

Lecture 12: Support Vector Machines

CS 475: Machine Learning

Raman Arora

March 8, 2017



Review

Review: optimal separating hyperplane

- Decision boundary parametrized as $\mathbf{w} \cdot \mathbf{x} + w_0 = 0$
- “confidence” = $y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0)$
- Distance from the hyperplane $\frac{1}{\|\mathbf{w}\|} y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0)$
- We seek $\operatorname{argmax}_{\mathbf{w}, w_0} \left\{ \frac{1}{\|\mathbf{w}\|} \min_i y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \right\}$
- Assuming $\forall i, y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) > 0$, we can rescale $\|\mathbf{w}\|, w_0$ so that

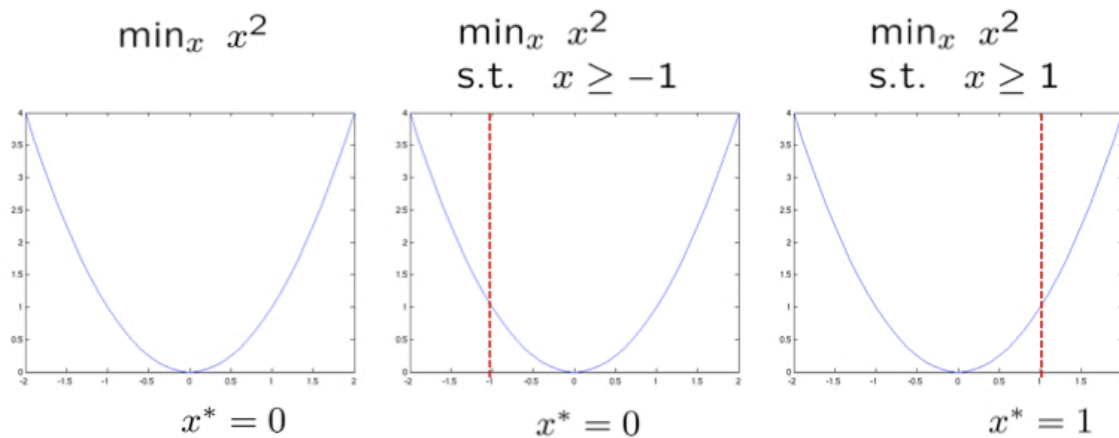
$$\min_i y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) = 1.$$

- Then, the optimization becomes:

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{w}, w_0} \frac{1}{\|\mathbf{w}\|} && \text{s.t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \forall i = 1, \dots, N. \\ \Rightarrow \operatorname{argmin}_{\mathbf{w}} && \|\mathbf{w}\|^2 && \text{s.t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \forall i = 1, \dots, N. \end{aligned}$$



Review: Constrained optimization



Review: large margin setup

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_{j=1}^d w_j^2,$$

subject to $y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i) - 1 \geq 0, \quad i = 1, \dots, N.$

- We will associate with each constraint the loss

$$\max_{\alpha_i \geq 0} \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)] = \begin{cases} 0, & \text{if } y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i) - 1 \geq 0, \\ \infty & \text{otherwise (constraint violated).} \end{cases}$$

- We can reformulate our problem now:

$$\min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \max_{\alpha_i \geq 0} \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)] \right\}$$

Max-margin optimization

- We want all the constraint terms to be zero:

$$\begin{aligned}
 & \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \max_{\alpha_i \geq 0} \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)] \right\} \\
 &= \min_{\mathbf{w}} \max_{\alpha \geq 0} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)] \right\} \\
 &= \max_{\alpha \geq 0} \min_{\mathbf{w}} \underbrace{\left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)] \right\}}_{J(\mathbf{w}, w_0; \alpha)}.
 \end{aligned}$$

- Why could we switch min and max? convexity!



Strategy for optimization

- We need to find

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} \underbrace{\left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)] \right\}}_{J(\mathbf{w}, w_0; \alpha)}.$$

- We will first fix α and treat $J(\mathbf{w}, w_0; \alpha)$ as a function of \mathbf{w}, w_0 .
 - Find *functions* $\mathbf{w}(\alpha), w_0(\alpha)$ that attain the minimum $\forall \alpha$.
- Next, maximize $J(\mathbf{w}(\alpha), w_0(\alpha); \alpha)$ as a function of α .
- In the end, the solution is given by α^* ;
find $\mathbf{w}(\alpha^*)$ and $w_0(\alpha^*)$ by substitution.



Minimizing J with respect to \mathbf{w}, w_0

- For fixed α we can minimize

$$J(\mathbf{w}, w_0; \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i [1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)]$$

by setting derivatives w.r.t. w_0, \mathbf{w} to zero:

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, w_0; \alpha) = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = 0,$$

$$\frac{\partial}{\partial w_0} J(\mathbf{w}, w_0; \alpha) = - \sum_{i=1}^N \alpha_i y_i = 0.$$

- Note that the bias term w_0 dropped out but has produced a “global” constraint on α .



Solving for α

$$\underbrace{\mathbf{w}(\alpha) = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i}_{\text{later: Representer theorem!}}, \quad \sum_{i=1}^N \alpha_i y_i = 0.$$

- Now can (with a bit of algebra) substitute this solution into

$$\begin{aligned} & \max_{\alpha \geq 0, \sum_i \alpha_i y_i = 0} \left\{ \frac{1}{2} \|\mathbf{w}(\alpha)\|^2 + \sum_{i=1}^N \alpha_i [1 - y_i(w_0(\alpha) + \mathbf{w}(\alpha) \cdot \mathbf{x}_i)] \right\} \\ &= \max_{\alpha \geq 0, \sum_i \alpha_i y_i = 0} \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \right\}. \end{aligned}$$



Max-margin and quadratic programming

- We started by writing down the max-margin problem and arrived at the *dual problem* in α :

$$\begin{aligned} & \max \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \right\} \\ & \text{subject to } \sum_{i=1}^N \alpha_i y_i = 0, \alpha_i \geq 0 \text{ for all } i = 1, \dots, N. \end{aligned}$$

- Solving this *quadratic program* with linear constraints yields α^* .
- We substitute α^* back to get \mathbf{w} :

$$\hat{\mathbf{w}} = \mathbf{w}(\alpha^*) = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$



Maximum margin decision boundary

$$\hat{\mathbf{w}} = \mathbf{w}(\alpha^*) = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$

- Suppose that, under the optimal solution, the margin (distance to the boundary) of a particular \mathbf{x}_i is

$$y_i (w_0 + \hat{\mathbf{w}} \cdot \mathbf{x}_i) > 1.$$

- Then, necessarily, $\alpha_i^* = 0 \Rightarrow$ not a support vector.
- The direction of the max-margin decision boundary is

$$\hat{\mathbf{w}} = \sum_{\alpha_i^* > 0} \alpha_i^* y_i \mathbf{x}_i.$$

- w_0 is set by making the margin equidistant to two classes.



Support vectors

$$\hat{\mathbf{w}} = \sum_{\alpha_i > 0} \alpha_i y_i \mathbf{x}_i.$$

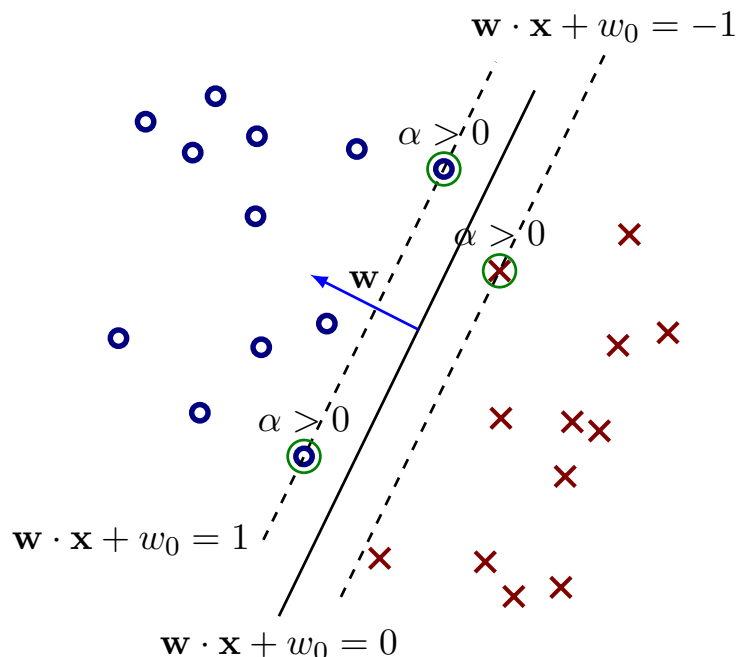
- Given a test example \mathbf{x} , it is classified by

$$\begin{aligned} \hat{y} &= \text{sign}(\hat{w}_0 + \hat{\mathbf{w}} \cdot \mathbf{x}) \\ &= \text{sign}\left(\hat{w}_0 + \left(\sum_{\alpha_i > 0} \alpha_i y_i \mathbf{x}_i\right) \cdot \mathbf{x}\right) \\ &= \text{sign}\left(\hat{w}_0 + \sum_{\alpha_i > 0} \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}\right) \end{aligned}$$

- The classifier is based on the expansion in terms of dot products of \mathbf{x} with support vectors.



SVM geometry



- Support vectors:

$$\alpha_i > 0$$

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) = 1$$

- Other examples:

$$\alpha_i = 0$$

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) > 1$$



Non-separable case

- Not linearly separable data: we can no longer satisfy $y_i (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1$ for all i .
- Recall the constraint-based terms in separable case:

$$\max_{\alpha \geq 0} \sum_i \alpha_i [1 - y_i (w_0 + \mathbf{w} \cdot \mathbf{x}_i)]$$

- We can no longer have $\alpha \geq 0$ if constraint violation is unavoidable; would yield $J = \infty$
- We will set maximum penalty on constraint violation:

$$\max_{0 \leq \alpha \leq C} \sum_i \alpha_i [1 - y_i (w_0 + \mathbf{w} \cdot \mathbf{x}_i)]$$



Slack variables

- We introduce *slack variables* to satisfy margin constraints

$$y_i (w_0 + \mathbf{w} \cdot \mathbf{x}_i) - 1 + \xi_i \geq 0, \quad \xi_i \geq 0.$$

- We want ξ_i to capture the *minimum* amount we need to fix:

$$\xi_i = \max \{0, 1 - y_i (w_0 + \mathbf{w} \cdot \mathbf{x}_i)\}$$

note: ξ_i is really a function of \mathbf{w}

- Our objective now:

$$\min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \right\}.$$



Non-separable case: solution

$$\min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \right\}.$$

- We can solve this using Lagrange multipliers
 - Introduce additional multipliers for the $\xi \geq 0$.
- The resulting dual problem:

$$\max \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \right\}$$

subject to $\sum_{i=1}^N \alpha_i y_i = 0, \quad 0 \leq \alpha \leq C.$

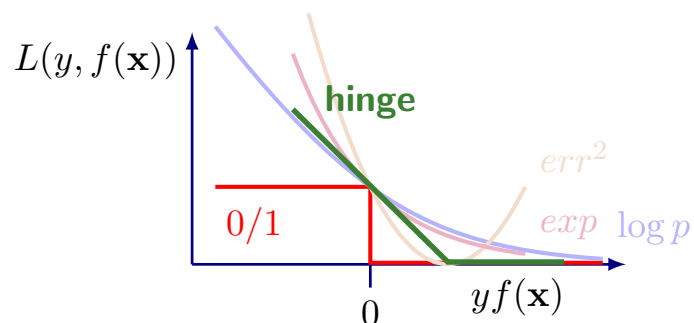
Loss in SVM

$$\min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \right\}$$

- L_2 -regularized loss, measured as

$$\sum_{i=1}^N \xi_i = \sum_{i=1}^N \max \{0, 1 - y_i(w_0 + \mathbf{w} \cdot \mathbf{x}_i)\}$$

- This surrogate loss is known as *hinge loss*



Solving SVM in the primal

- Setting $\lambda = 2/C$ we get

$$\text{primal:} \quad \min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \max\{0, 1 - y_i \mathbf{w} \cdot \mathbf{x}_i\}$$

- Traditional tactic: write the dual, solve using QP
- Alternative: optimize the primal directly using gradient descent
- Problem: hinge loss is not differentiable at $y\mathbf{w} \cdot \mathbf{x} = 1$
- Solution: *sub*gradient descent