

# Supervised Learning (COMP0078) – Coursework 2

Due : 03 January 2024.

## Submission

You may work in groups of up to two. You should produce a report (this should be in .pdf format) about your results. You will not only be assessed on the **correctness/quality** of your answers but also on **clarity of presentation**. Additionally make sure that your code is *well commented*. Please submit on moodle i) your report as well as a ii) zip file with your source code. Finally please ensure that if you are working in a group both of your student IDs are in the report. If you are working individually, you should include your student ID. Examples of lines such as for

**Note.** Each coursework part should be answered in a separate file. You are not allowed to assume the results of previous parts. If you are not able to prove a result, you should state this. In these cases, you are allowed to assume the results of previous parts. You should state this. You are not allowed to assume the results of previous parts. You should state this.

## 1 PART I [20%]

### Rademacher Complexity

In this problem we will find a bound on the Rademacher complexity of a set of hypotheses that depends only logarithmically on the size of the set. This will then be used to prove the bound on the generalization error seen in the previous part.

We will first show an intermediate result. Let  $X_1, \dots, X_m$  be independent random variables (namely  $\mathbb{E}X_i = 0$  for all  $i = 1, \dots, m$ ).

**1.1 [2 marks].** Let  $\bar{X} = \max_{i=1, \dots, m} X_i$ . Show that for any  $\lambda > 0$

$$\mathbb{E}\bar{X} \leq \frac{1}{\lambda} \log \mathbb{E}e^{\lambda\bar{X}}$$

**1.2 [5 marks].** Show that

$$\frac{1}{\lambda} \log \mathbb{E}e^{\lambda\bar{X}} \leq \frac{1}{\lambda} \log m + \lambda \frac{(b-a)^2}{8}$$

*Hint: use **Hoeffding's Lemma**: for any random variable  $X$  such that  $X - \mathbb{E}X \in [a, b]$  with  $a, b \in \mathbb{R}$ , and for any  $\lambda > 0$ , we have*

$$\mathbb{E}e^{\lambda(X - \mathbb{E}X)} \leq e^{\lambda^2(b-a)^2/8}$$

**1.3 [3 marks].** Conclude that by choosing  $\lambda$  appropriately,

$$\mathbb{E} \max_{i=1, \dots, m} X_i \leq \frac{b-a}{2} \sqrt{2 \log(m)}$$

as desired.

We are almost ready to provide the bound for the Rademacher complexity of a finite set of hypotheses. Let  $S$  a finite set of points in  $\mathbb{R}^n$  with cardinality  $|S| = m$ . We can define the Rademacher complexity of  $S$  similarly to how we have done for the Rademacher complexity of a space of hypotheses:

$$\mathcal{R}(S) = \mathbb{E}_{\sigma} \max_{x \in S} \frac{1}{n} \sum_{j=1}^n \sigma_j x_j,$$

With  $\sigma_1, \dots, \sigma_n$  Rademacher variables (independent and uniformly sampled from  $\{-1, 1\}$ ).

**1.4 [3 marks].** Show that

$$\mathcal{R}(S) \leq \max_{x \in S} \|x\|_2 \frac{\sqrt{2 \log(m)}}{n}$$

with  $\|\cdot\|_2$  denoting the Euclidean norm.

**1.5 [7 marks].** Let  $\mathcal{H}$  be a set of hypotheses with cardinality  $|\mathcal{H}| < +\infty$ . Let  $S = (x_i)_{i=1}^n$  be a set of points in  $\mathbb{R}^n$ . Let  $\hat{\mathcal{R}}(S, \mathcal{H})$  be the empirical Rademacher complexity of  $\mathcal{H}$  on  $S$ . Show that  $\hat{\mathcal{R}}(S, \mathcal{H})$  is an upper bound for  $\mathcal{R}(\mathcal{H})$  with probability at least  $1 - \delta$ .

## 2 PART II [40%]

### Bayes Decision Rule and Empirical Risk

In (binary) classification problems, we are given a set of training points  $S = \{(x_i, y_i)\}_{i=1}^n$  where  $\mathcal{Y} = \{1, -1\}$ . The quality of a classification rule  $c : \mathcal{X} \rightarrow \mathcal{Y}$ , is measured by the classification error  $R(c) = \mathbb{P}(c(x) \neq y)$ .

assuming to sample an input-output pair  $(x, y)$  from the joint distribution  $\rho$ .

**2.1 [4 marks].** Let  $\mathbf{1}_{y=y'}$  be the indicator function of the event  $y = y'$ . Show that the misclassification error can be written as  $R(c) = \int \mathbf{1}_{c(x) \neq y} d\rho(x, y)$ .

**(Surrogate Approaches)** Since minimizing the misclassification error is often difficult, we address the learning problem directly and in practice one usually looks for a real valued function  $f : \mathcal{X} \rightarrow \mathbb{R}$  solving a so-called *surrogate problem*

$$\mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \ell(f(x), y) d\rho(x, y)$$

where  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is a “suitable” convex loss function that makes the surrogate learning problem more amenable to computations. Given a function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , a classification rule  $c_f : \mathcal{X} \rightarrow \{-1, 1\}$  is given in terms of a “suitable” map  $d : \mathbb{R} \rightarrow \{-1, 1\}$  such that  $c_f(x) = d(f(x))$  for all  $x \in \mathcal{X}$ . Here we will look at some surrogate frameworks.

A good surrogate method satisfies the following two properties:

**(Fisher Consistency).** Let  $f_* : \mathcal{X} \rightarrow \mathbb{R}$  denote the expected risk minimizer for  $\mathcal{E}(f_*) = \inf_{f : \mathcal{X} \rightarrow \mathbb{R}} \mathcal{E}(f)$ , we say that the surrogate framework is Fisher consistent if

$$R(c_{f_*}) = \inf_{c : \mathcal{X} \rightarrow \{-1, 1\}} R(c)$$

**(Comparison Inequality).** The surrogate framework satisfies a *comparison inequality* if for any  $f : \mathcal{X} \rightarrow \mathbb{R}$

$$R(c_f) - R(c_{f_*}) \leq \sqrt{\mathcal{E}(f) - \mathcal{E}(f_*)}$$

In particular, if we have an algorithm producing estimators  $f_n$  for the surrogate problem such that  $\mathcal{E}(f_n) \rightarrow \mathcal{E}(f_*)$  for  $n \rightarrow +\infty$ , we automatically have  $R(c_{f_n}) \rightarrow R(c_{f_*})$ .

**2.2 [4 marks].** (Assuming to know  $\rho$ ), calculate the closed-form of the minimizer  $f_*$  of  $\mathcal{E}(f)$  for the:

- a) squared loss  $\ell(f(x), y) = (f(x) - y)^2$
- b) exponential loss  $\ell(f(x), y) = \exp(-yf(x))$ ,
- c) logistic loss  $\ell(f(x), y) = \log(1 + \exp(-yf(x)))$ ,
- d) hinge loss  $\ell(f(x), y) = \max(0, 1 - yf(x))$ .

(hint: recall that  $\rho(x, y) = \rho(y|x)\rho_{\mathcal{X}}(x)$  with  $\rho_{\mathcal{X}}$  the marginal distribution of  $\rho$  on  $\mathcal{X}$  and  $\rho(y|x)$  the corresponding conditional

$$\mathcal{E}(f) = \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(f(x), y) \rho(y|x) \rho_{\mathcal{X}}(x) dy dx$$

you can now solve the problem

**2.3 [4 marks].** The minimizer of  $\mathcal{E}(f)$  is called *Bayes decision rule*. Write explicitly the Bayes decision rule (i.e. the Bayes prior).

**2.4 [4 marks].** Are the surrogate risk  $R$  and the expected risk  $\mathcal{E}$  convex? If yes, can you find a minimizer of the surrogate risk  $R$ ? If it is the case, can you find a minimizer of the expected risk  $\mathcal{E}$ ?

**(Comparison Inequality for the sign function).** Let  $f_*$  be a minimizer of the expected risk  $\mathcal{E}$  for the surrogate risk  $R$  (2.2). Let  $\text{sign} : \mathbb{R} \rightarrow \{-1, 1\}$  denote the “sign” function.

Prove the following comparison inequality:

$$0 \leq R(\text{sign}(f)) - R(\text{sign}(f_*)) \leq \sqrt{\mathcal{E}(f) - \mathcal{E}(f_*)},$$

by showing the following intermediate steps:

**2.5.1 [8 marks].**  $|R(\text{sign}(f)) - R(\text{sign}(f_*))| = \int_{\mathcal{X}_f} |f_*(x)| d\rho_{\mathcal{X}}(x)$ ,  
Where  $\mathcal{X}_f = \{x \in \mathcal{X} \mid \text{sign}(f(x)) \neq \text{sign}(f_*(x))\}$ .

**2.5.2 [8 marks].**  $\int_{\mathcal{X}_f} |f_*(x)| d\rho_{\mathcal{X}}(x) \leq \int_{\mathcal{X}_f} |f_*(x) - f(x)| d\rho_{\mathcal{X}}(x) \leq \sqrt{\mathbb{E}(|f(x) - f_*(x)|^2)}$ .  
Where  $\mathbb{E}$  denotes the expectation with respect to  $\rho_{\mathcal{X}}$

**2.5.3 [8 marks].**  $\mathcal{E}(f) - \mathcal{E}(f_*) = \mathbb{E}(|f(x) - f_*(x)|^2)$ .

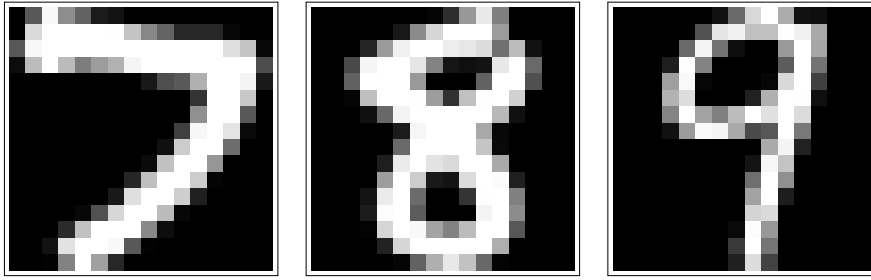


Figure 1: Scanned Digits

### 3 PART III [40%]

#### Kernel perceptron (Handwritten Digit Classification)

**Introduction:** In this exercise, the task is quasi-realistic and you will encounter data which you would not encounter in a standard machine learning course.

You may already be familiar with the perceptron algorithm. In this exercise, we first generalize the perceptron algorithm to a linear separating surface and second, we generalize it to a non-linear separating surface. That is, instead of separating only two classes we will separate multiple classes.

**Adding a kernel:** The kernel perceptron algorithm uses a set of basis functions so that classifying data is equivalent to finding a linear separating surface in a higher-dimensional space. We will consider a single type of kernel, the polynomial kernel, which is controlling the dimension of the space.

**Training and testing the kernel perceptron:** The kernel perceptron algorithm operates on a single example  $(\mathbf{x}_t, y_t)$  and a kernel function  $K(\mathbf{x}_t, \cdot)$  is added for each example. We repeatedly cycle through the training set, and the kernel function no longer changing when we cycle through the training set. On some datasets that the classifier may not converge, but on some datasets that the classifier may converge better if not trained to convergence. We will consider training a particular classifier algorithm to a batch algorithm. The perceptron algorithm describes training for a single epoch, however, explicit notation for multiple epochs is represented by repeating the dataset with the  $\mathbf{x}_i$ 's renumbered. I.e., suppose we have a 40 element training set  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_{40}, y_{40})\}$  to model additional epochs simply extend the data by duplication, hence an  $m$  epoch dataset is

$$\underbrace{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{40}, y_{40})}_{\text{epoch 1}}, \underbrace{(\mathbf{x}_{41}, y_{41}), \dots, (\mathbf{x}_{80}, y_{80})}_{\text{epoch 2}}, \dots, \underbrace{(\mathbf{x}_{(m-1) \times 40 + 1}, y_{(m-1) \times 40 + 1}), \dots, (\mathbf{x}_{(m-1) \times 40 + 40}, y_{(m-1) \times 40 + 40})}_{\text{epoch m}}$$

where  $\mathbf{x}_1 = \mathbf{x}_{41} = \mathbf{x}_{81} = \dots = \mathbf{x}_{(m-1) \times 40 + 1}$ , etc. Testing is performed as follows, once we have trained a classifier  $\mathbf{w}$  on the training set, we simply use the trained classifier with only the *prediction* step for each example in test set. It is a mistake when ever the prediction  $\hat{y}_t$  does not match the desired output  $y_t$ , thus the test error is simply the number of mistakes divided by test set size. Remember in testing the *update* step is never performed.



digits. The task is to classify these digits. The task is to classify these digits. The task is to classify these digits.

on in two ways, linear separating surface. That is, instead of separating only two classes we will separate multiple classes.

as we did with the perceptron algorithm. We will consider a single type of kernel, the polynomial kernel, which is controlling the dimension of the space.

algorithms operate on a single kernel function. In this exercise, we first generalize the perceptron algorithm to a linear separating surface and second, we generalize it to a non-linear separating surface. That is, instead of separating only two classes we will separate multiple classes. We will consider training a particular classifier algorithm to a batch algorithm. The perceptron algorithm describes training for a single epoch, however, explicit notation for multiple epochs is represented by repeating the dataset with the  $\mathbf{x}_i$ 's renumbered. I.e., suppose we have a 40 element training set  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_{40}, y_{40})\}$  to model additional epochs simply extend the data by duplication, hence an  $m$  epoch dataset is



	Two Class Kernel Perceptron (training)
<b>Input:</b>	$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \in (\mathbb{R}^n, \{-1, +1\})^m$
<b>Initialization:</b>	$\mathbf{w}_1 = \mathbf{0}$ ( $\alpha_0 = 0$ )
<b>Prediction:</b>	Upon receiving the $t$ th instance $\mathbf{x}_t$ , predict $\hat{y}_t = \text{sign}(\mathbf{w}_t(\mathbf{x}_t)) = \text{sign}(\sum_{i=0}^{t-1} \alpha_i K(\mathbf{x}_i, \mathbf{x}_t))$
<b>Update:</b>	<p>if <math>\hat{y}_t = y_t</math> then <math>\alpha_t = 0</math>  else <math>\alpha_t = y_t</math>  <math>\mathbf{w}_{t+1}(\cdot) = \mathbf{w}_t(\cdot) + \alpha_t K(\mathbf{x}_t, \cdot)</math></p>

**Generalizing to  $k$  classes:** Design a method (or research a method) to generalise your two-class classifier to  $k$  classes. The method should return a vector  $\boldsymbol{\kappa} \in \mathbb{R}^k$  where  $\kappa_i$  is the “confidence” in label  $i$ ; then you should predict either with a label that maximises confidence or alternately with a randomised scheme.

I’m providing you with mathematica code for a 3-classifier and a demonstration on a small subset of the data. First, however, my mathematica implementation is flawed and is relatively inefficient for large datasets. One aspect of your goals are to improve my code so that it can work on larger datasets. The mathematical logic of the algorithm should not change, however either the program logic and/or the data structures will need to change. You may find it useful to compare my code with a fast code in Python (or the latest version of R) to see how to implement sufficiently

**Files:** From

4. Within the dataset relative to your experiments there will be five hardest to predict correctly “pixelated images.” Print out the visualisation of these five digits along with their labels. Is it surprising that these are hard to predict? Explain why in your opinion that is the case.
5. Repeat 1 and 2 ( $d^*$  is now  $c$  and  $\{1, \dots, 7\}$  is now  $S$ ) above with a Gaussian kernel

$$K(\mathbf{p}, \mathbf{q}) = e^{-c\|\mathbf{p}-\mathbf{q}\|^2},$$

$c$  the width of the kernel is now a parameter which must be optimised during cross-validation however, you will also need to perform some initial experiments to decide a reasonable set  $S$  of values to cross-validate  $c$  over.

6. Choose (research) an alternate method to generalise the kernel perceptron to  $k$ -classes then repeat 1 and 2.

**Assessment:** In your report you will not only be assessed on the correctness/quality of your experiment (e.g., sound methods for choosing parameters, reasonable final test errors) but also on the clarity of presentation and the insightfulness of your observations. Thus the aim is that your report is sufficiently detailed so that the reader could largely reconstruct your experiment from the report alone. The report should also contain the following:

- A discussion of any parameters that were varied.
- A discussion of the two different methods of choosing  $d^*$  for the two different classifiers.
- A discussion comparing the two methods.
- A discussion of your implementation of the kernel perceptron. Discuss how the sum  $\mathbf{w}(\cdot) = \sum_{i=0}^m \alpha_i K(\mathbf{x}_i, \cdot)$  is calculated. Discuss how the sum during training.
- Any table produced in 1 and 2.

**Note: (further comments)** The report should be a qualitative assessment of the performance rather than a quantitative one. It should be clear that a report that is merely “correct/good as a baseline can be expected to receive 32-40 points. Regarding page limits the experiment should be no more of three pages of text (this does not include figures and tables). (The assignment). There is no strict limit, however, a report that is more than one page may be sufficient) and there are a number of things that you should do:

