

Machine Learning Theory

Linear classifiers

Ruth Urner

We have seen in the last lectures that, for binary hypothesis classes, learnability is equivalent to having bounded VC-dimension. Moreover, we have seen that, for classes of bounded VC-dimension, *any* ERM learning algorithm is a successful learner (that is, learns the class according to Definition 1, Lecture 2).

Note that, “ERM” for is not an learning algorithm. Rather it is a paradigm. Namely, it asks to output a classifier of minimal empirical error. Different learning methods can be ERM for some fixed hypothesis class. In this lecture, we will see a first concrete learning algorithm for the class of linear predictors: The perceptron algorithm is a very basic learner that, as we will see, successfully learns the class of linear predictors in the realizable case.

This chapter is based on the (first part of) Chapter 9 in the textbook [1].

1 Linear Classifiers

We consider the domain $\mathcal{X} = \mathbb{R}^d$ and label set $\mathcal{Y} = \{-1, 1\}$. The class $\mathcal{H}_{g\text{-lin}}$ of general linear classifiers (or the class of linear halfspaces) is defined as follows:

$$\mathcal{H}_{g\text{-lin}} = \{h_{\mathbf{w},b} : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\},$$

where

$$h_{\mathbf{w},b}(\mathbf{x}) = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle + b) = \text{sign}\left(\left(\sum_{i=1}^d x_i w_i\right) + b\right)$$

General linear halfspaces in \mathbb{R}^d can be viewed as *homogeneous* linear halfspaces in \mathbb{R}^{d+1} via a simple transformation. Recall that the class \mathcal{H}_{lin} of homogeneous linear halfspaces is defined as

$$\mathcal{H}_{\text{lin}} = \{h_{\mathbf{w}} : \mathbf{w} \in \mathbb{R}^d\},$$

where

$$h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle) = \text{sign}\left(\sum_{i=1}^d x_i w_i\right).$$

Let $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$, $\mathbf{w} = (w_1, w_2, \dots, w_d) \in \mathbb{R}^d$ and $b \in \mathbb{R}$ be given. We define

$$\mathbf{x}' = (1, x_1, x_2, \dots, x_d) \in \mathbb{R}^{d+1}$$

and

$$\mathbf{w}' = (1, w_1, w_2, \dots, w_d) \in \mathbb{R}^{d+1}.$$

Then we get $\langle \mathbf{x}, \mathbf{w} \rangle + b = \langle \mathbf{x}', \mathbf{w}' \rangle$ for all $\mathbf{x} \in \mathbb{R}^d$, and thus

$$h_{\mathbf{w},b}(\mathbf{x}) = h_{\mathbf{w}'}(\mathbf{x}').$$

We will thus now focus our analysis of learning linear classifiers to learning homogeneous linear classifiers. This is not a restriction of generality. For the problem of learning general linear classifiers we can employ the following simple reduction:

- Given a dataset $S = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)$ of data points \mathbf{x}_i in \mathbb{R}^d , transform them to data set $S' = (\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_m)$ of points in \mathbb{R}^{d+1} .
- Learn a homogeneous linear classifier $\mathbf{w}' \in \mathbb{R}^{d+1}$ from S' .
- Interpret the first component w'_1 in \mathbf{w}' as the term b and the remaining entries (w'_2, \dots, w'_d) as the vector \mathbf{w} defining the halfspace classifier $h_{\mathbf{w},b}$ in \mathbb{R}^d .

2 The VC-dimension of halfspaces

Theorem 1. *The VC-dimension of homogeneous linear halfspaces in \mathbb{R}^d is d .*

Proof. We first show that homogeneous linear halfspaces can shatter d points. Consider the set of vectors $\mathbf{e}_1, \dots, \mathbf{e}_d$, where \mathbf{e}_i is defined as the vector with entry 1 in the i -th component and entries 0 else. Let $U \subseteq [d]$ be a subset of indices and let y_1^U, \dots, y_d^U be the induced labeling, that is

$$y_i^U = \begin{cases} 1 & \text{if } i \in U \\ -1 & \text{else} \end{cases}$$

We now define \mathbf{w} by

$$w_i = \begin{cases} 1 & \text{if } i \in U \\ -1 & \text{else} \end{cases}$$

Then we get

$$h_{\mathbf{w}}(\mathbf{e}_i) = y_i$$

for all $i \in [d]$. Thus, the d unit vectors are shattered.

Next, we prove that homogeneous linear halfspaces can not shatter more than d points in \mathbb{R}^d . Let some points $\mathbf{x}_1, \dots, \mathbf{x}_{d+1} \in \mathbb{R}^d$ be given and assume by way of contradiction, that this set of $d+1$ points is shattered. Since they are also linearly dependent, there exists a set of scalars a_1, \dots, a_{d+1} such that $\sum_{i=1}^{d+1} a_i \mathbf{x}_i = \mathbf{0}$.

Let $I = \{i : a_i > 0\}$ and $J = \{i : a_i < 0\}$. At least one of these sets of indices is non-empty. We assume for now that both I and J are not empty. Then we get

$$\sum_{i \in I} a_i \mathbf{x}_i = \sum_{j \in J} |a_j| \mathbf{x}_j$$

Since the $\mathbf{x}_1, \dots, \mathbf{x}_{d+1}$ are shattered, there exists a vector \mathbf{w} with

$$\langle \mathbf{w}, \mathbf{x}_i \rangle > 0 \text{ for all } i \in I$$

and

$$\langle \mathbf{w}, \mathbf{x}_j \rangle < 0 \text{ for all } j \in J$$

This yields

$$0 < \sum_{i \in I} a_i \langle \mathbf{w}, \mathbf{x}_i \rangle = \left\langle \mathbf{w}, \sum_{i \in I} a_i \mathbf{x}_i \right\rangle = \left\langle \mathbf{w}, \sum_{j \in J} |a_j| \mathbf{x}_j \right\rangle = \sum_{j \in J} |a_j| \langle \mathbf{w}, \mathbf{x}_j \rangle < 0,$$

a contradiction. Finally note that if either $I = \emptyset$ or $J = \emptyset$, then one of the inequalities above is an equality, but the contradiction still follows. \square

Similarly, one can show that the VC-dimension of general linear classifiers in \mathbb{R}^d is $d + 1$.

Theorem 2. *The VC-dimension of general linear classifiers in \mathbb{R}^d is $d + 1$.*

Proof. Left as exercise. \square

The theory of learning classes of bounded VC-dimension (see Lectures 2 and 3) now tells us that there is a constant $C > 0$, such that general halfspaces in \mathbb{R}^d are learnable from samples of size larger than

$$m(\epsilon, \delta) = C \cdot \frac{(d + 1) + \ln(1/\delta)}{\epsilon^2},$$

in the agnostic case and samples of sizes at least

$$m(\epsilon, \delta) = C \cdot \frac{(d + 1) + \ln(1/\delta)}{\epsilon},$$

in the realizable case.

3 The Perceptron Algorithm

The perceptron algorithm is a simple learner that implements an ERM rule for \mathcal{H}_{lin} in the realizable case.

Note that a training example (\mathbf{x}_i, y_i) from S is misclassified by a (homogeneous) halfspace $h_{\mathbf{w}}$ if and only if

$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle < 0.$$

In the realizable case, a learner is an ERM learner if it always outputs a classifier that makes no errors on the training data. Thus an ERM for homogeneous halfspaces needs to output a vector \mathbf{w} such that

$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle > 0$$

for all $(\mathbf{x}_i, y_i) \in S$. The perceptron algorithm is defined as follows:

We will next show that, on realizable data, the perceptron will always find a classifier that makes no error (that is, the above algorithm will always halt). The following theorem furthermore gives a bound on the number of iterations the algorithm makes before finding such a classifier. It thus provides a bound on the computation time that the Perceptron algorithm requires.

Algorithm 1 Perceptron

```
1: Input  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$ 
2: Initialize:  $\mathbf{w}^1 = (0, 0, \dots, 0)$ 
3: for  $t = 1, 2, \dots$  do
4:   if there exists an  $i$  such that  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \leq 0$  then
5:      $\mathbf{w}^{t+1} = \mathbf{w}^t + y_i \mathbf{x}_i$ 
6:   else
7:     return  $\mathbf{w}^t$ 
8:   end if
9: end for
```

Theorem 3. Assume that the data generating distribution P is realizable by the class of homogeneous linear classifiers in \mathbb{R}^d . Let $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$ be a sample generated by P . Let

$$B = \min\{\|\mathbf{w}\| : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 \text{ for all } i \in [m]\}$$

and let

$$R = \max_{i \in [m]} \|\mathbf{x}_i\|.$$

Then the Perceptron algorithm stops after $t \leq (RB)^2$ iterations.

Note that, by the definition of the algorithm, we have

$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle > 0$$

for all $(\mathbf{x}_i, y_i) \in S$, after $t \leq (RB)^2$ iterations. That is, at that point the algorithm found a classifier that minimizes the empirical risk on the sample S .

Proof. The proof can be found as the proof of Theorem 9.1, Chapter 9, of the textbook [1]. \square

In the tutorials, it was shown that the bound $(RB)^2$ in the above theorem can be replaced with the quantity $\frac{R^2}{\gamma^2}$, where γ is the *margin* of the data. That is, γ is the distance of the closest point to the hyperplane defining the halfspace classifier, for the classifier that maximizes that distance. The margin, by which a dataset is separable, can be viewed as a parameter of niceness or easiness of the data. Theorem 3 quantifies how the runtime of the perceptron algorithm depends on this easiness parameter. The “easier” the dataset, that is, the larger the margin by which it is separable with a homogeneous linear halfspace, the faster the algorithm will converge to finding such a separating halfspace.

References

- [1] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning*. Cambridge University Press, 2014.