

Support Vector Machine

Wei-Ta Chu

Introduction – Linearly Separable Data

2

Consider the set of training data that consist of two classes $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$, $y_i \in \{-1, 1\}$, $\mathbf{x}_i \in \mathbf{R}^d$. A linear classifier able to separate the positive from the negative examples will be a hyperline in \mathbf{R}^d characterized by a normal \mathbf{w} and an offset b :

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

For a linearly separable data set S , there exists a hyperplane that satisfies all the points in S

$$\mathbf{w}^T \mathbf{x}_i + b \geq 0 \text{ for } y_i = +1$$

$$\mathbf{w}^T \mathbf{x}_i + b < 0 \text{ for } y_i = -1$$

$$\rightarrow y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 0 \quad \forall i$$

Equation of a plane

$$ax + by + cz + d = 0$$

Normal vector: (a,b,c)

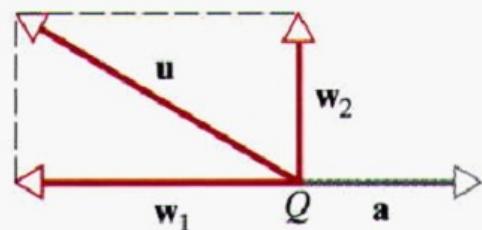
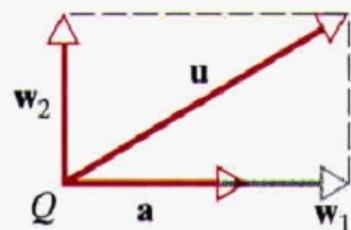
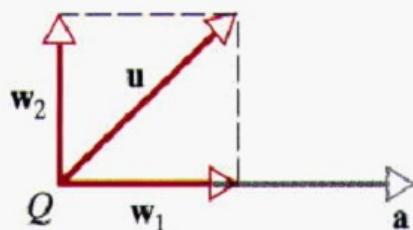
LA Recap: Orthogonal Projection

3

$$\text{proj}_{\mathbf{a}} \mathbf{u} = \frac{\mathbf{u} \cdot \mathbf{a}}{\|\mathbf{a}\|^2} \mathbf{a} \quad (\text{vector component of } \mathbf{u} \text{ along } \mathbf{a})$$

$$\mathbf{u} - \text{proj}_{\mathbf{a}} \mathbf{u} = \mathbf{u} - \frac{\mathbf{u} \cdot \mathbf{a}}{\|\mathbf{a}\|^2} \mathbf{a} \quad (\text{vector component of } \mathbf{u} \text{ orthogonal to } \mathbf{a})$$

$$\|\text{proj}_{\mathbf{a}} \mathbf{u}\| = \frac{|\mathbf{u} \cdot \mathbf{a}|}{\|\mathbf{a}\|} = \|\mathbf{u}\| \cos \theta$$

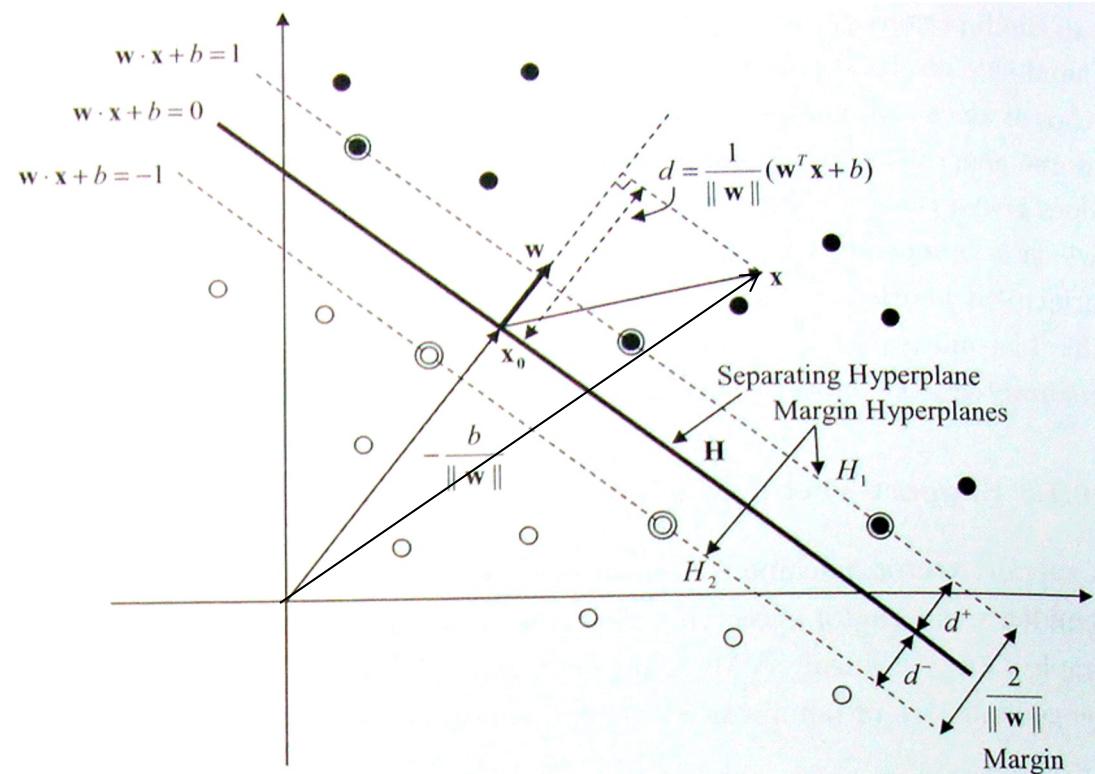


Introduction

4

- Rescale (\mathbf{w}, b) so that the closest points to the hyperplane satisfy $|\mathbf{w}^T \mathbf{x}_i + b| = 1$. This normalization leads to the canonical form for SVM

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$



Properties

5

1. For any two points \mathbf{x}_1 and \mathbf{x}_2 lying on the hyperplane, $\mathbf{w}^T(\mathbf{x}_1 - \mathbf{x}_2) = 0$.
Therefore, \mathbf{w} is the vector normal to the surface of the hyperplane.
2. For any point \mathbf{x}_h on the hyperplane, $\mathbf{w}^T \mathbf{x}_h = -b$.
3. The signed distance from a point \mathbf{x} to the hyperplane is given by

$$\vec{v} \cdot \vec{w} = \|\vec{v}\| \|\vec{w}\| \cos \theta$$

$$\|\vec{v}\| \cos \theta = \frac{\vec{v} \cdot \vec{w}}{\|\vec{w}\|}$$

$$\begin{aligned} d &= \frac{\mathbf{w}}{\|\mathbf{w}\|} (\mathbf{x} - \mathbf{x}_0) \\ &= \frac{1}{\|\mathbf{w}\|} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{x}_0) \\ &= \frac{1}{\|\mathbf{w}\|} (\mathbf{w}^T \mathbf{x} + b) \end{aligned} \tag{10.11}$$

where \mathbf{x}_0 is the intersection point between the normal vector \mathbf{w} and the hyperplane. Since \mathbf{x}_0 lies on the hyperplane, it satisfies the equality

Properties

6

$\mathbf{w}^T \mathbf{x}_0 = -b$ in item 2, which leads to the last equality in the above derivation.

4. The perpendicular distance from the hyperplane to the origin equals $|b|/\|\mathbf{w}\|$ (see Problem 10.2 at the end of the chapter).
5. The points for which the equality in (10.10) holds are those points that lie on the hyperplanes $\mathbf{w}^T \mathbf{x} + b = \pm 1$ (denoted as H_1, H_2 , respectively), and have the perpendicular distance $1/\|\mathbf{w}\|$ to the separating hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$.

$$d^+ + d^- = \frac{1}{\|\mathbf{w}\|} + \frac{1}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$

Substitute \mathbf{x} as $(0,0)$ in property 3

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i \quad (10.10)$$

- The goal of SVM is to find the pair of hyperplanes H1, H2 that maximize the margin, subject to the constraints

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

- This can be formulated as the constrained optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w},$$

$$\text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

- This is a convex optimization problem for which we are guaranteed to obtain its global optimal solution.

Introduction – Linearly Non-Separable Data

8

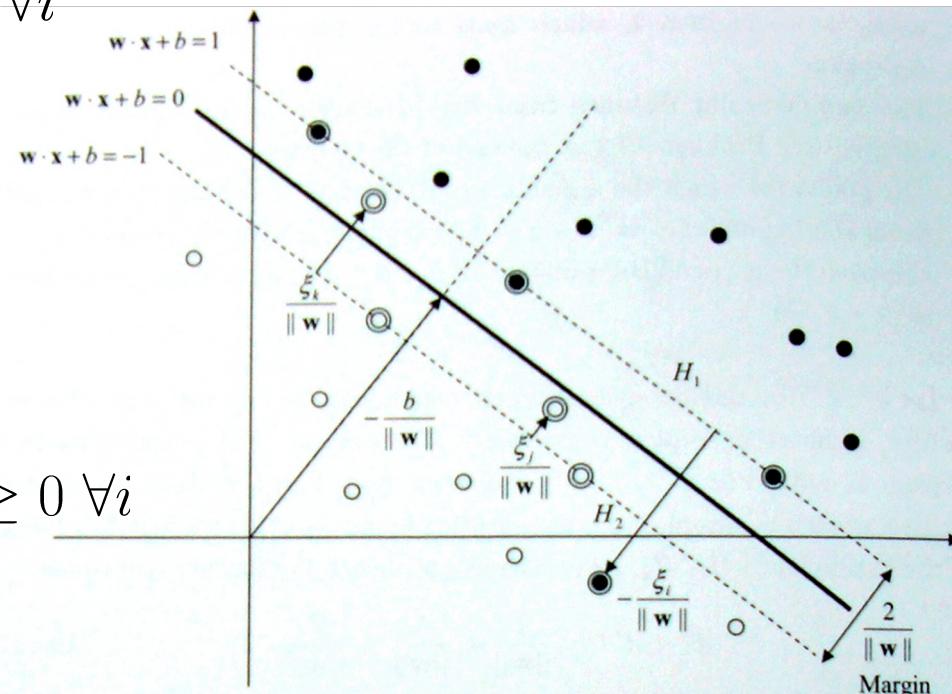
- When the SVM derived above is applied to non-separable data sets, some data points could be at a distance $\xi_i/\|\mathbf{w}\|$ on the wrong side.
- We transform the constraint to

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i$$

- The SVM for non-separable case can be casted as the optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{n} \sum_{i=1}^n \xi_i,$$

$$\text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i$$



Optimization Problem

9

- Use Lagrange multiplier method to find the optimal solution – transform the original problem to a dual problem

The Dual Problem

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{n} \sum_{i=1}^n \xi_i,$$

subject to $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \forall i$

10

- The Lagrange multiplier method first defines a Lagrange function using a set of non-negative Lagrangian multipliers

$$\boldsymbol{\alpha} = \{\alpha_i\} \quad \boldsymbol{\beta} = \{\beta_i\}$$

$$L_P = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{n} \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^b \beta_i \xi_i ,$$

- The unconstrained minimum of the Lagrangian function is computed with respect to \mathbf{w}, b , and ξ_i

$$\frac{\partial L_P}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \quad \rightarrow \quad \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0$$

$$\frac{\partial L_P}{\partial b} = - \sum_{i=1}^n \alpha_i y_i = 0 \quad \rightarrow \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial L_P}{\partial \xi_i} = \frac{C}{n} - \alpha_i - \beta_i = 0 \quad \rightarrow \quad \alpha_i = \frac{C}{n} - \beta_i$$

The Dual Problem

$$L_P = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{n} \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^b \beta_i \xi_i ,$$

11

- Substituting these solutions, we obtain the dual objective function

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

- Since the primal problem is convex and strictly feasible, its minimum solution can be obtained by equivalently maximizing the dual objective function

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

$$\text{subject to } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq \frac{C}{n}, \forall i$$

- Once the optimal $\boldsymbol{\alpha}$ set is obtained, the optimal classifier is then given by $f_{\boldsymbol{\alpha}}(\mathbf{x}) = \text{sign} (\sum_{i=1}^n \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b)$

Kernel Trick

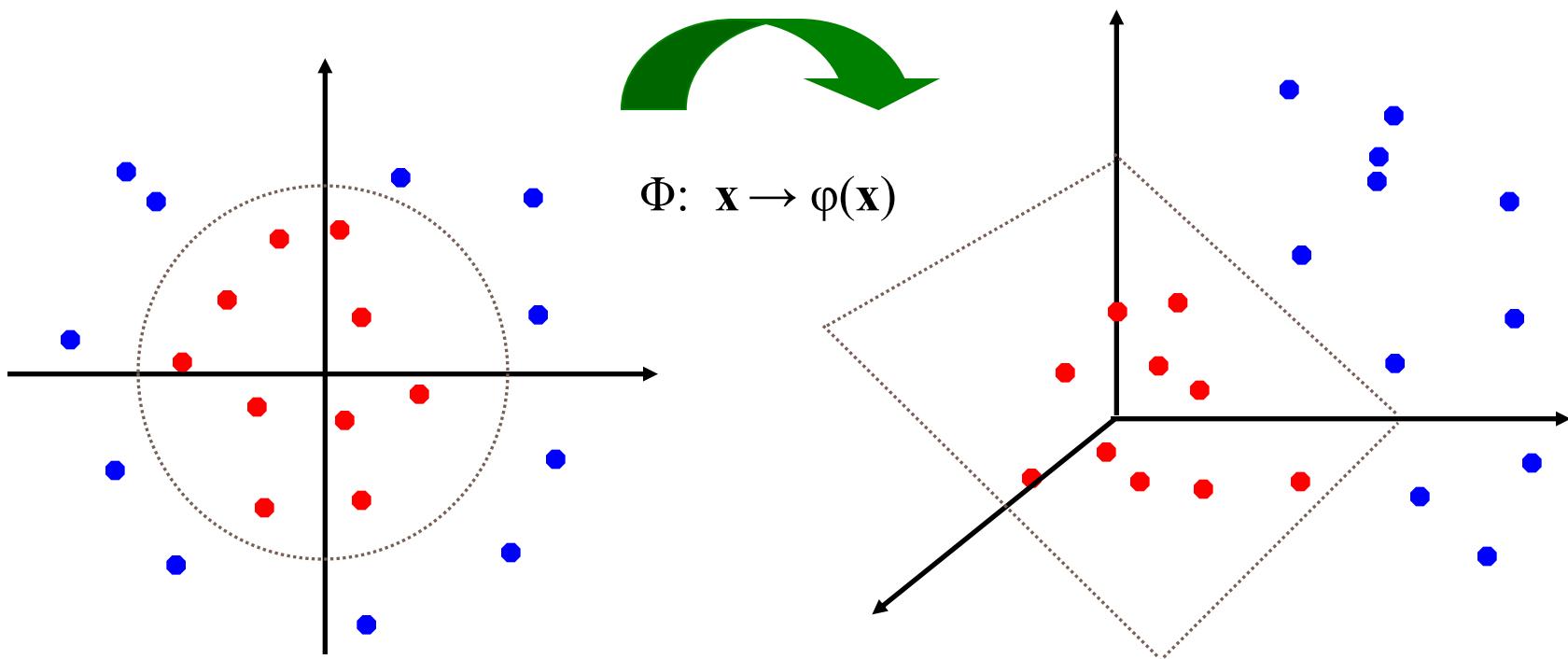
12

- So far, we discuss the case that finds a linear boundary in the input feature space which divides data into two classes.
- We now extend the ideas of linear SVM to enable the generation of a nonlinear classification boundary using the kernel trick techniques.
 - Implicitly map the original feature space into an enlarged feature space using a kernel function, and to use a linear classifier to conduct data classification in this enlarged space.

General Idea

13

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



Kernel Trick

$$f_{\alpha}(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b \right)$$

14

- Note that the input feature vectors appear only in the form of dot products. We do not need to care about individual components of the input vectors. All we need to know is dot products between the input vectors.
- Assume that we enlarge the original feature space by mapping the data into some high dimensional Euclidean space \mathcal{N} using a mapping function $\Phi : R_d \longmapsto \mathcal{N}$
- In the enlarged feature space, because the SVM depends on the data only through dot products $\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$, if there is a “kernel function” K such that $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$, then we only need K in computation rather than Φ

Kernel Trick

15

$$\begin{aligned}f_{\alpha}(\mathbf{x}) &= \text{sign} \left(\sum_{i=1}^n \alpha_i y_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b \right) \\&= \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right)\end{aligned}$$

Example: Consider a feature space with two inputs $\mathbf{x} = [x_1, x_2]$. Applying the degree-2 polynomial kernel $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^2$

$$\begin{aligned}K(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x} \cdot \mathbf{x}')^2 \\&= (1 + x_1 x'_1 + x_2 x'_2)^2 \\&= 1 + 2x_1 x'_1 + 2x_2 x'_2 + (x_1 x'_1)^2 + (x_2 x'_2)^2 + 2x_1 x'_1 x_2 x'_2\end{aligned}$$

If we choose the mapping function as

$$\Phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2]^T$$

Which is the function that maps the two dimensional feature space into a six dimensional one, then $K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$

Kernel Trick

16

- This example reveals that the degree-2 polynomial kernel function serves to map the original two dimensional feature space into a six dimensional one, and to conduct dot products in the it without the need to explicitly compute $\Phi(\mathbf{x})$ for each input data.
- Using kernel trick, we can generalize the linear SVM to a nonlinear one with a minimal change in the framework and a minimal computational cost.
- Examples of kernel functions:

Degree-d polynomial: $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^d$

Radial basis: $K(\mathbf{x}, \mathbf{x}') = \exp(-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{c})$

Neural network: $K(\mathbf{x}, \mathbf{x}') = \tanh(\kappa_1 \mathbf{x} \cdot \mathbf{x}' + \kappa_2)$

SVM Training/Testing

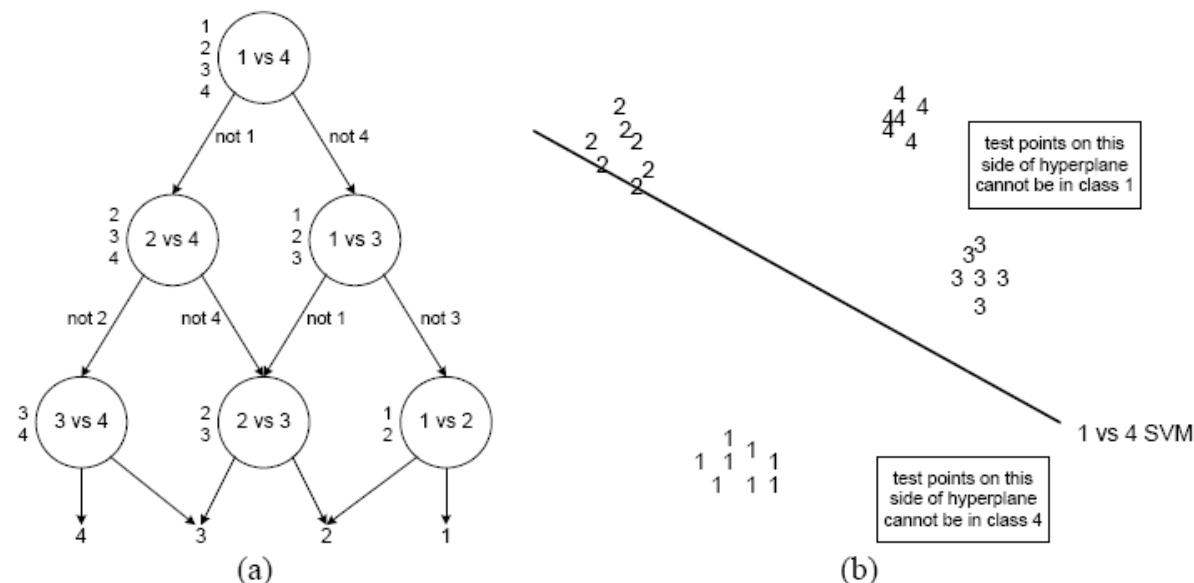
17

- Training with labeled data
- Training error vs. testing error
- k-fold cross validation
 - ▣ The original sample is partitioned into K subsamples. Of the K subsamples, a single subsample is retained as the validation data for testing the model, and the remaining K-1 subsamples are used as training data.
 - ▣ The cross-validation process is then repeated K times, with each of the K subsamples used exactly once as the validation data.
 - ▣ The K results from the folds then can be averaged (or otherwise combined) to produce a single estimation.

Multiclass SVMs

18

- The support vector machine is fundamentally a two-class classifier.
- For multiclass problems, combine multiple two-class SVMs.
 - “One against one” vs. “one against the rest”
 - E.g. DAGSVM (Decision Acyclic Graph SVM)



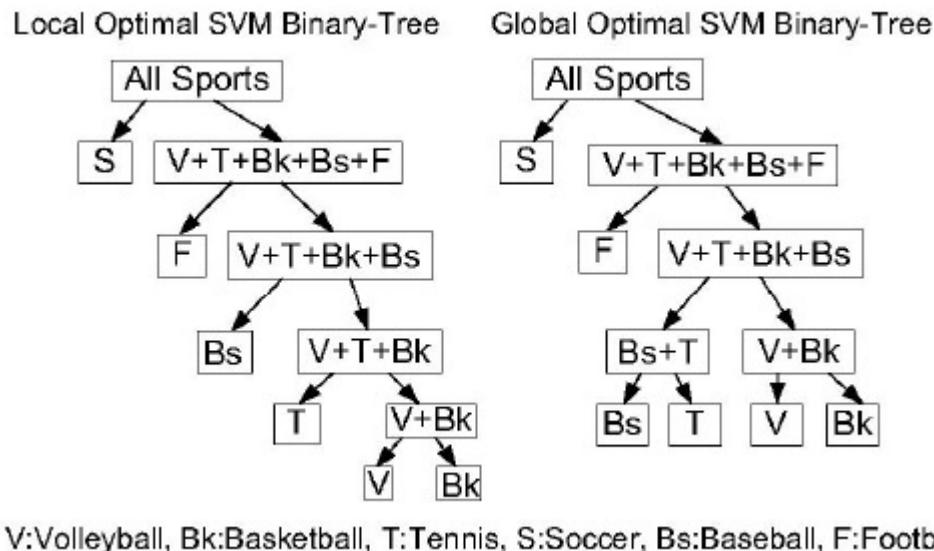
Platt et al. “Large margin DAGs for multiclass classification,” Advances in Neural Information Processing Systems, vol. 12, pp. 547-553, 2000

Figure 1: (a) The decision DAG for finding the best class out of four classes. The equivalent list state for each node is shown next to that node. (b) A diagram of the input space of a four-class problem. A 1-v-1 SVM can only exclude one class from consideration.

Case Study

19

- Automatic video genre categorization
- Temporal features
 - Shot length, cut percentage, average color diff., ...
- Spatial features
 - Face frames ratio, average brightness, average color entropy, ...



Yuan, et al. "Automatic video genre categorization using hierarchical SVM," Proc. of ICIP, pp. 2905-2908, 2006.

Related Resources

20

- LIBSVM
 - <http://www.csie.ntu.edu.tw/~cjlin/libsvm/index.html>

References

21

- L. Rabiner, “A tutorial on hidden Markov models and selected applications in speech recognition,” Proceedings of IEEE, vol. 77, no. 2, pp. 257-286, 1989.
- Statistical Data Mining Tutorials, Tutorial Slides by Andrew Moore:
<http://www.autonlab.org/tutorials/index.html>
- Jiang, et al., “A new method to segment playfield and its applications in match analysis in sports video,” In Proc. of ACM MM, pp. 292-295, 2004.
- Peng, et al., “Extract highlights from baseball game video with hidden Markov models,” In Proc. of ICIP, vol. 1, pp. 609-612, 2002.

References

22

- Platt et al. “Large margin DAGs for multiclass classification,” Advances in Neural Information Processing Systems, vol. 12, pp. 547-553, 2000
- Yuan, et al. “Automatic video genre categorization using hierarchical SVM,” Proc. of ICIP, pp. 2905-2908, 2006.
- LIBSVM, <http://www.csie.ntu.edu.tw/~cjlin/libsvm/index.html>
- Hidden Markov Model (HMM) Toolbox for Matlab
 - <http://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>
- The General Hidden Markov Model library (GHMM)
 - <http://ghmm.sourceforge.net/>
- HTK Speech Recognition Toolkit
 - <http://htk.eng.cam.ac.uk>