

Stability and risk control for SVM

ERM addresses underfitting by choosing the predictor minimizing the training error. Overfitting, instead, is controlled by restricting the class of predictors that ERM can choose from. Note that ERM must be sensitive to small changes in the training set, as these may change the predictor minimizing the training error. We now investigate stability, a technique that controls overfitting by preventing the learning algorithm to react to small changes in the training set.

Fix a training set S of examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ and use \mathbf{z}_t to denote the t -th example $(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \mathbb{R}$. Given a loss function ℓ and a predictor h , let $\ell(h, \mathbf{z}_t) = \ell(h(\mathbf{x}_t), y_t)$ and use

$$\ell_S(h) = \frac{1}{m} \sum_{t=1}^m \ell(h, \mathbf{z}_t)$$

to denote the training error of h . From now on, we assume S is a sample of examples $Z_t = (\mathbf{X}_t, Y_t)$ independently drawn from a distribution \mathcal{D} . We use $S^{(t)}$ to denote S where the t -th example (\mathbf{X}_t, Y_t) is replaced by $\mathbf{Z}'_t = (\mathbf{X}'_t, Y'_t)$, also drawn from \mathcal{D} independently of S .

Fix a learning algorithm A and let $h_S = A(S)$ and $h_{S^{(t)}} = A(S^{(t)})$. We say that A is ε -stable if, for each $t = 1, \dots, m$

$$\mathbb{E}[\ell(h_{S^{(t)}}) - \ell(h_S)] \leq \varepsilon$$

where the expected value is computed with respect to the random draw of S and \mathbf{Z}'_t . Note that, in general, $\ell(h_{S^{(t)}}) > \ell(h_S)$ because \mathbf{Z}_t is not in $S^{(t)}$. Stability demands that $\ell(h_{S^{(t)}})$ be not much bigger than $\ell(h_S)$ in expectation with respect to S and \mathbf{Z}'_t . In other words, an algorithm is stable if replacing a single example in the training set does not significantly increase the loss on that example. In our analysis, we use the following equivalent definition of stability:

$$\mathbb{E}[\ell(h_S) - \ell(h_{S^{(t)}})] \leq \varepsilon$$

The next two results show that a stable algorithm does not overfit.

Theorem 1. *If A is ε -stable, then $\mathbb{E}[\ell_{\mathcal{D}}(h_S) - \ell_S(h_S)] \leq \varepsilon$.*

PROOF. Let $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ and $S' = \{(\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_m, y'_m)\}$ be drawn i.i.d. from \mathcal{D} . Then

$$\mathbb{E}[\ell_S(h_S)] = \mathbb{E}\left[\frac{1}{m} \sum_{t=1}^m \ell(h_S, \mathbf{z}_t)\right] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}_t)] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_{S^{(t)}}) - \ell(h_S)] + \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_{S^{(t)}})] .$$

Moreover,

$$\ell_{\mathcal{D}}(h_S) = \mathbb{E}[\ell(h_S, \mathbf{z}'_t) | S] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}'_t) | S]$$

averaging with respect to the random draw of S , this implies

$$\mathbb{E}[\ell_{\mathcal{D}}(h_S)] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}'_t)] .$$

Therefore,

$$\mathbb{E}[\ell_{\mathcal{D}}(h_S) - \ell_S(h_S)] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}'_t) - \ell(h_{S^{(t)}}, \mathbf{z}'_t)] \leq \varepsilon$$

due to the stability assumption. □

Because stability and minimization of training error work against each other, ERM is typically not stable. On the other hand, a stable learning algorithm that outputs predictors with small empirical risk must have a small variance error.

Theorem 2. *If A is ε -stable and also approximately minimizes the empirical risk in a given class of predictors, that is*

$$\ell_S(h_S) \leq \inf_{h \in \mathcal{H}} \ell_S(h) + \gamma$$

for some $\gamma > 0$, then

$$\mathbb{E}[\ell_{\mathcal{D}}(h_S)] \leq \inf_{h \in \mathcal{H}} \ell_{\mathcal{D}}(h) + \varepsilon + \gamma .$$

PROOF. Let h^* be the predictor with smallest risk in \mathcal{H} . Then

$$\begin{aligned} \mathbb{E}[\ell_{\mathcal{D}}(h_S)] &= \mathbb{E}\left[\underbrace{\ell_{\mathcal{D}}(h_S) - \ell_S(h_S)}_{\leq \varepsilon \text{ (stability)}}\right] + \mathbb{E}\left[\underbrace{\ell_S(h_S) - \ell_S(h^*)}_{\leq \gamma \text{ (ERM approximation)}}\right] + \mathbb{E}[\ell_S(h^*)] \\ &\leq \varepsilon + \gamma + \mathbb{E}[\ell_S(h^*)] . \end{aligned}$$

The proof is concluded by observing that $\mathbb{E}[\ell_S(h^*)] = \ell_{\mathcal{D}}(h^*)$, namely the expected value of the empirical risk is the risk. □

In case of predictors parameterized by a vector $\mathbf{w} \in \mathbb{R}^d$ (like linear predictors), ERM can be made stable by adding to the empirical loss a so-called **regularization term**. We also need an additional condition, namely that the loss function ℓ is such that $\ell(\cdot, \mathbf{z})$ be convex and Lipschitz, where $\ell(\mathbf{w}, \mathbf{z})$ is the error of \mathbf{w} on the example \mathbf{z} . Recall that Lipschitz means that there exists a constant L such that $|\ell(\mathbf{w}, \mathbf{z}) - \ell(\mathbf{w}', \mathbf{z})| \leq L \|\mathbf{w} - \mathbf{w}'\|$ for all $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^d$ and for all $\mathbf{z} = (\mathbf{x}, y)$. No other assumptions on ℓ are required.

Theorem 3. *Let ℓ be a loss function such that $\ell(\cdot, \mathbf{z})$ is convex, differentiable¹ and Lipschitz with constant $L > 0$. Then the learning algorithm A such that*

$$A(S) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left(\ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right)$$

is $(2L)^2/(\lambda m)$ -stable for every $\lambda > 0$.

¹Some loss functions, notably the hinge loss, are not everywhere differentiable. However, they are everywhere subdifferentiable, which is sufficient for this theorem to hold.

PROOF. Introduce

$$F_S(\mathbf{w}) = \ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

and also

$$\mathbf{w}_S = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} F_S(\mathbf{w}) \quad \text{and} \quad \mathbf{w}_{S^{(t)}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} F_{S^{(t)}}(\mathbf{w}) .$$

In order to prove stability, we need to upper bound $\mathbb{E}[\ell(\mathbf{w}_S, \mathbf{z}'_t) - \ell(\mathbf{w}_{S^{(t)}}, \mathbf{z}'_t)]$. We actually prove a stronger result by bounding the quantity $\ell(\mathbf{w}_S, \mathbf{z}'_t) - \ell(\mathbf{w}_{S^{(t)}}, \mathbf{z}'_t)$ for all S and \mathbf{z}'_t . As a first step, we use the Lipschitz condition to write

$$\ell(\mathbf{w}_S, \mathbf{z}'_t) - \ell(\mathbf{w}_{S^{(t)}}, \mathbf{z}'_t) \leq L \|\mathbf{w}_S - \mathbf{w}_{S^{(t)}}\| . \quad (1)$$

Next, we upper bound $\|\mathbf{w}_S - \mathbf{w}_{S^{(t)}}\|$. Introduce the abbreviations $\mathbf{w} = \mathbf{w}_S$ e $\mathbf{w}' = \mathbf{w}_{S^{(t)}}$. Then

$$\begin{aligned} F_S(\mathbf{w}') - F_S(\mathbf{w}) &= \ell_S(\mathbf{w}') - \ell_S(\mathbf{w}) + \frac{\lambda}{2} (\|\mathbf{w}'\|^2 - \|\mathbf{w}\|^2) \\ &= \ell_{S^{(t)}}(\mathbf{w}') - \ell_{S^{(t)}}(\mathbf{w}) + \frac{\ell(\mathbf{w}', \mathbf{z}_t) - \ell(\mathbf{w}, \mathbf{z}_t)}{m} - \frac{\ell(\mathbf{w}', \mathbf{z}'_t) - \ell(\mathbf{w}, \mathbf{z}'_t)}{m} + \frac{\lambda}{2} (\|\mathbf{w}'\|^2 - \|\mathbf{w}\|^2) \\ &= F_{S^{(t)}}(\mathbf{w}') - F_{S^{(t)}}(\mathbf{w}) + \frac{\ell(\mathbf{w}', \mathbf{z}_t) - \ell(\mathbf{w}, \mathbf{z}_t)}{m} - \frac{\ell(\mathbf{w}', \mathbf{z}'_t) - \ell(\mathbf{w}, \mathbf{z}'_t)}{m} \\ &\leq \frac{|\ell(\mathbf{w}', \mathbf{z}_t) - \ell(\mathbf{w}, \mathbf{z}_t)|}{m} + \frac{|\ell(\mathbf{w}', \mathbf{z}'_t) - \ell(\mathbf{w}, \mathbf{z}'_t)|}{m} \\ &\leq \frac{2L}{m} \|\mathbf{w} - \mathbf{w}'\| \end{aligned}$$

where the first inequality holds because $\mathbf{w}' = \mathbf{w}_{S^{(t)}}$ minimizes $F_{S^{(t)}}$ and the second inequality holds because $\ell(\cdot, \mathbf{z})$ is Lipschitz.

We proceed by noting that the function F_S is λ -strongly convex: indeed, $\ell(\cdot, \mathbf{z})$ is convex, $\frac{\lambda}{2} \|\mathbf{w}\|^2$ is λ -strongly convex, which implies that their sum is also λ -strongly convex. Therefore, by definition of strongly convex function,

$$F_S(\mathbf{w}') \geq F_S(\mathbf{w}) + \nabla F_S(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}'\|^2 .$$

Because \mathbf{w} is the minimizer of F_S , $\nabla F_S(\mathbf{w}) = \mathbf{0}$ and so

$$F_S(\mathbf{w}') - F_S(\mathbf{w}) = \left(\ell_S(\mathbf{w}') + \frac{\lambda}{2} \|\mathbf{w}'\|^2 \right) - \left(\ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \geq \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}'\|^2$$

Combining the two inequalities we get

$$\frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}'\|^2 \leq \frac{2L}{m} \|\mathbf{w} - \mathbf{w}'\| \quad \text{ovvero} \quad \|\mathbf{w} - \mathbf{w}'\| \leq \frac{4L}{\lambda m}$$

which, together with (1) shows the stability of $\mathbf{w} = \mathbf{w}_S$. □

We now show how the notion of stability can be used to control the risk of the SVM predictor. First, recall that the hinge loss $\ell(\mathbf{w}, (\mathbf{x}, y)) = [1 - y \mathbf{w}^\top \mathbf{x}]_+$ is convex in \mathbf{w} for $(\mathbf{x}, y) \in \mathbb{R}^d \times \{-1, 1\}$. In

order to compute the Lipschitz constant L , observe that, using Taylor's theorem and the Cauchy-Schwartz inequality $\mathbf{u}^\top \mathbf{v} \leq \|\mathbf{u}\| \|\mathbf{v}\|$, we can write

$$\ell(\mathbf{w}', \mathbf{z}) - \ell(\mathbf{w}, \mathbf{z}) \leq \nabla \ell(\mathbf{w}, \mathbf{z})^\top (\mathbf{w}' - \mathbf{w}) \leq \|\nabla \ell(\mathbf{w}, \mathbf{z})\| \|\mathbf{w}' - \mathbf{w}\|$$

Now note that $\nabla \ell(\mathbf{w}, \mathbf{z}) = \mathbf{y} \mathbf{x} \mathbb{I}\{\mathbf{y} \mathbf{w}^\top \mathbf{x} \leq 1\}$, and thus $\|\nabla \ell(\mathbf{w}, \mathbf{z})\| \leq \|\mathbf{x}\|$. Assuming $\|\mathbf{x}_t\| \leq X$ for $t = 1, \dots, m$ we can then set $L = X$ and show the following result.

Theorem 4. *Given a training set S , the SVM solution*

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left(\ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \quad (2)$$

satisfies

$$\mathbb{E}[\ell_{\mathcal{D}}^{0-1}(\mathbf{w}_S)] \leq \min_{\mathbf{u} \in \mathbb{R}^d} \left(\ell_{\mathcal{D}}(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 \right) + \frac{4X^2}{\lambda m} .$$

where $\ell_{\mathcal{D}}^{0-1}(\mathbf{w}_S)$ is the risk of \mathbf{w}_S with respect to the zero-one loss.

PROOF. Clearly, for each $\mathbf{u} \in \mathbb{R}^d$ we have

$$\ell_S(\mathbf{w}_S) \leq \ell_S(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{w}_S\|^2 \leq \ell_S(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 . \quad (3)$$

Therefore, because Theorem 3 implies that \mathbf{w}_S is $(4L^2)/(\lambda m)$ -stable for $L = X$,

$$\begin{aligned} \mathbb{E}[\ell_{\mathcal{D}}(\mathbf{w}_S)] &\leq \mathbb{E}[\ell_S(\mathbf{w}_S)] + \frac{4X^2}{\lambda m} \quad \text{by Theorem 1} \\ &\leq \mathbb{E} \left[\ell_S(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 \right] + \frac{4X^2}{\lambda m} \quad \text{using (3)} \\ &= \ell_{\mathcal{D}}(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 + \frac{4X^2}{\lambda m} \end{aligned}$$

The proof is concluded by noting that, for any linear predictor \mathbf{w} , $\ell_{\mathcal{D}}^{0-1}(\mathbf{w}) = \mathbb{P}(Y \mathbf{w}^\top \mathbf{X} \leq 0) \leq \ell_{\mathcal{D}}(\mathbf{w})$. This is an immediate consequence of the fact that the hinge loss is a convex upper bound on the zero-one loss. \square

Consistency of SVM. In a kernel space \mathcal{H}_K the SVM objective function (2) becomes

$$g_S = \operatorname{argmin}_{g \in \mathcal{H}_K} \left(\ell_S(g) + \frac{\lambda}{2} \|g\|_K^2 \right)$$

If the kernel is Gaussian, then one can prove that SVM becomes a consistent learning algorithm (with respect to the zero-one loss) when the regularization parameter λ is chosen as a function λ_m of the training set size m . In particular, for $m \rightarrow \infty$, λ_m must satisfy the two following conditions: $\lambda_m = o(1)$ and $\lambda_m = \omega(m^{-1/2})$.