

Support Vector Machines

Lecture 4: Sept 20, 2013

CS886-2 Natural Language
Understanding
University of Waterloo

CS886 Lecture Slides (c) 2013 P. Poupart

1

Support Vector Machines

- Find linear separator that maximizes the distance (or margin) to closest data points
- Picture

CS886 Lecture Slides (c) 2013 P. Poupart

2

Notation

- Inputs: \mathbf{x} (column vector)
- Output: $y \in \{-1, 1\}$ (class)
- Weights: \mathbf{w} (vector perpendicular to linear separator)
- Constant: b (displacement of linear separator)
- Index: i (i^{th} instance in dataset)
 - \mathbf{x}_i and y_i : input and output of i^{th} instance

CS886 Lecture Slides (c) 2013 P. Poupart

3

Classification

- Linear separator: $\mathbf{w}^T \mathbf{x} + b = 0$

- Distance to linear separator:

$$\frac{y(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|} \quad \text{where } y \in \{-1, 1\}$$

- Classification:

$$y = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

$$\text{where } \text{sign}(\text{expr}) = \begin{cases} +1 & \text{if } \text{expr} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

CS886 Lecture Slides (c) 2013 P. Poupart

4

Margin

- Distance to linear separator:

$$\frac{y(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|} \text{ where } y \in \{-1, 1\}$$

- Margin: smallest distance to any data point

$$\frac{1}{\|\mathbf{w}\|} \left\{ \min_i y_i (\mathbf{w}^T \mathbf{x}_i + b) \right\}$$

- Maximum margin:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \left\{ \min_i y_i (\mathbf{w}^T \mathbf{x}_i + b) \right\}$$

CS886 Lecture Slides (c) 2013 P. Poupart

5

Maximum Margin

- Unique max margin linear separator

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \left\{ \min_i y_i (\mathbf{w}^T \mathbf{x}_i + b) \right\}$$

- Alternatively, we can fix the minimal distance to 1 and minimize $\|\mathbf{w}\|$

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

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

- This is a convex quadratic optimization problem that can easily be solved by many optimization packages

CS886 Lecture Slides (c) 2013 P. Poupart

6

Support Vectors

- Quadratic optimization problem

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

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

- Only the points where $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$ are necessary. These points define the active constraints and are known as the **support vectors**

CS886 Lecture Slides (c) 2013 P. Poupart

7

Linear Separability

- **What if the data is not linearly separable?**
- Two solutions:
 - Map data to new space (known as feature space)
 - Relax constraints with slack variables

CS886 Lecture Slides (c) 2013 P. Poupart

8

Feature Mapping

- Let $\phi(x)$ be a mapping that computes the feature representation of a point x
- When ϕ is non-linear, it transforms the data to become linearly separable.

CS886 Lecture Slides (c) 2013 P. Poupart

9

Dimensionality

- Feature spaces with high dimensionality
 - Data more likely to be linearly separable
 - If # of dimensions is higher than the number of data points, then the data is necessarily linearly separable
 - Problem: high dimensionality increases computation
- Solution: kernel trick
 - Computation depends on amount of data instead of dimensionality

CS886 Lecture Slides (c) 2013 P. Poupart

10

Kernel Function

- Let $\phi(x)$ be a set of basis functions that map inputs x to a feature space.
- In many algorithms, this feature space only appears in the dot product $\phi(x)^T \phi(x')$ of pairs inputs x, x' .
- Define the kernel function $k(x, x') = \phi(x)^T \phi(x')$ to be the dot product of any pair x, x' in feature space.
 - **We only need to know $k(x, x')$, not $\phi(x)$**

CS886 Lecture Slides (c) 2013 P. Poupart

11

Constructing Kernels

- Can we construct k directly without knowing ϕ ?
- Yes, any positive semi-definite k is fine since there is a corresponding implicit feature space. But positive semi-definiteness is not always easy to verify.
- Alternative, construct kernels from other kernels using rules that preserve positive semi-definiteness

CS886 Lecture Slides (c) 2013 P. Poupart

12

Rules to construct Kernels

- Let $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ be valid kernels
- The following kernels are also valid:
 1. $k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad \forall c > 0$
 2. $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad \forall f$
 3. $k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$ q is polynomial with coeffs ≥ 0
 4. $k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$
 5. $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
 6. $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$
 7. $k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}'))$
 8. $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}'$ \mathbf{A} is symmetric positive semi-definite
 9. $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b)$
 10. $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)$

CS886 Lecture Slides (c) 2013 P. Poupart

13

Common Kernels

- Polynomial kernel: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^d$
 - d is the degree
 - Feature space: all degree d products of entries in \mathbf{x}
 - Example: Let \mathbf{x} and \mathbf{x}' be two documents, then feature space could be all products of d word frequencies
- More general polynomial kernel:

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d \quad \text{with } c > 0$$
 - Feature space: all products of up to d entries in \mathbf{x}

CS886 Lecture Slides (c) 2013 P. Poupart

14

Common Kernels

- Gaussian Kernel: $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{2\sigma^2}\right)$
- Valid Kernel because:
- Implicit feature space is infinite!

CS886 Lecture Slides (c) 2013 P. Poupart

15

Dual representation

- Idea: reformulation where $\phi(\mathbf{x})$ always appears in a kernel $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$
- Approach: find the dual of the optimization problem
- Result: (sparse) kernel support vector machines

$$\begin{aligned} \max_{\mathbf{a}} \quad & \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_i a_i y_i = 0 \\ & a_i \geq 0 \quad \forall i \end{aligned}$$

CS886 Lecture Slides (c) 2013 P. Poupart

16

Dual derivation

- Transform constrained optimization

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t. } y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 \quad \forall n$$

into an unconstrained optimization problem

- Lagrangian

$$\max_{\mathbf{a}} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \mathbf{a}) \quad \text{s.t. } \mathbf{a} \geq 0$$

$$\text{where } L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i a_i \underbrace{[y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) - 1]}_{\text{penalty for violating the } i^{\text{th}} \text{ constraint}}$$

penalty for violating
the i^{th} constraint

CS886 Lecture Slides (c) 2013 P. Poupart

17

Dual derivation

- Solve inner minimization: $\min_{\mathbf{w}, b} L(\mathbf{w}, b, \mathbf{a})$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i a_i [y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) - 1]$$

- Set derivatives to 0

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_i a_i y_i \phi(\mathbf{x}_i)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow 0 = \sum_i a_i y_i$$

- Eliminate \mathbf{w} and b based on these conditions to obtain:

$$L(\mathbf{a}) = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

CS886 Lecture Slides (c) 2013 P. Poupart

18

Dual Problem

- We are then left with an optimization in a only known as the **dual problem**

$$\begin{aligned} \max_a L(\mathbf{a}) &= \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t. } \sum_i a_i y_i &= 0 \\ a_i &\geq 0 \quad \forall i \end{aligned}$$

- Sparse optimization:** many a_i 's are 0

CS886 Lecture Slides (c) 2013 P. Poupart

19

Classification

- Primal problem

$$y = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

- Dual problem

$$y = \text{sign}\left(\sum_i a_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b\right)$$

$$y = \text{sign}\left(\sum_i a_i y_i k(\mathbf{x}_i, \mathbf{x}) + b\right)$$

CS886 Lecture Slides (c) 2013 P. Poupart

20

Generalization

- Support vector machines generalize quite well
 - i.e., overfitting is rare
- Reason: maximizing the margin is equivalent to minimizing an upper bound on the worst case loss (worst loss for any underlying input distribution).

Case Study: Text Categorization

- T. Joachims, *Text Categorization with Support Vector Machines: Learning with Many Relevant Features*. Proceedings of the European Conference on Machine Learning (ECML), Springer, 1998.
- Early success that helped SVMs become popular

Text Categorization

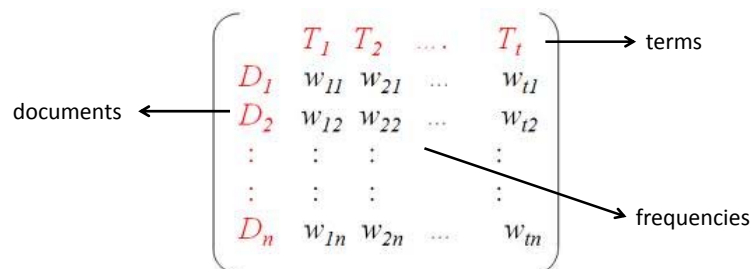
- **Problem:** how to categorize a new article as finance, sports, politics, science, health, etc.?
- **Idea:** train a classifier with archives of news articles that have already been classified

CS886 Lecture Slides (c) 2013 P. Poupart

23

Representation

- How should we represent a document?
- Idea: vector of word counts (vector space model)



CS886 Lecture Slides (c) 2013 P. Poupart

24

Challenges

- **High dimensional input space:**
 - Length of vector is # of words in dictionary (e.g., 10,000)
- **Few irrelevant features:**
 - Most words carry some information that reflect their meaning
- Need an approach that scales well with input dimensionality: **support vector machines**

CS886 Lecture Slides (c) 2013 P. Poupart

25

Experiment

- [Joachim 98]
 - Data: Reuters dataset
 - Compare precision/recall breakeven point
 - i.e., precision = recall
 - Precision: $\frac{|{\text{relevant docs}}|n|{\text{retrieved docs}}|}{|{\text{relevant docs}}|}$
 - Recall: $\frac{|{\text{relevant docs}}|n|{\text{retrieved docs}}|}{|{\text{retrieved docs}}|}$
 - Algorithms
 - Naïve Bayes: 72.0%
 - Decision trees: 79.4%
 - Rochio: 79.9%
 - K-Nearest Neighbors: 82.3%
 - SVMs: 86.0% (polynomial kernel), 86.4% (Gaussian kernel)

CS886 Lecture Slides (c) 2013 P. Poupart

26

SVM summary

- Find (generalized) linear separator
 - Dual representation (kernel): non-linear separator
- Unique max-margin separator
 - Good generalization
- Convex quadratic optimization
 - Polynomial complexity
 - Global optimality
- Sparse optimization
 - many variables are 0