

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 in a group, you should implement regression using multiple regression lines such as for example cross-validation should be used.

Note. Each coursework part should be an individual sub-part (e.g. 1.1, 1.2, 1.3). If you are not able to prove a result in these cases, you are allowed to assume the results in the previous parts to prove them.

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 number of hypotheses. This will improve 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 with $\mathbb{E}X_i = 0$ for all $i = 1, \dots, m$.

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

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

1.2 [5 marks]. Show that

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

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

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

1.3 [3 marks]. Conclude that by choosing ϵ 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:

$$R(S) = \mathbb{E} \max_{x \in S} \frac{1}{n} \sum_{j=1}^n \epsilon_j x_j$$

With $\epsilon_1, \dots, \epsilon_n$ Rademacher variables (independent and uniformly sampled from $\{-1, 1\}$).

1.4 [3 marks]. Show that

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

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

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

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 $Y = \{-1, 1\}$. The quality of a hypothesis $h: X \rightarrow Y$ is measured by the classification error $E(h) = \mathbb{P}(h(x) \neq y)$.

assuming to sample an input-output pair (x, y) from the joint distribution \mathbb{P} .

2.1 [4 marks]. Let $\mathbb{P}_{y=y^0}$ be the marginal distribution of y when $y = y^0$. Show that the classification error of a hypothesis h can be written as $E(h) = \mathbb{P}_{y=y^0}(h(x) \neq y)$.

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

$$E(f) = \int_{X \times Y} \eta(f(x); y) d\mathbb{P}(x, y)$$

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

A good surrogate method satisfies the following two properties:

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

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

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

$$R(c_f) - R(f) \leq \rho \frac{E(f) - E(f^*)}{\rho}$$

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

2.2 [4 marks]. (Assuming to know), calculate the closed-form of the minimizer f^* of $E(f)$ 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 $(x; y) = (y|x) \cdot x(x)$ with x the marginal distribution of x on X and $(y|x)$ the corresponding conditional

$$E(f) = \int f(x) d\mu(x)$$

you can now solve the problem

2.3 [4 marks]. The minimizer of the surrogate risk is called *Bayes decision rule*. Write explicitly the Bayes decision rule (in terms of the Bayes prior).

2.4 [4 marks]. Are the surrogate risk and the expected risk $E(f)$ such that $E(f) \leq R(f)$ for all f ? If it is true, can you find a minimizer of the surrogate risk E ? If it is false, provide a counterexample.

(Comparison Inequality for the surrogate risk) Let f^* be a minimizer of the expected risk for the surrogate loss ℓ (see (2.2)). Let $\text{sign} : \mathbb{R} \rightarrow \{-1, 1\}$ denote the "sign" function defined by

Prove the following comparison inequality:

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

by showing the following intermediate steps:

2.5.1 [8 marks]. $R(\text{sign}(f)) - R(\text{sign}(f^*)) = \int_{X_f} |f(x) - f^*(x)| d\mu(x)$,
Where $X_f = \{x \in X : \text{sign}(f(x)) \neq \text{sign}(f^*(x))\}$:

2.5.2 [8 marks]. $\int_{X_f} |f(x) - f^*(x)| d\mu(x) = \int_{X_f} |f(x) - f^*(x)|^2 d\mu(x) \leq E((f(x) - f^*(x))^2)$.
Where E denotes the expectation with respect to μ .

2.5.3 [8 marks]. $E(f) - E(f^*) = E((f(x) - f^*(x))^2)$.

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 be working with a dataset which you would not encounter in a standard machine learning course.

You may already be familiar with the perceptron algorithm. First we generalize the perceptron to a linear classifier, then we generalize the linear classifier to a non-linear classifier by using a kernel function. Second, we generalize the perceptron to a non-linear classifier by using a kernel function. Finally, we generalize the perceptron to a non-linear classifier by using a kernel function.

Adding a kernel: The kernel function is a function that takes two inputs and returns a scalar value. It is used to map the input data into a higher-dimensional space where it is easier to separate the classes. The kernel function is used to map the input data into a higher-dimensional space where it is easier to separate the classes.

Training and testing the kernel perceptron: The kernel perceptron is trained on a single example (\mathbf{x}_t, y_t) and the kernel function $K(\mathbf{x}_t, \cdot)$ is added for each example. The training process involves repeatedly cycling through the training set and updating the weights. The testing process involves using the trained classifier to predict the class of new examples. The testing process involves using the trained classifier to predict the class of new examples.

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

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

where $\mathbf{x}_1 = \mathbf{x}_{41} = \mathbf{x}_{81} = \dots = \mathbf{x}_{(m-1) \cdot 40 + 1}$, etc. Testing is performed as follows, once we have trained a classifier 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 into one of the 10 classes (0-9).

The task is to classify these digits into one of the 10 classes (0-9). The task is to classify these digits into one of the 10 classes (0-9).

The task is to classify these digits into one of the 10 classes (0-9). The task is to classify these digits into one of the 10 classes (0-9).

The task is to classify these digits into one of the 10 classes (0-9). The task is to classify these digits into one of the 10 classes (0-9). The task is to classify these digits into one of the 10 classes (0-9).

	Two Class Kernel Perceptron (training)
Input:	$f(\mathbf{x}_1; y_1), \dots, (\mathbf{x}_m; y_m) g 2 (<^n; f^{-1}; +1g)^m$
Initialization:	$\mathbf{w}_1 = \mathbf{0} \quad (\phi_0 = 0)$
Prediction:	Upon receiving the t th instance \mathbf{x}_t ; predict $\hat{y}_t = \text{sign}(\mathbf{w}_t(\mathbf{x}_t)) = \text{sign}(\sum_{i=0}^t \frac{1}{i} K(\mathbf{x}_i; \mathbf{x}_t))$
Update:	<p>if $\hat{y}_t = y_t$ then $\phi_t = 0$</p> <p>else $\phi_t = y_t$</p> <p>$\mathbf{w}_{t+1}(\cdot) = \mathbf{w}_t(\cdot) + \phi_t K(\mathbf{x}_t; \cdot)$</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 $\hat{\mathbf{y}} \in \mathbb{R}^k$ where \hat{y}_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. I'm providing you with the fast code in Python (or the last version of Mathematica) to implement sufficiently.

Files: From <http://www0.cs.cmu.edu/~paulsm/10701/assignments/assignment3/> files relevant to

this assignment:
 poorCodeDemoD
 dtrain123.dat
 dtest123.dat
 zipcombo.dat

each of the data files consists of a 1000x256 matrix where the first value is the digit, the remaining 256 values are the pixel intensities. In attempting to understand the data, I found that the first value is between -1 and 1. However, remember the demonstration that the Mathematica code, though less efficient, requires thought and observation of behaviour on the perceptron.

Experimental Protocol: You should be percentages not raw counts. (errors reported)

1. Basic Results: Perform 20 runs : when using the 80% training data split from within to perform 5-fold cross-validation to select the "best" parameter d then retrain on full 80% training set using d and then record the test errors on the remaining 20%. Thus you will find 20 d and 20 test errors. Your final result will consist of a mean test error std and a mean d with std.
2. Cross-validation: Perform 20 runs : when using the 80% training data split from within to perform 5-fold cross-validation to select the "best" parameter d then retrain on full 80% training set using d and then record the test errors on the remaining 20%. Thus you will find 20 d and 20 test errors. Your final result will consist of a mean test error std and a mean d with std.
3. Confusion matrix: Perform 20 runs : when using the 80% training data split that further to perform 5-fold cross-validation to select the "best" parameter d retrain on the full "80%" training set using d and then produce a *confusion matrix*. Here the goal is to find "confusions" thus if the true label (on the test set) was "7" and "2" was predicted then a "error" should be recorded for "(7,2)"; the final output will be a 10x10 matrix where each cell contains a confusion error *rate* and its standard deviation (here you will have averaged over the 20 runs). Note the diagonal will be 0. In computing the error rate for a cell use

$$\frac{\text{"Number of times digit } a \text{ was mistaken for digit } b \text{ (test set)"}}{\text{"Number of digit } a \text{ points (test set)"}};$$

4. Within the dataset relative to your experiments there will be "very hard to predict correctly" pixelated images." Print out the visualisation of these "very hard to predict" 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 $f1, \dots, 7g$ is now S) above with a Gaussian kernel

$$K(p; q) = e^{-ckp - qk^2};$$

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

- A discussion of any parameters.
- A discussion of the two models.
- A discussion comparing the two models.
- A discussion of your implementation.
 - sum $w() = \sum_{i=0}^m i K(x_i)$
 - sum during training.
- Any table produced in 1.

Note: (further comments rather it will be a qualitative correct/good as a baseline can Regarding page limits the exp pages of text (this does not incl There is no strict limit, howev be su cient) and there are a

