
Machine Learning 2020-2021

Home Assignment 3

Yevgeny Seldin Christian Igel

Department of Computer Science
University of Copenhagen

The deadline for this assignment is **8 December 2020, 22:00**. You must submit your *individual* solution electronically via the Absalon home page.

A solution consists of:

- A PDF file with detailed answers to the questions, which may include graphs and tables if needed. Do *not* include your source code in the PDF file.
- A .zip file with all your solution source code with comments about the major steps involved in each question (see below). Source code must be submitted in the original file format, not as PDF. The programming language of the course is Python.
- **IMPORTANT: Do NOT zip the PDF file**, since zipped files cannot be opened in speed grader. Zipped pdf submissions will not be graded.
- Your PDF report should be self-sufficient. I.e., it should be possible to grade it without opening the .zip file. We do not guarantee opening the .zip file when grading.
- Your code should be structured such that there is one main file (or one main file per question) that we can run to reproduce all the results presented in your report. This main file can, if you like, call other files with functions, classes, etc.
- Handwritten solutions will not be accepted, please use the provided latex template to write your report.

1 The Role of Independence (5 points)

Design an example of identically distributed, but *dependent* Bernoulli random variables X_1, \dots, X_n (i.e., $X_i \in \{0, 1\}$), such that

$$\mathbb{P}\left(\left|\mu - \frac{1}{n} \sum_{i=1}^n X_i\right| \geq \frac{1}{2}\right) = 1,$$

where $\mu = \mathbb{E}[X_i]$.

Note that in this case $\frac{1}{n} \sum_{i=1}^n X_i$ does not converge to μ as n goes to infinity. The example shows that independence is crucial for convergence of empirical means to the expected values.

2 How to Split a Sample into Training and Test Set (25 points)

In this question you will analyze one possible approach to the question of how to split a dataset S into training and test sets, S^{train} and S^{test} . As we have already discussed, overly small test sets lead to unreliable loss estimates, whereas overly large test sets leave too little data for training, thus producing poor prediction models. The optimal trade-off depends on the data and the prediction model. So can we let the data speak for itself? We will give it a try.

1. To warm up: assume that you have a fixed split of S into S^{train} and S^{test} , where the size of S^{test} is n^{test} . You train a model $\hat{h}_{S^{\text{train}}}^*$ on S^{train} using whatever procedure you want. Then you compute the test loss $\hat{L}(\hat{h}_{S^{\text{train}}}^*, S^{\text{test}})$. Derive a bound on $L(\hat{h}_{S^{\text{train}}}^*)$ in terms of $\hat{L}(\hat{h}_{S^{\text{train}}}^*, S^{\text{test}})$ and n^{test} that holds with probability at least $1 - \delta$.
2. Now we want to find a good balance between the sizes of S^{train} and S^{test} . We consider m possible splits $\{(S_1^{\text{train}}, S_1^{\text{test}}), \dots, (S_m^{\text{train}}, S_m^{\text{test}})\}$, where the sizes of the test sets are n_1, \dots, n_m , correspondingly. For example, it could be (10%, 90%), (20%, 80%), \dots , (90%, 10%) splits or anything else with a reasonable coverage of the possible options. We train m prediction models $\hat{h}_1^*, \dots, \hat{h}_m^*$, where \hat{h}_i^* is trained on S_i^{train} . We calculate the test loss of the i -th model on the i -th test set $\hat{L}(\hat{h}_i^*, S_i^{\text{test}})$. Derive a bound on $L(\hat{h}_i^*)$ in terms of $\hat{L}(\hat{h}_i^*, S_i^{\text{test}})$ and n_i that holds for all \hat{h}_i^* simultaneously with probability at least $1 - \delta$.

Comment: No theorem from the lecture notes applies directly to this setting, because they all have a fixed sample size n , whereas here the sample sizes vary, n_1, \dots, n_m . You have to provide a complete derivation.

3 Occam's Razor (20 points)

We want to design an application for bilingual users. The application should detect the language in which the person is typing based on the first few letters typed. In other words, we want to design a classifier that takes a short string (that may be less than a full word) as input and predicts one of two languages, say, Danish or English. For simplicity we will assume that the alphabet is restricted to a set Σ of 26 letters of the Latin alphabet plus the white space symbol (so in total $|\Sigma| = 27$). Let Σ^d be the space of strings of length d . Let \mathcal{H}_d be the space of functions from Σ^d to $\{0, 1\}$, where Σ^d is the input string and $\{0, 1\}$ is the prediction (Danish or English). Let $\mathcal{H} = \bigcup_{d=0}^{\infty} \mathcal{H}_d$ be the union of \mathcal{H}_d -s.

1. Derive a high-probability bound¹ for $L(h)$ that holds for all $h \in \mathcal{H}_d$.
2. Derive a high-probability bound for $L(h)$ that holds for all $h \in \mathcal{H}$.
3. Explain the trade-off between picking short strings (small d) and long strings (large d). Which terms in the bound favor small d (i.e., they increase with d) and which terms in the bound favor large d (i.e., they decrease with d)?
4. We have presented a lower bound, where we constructed an example of a problem with a large hypothesis class (of size 2^{2^n}), where the empirical loss of the empirically best hypothesis was always zero, but the expected loss of the empirically best hypothesis was at least $1/4$. The hypothesis class \mathcal{H} in this question is obviously infinite. Explain why there is no contradiction between the bound in Point 2 and the lower bound.

Optional, not for submission You are very welcome to experiment with the influence of the string length d on the performance. You can find a lot of texts in different languages at <http://www.gutenberg.org/catalog/>. Do you observe the effect of initial improvement followed by overfitting as you increase d ?

4 Kernels (50 points)

The first question should improve the understanding of the geometry of the kernel-induced feature space. You can directly use the result to implement a kernel nearest-neighbor algorithm.

¹A bound that holds with probability at least $1 - \delta$.

The second question should make you more familiar with the basic definition of the important concept of positive definiteness.

The third question is important to understand the real dimensionality of learning problems using a linear kernel – one reason why linear kernels are often treated differently in efficient implementations.

The fourth question derives the general Cauchy-Schwarz inequality, which is of general importance. In particular, it needed to close a gap in the lecture. When we defined the scalar product $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}$ in the RKHS \mathcal{F} , we showed that $\forall \mathbf{x} \in \mathcal{F} : \langle \mathbf{x}, \mathbf{x} \rangle \geq 0$. But we also have to show that $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ implies $\mathbf{x} = \mathbf{0}$ – and this requires the general Cauchy-Schwarz inequality.

4.1 Distance in feature space

Given a kernel k on input space \mathcal{X} defining RKHS \mathcal{H} . Let $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ denote the corresponding feature map (think of $\Phi(x) = k(x, \cdot)$). Let $x, z \in \mathcal{X}$. Show that the distance of $\Phi(x)$ and $\Phi(z)$ in \mathcal{H} is given by

$$\|\Phi(x) - \Phi(z)\| = \sqrt{k(x, x) - 2k(x, z) + k(z, z)}$$

(if distance is measured by the canonical metric induced by k).

4.2 Sum of kernels

Let $k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be positive-definite kernels.

Prove that $k(x, z) = k_1(x, z) + k_2(x, z)$ is also positive-definite.

4.3 Rank of Gram matrix

Let the input space be $\mathcal{X} = \mathbb{R}^d$. Assume a linear kernel, $k(x, z) = x^T z$ for $x, z \in \mathbb{R}^d$ (i.e., the feature map Φ is the identity) and m input patterns $x_1, \dots, x_m \in \mathbb{R}^d$.

Prove an upper bound on the rank of the Gram matrix from the m input patterns in terms of d and m .

4.4 Cauchy-Schwarz

Solve exercise (4.3) from the textbook by Steinwart and Christmann (2008). Given a vector space \mathcal{F} and a positive semi-definite symmetric bilinear form $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}$, that is,

1. $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$

$$2. \langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$$

$$3. \langle \alpha \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \alpha \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$$

for $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{F}$, show the Cauchy-Schwarz inequality

$$|\langle \mathbf{x}, \mathbf{y} \rangle|^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \cdot \langle \mathbf{y}, \mathbf{y} \rangle$$

for all $\mathbf{x}, \mathbf{y} \in \mathcal{F}$.

Hint: Start with $0 \leq \langle \mathbf{x} + \alpha \mathbf{y}, \mathbf{x} + \alpha \mathbf{y} \rangle$ and consider the case $\alpha = 1$ and $\alpha = -1$ if $\langle \mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{y} \rangle = 0$. Otherwise, if, e.g., $\langle \mathbf{y}, \mathbf{y} \rangle \neq 0$, use $\alpha = -\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{y}, \mathbf{y} \rangle}$.

Why is this relevant? The Cauchy-Schwarz inequality is of general importance. For example, it is needed for proving the following property completing the argument in the lecture.

Recall from the lecture slides that a *dot product* on a vector space \mathcal{F} is a symmetric bilinear form $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}$ that is *strictly* positive definite, i.e., $\forall \mathbf{x} \in \mathcal{F} : \langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ with equality only for $\mathbf{x} = \mathbf{0}$.

Given a kernel k and

$$f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i) \qquad g(\cdot) = \sum_{j=1}^{m'} \beta_j k(\cdot, x'_j)$$

(see lecture slides for details and context) we defined the scalar product

$$\langle f, g \rangle := \sum_{i=1}^m \sum_{j=1}^{m'} \alpha_i \beta_j k(x_i, x'_j) \quad .$$

In the lecture, we stated

$$\langle f, f \rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j k(x_i, x_j) \geq 0 \quad ,$$

which follows from the kernel being positive semi-definite. However, we did not prove that $\langle f, f \rangle = 0$ implies $f = \mathbf{0}$, but this is necessary for having a scalar product. A proof can be found in Schölkopf and Smola (2002), which is, however, not fully correct. It argues with a Cauchy-Schwarz inequality for kernels, but a Cauchy Schwarz inequality for a positive symmetric bilinear form would be needed – and this is what you are supposed to prove in this assignment.

References

- B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- I. Steinwart and A. Christmann. *Support Vector Machines*. Information Science and Statistics. Springer-Verlag, 2008.