



EE 6885 Statistical Pattern Recognition

Fall 2005
Prof. Shih-Fu Chang
<http://www.ee.columbia.edu/~sfchang>

Lecture 15 (11/16/05)

EE6887-Chang

15-1

- Today's lecture
 - SVM with kernel, error bounds
 - Paper:
Christopher J.C. Burges, "A Tutorial on Support Vector Machines for Pattern Recognition," Data Mining and Knowledge Discovery 2, 121-167, 1998.
- Guest Lecture: Application of SVM in video stream concept detection
- Project data and guideline announced. Due Dec 12.
- Next Lecture: Nov. 21st Monday, Long lecture, starts at 12:20pm
- Final Exam
 - Dec. 16th Friday 1:10-4pm, Mudd Rm 644

EE6887-Chang

15-2

Support Vector Machine (tutorial by Burges '98)

- Look for separation plane with the highest margin

Decision boundary

$$H_0: \mathbf{w}'\mathbf{x} + b = 0$$

- Linearly separable

$$\mathbf{w}'\mathbf{x}_i + b \geq +1 \quad \forall \mathbf{x}_i \text{ in class } \omega_1 \text{ i.e. } y_i = +1$$

$$\mathbf{w}'\mathbf{x}_i + b \leq -1 \quad \forall \mathbf{x}_i \text{ in class } \omega_2 \text{ i.e. } y_i = -1$$

$$\text{Inequality constraints: } y_i(\mathbf{w}'\mathbf{x}_i + b) - 1 \geq 0, \quad \forall i$$

- Two parallel hyperplanes defining the margin

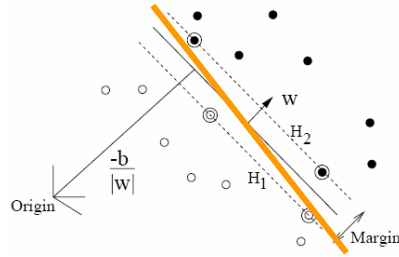
$$\text{hyperplane } H_1(H_+): \mathbf{w}'\mathbf{x}_i + b = +1$$

$$\text{hyperplane } H_2(H_-): \mathbf{w}'\mathbf{x}_i + b = -1$$

- Margin: sum of distances of the closest points to the separation plane

$$\text{margin} = 2 / \|\mathbf{w}\|$$

- Best plane defined by \mathbf{w} and b



EE6887-Chang

15-3

KKT conditions (iff) for separable case

$$\frac{\partial}{\partial w_\nu} L_P = w_\nu - \sum_i \alpha_i y_i x_{i\nu} = 0 \quad \nu = 1, \dots, d \rightarrow \mathbf{w}^* = \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i$$

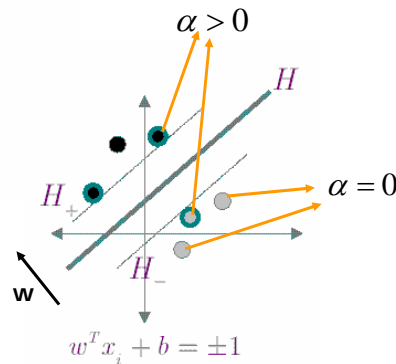
$$\frac{\partial}{\partial b} L_P = -\sum_i \alpha_i y_i = 0$$

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad i = 1, \dots, l$$

$$\alpha_i \geq 0 \quad \forall i$$

$$\alpha_i(y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0 \quad \forall i$$

- Weight sum from positive class = Weight sum from negative class
- Direction of \mathbf{w} : roughly from negative support vectors to positive ones



if $\alpha_i > 0$, \mathbf{x}_i is on H_+ or H_- and is a support vector

- How to compute \mathbf{w} and b ?
- How to classify new data?

EE6887-Chang

15-4

Non-separable

- Add slack variables ξ_i

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 - \xi_i \quad \text{for } y_i = +1$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 + \xi_i \quad \text{for } y_i = -1$$

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

if $\xi_i > 1$, then \mathbf{x}_i is misclassified (i.e. training error)

Lagrange multiplier: minimize

$$L_P = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i}_{\text{New objective function}} - \sum_i \alpha_i \{y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i\} - \sum_i \mu_i \xi_i$$

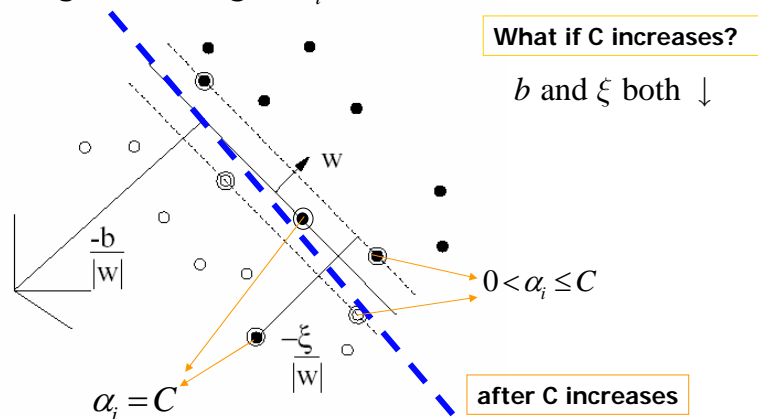
New objective function

Ensure positivity

EE6887-Chang

15-5

- All the points located in the margin gap or the wrong side will get $\alpha_i = C$



- When C increases, samples with errors get more weights
 - better training accuracy, but smaller margin
 - less generalization performance

EE6887-Chang

15-6

Generalized Linear Discriminant Functions

- Include more than just the linear terms

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j = w_0 + \mathbf{w}'\mathbf{x} + \mathbf{x}'\mathbf{W}\mathbf{x}$$

- Shape of decision boundary
 - ellipsoid, hyperhyperboloid, lines etc

- In general $g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x}) = \mathbf{a}'\mathbf{y}$

- Example

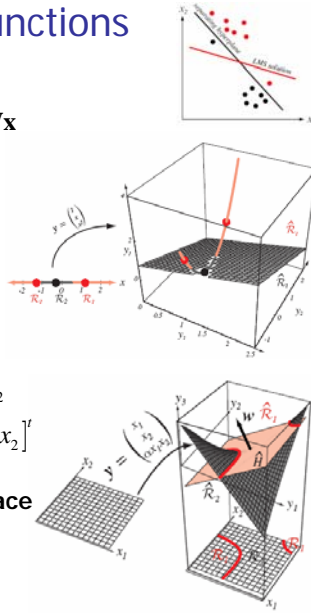
$$g(x) = a_1 + a_2 x + a_3 x^2 \quad g(x) = a_1 x_1 + a_2 x_2 + a_3 x_1 x_2$$

$$= [a_1 \ a_2 \ a_3] [1 \ x \ x^2]^T \quad = [a_1 \ a_2 \ a_3] [x_1 \ x_2 \ x_1 x_2]^T$$

- Data become separable in higher-dimensional space
 - learning parameters in high dimension is hard (curse of dim.)
 - instead, try to maximize margins \rightarrow SVM

EE6887-Chang

15-7



Non-Linear Space

$\Phi: \mathbf{R}^d \mapsto \mathcal{H}$. Map to a high dimensional space, to make the data separable

$$\Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix}$$

- Find the SVM in the high-dim space (embedding space)

$$g(\mathbf{x}) = \sum_{i=1}^{N_s} \alpha_i y_i \underbrace{\Phi(\mathbf{s}_i)}_{\mathbf{w}} \cdot \Phi(\mathbf{x}) + b$$

- Luckily, we don't have to find $\Phi(\mathbf{s}_i)$ nor $\sum_{i=1}^l \alpha_i y_i \Phi(\mathbf{s}_i)$

- Instead, we define kernel $K(\mathbf{s}_i, \mathbf{x}) = \Phi(\mathbf{s}_i) \cdot \Phi(\mathbf{x})$

$$\Rightarrow g(\mathbf{x}) = \sum_{i=1}^{N_s} \alpha_i y_i K(\mathbf{s}_i, \mathbf{x}) + b$$

- We can use the same method to maximize L_D to find α_i

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

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

EE6887-Chang

15-8

Some popular kernels

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^p \quad \text{polynomial}$$

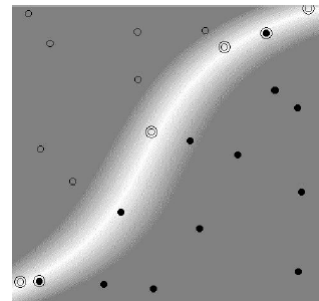
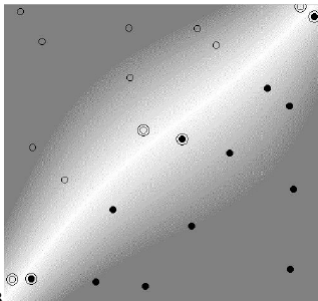
$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2} \quad \text{Gaussian Radial Basis Function (RBF)}$$

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\kappa \mathbf{x} \cdot \mathbf{y} - \delta) \quad \text{sigmoidal neural network}$$

separable

Cubic polynomial

non-separable



EE6887-Chang

15-9

Error bound based on VC dimension (Vapnik '95)

- Risk: expectation of loss → loss

$$R(\alpha) = E \left[\frac{1}{2} |y_i - f(\mathbf{x}_i, \alpha)| \right], \quad \text{where } \alpha \text{ is the classifier model parameter}$$

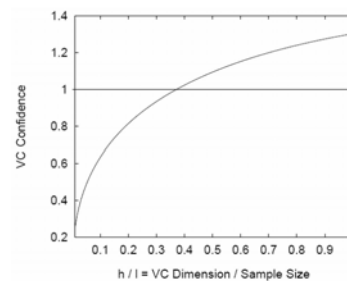
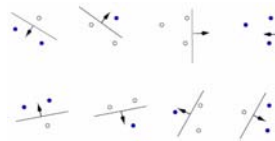
$R_{emp}(\alpha)$: empirical risk for a specific classifier over a training/test set

$$R(\alpha) \leq R_{emp}(\alpha) + \underbrace{\sqrt{\left(\frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l} \right)}}_{\text{VC confidence}} \quad \text{with probability } (1-\eta)$$

h : VC dimension, capacity

l : size of the training set

- VC dim of hyperplane in \mathbb{R}^n is $(n+1)$

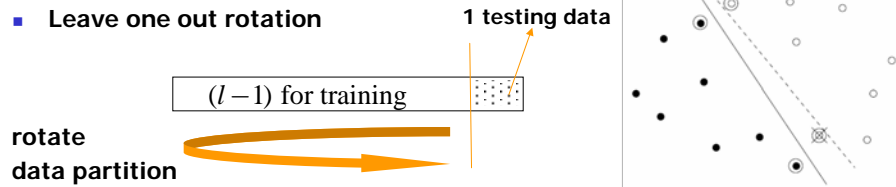


EE6887-Chang

15-10

A tighter error bound of SVM

- Leave one out rotation



- First, train a SVM over l samples
- In each rotation, re-train the SVM over the $l-1$ samples, test on the remaining data
- if the test sample is not SV, then SVM does not change and there is no error. Otherwise, there might be an error.

$$E[P(\text{error})] \leq \frac{E[\# \text{ of support vectors}]}{\text{number of training samples}} = \frac{E[\# \text{ of support vectors}]}{(l-1)}$$

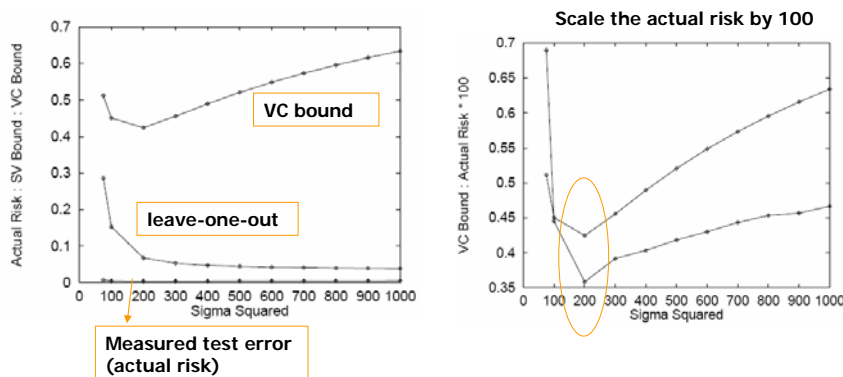
$P(\text{error})$: risk (expected test error) for a learned classifier trained on $l-1$ samples

$E[P(\text{error})]$: expected risk over all choices of training set of $l-1$ samples

$E[\# \text{ of s.v.}]$: expected # of s.v. over all choices of training set of size l

Potential use and issue of the error bound

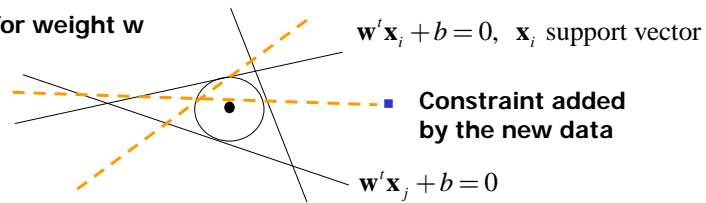
$$K(\mathbf{x}, \mathbf{y}) = e^{-(\mathbf{x} \cdot \mathbf{y})^2 / 2\sigma^2} \quad \bullet \quad \text{How to determine the best sigma value?}$$



- Leave-one-out bound is tighter than VC bound in this experiment (NIST digit classification)
- But VC bound has better predictive power for selecting a good classifier (machine)

Applications (Active SVM)

- Space for weight w



- In image retrieval

- first train a SVM from labeled data
- now in interactive retrieval
- select a new sample and present it to user
- user label the new data
- use the new label to re-train the weight w
- which sample to choose?

- Choose the un-labeled sample that is closest to the current separation plane. Why?