left(\sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$

Output (decision function)

$$y = \text{sign} \left(\sum_{i=1}^s \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right)$$

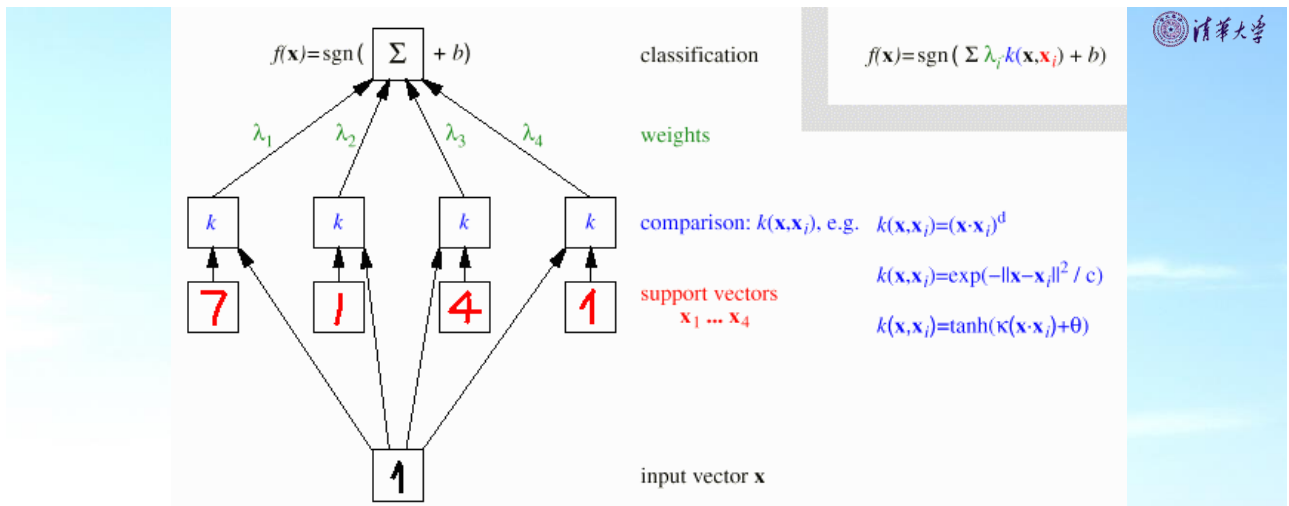
weights $w_i = \alpha_i y_i$

inner-products with s support vectors $\mathbf{x}_1, \dots, \mathbf{x}_s$

Input vector $\mathbf{x} = (x^1, x^2, \dots, x^d)$

Xuefeng Zhang

42



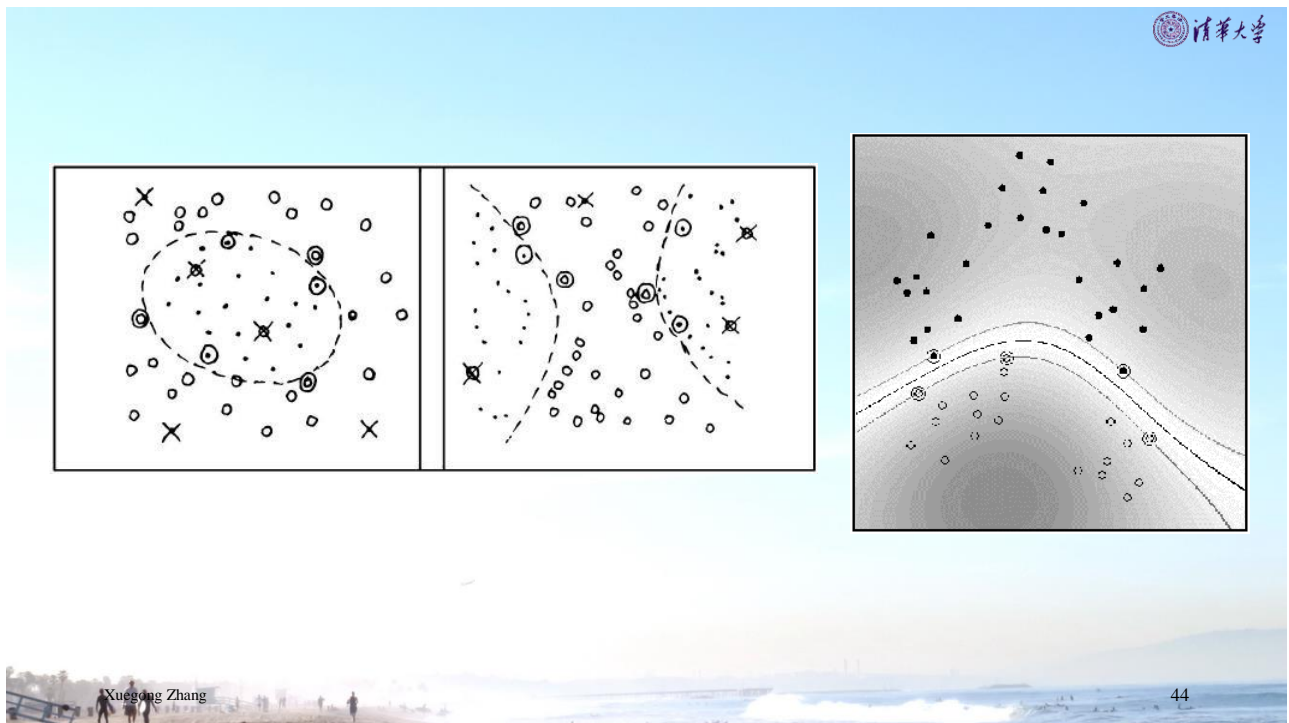
An intuitive understanding

Templates \rightarrow Similarity \rightarrow Weighted Voting

Distance/correlation... \rightarrow Kernel

Xuegong Zhang

43



Former Bell Lab's experiments on USPS data



2601446357146371037114497
1105711129981102860028870
3301033010290602840029012
9405290672980129650299055
510129201803270124431064
1161176057188600158701899
1157557212570688327499516
9950512001536272203242370
355121272315393053880319
1371914119129192511917014
1011919485726803226414186
635972029929972251004670
3084114591010615406103631
1064111030425262009979966
8912056708553131427955460
1018730187112991089970984
0109707597331972015519065
1075318255182814380109463
1787521655460544603546055
18255108503047520439401

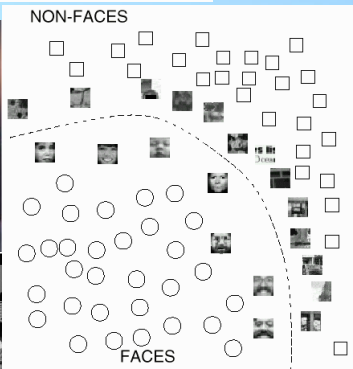
Method	Test Error
Human	2.5%
Decision Tree	16.2%
2-layer MLP	5.9%
5-layer MLP	5.1%
SVM with 3 types of kernels	4.0% 4.1% 4.1%

Xuegong

45



Minnesota Vikings Cheerleaders www.cheerleaders.com



Xuegong Zhang

46



5-second break



Xuegong Zhang

47



6.4 SVM Training Algorithms

Xuegong Zhang

48



Basic concepts of ML: Perceptron

- How can we make a learning machine?
 - It needs a teacher.
 - The model: $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$
 - We need materials to train it. / It needs materials to learn from.
 - Training data: $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in \{-1, 1\}$
 - We need to tell what is the goal of the learning.
 - Objective function: $\min J_P(\alpha) = \sum_{y_j \in Y^k} (-\alpha^T y_j)$
 - We need to tell it how to learn.
 - Learning algorithm: $\alpha(k+1) = \alpha(k) - \rho_k \nabla J = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$

Basic concepts of ML: Logistic Regression

- How can we make a learning machine?
 - It needs a teacher.
 - The model: $h(x) = \theta(w^T x)$
 - We need materials to train it. / It needs materials to learn from.
 - Training data: $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in \{-1, 1\}$
 - We need to tell what is the goal of the learning.
 - Objective function: $\min E(w) = \frac{1}{N} \sum_{j=1}^N \ln(1 + e^{-y_j w^T x_j})$
 - We need to tell it how to learn.
 - Learning algorithm: $w(k+1) = w(k) - \rho_k \nabla E$

Basic concepts of ML: Linear Regression

- How can we make a learning machine?
 - It needs a teacher.
 - The model: $f(x) = \sum_{i=0}^d w_i x_i = w^T x$
 - We need materials to train it. / It needs materials to learn from.
 - Training data: $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in R$
 - We need to tell what is the goal of the learning.
 - Objective function: $\min E = \frac{1}{N} \sum_{j=1}^N (f(x_j) - y_j)^2$
 - We need to tell it how to learn.
 - Learning algorithm: $w(k+1) = w(k) - \rho_k \nabla E$

Basic concepts of ML: MLP

- How can we make a learning machine?
 - It needs a teacher.
 - The model: $g(x) = f(\sum_j w_{jk} f(\sum_i w_{ij} x_i + w_{j0}) + w_{k0})$
 - We need materials to train it. / It needs materials to learn from.
 - Training data: $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in R$
 - We need to tell what is the goal of the learning.
 - Objective function: $\min E = \frac{1}{2} \sum_{j=1}^N (g(x_j) - y_j)^2$
 - We need to tell it how to learn.
 - Learning algorithm: $w(k+1) = w(k) - \rho_k \nabla E$ via the BP algorithm

How to train a SVM?

Xuegong Zhang

49

The optimization problem of SVM



$$f(x) = \text{sign} \left(\sum y_i \alpha_i K(x_i, x) - b \right)$$

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad & \sum_{i=1}^l y_i \alpha_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l \end{aligned}$$

A quadratic programming problem with equality and inequality constraints.

→ Convex Optimization

https://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf



Xuegong Zhang

50

Exterior Point Method for nonlinear programming



$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s. t.} & g_i(\mathbf{x}) \geq 0 \quad i \in I \\ & h_i(\mathbf{x}) = 0 \quad i \in E \end{array}$$



Define a penalty function

$$T(\mathbf{x}, M_k) = f(\mathbf{x}) + M_k \left(\sum_{i \in E} [h_i(\mathbf{x})]^2 \right) + M_k \left(\sum_{i \in I} [\min(0, g_i(\mathbf{x}))]^2 \right)$$

- Increasing the penalty factor M_k during optimization
- Using [the conjugate gradient method](#) to solve the quadratic programming problem.

Xuegong Zhang

51

Algorithms for SVM



$$\max_{\alpha} W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

- The complexity of the problem
 - does not depend on the dimensionality of samples
 - key component: $K(\mathbf{x}_i, \mathbf{x}_j)$, a $l \times l$ matrix
 - huge requirement on memory and computation for large sample sizes
- Key ideas of solutions
 - Dividing the problem on all samples to sub-problems with fewer samples
- An old-styled resource for SVM
 - www.kernel-machine.org or www.kernel-machines.com



Xuegong Zhang

52



Some early references on the implementation

- Osuna, E., Freund, R., Girosi, F., Support Vector Machines: Training and Applications. *MIT AI Memo* 1602, March, 1997
- Edgar Osuna et al. Training Support Vector Machines: an Application to Face Detection, *CVPR'97*
- John C. Platt. Sequential minimal optimization: A fast algorithm for training support vector machines. *Advances in Kernel Methods-Support Vector Learning*. Cambridge, MA: MIT Press, 1999, pp.185-208



Xuegong Zhang

53

Early Popular SVM Packages

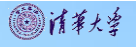


- **SVM^{light}**
 - Thorsten Joachims, Making large-scale SVM learning practical, in B. Schoekopf et al eds. *Advances in Kernel Methods- Support Vector Learning*, MIT Press, 1998
 - <http://svmlight.joachims.org/> (version 6.02, 14.08.2008)
 - https://www.cs.cornell.edu/people/tj/svm_light/
- **SVM^{Torch}**
 - Ronan Collobert & Samy Bengio, SVM^{Torch}: support vector machines for large-scale regression problems, *JMLR*, 1:143-160, 2001
 - <http://bengio.abracadoudou.com/SVMTorch.html>
 - <http://www.torch.ch/>
- **LibSVM**
 - Chih-Chung Chang and Chih-Jen Lin, LIBSVM : a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1--27:27, 2011
 - <http://www.csie.ntu.edu.tw/~cjlin/libsvm> (Version 3.23 released on July 15, 2018)



Xuegong Zhang

54



5-second break



Xuegong Zhang

55



6.5 SVM for Classification

Xuegong Zhang

56

Support Vector Machine (SVM)



Training samples $(y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l), y \in \{-1, 1\}$

The Dual Problem

$$\begin{aligned} \max_{\alpha} \quad Q(\alpha) &= \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s. t.} \quad &\sum_{i=1}^l y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l \end{aligned}$$

Decision Function:

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^l \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + b^* \right)$$

Xuegong Zhang

57

Questions



- Why can we call it “optimal”?
- What if the samples are not separable?
- Can we build optimal nonlinear machines like this?

Xuegong Zhang

58



Why largest margin is optimal?

- Generalization: the expected performance of a machine on future samples after being trained on limited samples
 - The difference between the expected risk and empirical risk

- **Statistical Learning Theory**

- Large margin

→ Low VC dimension

→ Low complexity

→ High generalization ability

$$R(w) \leq R_{emp}(w) + \Phi\left(\frac{h}{n}\right)$$

Let's leave it for future classes.

Xugong Zhang

59

Example:

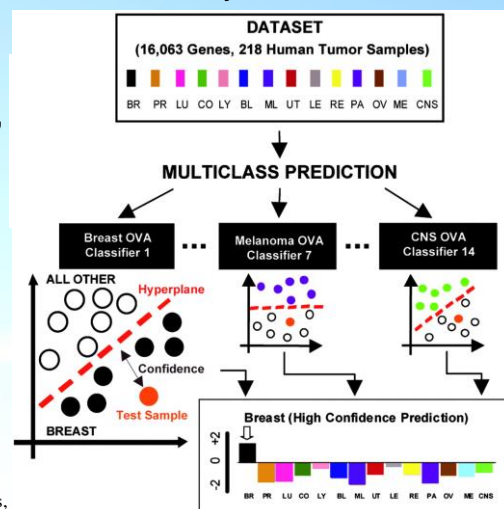
SVM in **multi-class classification** of cancers with microarray data



- 14 tumor classes
- Methods: SVM, Recursive Feature Elimination, etc.
 - they concluded SVM performs the best

- Multi-class problem:

multiple one-over-all (OVA) binary classifiers



S. Ramaswamy et. al. Multiclass cancer diagnosis using tumor gene expression signatures, *PNAS*, **98**(26): 15149-15154, 2001

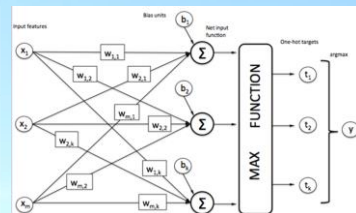
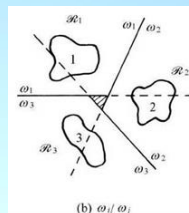
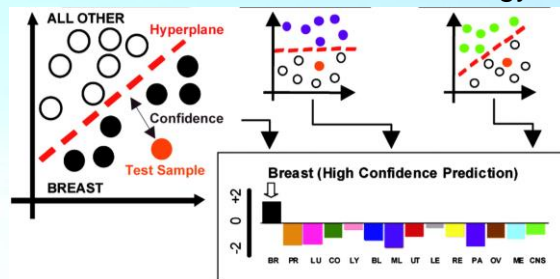
Xugong Zhang

60

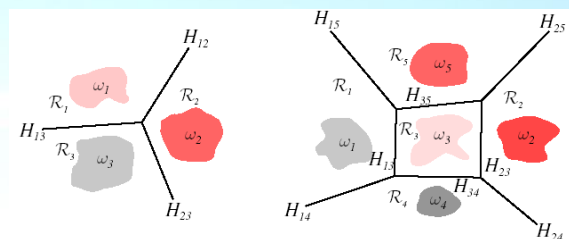
Multi-class classification



- One-vs-all or one-over-all strategy



Can we build a multi-class SVM?



- Question

- Since each SVM classifier is built independently, with its own scaling factor, are the output values of different SVM classifiers comparable?

Xuegong Zhang

61

Multicategory SVM



- Yoonkyung Lee, Yi Lin & Grace Wahba, Multicategory Support Vector Machines, *Technical Report No. 1043, Dept of Statistics, Univ of Wisconsin, Madison*, Sept. 29, 2001
- slides adopted from Lee's presentations

SVM in Regularization Framework

Wahba (1998)

- $f(\mathbf{x}) = h(\mathbf{x}) + b$ with $h \in \mathcal{H}_K$ a Reproducing Kernel Hilbert space (RKHS) with reproducing kernel K .
- Classification rule : $\phi(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$
- SVM scheme : to find f minimizing

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i f(\mathbf{x}_i))_+ + \lambda \|h\|_{\mathcal{H}_K}^2$$

λ is a tuning parameter which balances the data fit and the complexity of $f(\mathbf{x})$.

Xuegong Zhang

62

Multicategory SVM

- Class codes :
 \mathbf{y}_i is represented by a k -dimensional vector with 1 in the j th coordinate and $-\frac{1}{k-1}$ elsewhere if example i falls into class j .
For example, when $k = 3$,
$$\mathbf{y}_i = \begin{cases} (1, -\frac{1}{2}, -\frac{1}{2}) & \text{for class 1} \\ (-\frac{1}{2}, 1, -\frac{1}{2}) & \text{for class 2} \\ (-\frac{1}{2}, -\frac{1}{2}, 1) & \text{for class 3} \end{cases}$$
- Separating functions :
 $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_k(\mathbf{x}))$ with sum-to-zero constraint,
 $\sum_{j=1}^k f_j(\mathbf{x}) = 0$ for any $\mathbf{x} \in \mathbb{R}^d$, and
 $f_j(\mathbf{x}) = h_j(\mathbf{x}) + b_j$ with $h_j \in \mathcal{H}_K$.

• Multicategory SVM formulation :

Find $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_k(\mathbf{x}))$, with sum-to-zero constraint, minimizing

$$\frac{1}{n} \sum_{i=1}^n L(\mathbf{y}_i) \cdot (\mathbf{f}(\mathbf{x}_i) - \mathbf{y}_i)_+ + \frac{\lambda}{2} \sum_{j=1}^k \|h_j\|_{\mathcal{H}_K}^2$$

where $L(\mathbf{y}_i)$ = the j th row of the cost matrix C if \mathbf{y}_i indicates class j .
 C is defined as k by k matrix with 0 on the diagonal, and 1 elsewhere.
For example, when $k = 3$,

$$C = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$$

Figure 4: Predicted decision vectors (f_1, f_2, f_3, f_4) at the test samples. EW: $(1, -1/3, -1/3, -1/3)$. BL: $(-1/3, 1, -1/3, -1/3)$. NB: $(-1/3, -1/3, 1, -1/3)$. and RMS: $(-1/3, -1/3, -1/3, 1)$. The colors indicate the true class identities of the test samples.

Xuegong Zhang

63

单选题 1分

设置

5-min break

A Yes, I'm back to my computer.

B Sorry, not yet.

Xuegong Zhang

提交

64

32



6.6

SVM for Regression (Support Vector Regression)

Xuegong Zhang

65

The primal problem of linear SVR



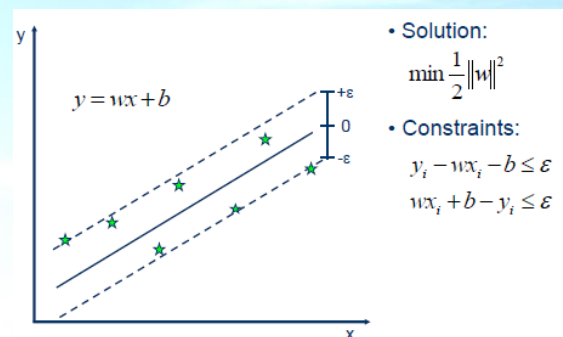
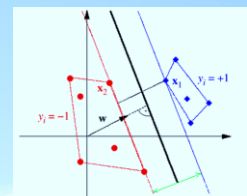
- The task
 - Using function $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ to fit samples $\{\mathbf{x}_i, y_i\}$, $\mathbf{x}_i \in R^d, y_i \in R, i = 1, \dots, l$
- If all samples that can be fit within precision of ε

$$\begin{cases} y_i - \mathbf{w} \cdot \mathbf{x}_i - b \leq \varepsilon, \\ \mathbf{w} \cdot \mathbf{x}_i + b - y_i \leq \varepsilon, \end{cases} \quad i = 1, \dots, l,$$

minimize

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

to obtain the optimal regression.


https://www.saedsayad.com/support_vector_machine_reg.htm

Xuegong Zhang

66



The primal problem of linear SVR

Given samples $\{\mathbf{x}_i, y_i\}, \mathbf{x}_i \in R^d, y_i \in R, i = 1, \dots, l$

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \begin{cases} y_i - \mathbf{w} \cdot \mathbf{x}_i - b \leq \varepsilon, \\ \mathbf{w} \cdot \mathbf{x}_i + b - y_i \leq \varepsilon, \end{cases} \quad i = 1, \dots, l \end{aligned}$$

The primal problem of SVM

$$\begin{aligned} & (y_1, x_1), \dots, (y_l, x_l), y \in \{-1, 1\} \\ \min \quad & \Phi(w) = \frac{1}{2} (w \cdot w) \\ \text{s.t.} \quad & y_i [(x_i \cdot w) - b] \geq 1, \quad i = 1, 2, \dots, l \end{aligned}$$

Xu, Ge, Zhang

67

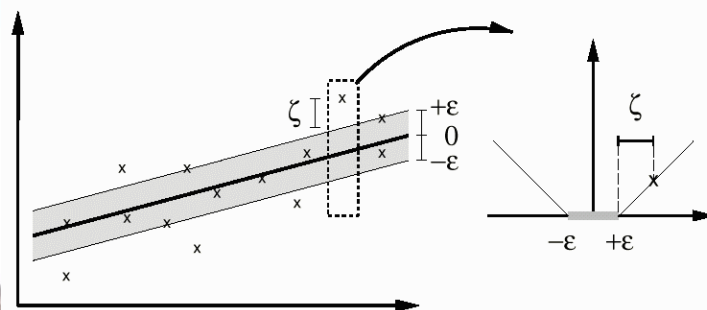


To allow for errors:

- Introducing slack variables $\xi_i \geq 0, \xi_i^* \geq 0$,

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \\ \text{s.t.} \quad & \begin{cases} y_i - \mathbf{w} \cdot \mathbf{x}_i - b \leq \varepsilon + \xi_i^*, \\ \mathbf{w} \cdot \mathbf{x}_i + b - y_i \leq \varepsilon + \xi_i, \end{cases} \quad i = 1, \dots, l \end{aligned}$$

where parameter C controls the penalty to samples that exceed the accuracy threshold ε .



Xu, Ge, Zhang

68

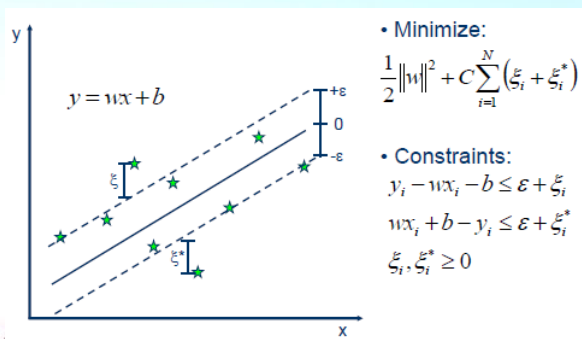


The Primal Problem of SVR

$$\min \quad \Phi(\mathbf{w}, \xi^*, \xi) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \left(\sum_{i=1}^l \xi_i^* + \sum_{i=1}^l \xi_i \right)$$

s.t.

$$\begin{aligned} y_i - (\mathbf{w} \cdot \mathbf{x}_i) - b &\leq \varepsilon + \xi_i^*, & i = 1, \dots, l \\ (\mathbf{w} \cdot \mathbf{x}_i) + b - y_i &\leq \varepsilon + \xi_i, & i = 1, \dots, l \\ \xi_i^* &\geq 0, & \xi_i &\geq 0, & i = 1, \dots, l \end{aligned}$$



https://www.saeedasad.com/support_vector_machine_reg.htm

ε -insensitive loss function

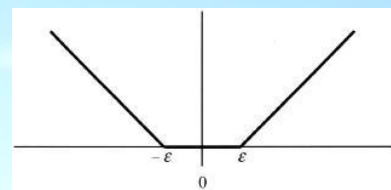
$$L(y, f(x, \alpha)) = |y - f(x, \alpha)|_\varepsilon$$

$$|y - f(x, \alpha)|_\varepsilon = \begin{cases} 0, & \text{if } |y - f(x, \alpha)| \leq \varepsilon \\ |y - f(x, \alpha)| - \varepsilon, & \text{else} \end{cases}$$

The objective function of SVR can be re-written as:

$$\min \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \cdot R_{emp}^\varepsilon(f)$$

$$R_{emp}^\varepsilon(f) := \frac{1}{l} \sum_{i=1}^l |y_i - f(\mathbf{x}_i)|_\varepsilon$$



Xu, Ge, Zhang

70

The Lagrangian

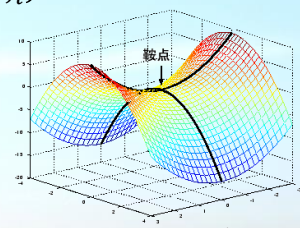
$$L(\mathbf{w}, b, \alpha) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) - b]y_i - 1\}$$



- The solution to the SVR primal problem is the saddle point of the Lagrangian:

$$\begin{aligned} & L(\mathbf{w}, b, \xi^*, \xi; \alpha^*, \alpha, \gamma, \gamma^*) \\ &= \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C \sum_{i=1}^l (\xi_i^* + \xi_i) - \sum_{i=1}^l \alpha_i [y_i - (\mathbf{w} \cdot \mathbf{x}_i) - b + \varepsilon + \xi_i] \\ &\quad - \sum_{i=1}^l \alpha_i^* [(\mathbf{w} \cdot \mathbf{x}_i) + b - y_i + \varepsilon + \xi_i^*] - \sum_{i=1}^l (\gamma_i^* \xi_i^* + \gamma_i \xi_i) \end{aligned}$$

- The minimal point w.r.t. $\mathbf{w}, b, \xi^*, \xi$ and the maximal point w.r.t. $\alpha_i^* \geq 0, \alpha_i \geq 0, \gamma_i^* \geq 0, \gamma_i \geq 0, i = 1, \dots, l$



Xuebin Zhang

71

The Dual Problem of SVR



Following the same procedure as for SVM classification, we can get

$$\mathbf{w} = \sum_{i=1}^l (\alpha_i^* - \alpha_i) \mathbf{x}_i$$

where α_i^*, α_i are the solution of

$$\max \quad W(\alpha, \alpha^*) = -\varepsilon \sum_{i=1}^l (\alpha_i^* + \alpha_i) + \sum_{i=1}^l y_i (\alpha_i^* - \alpha_i) - \frac{1}{2} \sum_{i,j=1}^l (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) (\mathbf{x}_i \cdot \mathbf{x}_j)$$

s.t.

$$\begin{aligned} \sum_{i=1}^l \alpha_i^* &= \sum_{i=1}^l \alpha_i \\ 0 &\leq \alpha_i^* \leq C, \quad i = 1, \dots, l \\ 0 &\leq \alpha_i \leq C, \quad i = 1, \dots, l \end{aligned}$$

And the regression function is

$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b = \sum_{i=1}^l (\alpha_i^* - \alpha_i) (\mathbf{x}_i \cdot \mathbf{x}) + b^*$$

Xuebin Zhang

72



Support Vectors in SVR

- $\alpha_i^* \alpha_i = 0$ for any sample
- Only samples with $\alpha_i^{(*)} = C$ falls outside the ε -insensitive tube
- Samples with $\alpha_i^{(*)} \in (0, C)$ have $\xi_i^{(*)} = 0$ and $|y_i - f(x_i)| = \varepsilon$.

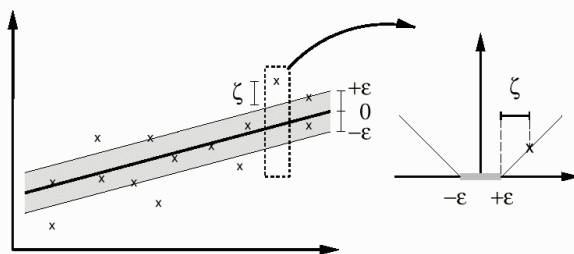
The b can be obtained on these samples:

$$b = y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle - \varepsilon \quad \text{for } \alpha_i \in (0, C)$$

$$b = y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle + \varepsilon \quad \text{for } \alpha_i^* \in (0, C)$$

Corresponds to error samples
in classification

Corresponds to
boundary samples
in classification

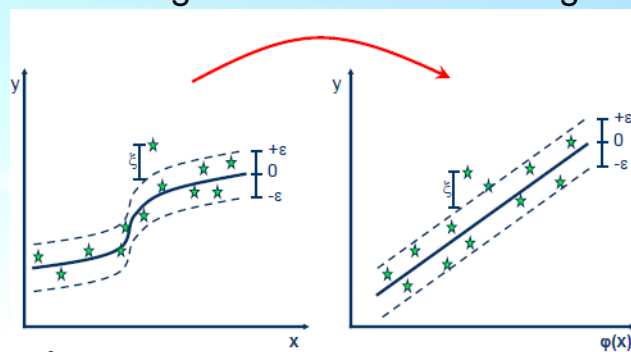


Xu, Ge, Zhang

73

Nonlinear SVR with the kernel trick

- Construct nonlinear regression with SVM using kernels



$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b = \sum_{i=1}^l (\alpha_i^* - \alpha_i) (\mathbf{x}_i \cdot \mathbf{x}) + b^* \quad \Rightarrow \quad f(\mathbf{x}; \mathbf{v}, \boldsymbol{\beta}) = \sum_{i=1}^N \beta_i K(\mathbf{x}, \mathbf{v}_i) + b$$



Xu, Ge, Zhang

https://www.saedsayad.com/support_vector_machine_reg.htm

74



$$f(x) = \sum_{i=1}^l \beta_i K(x, x_i) + b, \quad \beta_i = \alpha_i^* - \alpha_i, \quad i = 1, \dots, l$$



The dual problem of SVR

$$\max W(\alpha, \alpha^*) = -\varepsilon \sum_{i=1}^l (\alpha_i^* + \alpha_i) + \sum_{i=1}^l y_i (\alpha_i^* - \alpha_i) - \frac{1}{2} \sum_{i,j=1}^l (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) K(x_i, x_j)$$

s.t.

$$\begin{aligned} \sum_{i=1}^l \alpha_i^* &= \sum_{i=1}^l \alpha_i \\ 0 &\leq \alpha_i^* \leq C, \quad i = 1, \dots, l \\ 0 &\leq \alpha_i \leq C, \quad i = 1, \dots, l \end{aligned}$$

Xu, Gao, Zhang

75

SVR examples on toy data (1)



- Data from 1D sinc function $f(x) = \frac{\sin(x-10)}{x-10}$
- Regression function $y = \sum_{i=1}^N (\alpha_i^* - \alpha_i) K_1(x, x_i) + b$
- The kernel $K_1(x, x_i) = 1 + x_i x + \frac{1}{2} |x - x_i| (x \wedge x_i)^2 + \frac{(x \wedge x_i)^3}{3}$

where $(x \wedge x_i) := \min(x, x_i)$

Xu, Gao, Zhang

76

The number of SVs decreases with relaxed requirement on accuracy.

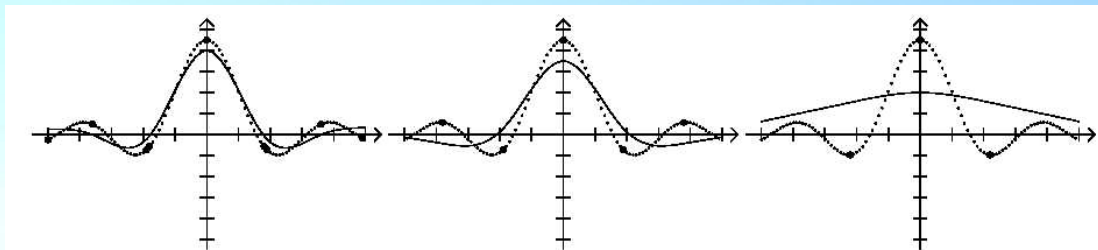
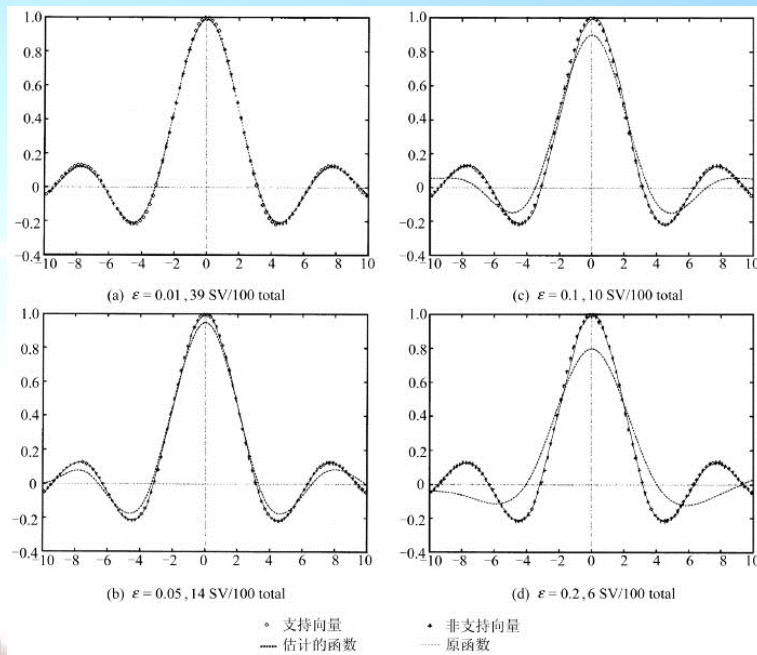


Figure 9.4 Left to right: regression (solid line), data points (small dots) and SVs (big dots) for an approximation of $\text{sinc } x$ (dotted line) with $\varepsilon = 0.1, 0.2$, and 0.5 . Note the decrease in the number of SVs.





SVR examples on toy data (2)

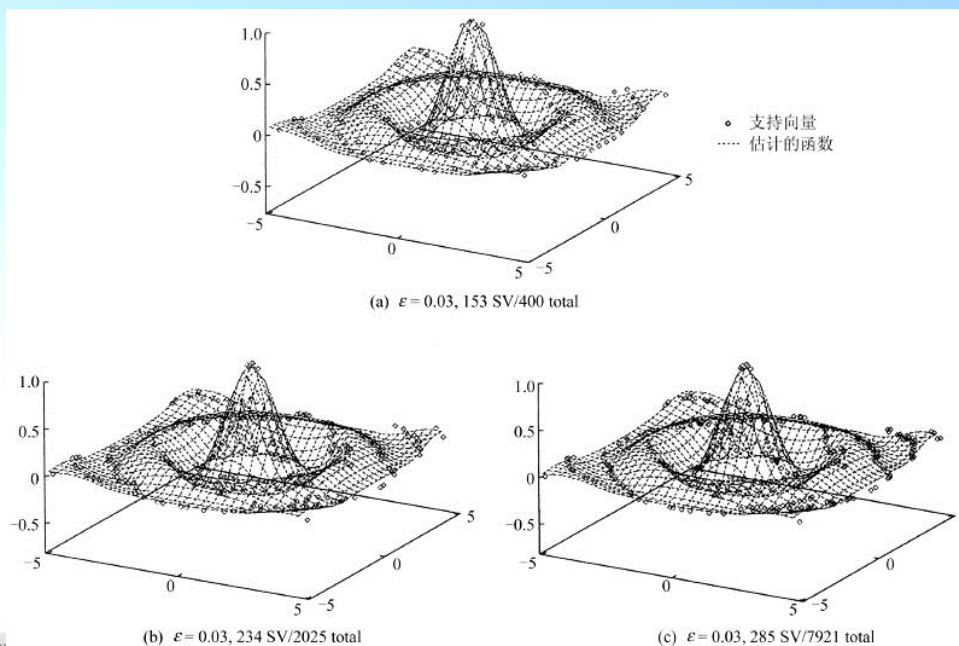
- Data from 2D sinc function $f(x) = \frac{\sin \sqrt{(x-10)^2 + (y-10)^2}}{\sqrt{(x-10)^2 + (y-10)^2}}$
- Regression function $y = \sum_{i=1}^N (\alpha_i^* - \alpha_i) K(x, x_i) K(y, y_i) + b$
- The kernel $K(x, y; x_i, y_i) = K(x, x_i) K(y, y_i)$

$$= \left(1 + xx_i + \frac{1}{2} |x - x_i| (x \wedge x_i)^2 + \frac{(x \wedge x_i)^3}{3} \right) \times \left(1 + yy_i + \frac{1}{2} |y - y_i| (y \wedge y_i)^2 + \frac{(y \wedge y_i)^3}{3} \right)$$

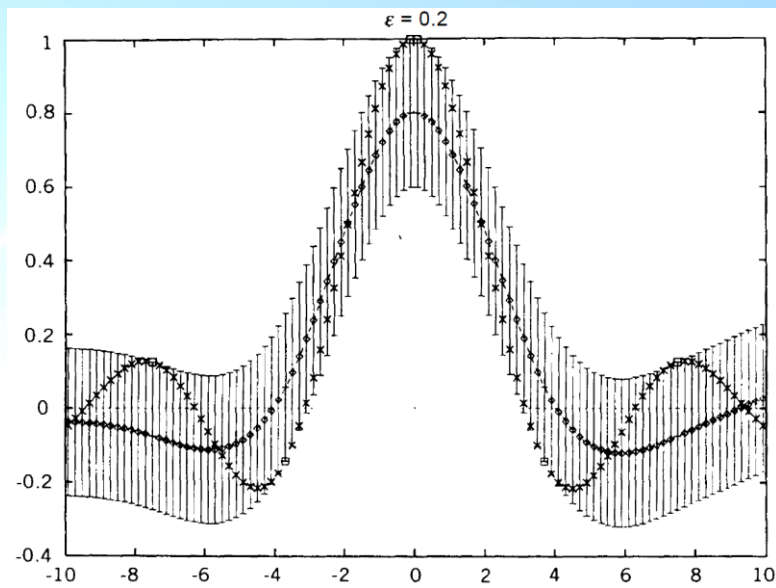
where $(x \wedge x_i) := \min(x, x_i)$



number of SVs increased <2 fold when number of training samples increased ~20 fold



The ε -tube model of function approximation



Xuegong Zhang

81

5-second break



Xuegong Zhang

82



6.7 Kernel Machines

Xuegong Zhang

83

The idea of kernel machines



- Any linear method, if it only involves the inner-product of its input vectors, it can be nonlinearized by adopting the kernel trick.
- Large margin is sometimes desirable for good generalization as the kernel trick raises the dimensionality. → **“Large Margin Machines”**
- **Kernel Fisher’s Discriminant (KFD):** S. Mika, G. Ratsch, J. Weston, B. Scholkopf and K.-R. Muller, Fisher discriminant analysis with kernels, *Neural Networks for Signal Processing IX*, pp.41-48, IEEE, 1999
- **Kernel MSE:** Jianhua Xu, Xuegong Zhang, Yanda Li, Kernel MSE algorithm: a unified framework for KFD, LS-SVM and KRR, *Proceedings of IJCNN’01*, pp.1486-1491
- **Kernel Pocket:** Jianhua Xu, Xuegong Zhang, Yanda Li, Large margin kernel pocket algorithm, *ibid.*, pp. 1480-1485
- **Kernel Nearest Neighbor:** Kai Yu, Liang Ji, Xuegong Zhang, Kernel nearest-neighbor algorithm, *Neural Processing Letters*, 15: 147-156, 2002



Xuegong Zhang

84

6.8 Variations of SVM

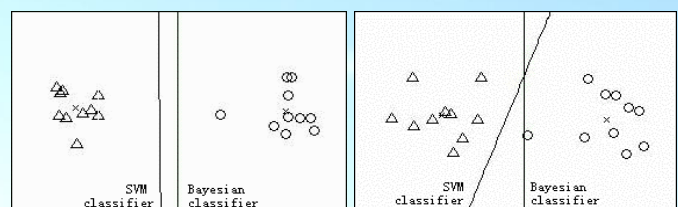
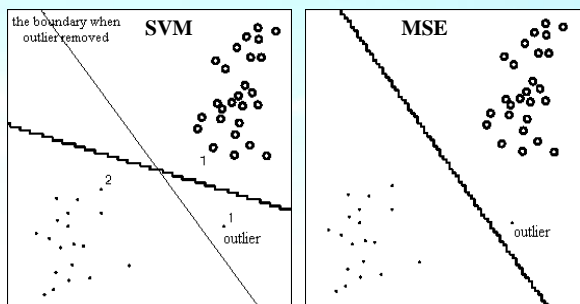
Xuegong Zhang

85

6.8.1 Central Support Vector Machines (CSVM) 清华大学

X. Zhang, Using class-center vectors to build support vector machines, *Neural Networks for Signal Processing, IX*, pp. 3-11, 1999

- Observation: SVM can be sensitive to noise when sample size is very small



Examples of synthetic samples from two normal distributions

Xuegong Zhang

86



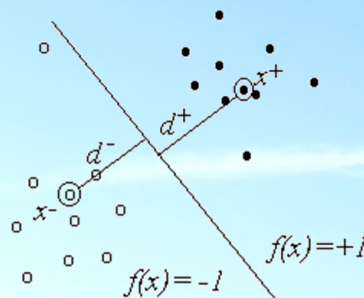
Central Margin

$$d = d^+ + d^- = \frac{\sum_{i=1}^n l_i y_i (\mathbf{w} \cdot \mathbf{x}_i)}{\|\mathbf{w}\|}$$

$$l_i = 1/n^- \quad \text{if } y_i = -1$$

$$l_i = 1/n^+ \quad \text{if } y_i = +1$$

$$\text{normalization: } \sum_{i=1}^n l_i y_i (\mathbf{w} \cdot \mathbf{x}_i) = 1$$



Xuegong Zhang

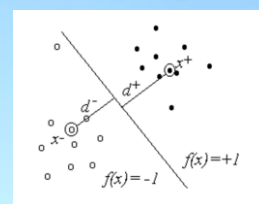
87

CSVM (Linearly separable case)

$$\min \quad \psi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} (\mathbf{w} \cdot \mathbf{w})$$

$$\text{subject to } \sum_{i=1}^n l_i y_i (\mathbf{w} \cdot \mathbf{x}_i) = 1$$

$$\text{and } y_i [(\mathbf{w} \cdot \mathbf{x}_i) + b] - \varepsilon \geq 0, \quad i = 1, \dots, n$$



where $\varepsilon \geq 0$ is a constant indicating the least distance that the samples should be away from the separation hyperplane.

Solution of the dual problem:

$$\mathbf{w}^* = \sum_{i=1}^n (\alpha_i^* + \beta^* l_i) y_i \mathbf{x}_i = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i + \beta^* (\mathbf{x}^+ - \mathbf{x}^-)$$

- Intuition: maximizing the central margin while keeping all the training data not only correctly classified but also be away from the separation hyperplane by at least certain small distance.

Xuegong Zhang

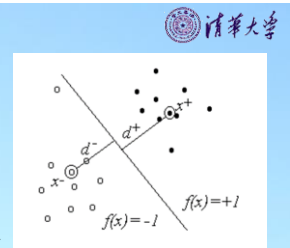
88

CSVM (Linearly non-separable case)

$$\min \psi(\mathbf{w}, \xi) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + C(\sum_{i=1}^n \xi_i)$$

$$\text{subject to } \sum_{i=1}^n l_i y_i (\mathbf{w} \cdot \mathbf{x}_i) = 1$$

$$\text{and } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) + \xi_i \geq \varepsilon > 0, \quad \xi_i \geq 0, \quad i = 1, \dots, n$$



where $\varepsilon \geq 0$ is the constant indicating the least distance that the correctly-classified samples should be away from the separation hyperplane.

Solution to the dual problem :

$$\mathbf{w}^* = \sum_{i=1}^n (\alpha_i^* + \beta^* l_i) y_i \mathbf{x}_i = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i + \beta^* (\mathbf{x}^+ - \mathbf{x}^-)$$

- Intuition: maximizing the central margin while keeping smallest possible mistakes (samples that are misclassified or within the boundary zone are regarded as possible mistakes).

Xuegong Zhang

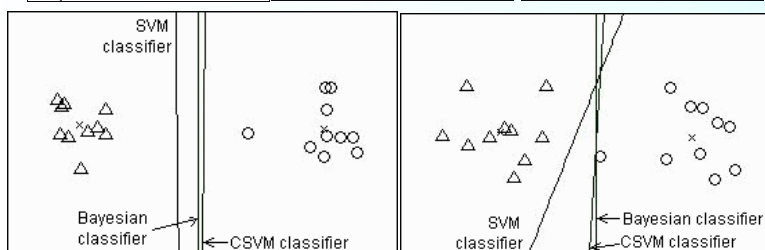
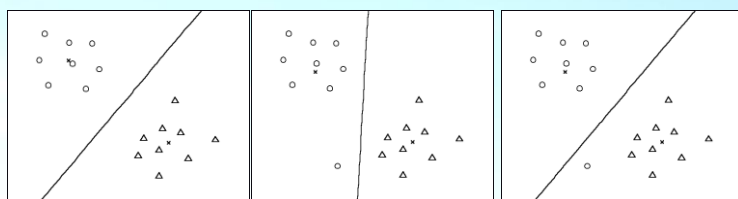
89

CSVM implementation

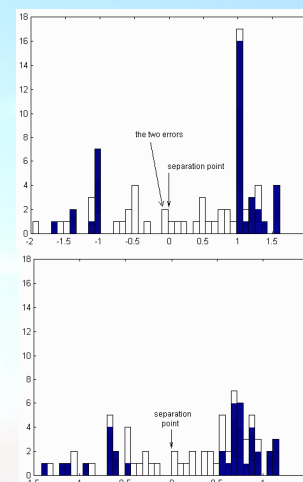
Simplified calculation: $\mathbf{w}^{\text{CSVM}} = (1 - \lambda) \mathbf{w}^{\text{SVM}} + \lambda(\mathbf{x}^+ - \mathbf{x}^-)$

Kernel version: $f(\mathbf{x}) = \text{sgn}((1 - \lambda) \sum_{i=1}^n \alpha_i^{\text{SVM}} y_i K(\mathbf{x}_i, \mathbf{x}) + \lambda \sum_{i=1}^n l_i y_i K(\mathbf{x}_i, \mathbf{x}) + b^*)$

Experiments:



Xuegong Zhang



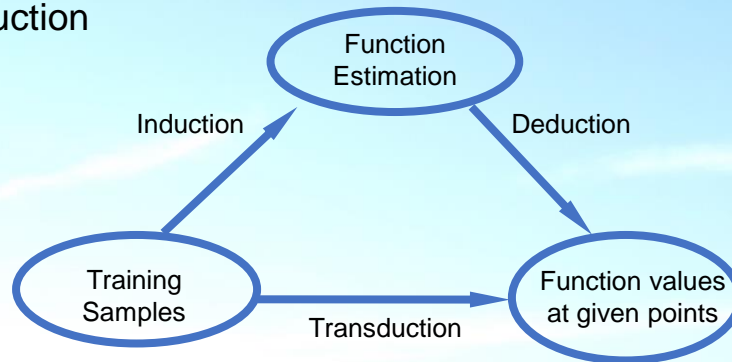
Experiment on a small sample high-dimensional gene expression data

90

6.8.2 Transduction and Semi-supervised SVM



- Transduction

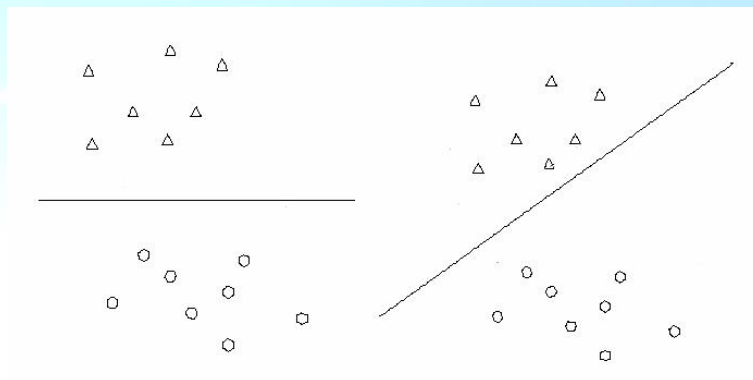


Xuegong Zhang

91

- Semi-supervised SVM (S^3VM)
- Idea: to consider the unknown test samples in the training

Bennett, K. & Demiriz, A., Semi-supervised support vector machines. in: M. S. Kearns, S. A. Solla, and D. A. Cohn, editors, Advances in Neural Information Processing Systems 11, MIT Press, Cambridge, MA, 1999



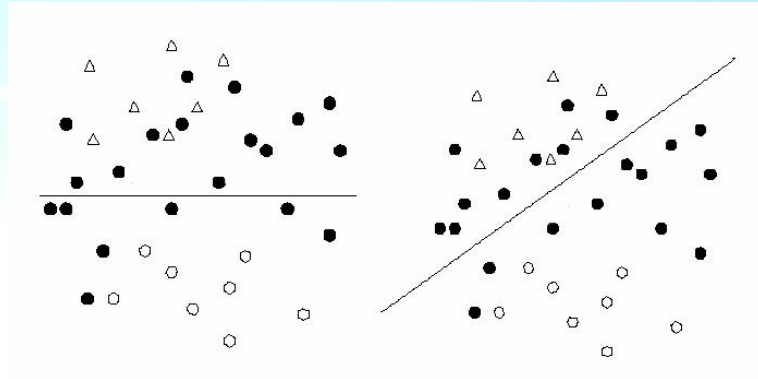
Xuegong Zhang

92



- Semi-supervised SVM (S^3VM)
- Idea: to consider the unknown test samples in the training

Bennett, K. & Demiriz, A., Semi-supervised support vector machines. in: M. S. Kearns, S. A. Solla, and D. A. Cohn, editors, Advances in Neural Information Processing Systems 11, MIT Press, Cambridge, MA, 1999



Xuegong Zhang

93

Homework



- Problems (Pr4)
 1. Derive the SVM dual problem for linearly non-separable cases.
 2. Practice with SVM on a 2-D toy example.
- Deadline:
 - Oct. 20 (Wednesday), 23:00
- Computer exercises (Ex3)
 - Study the SVM package in scikit-learn.
 - Experiment on the medical dataset
- Deadline:
 - Oct. 27 (Wednesday), 23:00



Xuegong Zhang

94

See you next week
for
Statistical Learning Theory



Xuegong Zhang

95