

We are now in a position to prove a learning estimate for SVM.

Defn (leave-one-out error): Let h_S be the hypothesis returned by an algorithm \mathcal{A} when trained on a finite sample S . The leave-one-out error of \mathcal{A} on S is

$$\hat{R}_{\text{Loo}}(\mathcal{A}) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{h_{S \setminus \{x_i\}}(x_i) \neq y_i}$$

In other words, for each x_i , we train ~~algorithm~~ on $S \setminus \{x_i\}$, then compute the error of predicting on x_i . We then average these errors.

Lemma: The average Loo error for samples of size $m \geq 2$ is an unbiased estimate of the average generalization error for samples of size $m-1$:

$$\mathbb{E}_{S \sim D^m} (\hat{R}_{\text{Loo}}(\mathcal{A})) = \mathbb{E}_{S' \sim D^{m-1}} (R(h_{S'}))$$

Proof: By linearity of expectation, we compute

$$\mathbb{E}_{S \sim D^m} (\hat{R}_{\text{Loo}}(\mathcal{A})) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{S \sim D^m} (\mathbb{1}_{h_{S \setminus \{x_i\}}(x_i) \neq y_i})$$

$$\xrightarrow{\text{iid samples}} \mathbb{E}_{S \sim D^m} (\mathbb{1}_{h_{S \setminus \{x_1\}}(x_1) \neq y_1})$$

$$= \mathbb{E}_{\substack{S' \sim D^{m-1} \\ x_1 \sim D}} (\mathbb{1}_{h_{S'}(x_1) \neq y_1})$$

$$= \mathbb{E}_{S \sim D^{m+1}} \left(\mathbb{E}_{x_1 \sim D} \left(\mathbb{1}_{h_S(x_1) \neq y_1} \right) \right)$$

$$= \mathbb{E}_{S \sim D^{m+1}} \left(R(h_S) \right) . \blacksquare$$

Theorem (learning bound for SVM): Let h_S be the hypothesis returned by SVM for a sample S , and let $N_{SV}(S)$ be the number of support vectors that define h_S .

Then:
$$\mathbb{E}_{S \sim D^m} \left(R(h_S) \right) \leq \mathbb{E}_{S \sim D^{m+1}} \left(\frac{N_{SV}(S)}{m+1} \right)$$

Proof: Let S be a linearly separable sample of size $m+1$. If x is not a support vector, then removing it from the sample does not change the SVM solution. Hence, $h_{S \setminus \{x\}} = h_S$ and $h_{S \setminus \{x\}}$ correctly classifies x . Hence, if h_S misclassifies x , then x must be a support vector. Thus,

$$\hat{R}_{\text{Loo}}(\text{SVM}) \leq \frac{N_{SV}(S)}{m+1}$$

$$\Rightarrow \mathbb{E}_{S \sim D^m} \left(\hat{R}_{\text{Loo}}(\text{SVM}) \right) \leq \mathbb{E}_{S \sim D^m} \left(\frac{N_{SV}(S)}{m+1} \right)$$

$$\mathbb{E}_{S \sim D^m} \left(R(h_S) \right) . \blacksquare$$

• So far, we have assumed that the data is linearly separable, i.e. that there exists $\textcircled{3}$ some $(w, b) \in \mathbb{R}^{D+1}$ such that $y_i (w \cdot x_i + b) \geq 1 \quad \forall i=1, \dots, m$.

• This is not really very realistic, so we consider relaxing these conditions as follows. For each $i=1, \dots, m$, suppose $\exists z_i \geq 0$ s.t.

$$y_i [w \cdot x_i + b] \geq 1 - z_i$$

We call the $\{z_i\}_{i=1}^m$ slack variables.

• If $z_i > 0$, this means x_i is on the "wrong side" of the hyperplane, and may be considered an outlier.

• We would like a large margin classifier still, but we would also like the total amount of slack (e.g. $\sum_{i=1}^m z_i$) to be small. These are in tension.

• Consider the following constrained optimization problem: for $\lambda > 0, p \geq 1$,

$$\min_{\substack{(w,b) \in \mathbb{R}^{D+1} \\ z \in \mathbb{R}^m}} \frac{1}{2} \|w\|_2^2 + \underbrace{\lambda \sum_{i=1}^m z_i^p}_{\lambda \|z\|_p^p} \quad \text{subject to } \begin{cases} y_i (w \cdot x_i + b) \geq 1 - z_i \\ z_i \geq 0 \end{cases} \quad \forall i=1, \dots, m$$



• The p -norms most commonly considered for the slack penalty are $p=1$ and $p=2$. In what follows, we will analyze the $p=1$ case.

Remark: This optimization $\textcircled{★}$ is convex, since the objective function is convex, and the constraints are affine.

As in the separable case, we may introduce the Lagrangian and examine the KKT conditions.

We have $2m$ constraints. Let $\{\alpha_i\}_{i=1}^m$ be the multipliers associated to $y_i(w \cdot x_i + b) \geq 1 - \zeta_i$ and let $\{\beta_i\}_{i=1}^m$ be the multipliers associated to $\zeta_i \geq 0$. This yields a Lagrangian

$$\mathcal{L}(w, b, \zeta, \alpha, \beta) = \frac{1}{2} \|w\|_2^2 + \lambda \sum_{i=1}^m \zeta_i - \sum_{i=1}^m \alpha_i (y_i (w \cdot x_i + b) - 1 + \zeta_i) - \sum_{i=1}^m \beta_i \zeta_i$$

Setting $\nabla_w \mathcal{L}$, $\nabla_b \mathcal{L}$, $\nabla_{\zeta} \mathcal{L}$ equal to 0 yields

$$\nabla_w \mathcal{L} = w - \sum_{i=1}^m \alpha_i y_i x_i = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\nabla_b \mathcal{L} = - \sum_{i=1}^m \alpha_i y_i = 0 \Rightarrow 0 = \sum_{i=1}^m \alpha_i y_i$$

$$\nabla_{\zeta} \mathcal{L} = \lambda - (\alpha_i + \beta_i) = 0 \Rightarrow \lambda = \alpha_i + \beta_i \quad \forall i$$

Moreover, $\forall i \quad \alpha_i (y_i (w \cdot x_i + b) - 1 + \zeta_i) = 0 \Rightarrow \alpha_i = 0$
or $y_i (w \cdot x_i + b) = 1 - \zeta_i$

$$\forall i \quad \beta_i \zeta_i = 0 \Rightarrow \beta_i = 0 \text{ or } \zeta_i = 0.$$

As before, x_i contributes to $w \Leftrightarrow \alpha_i \neq 0$. But, $\alpha_i \neq 0 \Leftrightarrow y_i (w \cdot x_i + b) = 1 - \zeta_i$

If $\zeta_i = 0$, then $y_i (w \cdot x_i + b) = 1$ and x_i is on a marginal hyperplane, as in the case when the data is linearly separable. Otherwise if $\zeta_i > 0$, then $\beta_i = 0 \Rightarrow \alpha_i = \lambda$

To summarize, either support vectors lie on the marginal hyperplanes ($\zeta_i = 0$) or else they are outliers and $\alpha_i = \lambda$ ($\zeta_i > 0$).