

---

# *Machine Learning A*

2022-2023

## Home Assignment 2

---

Sadegh Talebi      Christian Igel

Department of Computer Science

University of Copenhagen

The deadline for this assignment is **20 September 2022, 18:00**. 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 full source code in the PDF file, only selected lines if you are asked to do so.
- 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 Illustration of Markov's, Chebyshev's, and Hoeffding's Inequalities (23 points)

**2.a** Make 1,000,000 repetitions of the experiment of drawing 20 i.i.d. Bernoulli random variables  $X_1, \dots, X_{20}$  (20 coins) with bias  $\frac{1}{2}$  and answer the following questions.

Empirical probability is the number of times some event was observed in the data you have divided by the total sample size.

empirical freq = count/sample size

1. Plot the empirical frequency of observing  $\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha$  for  $\alpha \in \{0.5, 0.55, 0.6, \dots, 0.95, 1\}$ .
2. Explain why the above granularity of  $\alpha$  is sufficient. I.e., why, for example, taking  $\alpha = 0.51$  will not provide any extra information about the experiment.   
because we only have ranges from 0.05 so it doesnt add information
3. In the same figure plot the Markov's bound<sup>1</sup> on  $\mathbb{P}(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$ .
4. In the same figure plot the Chebyshev's bound<sup>2</sup> on  $\mathbb{P}(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$ . (You may have a problem calculating the bound for some values of  $\alpha$ . In that case and whenever the bound exceeds 1, replace it with the trivial bound of 1, because we know that probabilities are always bounded by 1.)
5. In the same figure plot the Hoeffding's bound<sup>3</sup> on  $\mathbb{P}(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$ .
6. Compare the four plots.
7. For  $\alpha = 1$  and  $\alpha = 0.95$  calculate the exact probability  $\mathbb{P}(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$ . (No need to add this one to the plot.)

**2.b** Repeat the question with  $X_1, \dots, X_{20}$  with bias 0.1 (i.e.,  $\mathbb{E}[X_1] = 0.1$ ) and  $\alpha \in \{0.1, 0.15, \dots, 1\}$ .

**2.c** Discuss the results.

Do not forget to put axis labels and a legend in your plot!

---

<sup>1</sup>Markov's bound is the right hand side of Markov's inequality.

<sup>2</sup>Chebyshev's bound is the right hand side of Chebyshev's inequality.

<sup>3</sup>Hoeffding's bound is the right hand side of Hoeffding's inequality.

## 2 The Role of Independence (14 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

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

x1 = 1, x2 = 1 if x1 = 0 else 1 ecc  
need to use n that goes to inf.

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

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

## 3 Tightness of Markov's Inequality (14 points)

In the previous question you have seen that Markov's inequality may be quite loose. In this question we will show that in some situations it is actually tight. Let  $\varepsilon^*$  be fixed. Design an example of a random variable  $X$  for which

$$\mathbb{P}(X \geq \varepsilon^*) = \frac{\mathbb{E}[X]}{\varepsilon^*}.$$

Prove that the above equality holds for your random variable.

*Hint:* It is possible to design an example satisfying the above requirement with a random variable  $X$  that accepts just two possible values. What should be the values and the probabilities that  $X$  accepts these values?

## 4 The effect of scale (range) and normalization of random variables in Hoeffding's inequality (14 points)

Prove that Corollary 2.5 in Yevgeny's lecture notes (simplified Hoeffding's inequality for random variables in the  $[0, 1]$  interval) follows from Theorem 2.3 (general Hoeffding's inequality). [Showing this for one of the two inequalities is sufficient.]

## 5 Logistic Regression

### 5.1 Cross-entropy error measure (11 points)

Read section 3.3 in the course textbook (Abu-Mostafa et al., 2012). You can also find a scanned version of the chapter on Absalon. Solve exercise 3.6 on page 92 in the course textbook. The *in-sample error*  $E_{\text{in}}$  corresponds to what we call the empirical risk (or training error).

### 5.2 Logistic regression loss gradient (13 points)

Solve exercise 3.7 on page 92 in the course textbook (Abu-Mostafa et al., 2012).

The book assumes labels in  $\{-1, 1\}$ . Solve exercise 3.7 again assuming the labels  $\{0, 1\}$ , which leads to

$$\nabla E_{\text{in}}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N [y_n - \theta(\mathbf{w}^T \mathbf{x})] \mathbf{x}_n .$$

hai 1 e predicti prob alta , l'errore è basso, se predicti 0 hai un errore alto, se hai 0 e predicti 1 quindi prob alta, hai una quantità negativa?? se hai 0 e predicti poco ha un errore basso o negativo???

Hints: Do not forget the “Argue ... one.” part in the exercise for both parts of this question. For the  $\{0, 1\}$  case the slides provide the answer, you just need to add an explanation and intermediate steps.

### 5.3 Log-odds (11 points)

We consider binary logistic regression. Let the input space be  $\mathbb{R}^d$  and the label space be  $\{0, 1\}$ . Let our model  $f$  with parameters  $\mathbf{w} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  model:

$$f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = P(Y = 1 | X = \mathbf{x}) \quad (1)$$

Prove that if the (affine) linear part of the model encodes the log-odds, that is, if

$$\mathbf{w}^T \mathbf{x} + b = \ln \frac{P(Y = 1 | X = \mathbf{x})}{P(Y = 0 | X = \mathbf{x})} , \quad (2)$$

then  $\sigma$  is the logistic function. That is, if  $\mathbf{w}^T \mathbf{x} + b$  encodes on log-scale how frequent class 1 occurs relative to class 0, then  $\sigma$  is the logistic function.

## References

Y. S. Abu-Mostafa, M. Magdon-Ismael, and H.-T. Lin. *Learning from Data*. AMLbook, 2012.