

Text Categorization: Discriminative Classifiers

Part 2

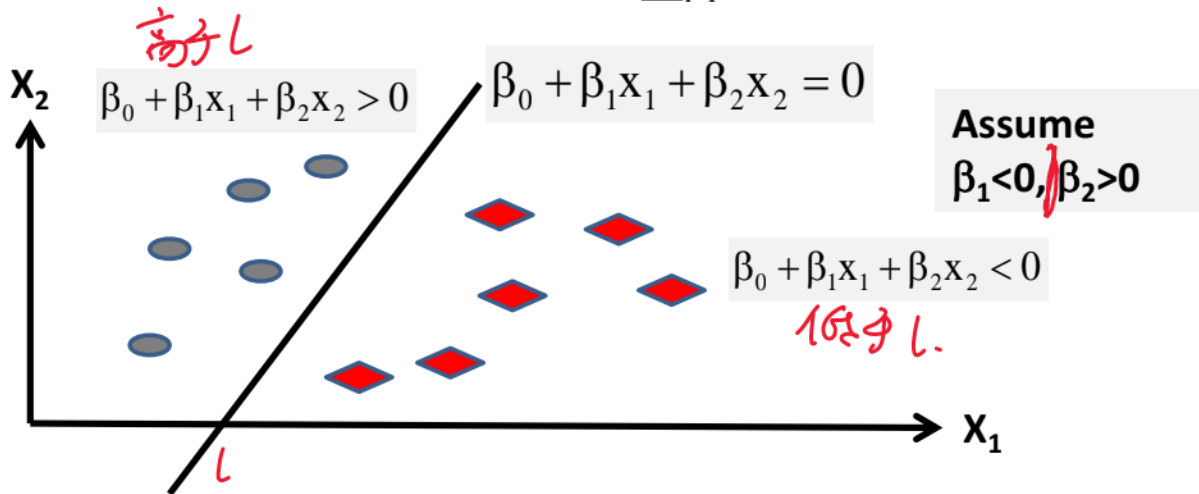
ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

Discriminative Classifier 3: Support Vector Machine (SVM)

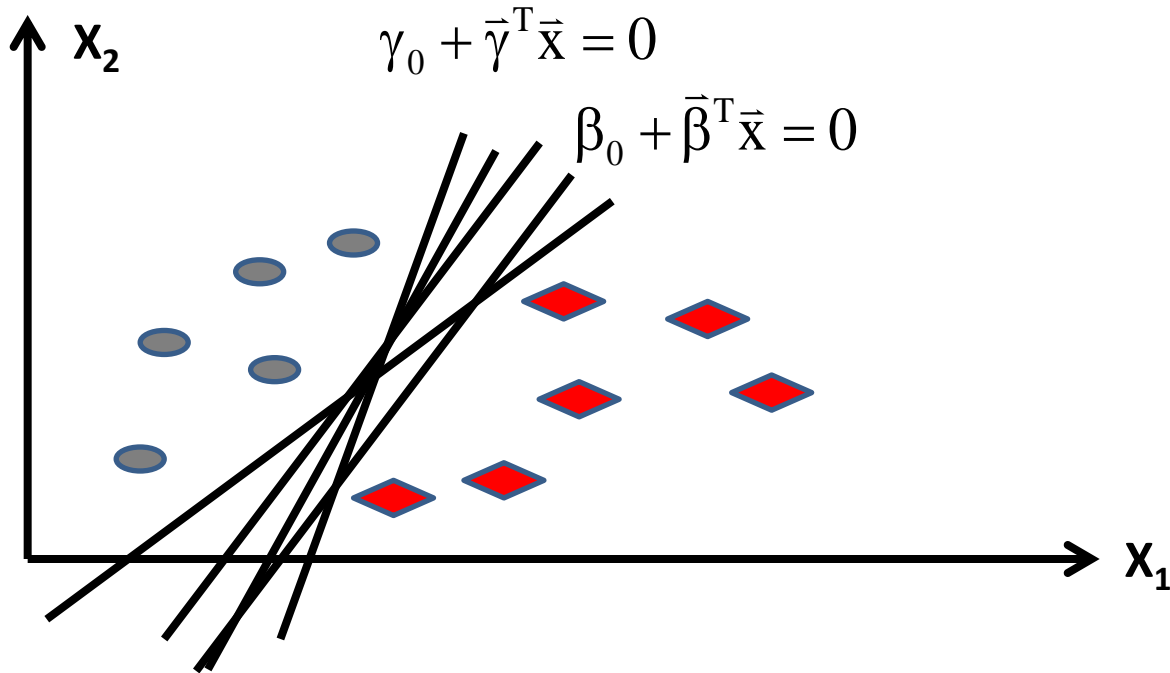
- Consider two categories: $\{\theta_1, \theta_2\}$

$f(X) \geq 0 \Rightarrow X$ is in category θ_1
 $f(X) < 0 \Rightarrow X$ is in category θ_2

- Use a linear separator $f(X) = \beta_0 + \sum_{i=1}^M x_i \beta_i \quad \beta_i \in \mathbb{R}$



Which Linear Separator Is the Best?



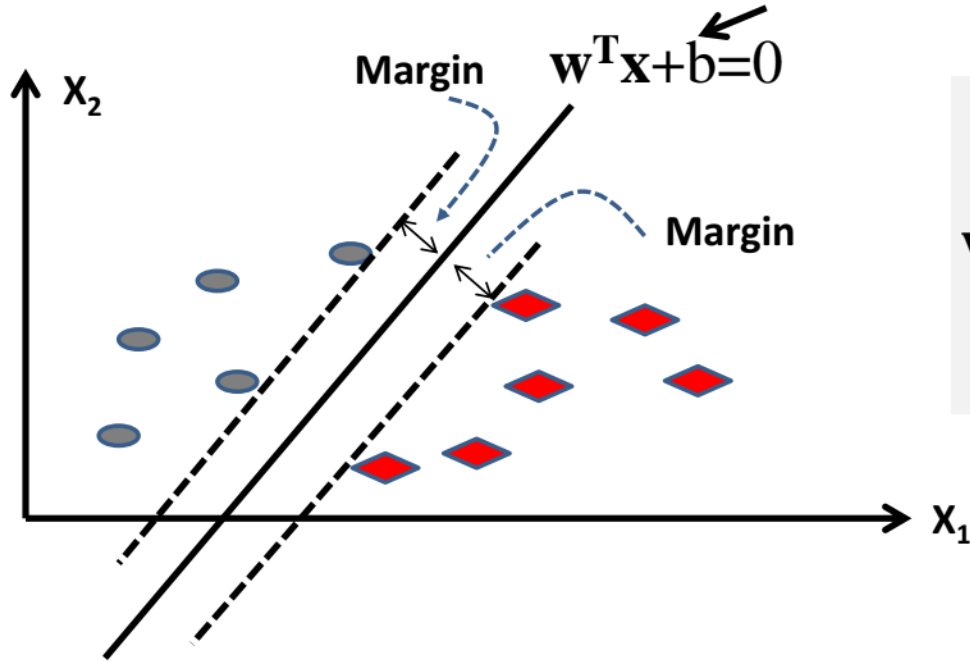
Best Separator = Maximize the Margin

Notation Change: $\beta \rightarrow w$; $\beta_0 \rightarrow b$

Bias constant

Feature Weights

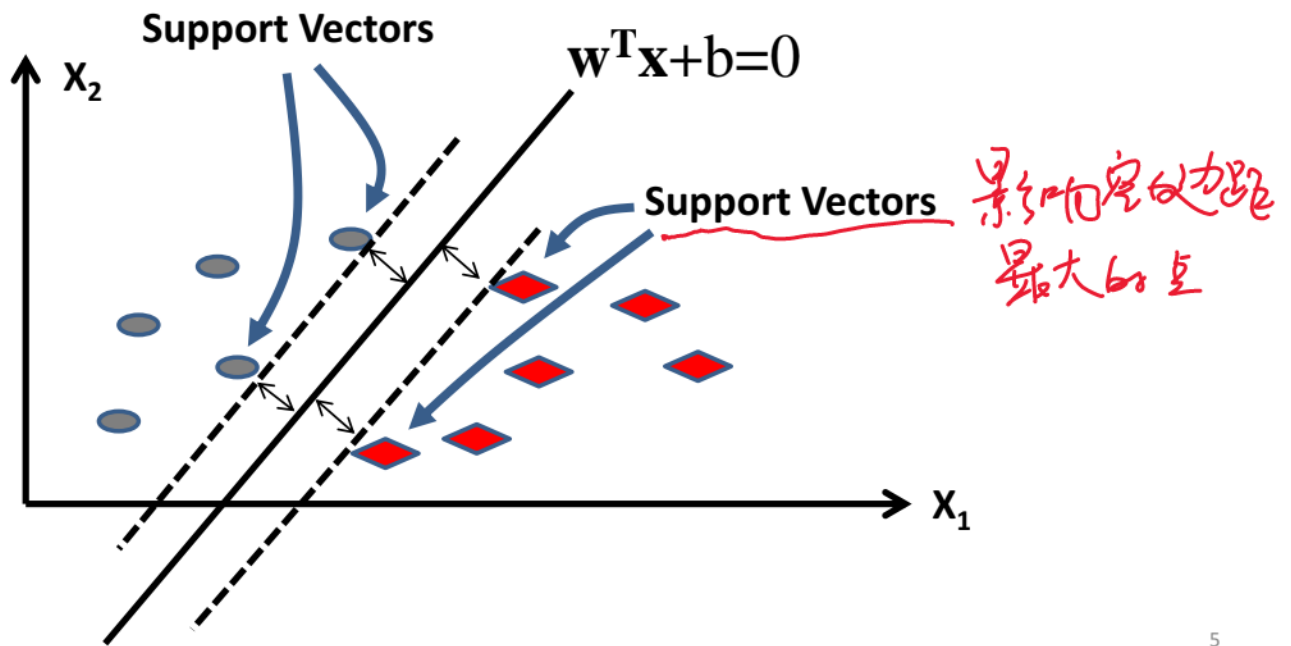
Feature Vector
(e.g., word counts)



$$\mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_M \end{pmatrix}$$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_M \end{pmatrix}$$

Only the Support Vectors Matter



Linear SVM

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

Parameters: \mathbf{w}, b

$f(\mathbf{X}) \geq 0 \Rightarrow \mathbf{X}$ is in category θ_1

$f(\mathbf{X}) < 0 \Rightarrow \mathbf{X}$ is in category θ_2

Training Data: $T = \{(\mathbf{x}_i, y_i)\}, i=1, \dots, |T|$. \mathbf{x}_i is a feature vector; $y_i \in \{-1, 1\}$

Goal 1: Correct labeling on training data:

If $y_i = 1 \Rightarrow \mathbf{w}^T \mathbf{x}_i + b \geq 1$

If $y_i = -1 \Rightarrow \mathbf{w}^T \mathbf{x}_i + b \leq -1$

} 结合两者为一个式子

Constraint

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

Goal 2: Maximize margin

所有权重的平方和

Large margin \Leftrightarrow Small $\mathbf{w}^T \mathbf{w}$

Objective

$$\text{Minimize } \Phi(\mathbf{w}) = \mathbf{w}^T \mathbf{w}$$

The optimization problem is quadratic programming with linear constraints

Linear SVM with Soft Margin

Classifier: $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b > 0$?

Parameters: \mathbf{w} , b

Training Data: $T = \{(\mathbf{x}_i, y_i)\}, i=1, \dots, |T|$.

Find \mathbf{w} , b , and ξ_i to minimize $\Phi(\mathbf{w}) = \mathbf{w}^T \mathbf{w} + C \sum_{i \in [1, |T|]} \xi_i$

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

Added to allow training errors

$\xi_i > 0$: 允许的偏差.
 $\xi_i = 0$: 不允许有偏差.



$C > 0$ is a parameter to control the trade-off between minimizing the errors and maximizing the margin

将减小误差和最大化 margin 之间作权衡.

The optimization problem is still quadratic programming with linear constraints

Summary of Text Categorization Methods

- Many methods are available, but no clear winner
 - All require effective feature representation (need domain knowledge)
 - It is useful to compare/combine multiple methods for a particular problem
- Most techniques rely on supervised machine learning and thus can be applied to **any** text categorization problem!
 - Humans annotate training data and design features
 - Computer optimizes the combination of features
 - Good performance requires 1) effective features and 2) plenty of training data
 - Performance is generally (much) more affected by the effectiveness of features than by the choice of a specific classifier

Summary of Text Categorization Methods (cont.)

- How to design effective features? (application-specific)
 - Analyze the categorization problem and exploit domain knowledge
 - Perform error analysis to obtain insights
 - Leverage machine learning techniques (e.g., feature selection, dimension reduction, deep learning)
- How to obtain “enough” training examples?
 - Low-quality (“pseudo”) training examples may be leveraged 半监督学习
 - Exploit unlabeled data (using semi-supervised learning techniques)
 - Domain adaptation/transfer learning (“borrow” training examples from a related domain/problem)

Suggested Reading

Manning, Chris D., Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge: Cambridge University Press, 2007.
(Chapters 13-15)