

Foundations of data science, summer 2020  
JONATHAN LENNARTZ, MICHAEL NÜSKEN, ANNIKA TARNOVSKI

**11. Exercise sheet**

**Hand in solutions until Thursday, 2 July 2020, 12:00**

**Exercise 11.1** (Mixture of densities).

(8 points)

Suppose you are given some random variables  $X^{(i)} \stackrel{\text{iid}}{\leftarrow} \mathbb{R}$  with density  $p_i$ . For the computer scientist: some routine `xi` produces samples of  $X^{(i)}$ . How do you construct a routine `x` that samples acc. to the overlayed density  $\sum_{i < k} w_i p_i$ ? Prove correctness: 8

**Theorem.** Consider  $X^{(i)} \stackrel{\text{iid}}{\leftarrow} p_i$  for  $i < k$  and  $\hat{i} \stackrel{\text{iid}}{\leftarrow} w$ , reading  $w$  as a distribution on  $\mathbb{N}_{<k}$ . Finally, let  $X \leftarrow X^{(\hat{i})}$ . Then  $X \sim p$ .

*Hint:*  $X \sim p$  means that  $p$  is the density of  $X$ , ie.  $\text{prob}(X \in [a, b]) = \int_a^b p(x) \, dx$  for all  $a < b$ .

**Remark:** This generalizes to random variables with other outputs instead of values in  $\mathbb{R}$ .

**Solution.** To show that the density of  $X$  is  $p$  means that we have to show

$$\text{prob}(X \in [a, b]) = \int_a^b p(x) \, dx$$

by definition. We have

$$X \in [a, b] \Leftrightarrow \exists i < k: X^{(i)} \in [a, b] \wedge w = i.$$

This directly follows from the construction of  $X$ . So

$$\begin{aligned} \text{prob}(X \in [a, b]) &= \sum_{i < k} \text{prob}(X \in [a, b] \mid \hat{i} = i) \cdot \text{prob}(\hat{i} = i) \\ &= \sum_{i < k} \text{prob}(X^{(i)} \in [a, b] \mid \hat{i} = i) \cdot \text{prob}(\hat{i} = i). \end{aligned}$$

Since  $X^{(i)}$  and  $\hat{i}$  are independent, we have

$$\text{prob}(X^{(i)} \in [a, b] \mid \hat{i} = i) = \text{prob}(X^{(i)} \in [a, b]).$$

Thus

$$\begin{aligned}\text{prob}(X \in [a, b]) &= \sum_{i < k} \text{prob}(X^{(i)} \in [a, b]) \cdot w_i \\ &= \sum_{i < k} w_i \cdot \int_a^b p_i(x) \, dx \\ &= \int_a^b \sum_{i < k} w_i p_i(x) \, dx \\ &= \int_a^b p(x) \, dx ,\end{aligned}$$

which was what we wanted. ○

**Exercise 11.2** (Application of the SVD).

(0+13 points)

In this exercise you shall play with the example from

Alex Thomo (2009). Latent Semantic Analysis (Tutorial).

+3

(i) Reprogram it, denote by  $k$  the used dimension.

+5

(ii) Examine the resulting ranking if...

- (a) ...you modify  $k \in \{2, 3, 4, 5\}$ .
- (b) ...you omit the scaling step.
- (c) ...you change the selection of words by omitting words that only occur in a single document or by adding more words.
- (d) ...you use the Euclidean metric instead of the angle metric.

That's a total of at least 24 cases. You need a careful analysis to isolate important insights.

+5

(iii) Redo similar analysis with a larger dataset: You will find documents 11-document\*.txt in the exercises folder, which contain (parts of) the short overviews of some Wikipedia articles.

*Hint:* We expect you to present an analysis with insights, explanations and arguments. So, no large tables or thelike.

**Solution** (Hints). There is no optimal solution to this exercise, because your choices in the implementation and the sheer possibilities in experimenting can lead to a lot of different results, that might not necessarily be always true. The important aspect in a solution are well-planned experiments and a critical treatment of the result.

In general, a higher  $k$  leads to more accurate results, but for some queries, even  $k$  high did not always lead to the expected solution, as our sample size was quite small. ○

