

CS485/685

Lecture 12: Feb 9, 2012

Support Vector Machines (continued)
[B] Section 7.1

CS485/685 (c) 2012 P. Poupart

1

Overlapping Class Distributions

- So far we assumed that data is linearly separable
 - High dimensions help for linear separability, but may hurt for generalization
- But what if the data is noisy (mistakes or outliers)
 - Constraints should allow misclassifications
- Picture

CS485/685 (c) 2012 P. Poupart

2

Soft margin

- Idea: relax constraints by introducing slack variables
 $\xi_n \geq 0$

$$y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n \quad \forall n$$
- Picture:

CS485/685 (c) 2012 P. Poupart

3

Soft margin classifier

- New optimization problem:

$$\min_{\mathbf{w}, b, \xi} \quad C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

s.t. $y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n$
and $\xi_n \geq 0 \quad \forall n$

- where $C > 0$ controls the trade-off between the slack variable penalty and the margin

CS485/685 (c) 2012 P. Poupart

4

Soft margin classifier

- Notes:
 1. Since $\sum_n \xi_n$ is an upper bound on the # of misclassifications, C can also be thought as a regularization coefficient that controls the trade-off between error minimization and model complexity
 2. When $C \rightarrow \infty$, then we recover the original hard margin classifier
 3. Soft margins handle minor misclassifications, but the classifier is still very sensitive to outliers

CS485/685 (c) 2012 P. Poupart

5

Support Vectors

- As before support vectors correspond to active constraints

$$y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1 - \xi_n$$

– i.e., all points that are in the margin or misclassified

- Picture:

CS485/685 (c) 2012 P. Poupart

6

Dual derivation

- Transform constrained optimization

$$\min_{\mathbf{w}, b, \xi} C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n \text{ and } \xi_n \geq 0 \quad \forall n$$

into an unconstrained optimization problem

- Lagrangian

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

$$\text{where } L(\mathbf{w}, b, \xi, \mathbf{a}, \mu) = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

$$- \sum_n a_n [y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1 + \xi_n] - \sum_n \mu_n \xi_n$$

CS485/685 (c) 2012 P. Poupart

7

Dual derivation

- Solve $\min_{\mathbf{w}, b, \xi} L(\mathbf{w}, b, \xi, \mathbf{a}, \mu)$ by setting derivatives to 0

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_n a_n y_n \phi(\mathbf{x}_n)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow 0 = \sum_n a_n y_n$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n$$

- Eliminate \mathbf{w}, b, ξ and μ based on these conditions:

$$L(\mathbf{a}) = \sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m y_n y_m k(\mathbf{x}_n, \mathbf{x}_m)$$

CS485/685 (c) 2012 P. Poupart

8

Dual Problem

- The resulting **dual problem** is

$$\begin{aligned} \max_{\mathbf{a}} L(\mathbf{a}) \\ \text{s.t. } \sum_n a_n y_n = 0 \\ 0 \leq \underbrace{a_n}_{\substack{\uparrow \\ \text{NB: Same optimization problem as for hard margins,} \\ \text{except for the upper bound on } a_n}} \leq C \quad \forall n \end{aligned}$$

- NB: Same optimization problem as for hard margins, except for the upper bound on a_n

CS485/685 (c) 2012 P. Poupart

9

Dual Problem

- Notes:

- $a_n = 0 \Rightarrow$ irrelevant point
- $a_n > 0 \Rightarrow$ support vector
 - $a_n < C \Rightarrow$ point on the margin
 - $a_n = C \Rightarrow$ point inside the margin or misclassified

CS485/685 (c) 2012 P. Poupart

10

Classification

- Same as for hard margins

- Primal problem

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

- Dual problem

$$y = \text{sign}\left(\sum_n a_n y_n k(\mathbf{x}_n, \mathbf{x}) + b\right)$$

CS485/685 (c) 2012 P. Poupart

11

Multiclass SVMs

- Three methods:
 1. One-against-all: for R classes, train R SVMs to distinguish each class from the rest
 2. Continuous ranking: single SVM that returns a continuous value to rank all classes
 3. Pairwise comparison: train $O(R^2)$ SVMs to compare each pair of classes

CS485/685 (c) 2012 P. Poupart

12

One-Against-All

- For R classes, train R SVMs to distinguish each class from the rest
- Picture:
- Problem: what if different classes are returned by different SVMs?

CS485/685 (c) 2012 P. Poupart

13

Continuous Ranking

- Single SVM that returns a continuous value to rank all classes
- Picture:
- Most popular approach today

CS485/685 (c) 2012 P. Poupart

14

Pairwise Comparison

- Train $O(R^2)$ SVMs to compare each pair of classes
- Picture:
- Problem: how do we pick the best class?

CS485/685 (c) 2012 P. Poupart

15

Continuous Ranking

- Idea: instead of computing the sign of a linear separator, compare the values of linear functions for each class r

$$y = \operatorname{argmax}_r \mathbf{w}_r^T \phi(\mathbf{x}) + b_r$$

CS485/685 (c) 2012 P. Poupart

16

Multiclass Margin

- For each class r define a linear constraint:

$$\mathbf{w}_y^T \phi(\mathbf{x}) + b_y + \delta_{yr} - \mathbf{w}_r^T \phi(\mathbf{x}) - b_r \geq 1 \quad \forall r$$
 where $\delta_{yr} = \begin{cases} 1 & y = r \\ 0 & \text{otherwise} \end{cases}$
- This guarantees a margin of at least 1

CS485/685 (c) 2012 P. Poupart

17

Multiclass Classification

- Optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \sum_r \|\mathbf{w}_r\|^2$$
 s.t. $\mathbf{w}_{y_n}^T \phi(\mathbf{x}_n) + b_{y_n} + \delta_{y_n r} - \mathbf{w}_r^T \phi(\mathbf{x}) - b_r \geq 1 \quad \forall n, r$
- Equivalent to binary SVM when we have only two classes

CS485/685 (c) 2012 P. Poupart

18

Overlapping classes

- Add slack variables:

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{b}, \xi} \quad & C \sum_n \xi_n + \frac{1}{2} \sum_r \|\mathbf{w}_r\|^2 \\ \text{s.t.} \quad & \mathbf{w}_{y_n}^T \phi(\mathbf{x}_n) + b_{y_n} + \delta_{y_n r} - \mathbf{w}_r^T \phi(\mathbf{x}) - b_r \geq 1 - \xi_n \quad \forall n, r \end{aligned}$$

- Equivalent to binary SVM when we have only two classes

CS485/685 (c) 2012 P. Poupart

19

Dual representation

- Kernelized form

$$\begin{aligned} \max_{\{\mathbf{a}_n\}} \quad & \sum_n \boldsymbol{\tau}_n^T \mathbf{e}_{y_n} - C \frac{1}{2} \sum_n \sum_m \boldsymbol{\tau}_n^T \boldsymbol{\tau}_m k(\mathbf{x}_n, \mathbf{x}_m) \\ \text{s.t.} \quad & \boldsymbol{\tau}_n \leq \mathbf{e}_{y_n} \quad \text{and} \quad \sum_r \tau_{nr} = 0 \end{aligned}$$

$$\text{where } e_{y_n r} = \begin{cases} 1 & \text{when } y_n = r \\ 0 & \text{otherwise} \end{cases}$$

$$\text{and } \boldsymbol{\tau}_n = \mathbf{e}_{y_n} - \mathbf{a}_n$$

CS485/685 (c) 2012 P. Poupart

20

Classification

- Primal

$$y = \operatorname{argmax}_r \mathbf{w}_r^T \phi(\mathbf{x}) + b_r$$

- Dual

$$y = \operatorname{argmax}_r \sum_n \tau_{y_n r} k(\mathbf{x}_n, \mathbf{x})$$