

# FROM BASIC MACHINE LEARNING MODELS TO ADVANCED KERNEL LEARNING

## HOME ASSIGNMENT 1

This homework is due by **November 14, 2022**. It is to be returned by email to pierre.gaillard@inria.fr as a **pdf** report together with the ipython notebook (or python file) used for the code. The results and the figures must be included into the pdf report but not the code.

The goal of this project is to automatically classify letters from different computer fonts. An example of samples of the letter “A” can be seen below.



The data comes from the `notMNIST` dataset and can be downloaded at <https://kernel-learning.github.io/docs/data1.zip>. The zip archive contains two folders:

- train: contains  $n = 6\,000$  labelled images of three classes “A”, “B” and “C” (2000 each)
- test: contains  $n_1 = 750$  labelled images (250 for each of the three classes).

The train folder will be used to train the forecasting methods. The test folder will be used to assess their performance. If for some reasons, the datasets are too large to be used on your computer, you can use subsets of with  $n$  and  $n_1$  sufficiently small to be computable but large enough to get prediction accuracy.

The goal is to classify if an image  $X_i$  corresponds to the letter “A”: i.e., the output is  $Y_i = 1$  if image  $i$  is “A” and  $-1$  otherwise (if the image is “B” or “C”).

1. Formalize the problem by defining the input space  $\mathcal{X}$ , the output space  $\mathcal{Y}$  and the training data set. What are their dimension?
2. If  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is a predictor from images to  $\mathcal{Y} = \{-1, 1\}$ , we define for a couple image/label  $(X_i, Y_i)$ :
  - the 0-1 loss:  $\ell_1(f_\theta(X_i), Y_i) = \mathbb{1}_{f_\theta(X_i) \neq Y_i}$
  - the square loss:  $\ell_2(f_\theta(X_i), Y_i) = (f_\theta(X_i) - Y_i)^2$
  - the logistic loss:  $\ell_3(f_\theta(X_i), Y_i) = \log(1 + e^{-Y_i f_\theta(X_i)})$ .
  - (a) What are the empirical risk (training error) associated with the 0-1 loss and the true risk? Why is it complicated to minimize the empirical risk in this case?
  - (b) Why should we use the test data to assess the performance?
  - (c) Recall the definition of the optimization problems associated with the linear least square regression, the linear logistic regression.
  - (d) What is the probability of  $\mathbb{P}(Y = 1|X)$  under the logistic model?

- (e) Explain the link between logistic regression and maximum likelihood (justify the choice of the logistic loss above) and the practical reasons why we work with the log of the likelihood.
3. You are performing least-squares polynomial regression. As the degree of your polynomials increases, which of the following is commonly seen to go down at first but then go up?
- (a) Training error (c) Variance  
(b) Test error (d) Bias
4. Assuming that  $f_\theta(x) = \langle \theta, x \rangle$ . Write the update rule of SGD for
- (a) the linear least-squares regression (c) the perceptron algorithm.  
(b) the logistic regression
5. Implement the stochastic gradient descent algorithm (SGD) to solve these problems.
- (a) Consider the logistic regression minimization problem. Plot the training errors and the test errors as functions of the number of access to the data points of SGD for well-chosen (by hand) values of the step sizes.
- (b) Denote by  $\hat{\beta}_n^{\text{logist}}(t) \in \mathbb{R}^{28 \times 28}$  the estimator of logistic regression after  $t$  gradient iterations of SGD. Plot as images the estimators  $\hat{\beta}_n^{\text{logist}}(t) \in \mathbb{R}^{28 \times 28}$  for  $t \in \{10, 100, 1000, 10000\}$ . Repeat for the linear least squares regression (OLS) and perceptron.
6. k-Nearest Neighbors (KNN).
- (a) Recall the definition of the  $k$ -nearest neighbors classification rule with  $\ell_2$  metric.  
(b) Implement it and plot as a function of  $k$ , its training and test errors.  
(c) Calibrate  $k$  using  $K$ -fold cross-validation with  $K = 5$ .
7. Multi-layer Perceptron (MLP). Given a multi-layer perceptron with 1 input layer containing  $28 \times 28$  neurons, 1 hidden layer containing 32 neurons, and 1 output layer containing 3 neurons (one for each class), how many parameters need to be trained? Implement it with ReLu activation function by using your favorite library.
8. Fill the following table and comment:

	Logistic regression	OLS	Perceptron	KNN	MLP
Empirical error (0-1 loss)					
Test error (0-1 loss)					

9. Why is it often important to regularize? What would be the updates of the three models of question 5) with
- (a)  $\ell_2$  regularization  $\lambda \|\theta\|_2^2$  (b)  $\ell_1$  regularization  $\lambda \|\theta\|_1$ ?
10. (Optional) Regularization of logistic regression. Using SGD, solve the following minimization problem (up to small enough precision):

$$\hat{\theta}_\lambda \in \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-Y_i \langle \theta, X_i \rangle}) + \lambda \|\theta\|_1.$$

- (a) Plot the test and training classification errors with the 0-1 loss associated with  $\hat{\theta}_\lambda$  as a function of  $\lambda$ .  
(b) What would be the best value for  $\lambda$ ? How would you tune it?  
(c) Plot as images the estimators  $\hat{\theta}_\lambda$  for four values of  $\lambda$ .  
(d) Repeat questions 10.a-c) by replacing the  $\ell_1$  regularization with a  $\ell_2$  regularization  $\lambda \|\theta\|_2^2$ .