

SVM as a Convex Optimization Problem

Leon Gu

CSD, CMU

Convex Optimization

- ▶ *Convex set*: the line segment between any two points lies in the set.
- ▶ *Convex function*: the line segment between any two points $(x, f(x))$ and $(y, f(y))$ lies on or above the graph of f .
- ▶ Convex optimization

$$\text{minimize} \quad f_0(x) \quad (1)$$

$$\text{s.t.} \quad f_i(x) \leq 0 \quad i = 1, \dots, m \quad (2)$$

$$h_i(x) = 0 \quad i = 1, \dots, p \quad (3)$$

- ▶ f_0 and f_i convex, h_i linear.
 - ▶ convex objective function, convex domain (feasible set).
 - ▶ any local optimum is also a global optimum.
- ▶ Operations preserve convexity
 - ▶ (for convex sets) **intersection**, affine transformation, perspective transformation, ...
 - ▶ (for convex functions) nonnegative weighted sum, **maximum and supremum**, composition with affine functions, composition with monotonic convex/concave functions, ...

Optimal Separating Hyperplane

Suppose that our data set $\{x_i, y_i\}_{i=1}^N$ is *linear separable*. Define a hyperplane by

$$\{x : f(x) = \beta^T x + \beta_0 = \beta^T (x - x_0) = 0\} \text{ where } \|\beta\| = 1.$$

- ▶ $f(x)$ is the sign distance to the hyperplane.
- ▶ we can define a classification rule induced by $f(x)$: $\text{sgn}[\beta^T (x - x_0)]$;

Define the **margin** of $f(x)$ to be the minimal $y f(x)$ through the data

$$C = \min_i y_i f(x_i)$$

A *optimal separating hyperplane* is the hyperplane that maximizes the margin,

$$\max_{\beta, \beta_0, \|\beta\|=1} C, \text{ s.t. } y_i(\beta^T x_i + \beta_0) \geq C, \quad i = 1, \dots, N$$

We can get rid of the norm constraint on β ,

$$\frac{1}{\|\beta\|} y_i (\beta^T x_i + \beta_0) \geq C$$

and arbitrarily set $\|\beta\| = 1/C$, then we can rephrase the problem as

$$\min_{\beta, \beta_0} \|\beta\|, \quad \text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1, \quad i = 1, \dots, N$$

This is a **convex optimization** problem.

Soft Margin SVM

The data is not always perfect. We need to extend optimal separating hyperplane to non-separable cases. The trick is to relax the margin constraints by introducing some “slack” variables.

$$\text{minimize} \quad \|\beta\| \quad \text{over } \beta, \beta_0 \quad (4)$$

$$\text{s.t.} \quad y_i(\beta^T x_i + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, N \quad (5)$$

$$\xi_i \geq 0; \quad \sum_{i=1}^N \xi_i \leq Z \quad (6)$$

- ▶ still convex.
- ▶ $\xi_i > 1$ – misclassification
 $\xi_i > 0$ – the data is correctly classified but lies in the margin.
- ▶ Z is a tuning parameter.

How to solve it? Use Lagrange/duality theory.

Lagrangian Theory

Lagrangian theory characterizes the solution of a constrained optimization problem. Recall the *primal* problem:

$$\text{minimize} \quad f_0(x) \quad (7)$$

$$\text{s.t.} \quad f_i(x) \leq 0 \quad i = 1, \dots, m \quad (8)$$

$$h_i(x) = 0 \quad i = 1, \dots, p \quad (9)$$

The stationary points are given by

$$\frac{df_0(x)}{dx} + \sum_{i=1}^m \lambda_i \frac{df_i(x)}{dx} + \sum_{i=1}^p \nu_i \frac{dh_i(x)}{dx} = 0$$

where λ, ν are free parameters called *Lagrange multipliers*. Accordingly, we define *Lagrangian prime function* (or *Lagrangian*) as

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x).$$

We define *Lagrangian dual function* $g(\lambda, \nu)$ as

$$g(\lambda, \nu) = \inf_{x \in \mathcal{X}} L(x, \lambda, \nu).$$

The so-called *Lagrangian dual problem* is the following:

$$\text{maximize} \quad g(\lambda, \nu) \tag{10}$$

$$\text{s.t.} \quad \lambda > 0. \tag{11}$$

The *weak duality theorem* says

$$g(\lambda, \nu) \leq f_0(x^*) \quad \text{for all } \lambda \text{ and } \nu$$

In other words, maximizing $g(\lambda, \nu)$ over λ and ν produce a bound on $f_0(x^*)$ (Note that $g(\lambda, \nu)$ is piecewise linear and convex). The difference between $g(\lambda^*, \nu^*)$ and $f_0(x^*)$ is called the “duality gap”.

K.K.T. Conditions

Slater's Theorem (Strong Duality Theorem) says: if the constraint functions are affine, the duality gap is zero.

Then, *K.K.T. conditions* provide **necessary and sufficient** conditions for a point x^* to be an optimum

$$\begin{array}{ll} \left. \frac{\partial L(x, \lambda^*, \nu^*)}{\partial x} \right|_{x^*} = 0 & \text{first-order derivative of optimality} \\ \lambda_i^* f_i(x^*) = 0 & \text{complementary slackness conditions} \\ \lambda_i^* \geq 0 & \text{dual constraints} \\ f_i(x^*) \leq 0 & \text{prime constraints} \\ h_i(x^*) = 0 & \text{prime constraints} \end{array}$$

Remarks: complementary slackness conditions are directly related to support vectors.

The Dual Problem

Recall the prime problem (soft-margin SVM)

$$\text{minimize} \quad \|\beta\| \quad \text{over } \beta, \beta_0 \quad (12)$$

$$\text{s.t.} \quad y_i(\beta^T x_i + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, N \quad (13)$$

$$\xi_i \geq 0; \quad \sum_{i=1}^N \xi_i \leq Z \quad (14)$$

Obviously strong duality holds. So we can find its dual problem by the following steps

1. Define Lagrange primal function (and Lagrange multipliers).
2. Take the first-order derivatives w.r.t. β , β_0 and ξ_i , and set to zero.
3. Substitute the results into the primal function.

$$\text{Maximize} \quad L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} \langle x_i, x_{i'} \rangle \quad (15)$$

$$\text{s.t.} \quad 0 \leq \alpha_i \leq \gamma, \quad i = 1, \dots, N \quad (16)$$

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad (17)$$

$$\text{Solution :} \quad (18)$$

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i \quad (19)$$

$$f(x) = \beta^T x + \beta_0 = \sum_{i=1}^N \hat{\alpha}_i y_i \langle x_i, x \rangle + \beta_0 \quad (20)$$

- ▶ Sparse representation: the separating hyperplane $f(x)$ is spanned those data points i where $\alpha_i \neq 0$, called **Support Vectors**.

- ▶ follows directly from complementary slackness conditions:

$$\alpha_i [y_i (\beta^T x_i + \beta_0) + \xi_i - 1] = 0$$

- ▶ Both the estimation and the evaluation of $f(x)$ only involve **dot product**.