Support Vector Machine

Mohammad Emtiyaz Khan EPFL

Oct 29, 2015



©Mohammad Emtiyaz Khan 2015

Motivation

By changing the cost function from Logistic to Hinge, we obtain SVMs. However, the resulting cost is difficult to optimize. We will use duality (similar to kernelized ridge) to show a surprising result: the solution to the dual problem is sparse. The non-zero entries will be our support vectors.

Support vector machine $[1+\varphi\beta]_{+}$

Throughout, we will work with a classification problem and assume that $y_n \in \{-1, +1\}$ (instead of $\in \{0,1\}$). Also, we will work with $\phi(\mathbf{x})$ instead of $\mathbf{\widetilde{x}}$ (bias included). A

SVM optimizes the following cost:

$$\underset{\beta}{\rightarrow} \min_{\beta} \sum_{n=1}^{N} [1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta}]_+ + \frac{\lambda}{2} \sum_{j=1}^{M} \beta_j^2 \times \mathbf{y}_j^{M} \mathbf{\beta}_j^2$$

where the first term is the Hinge loss defined as $[t]_{+} = \max(0, t)$. A "conventional" definition is shown below:

$$\min_{\beta} \sum_{n=1}^{N} \widetilde{C}[1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta}]_+ + \frac{1}{2} \sum_{j=1}^{M} \beta_j^2.$$

Logistic likelihood (recap)

$$\frac{1}{N}\sum_{n\geq 1}^{N} -\log \frac{1}{2} \left(\frac{y_n}{\varphi_n}\right) + \frac{1}{2}\sum_{j=1}^{N} \beta_j^2$$

$$\int_{=}^{\infty} \left(\widetilde{p_n} p + \log(1 + e^{\widetilde{p_n}} p) + \frac{y_{n=1}}{y_{n=1}}\right) + \log(1 + e^{\widetilde{p_n}} p) + \log(1 + e^{\widetilde{p_n}} p)$$

$$\frac{\text{Hinge loss}}{[t]_{+}} = \max(o_{1}t)$$

O class -1

1.1

$$\min \sum_{\underline{P}} \left[M - M_{n} \widetilde{P}_{n} (M \underline{P}) \right]$$

 $\left[M-y_{n}\widetilde{\beta}_{n}\widetilde{\beta}\right]$

1



Notice the margin in the Hinge loss. SVM is a maximum margin method.



See section 14.5.2.2 of KPM book.

Issues with optimization

Is this function convex? Is it differentiable?

$$\min_{\beta} \sum_{n=1}^{N} C[1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta}]_+ + \frac{1}{2} \sum_{j=1}^{M} \beta_j^2 := g(\boldsymbol{\beta})$$

Duality: the big picture

Let us say that we are interested in optimizing a function $g(\boldsymbol{\beta})$ and it is a difficult problem. Define an auxiliary function $G(\boldsymbol{\beta}, \boldsymbol{\alpha})$ as follows:

$$g(\boldsymbol{\beta}) = \max_{\alpha} G(\boldsymbol{\beta}, \boldsymbol{\alpha}).$$

Three questions.

- 1. How do you set $G(\boldsymbol{\alpha}, \boldsymbol{\beta})$?
- 2. When is it OK to switch max and min?
- 3. When is the dual better than the primal, and why?
- **Q1:** How do you set $G(\boldsymbol{\alpha}, \boldsymbol{\beta})$? $C[v_n]_+ = \max(0, Cv_n) = \max_{\alpha_n} \alpha_n v_n$ where $\alpha_n \in [0, C]$ $C[1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta}]_+ = \max_{\alpha_n \in [0, C]} \alpha_n (1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta})$

()

Max

x

We can rewrite the problem as:

$$\min_{\boldsymbol{\beta}} \max_{\boldsymbol{\alpha} \in [0,C]^N} \sum_{n=1}^N \alpha_n (1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta}) + \frac{1}{2} \sum_{j=1}^M \beta_j^2$$

This is differentiable, convex in β and concave in α .

Q2: When is it OK to switch max and min? Using a minimax theorem, it is OK to do so when $G(\alpha, \beta)$ is convex in β and concave in α , and the sets over which α and β are optimized are convex. In this case, we have:

 $\min_{\beta} \max_{\alpha} G(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \max_{\alpha} \min_{\beta} G(\boldsymbol{\beta}, \boldsymbol{\alpha})$

See Bertsekas' "Nonlinear Programming" for many more variants of this type of duality.



Switch the min and max and optimize w.r.t. β to get the dual.

$$\max_{\alpha \in [0,C]^N} \min_{\underline{\beta}} \sum_{n=1}^N \alpha_n (1 - y_n \widetilde{\boldsymbol{\phi}}_n^T \boldsymbol{\beta}) + \frac{1}{2} \sum_{j=1}^M \beta_j^2$$

Take derivative w.r.t. $\boldsymbol{\beta}$:

Derive this at home !

$$\frac{\partial G}{\partial \boldsymbol{\beta}} = -\left[\sum_{n=1}^{N} \alpha_n y_n \widetilde{\boldsymbol{\phi}}_n\right] + \begin{bmatrix} 0\\ \boldsymbol{\beta}_{1:M} \end{bmatrix} \checkmark$$

where $\boldsymbol{\beta}_{1:M}$ is a vector of all β_j except β_0 .

Equating this to $\underline{0}$, we get:

$$\boldsymbol{\beta}_{1:M}^* = \sum_{n=1}^{N} \alpha_n y_n \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T \operatorname{diag}(\mathbf{y}) \boldsymbol{\alpha} = \boldsymbol{\Phi}^T \mathbf{Y} \boldsymbol{\alpha}$$
$$\boldsymbol{\alpha}^T \mathbf{y} = 0$$

where $\mathbf{Y} := \operatorname{diag}(\mathbf{y})$.

Plugging $\boldsymbol{\beta}^*$ back in, we get the dual problem: $\max_{\alpha \in [0,C]^N} \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Y} \boldsymbol{\Phi} \boldsymbol{\Phi}^T \mathbf{Y} \boldsymbol{\alpha}$ subject to $\boldsymbol{\alpha}^T \mathbf{y} = 0$ **Q3:** When is the dual better than the primal and why?

(1) The dual is a differentiable (but constrained) least-squares problem.

$$\max_{\alpha \in [0,C]^N} \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha},$$

where $\mathbf{Q} := \operatorname{diag}(\mathbf{y}) \mathbf{\Phi} \mathbf{\Phi}^T \operatorname{diag}(\mathbf{y})$. Optimization is super easy using Sequential Minimal Optimization (SMO). See Wikipedia for details. Read the details on your own

Summary: Take two variables α_1 and α_2 and fix others. This gives rise to a 1-D quadratic problem. Minimize and repeat by choosing two different elements of α .

(2) The dual is naturally kernelized (just like the kernelized ridge) with $\mathbf{K} := \mathbf{\Phi} \mathbf{\Phi}^T$.

(3) The solution α is sparse, and is non-zero only for the training examples that are instrumental in determining the decision boundary. Recall that α_n is the slope of lines that are lower bounds to the Hinge loss.

$$C[1 - y_n f_n]_+ = \max_{\alpha_n \in [0,C]} \alpha_n (1 - y_n f_n)$$

There are 3 kinds of data vectors $\boldsymbol{\phi}_n$.

- 1. Not support vectors. Examples that lie outside the margin, therefore $\alpha_n = 0$.
- 2. Essential support vectors. Examples that lie right on the margin, therefore $\alpha_n \in (0, C)$.
- 3. Bound support vectors. Examples that lie inside the margin, therefore $\alpha_n = C$.





Issues with SVM

- There are no obvious probabilistic interpretation of SVM.
- Extension to multiclass is difficult (see Section 14.5.2.4 of KPM book).
- Choosing C is difficult in the presence of Kernels.
- The method does not work for positive semidefinite Kernels.

To do

- 1. Understand and visualize hinge loss and the margin.
- 2. Get comfortable with duality. Work out the derivation for SVM.
- 3. Clearly understand the reasons why dual is better than the primal.
- 4. What does "support vector" mean? Why do they arise? Where do they lie in the data space?
- 5. Read about SMO algorithm from Wikipedia and implement it.
- 6. Read about SVM for regression (section 14.5.1 of KPM).
- 7. Read Section 14.5.2.4 of KPM book and understand why extension of SVM to multiclass is difficult.
- 8. Read about maximum-margin methods in section 14.5.2.2 of KPM book.
- 9. Resource: SVM tutorial by Christopher J.C. Burges at http:// research.microsoft.com/pubs/67119/svmtutorial.pdf
- 10. Read SVM from HTF.

$$\sum_{n=1}^{N} \left[1 - y_n \tilde{y}_n^T \beta \right]_{+} + \frac{\lambda}{2} \sum_{j=1}^{M} \beta_j^2$$
For linearly separable case, all n are correctly designed
$$\Rightarrow y_n \tilde{y}_n^T \beta \ge 1 \text{ str} \qquad -\frac{1}{1 + 1 + 1}$$

$$\Rightarrow An \text{ equivalent problem } y_{n=-1} \qquad y_n = +1$$

$$\begin{bmatrix} \min_{\beta} \lambda \sum_{j=1}^{2} \beta_j^2, \text{ st. } y_n (\tilde{p}_n^T \beta) \ge 1 \end{bmatrix}^{p_n^T \beta \ge 0} - (1 \qquad \tilde{p}_n^T \beta \ge 0)$$

$$x^T \beta + \beta_0 = 0 \qquad + 1$$

$$M = \frac{1}{\|\beta\|}$$

$$x^T \beta + \beta_0 = 0 \qquad + 1$$

$$M = \frac{1}{\|\beta\|}$$

$$M = \frac{1}{\|\beta\|}$$

FIGURE 12.1. Support vector classifiers. The left panel shows the separable case. The decision boundary is the solid line, while broken lines bound the shaded maximal margin of width $2M = 2/||\beta||$. The right panel shows the nonseparable (overlap) case. The points labeled ξ_j^* are on the wrong side of their margin by an amount $\xi_j^* = M\xi_j$; points on the correct side have $\xi_j^* = 0$. The margin is maximized subject to a total budget $\sum \xi_i \leq \text{constant. Hence } \sum \xi_j^*$ is the total distance of points on the wrong side of their margin.

When the problem is not linearly separable,
min
$$\frac{\lambda}{2} \sum_{j} \beta_{j}^{2}$$
, s.t. $y_{m}(\rho_{m}^{T}\beta_{j}) \ge 1 - e_{n}$
 $+ \sum_{j} e_{n}^{2}$