Introduction to Machine Learning

Linear Support Vector Machines Hard-Margin SVM Dual

× 0 0 × 0 × ×



Learning goals

 Know how to derive the SVM dual problem

We before derived the primal quadratic program for the hard margin SVM. We could directly solve this, but traditionally the SVM is solved in the dual and this has some advantages. In any case, many algorithms and derivations are based on it, so we need to know it.

$$\begin{split} \min_{\boldsymbol{\theta}, \theta_0} & \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} & \boldsymbol{y}^{(i)} \left(\left\langle \boldsymbol{\theta}, \boldsymbol{x}^{(i)} \right\rangle + \theta_0 \right) \geq 1 \quad \forall \, i \in \{1, \dots, n\}. \end{split}$$

The Lagrange function of the SVM optimization problem is

$$L(\boldsymbol{\theta}, \theta_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{i=1}^n \alpha_i \left[\boldsymbol{y}^{(i)} \left(\left\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \right\rangle + \theta_0 \right) - 1 \right]$$

s.t. $\alpha_i \ge 0 \quad \forall i \in \{1, \dots, n\}.$

The **dual** form of this problem is

$$\max_{\alpha} \min_{\boldsymbol{\theta}, \theta_0} L(\boldsymbol{\theta}, \theta_0, \boldsymbol{\alpha}).$$



Notice how the (p+1) decision variables (θ, θ_0) have become *n* decisions variables α , as constraints turned into variables and vice versa. Now every data point has an associated non-negative weight.

$$L(\boldsymbol{\theta}, \theta_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{i=1}^n \alpha_i \left[y^{(i)} \left(\left\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \right\rangle + \theta_0 \right) - 1 \right]$$

s.t. $\alpha_i \ge 0 \quad \forall i \in \{1, \dots, n\}.$

We find the stationary point of $L(\theta, \theta_0, \alpha)$ w.r.t. θ, θ_0 and obtain

$$\boldsymbol{\theta} = \sum_{i=1}^{n} \alpha_i \boldsymbol{y}^{(i)} \mathbf{x}^{(i)},$$

$$\boldsymbol{0} = \sum_{i=1}^{n} \alpha_i \boldsymbol{y}^{(i)} \quad \forall i \in \{1, \dots, n\}.$$

× × 0 × × ×

By inserting these expressions & simplifying we obtain the dual problem

$$\begin{aligned} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \boldsymbol{y}^{(i)} \boldsymbol{y}^{(j)} \left\langle \boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)} \right\rangle \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i \boldsymbol{y}^{(i)} = \boldsymbol{0}, \\ & \alpha_i \ge \boldsymbol{0} \ \forall i \in \{1, \dots, n\}, \end{aligned}$$

× × ×

or, equivalently, in matrix notation:

$$\begin{split} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} & \mathbf{1}^T \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^T \operatorname{diag}(\mathbf{y}) \mathbf{K} \operatorname{diag}(\mathbf{y}) \boldsymbol{\alpha} \\ \text{s.t.} & \boldsymbol{\alpha}^T \mathbf{y} = \mathbf{0}, \\ & \boldsymbol{\alpha} \geq \mathbf{0}, \end{split}$$

with $\boldsymbol{K} := \boldsymbol{X} \boldsymbol{X}^{T}$.

If $(\theta, \theta_0, \alpha)$ fulfills the KKT conditions (stationarity, primal/dual feasibility, complementary slackness), it solves both the primal and dual problem (strong duality).

Under these conditions, and if we solve the dual problem and obtain $\hat{\alpha}$, we know that θ is a linear combination of our data points:

$$\hat{\theta} = \sum_{i=1}^{n} \hat{\alpha}_i \boldsymbol{y}^{(i)} \boldsymbol{x}^{(i)}$$

Complementary slackness means:

$$\hat{\alpha}_i \left[\boldsymbol{y}^{(i)} \left(\left\langle \boldsymbol{\theta}, \boldsymbol{x}^{(i)} \right\rangle + \theta_0 \right) - 1 \right] = 0 \quad \forall \ i \in \{1, ..., n\}.$$

× < 0 × × ×

$$\hat{\theta} = \sum_{i=1}^{n} \hat{\alpha}_{i} \boldsymbol{y}^{(i)} \boldsymbol{x}^{(i)}$$
$$\hat{\alpha}_{i} \left[\boldsymbol{y}^{(i)} \left(\left\langle \boldsymbol{\theta}, \boldsymbol{x}^{(i)} \right\rangle + \theta_{0} \right) - 1 \right] = \boldsymbol{0} \quad \forall i \in \{1, ..., n\}.$$

- So either â_i = 0, and is not active in the linear combination, or â_i > 0, then y⁽ⁱ⁾ (⟨θ, x⁽ⁱ⁾⟩ + θ₀) = 1, and (x⁽ⁱ⁾, y⁽ⁱ⁾) has minimal margin and is a support vector!
- We see that we can directly extract the support vectors from the dual variables and the θ solution only depends on them.
- We can reconstruct the bias term θ_0 from any support vector:

$$heta_0 = \mathbf{y}^{(i)} - \left\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \right\rangle.$$

× < 0 × × ×

DUAL VARIABLE AND SUPPORT VECTORS

- SVs are defined to be points with
 â_i > 0. In the case of hard margin linear SVM, the SVs are on the edge of margin.
- However, not all points on edge of margin are necessarily SVs.
- In other words, it is possible that both $\hat{\alpha}_i = 0$ and $y^{(i)} \left(\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \rangle \right) 1 = 0$ hold.



