

STATISTICAL MACHINE LEARNING (WS2020)
SEMINAR 4

Remark: Assignments 3 and 4 are related to the upcoming lecture on Tuesday, September 20. So, please work on them later.

Assignment 1. Let us consider the space of linear classifiers mapping $\mathbf{x} \in \mathbb{R}^n$ to $\{-1, +1\}$, that is

$$\mathcal{H} = \{h(\mathbf{x}; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b) \mid (\mathbf{w}, b) \in (\mathbb{R}^d \times \mathbb{R})\}.$$

Show that the VC dimension of \mathcal{H} is $n + 1$.

Hint: The proof has two steps:

- (1) Show that the VC dimension is at least $n + 1$ by constructing $n + 1$ points that are shattered by \mathcal{H} .
- (2) Show that the VC dimension is less than $n + 2$ by proving that $n + 2$ points cannot be shattered by \mathcal{H} .

The step 1 should be easy. If you find step 2 difficult, skip it for sake of the other assignments.

Assignment 2. Let $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$ be a hypothesis class with VC dimension $d < \infty$ and $\mathcal{T}^m = \{(x^1, y^1), \dots, (x^m, y^m)\} \in (\mathcal{X} \times \mathcal{Y})^m$ a training set drawn from i.i.d. random variables with distribution $p(x, y)$. Then, the following inequality holds for any $\varepsilon > 0$,

$$\mathbb{P}\left(\sup_{h \in \mathcal{H}} \left| R^{0/1}(h) - R_{\mathcal{T}^m}^{0/1}(h) \right| \geq \varepsilon\right) \leq 4 \left(\frac{2em}{d}\right)^d e^{-\frac{m\varepsilon^2}{8}},$$

where $R^{0/1}(h) = \mathbb{E}_{(x,y) \sim p}(\mathbb{1}[y \neq h(x)])$ and $R_{\mathcal{T}^m}^{0/1}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}[y^i \neq h(x^i)]$.

Show that this implies the ULLN for the class of strategies \mathcal{H} .

Assignment 3. Assume we are given a training set of examples $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \{+1, -1\}) \mid i = 1, \dots, m\}$ which is known to be linearly separable with respect to a feature map $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$. In this case, we can find parameters $(\mathbf{w}, b) \in \mathbb{R}^{n+1}$ of a linear classifier $h(x; \mathbf{w}, b) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$ which has zero training error by the Perceptron algorithm:

- (1) $\mathbf{w} \leftarrow 0, b \leftarrow 0$
- (2) Find an example $(x^u, y^u) \in \mathcal{T}^m$ whose label is incorrectly predicted by the current classifier, that is $h(x^u; \mathbf{w}, b) \neq y^u$.
- (3) If all examples are classified correctly exit the algorithm. Otherwise update the parameters by

$$\mathbf{w} \leftarrow \mathbf{w} + y^u \phi(x^u) \quad \text{and} \quad b \leftarrow b + y^u$$

and go to Step 2.

Assume that you cannot evaluate the feature map $\phi(x)$ because it is either unknown or its evaluation is expensive. However, you know how to cheaply evaluate a kernel function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that $k(x, x') = \langle \phi(x), \phi(x') \rangle$, $\forall x, x' \in \mathcal{X}$. Show that you can still use the Perceptron algorithm to find a linear classifier with zero training error and that you can evaluate this classifier on any $x \in \mathcal{X}$.

Assignment 4. Let the input observation be a vector $\mathbf{x} \in \mathbb{R}^d$. Let us consider a feature map $\phi_q: \mathbb{R}^d \rightarrow \mathbb{R}^n$, $n = d^q$, whose entries are all possible q -th degree ordered products of the entries of \mathbf{x} . For example, if $\mathbf{x} = (x_1, x_2, x_3)^T \in \mathbb{R}^3$ and $q = 2$ then

$$\phi_q(\mathbf{x}) = \begin{pmatrix} x_1x_1 \\ x_2x_1 \\ x_3x_1 \\ x_1x_2 \\ x_2x_2 \\ x_3x_2 \\ x_1x_3 \\ x_2x_3 \\ x_3x_3 \end{pmatrix}$$

a) Show that for any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ we can compute the dot product between $\phi_q(\mathbf{x})$ and $\phi_q(\mathbf{x}')$ as

$$\langle \phi_q(\mathbf{x}), \phi_q(\mathbf{x}') \rangle = \langle \mathbf{x}, \mathbf{x}' \rangle^q,$$

that is, as the dot product of the original vectors \mathbf{x} and \mathbf{x}' powered to q .

b) Consider a slightly different feature map $\phi': \mathbb{R}^d \rightarrow \mathbb{R}^{d(d+1)/2}$ whose entries are

$$\phi'(\mathbf{x}) = \begin{pmatrix} x_1^2, & \sqrt{2}x_1x_2, & \sqrt{2}x_1x_3, & \dots, & \sqrt{2}x_1x_d, \\ & x_2^2, & \sqrt{2}x_2x_3, & \dots, & \sqrt{2}x_2x_d, \\ & & & & \vdots \\ & & & & x_d^2 \end{pmatrix}^T,$$

so that the features correspond to all possible products of unordered pairs of entries from \mathbf{x} , and the products of different entries are multiplied by a constant factor $\sqrt{2}$. For example, if $\mathbf{x} = (x_1, x_2, x_3)^T \in \mathbb{R}^3$ then

$$\phi'(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, x_2^2, \sqrt{2}x_2x_3, x_3^2)^T.$$

This feature map defines a kernel $k(\mathbf{x}, \mathbf{x}') = \langle \phi'(\mathbf{x}), \phi'(\mathbf{x}') \rangle$ referred to as the homogeneous polynomial kernel of degree 2. Show that the kernel value equals to the square of the dot product of the input vectors, that is prove the identity

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi'(\mathbf{x}), \phi'(\mathbf{x}') \rangle = \langle \mathbf{x}, \mathbf{x}' \rangle^2, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d.$$

Hint: Exploit the relation between $\phi(\mathbf{x})$ and $\phi'(\mathbf{x})$.