

Introduction to machine learning

Exercise 2

Winter 2024

Submission guidelines, **read and follow carefully**:

- The exercise **must** be submitted in pairs.
- Submit via Moodle.
- The submission should include two separate files:
 1. A pdf file that includes your answers to all the questions.
 2. The code files for the python question. You must submit a copy of the shell python file provided for this exercise in Moodle, with the required functions implemented by you. **Do not change the name of this file!** In addition, you can also submit other code files that are used by the shell file.
- Your python code should follow the course python guidelines. See the Moodle website for guidelines and python resources.
- Before you submit, **make sure that your code works in the course environment**, as explained in the guidelines. Specifically, **make sure that the test `simple_test` provided in the shell file works**.
- You may only use python modules that are explicitly allowed in the exercise or in the guidelines. If you are wondering whether you can use another module, ask a question in the exercise forum. No module containing machine learning algorithms will be allowed.
- For questions, use the exercise forum, or if they are not of public interest, send them via the course requests system.
- Grading: Q.1 (python code): 10 points, Q.2: 20 points, Q.3: 15 points, Q.4: 19 points, Q.5: 18 points, Q.6: 18 points

Question 1. Implement the soft-SVM algorithm that we learned in class in python. The shell file “softsvm.py” is provided for this exercise in Moodle. It contains an empty implementation of the function required below. You should implement it and submit according to the submission instructions.

```
def softsvm(l, trainX, trainy)
```

The input parameters are:

- λ - the parameter λ of the soft SVM algorithm.
- `trainX` - a 2-D matrix of size $m \times d$, where m is the sample size and d is the dimension of the examples. Row i in this matrix is a vector with d coordinates that describes an example x_i from the training sample.
- `trainy` - a column vector of length m . The i 's number in this vector is the label $y_i \in \{-1, 1\}$ from the training sample.

The function returns the linear predictor w which is a column vector in \mathbb{R}^d .

- You may assume all the input parameters are legal.
- We will use the library `cvxopt` for our Quadratic Program solver.

Instructions for using `cvxopt`:

- First, you will need to define the matrices H , u , A , and v which correspond to the vectors and matrices with the same names in the quadratic programming problem you learned in class. Those matrices should be `cvxopt` matrices, check how to create `cvxopt` matrices or convert `numpy` arrays to `cvxopt` matrices here: <http://cvxopt.org/userguide/matrices.html>.
- In order to conserve memory, use sparse matrices when possible.
- Run `sol = cvxopt.solvers.qp(H, u, -A, -v)` to solve the quadratic programming problem. Here, we pass A and v with a minus sign, since this solver assumes the constraints are $Az \leq v$, while in class we assumed they were $Az \geq v$. The solution of the quadratic program is provided in `sol["x"]`.
- See the note at the end of the exercise regarding a possible error and how to solve it.

Question 2. In this question, you will run your soft SVM implementation on data from the MNIST dataset you saw in exercise 1. For this task, we took a subset of this dataset which include and digits 3 and 7, and the goal of the predictor is to distinguish between the two digits. You can load the dataset, which is already divided to train and test, from the file `EX2q2_mnist.npz` on the course website.

Run two experiments on this data set. In the first experiment, use a sample size of 100. To generate this small sample, draw it randomly from the provided training sample. Repeat the “small sample” experiment 10 times, and when you report the results, average over these 10 experiments, and plot also error bars which show the maximum and minimum values you got over all experiments. Run your soft-SVM implementation with each of the following values of λ : $\lambda = 10^n$, for $n \in \{1, \dots, 10\}$.

In the second experiment, use a sample size of 1000, which you should also draw randomly from the training set. Run your soft-SVM implementation with each of the following values of λ : $\lambda = 10^n$, for $n \in \{1, 3, 5, 8\}$. To make the running time feasible, you should run this experiment only once for each value of λ .

- Submit a plot of the training error and test error of the small sample size results as a function of λ (plot λ on a logarithmic scale), with one line for the train error and another line for the test error. Each line should show an average of the 10 experiments, and error bars which show the maximum and minimum values you got over all experiments.
- Add to the plot the points describing the training error and test error of the large sample size. For this part, don't draw lines between the points in this case, only show each point individually, since you tested values of λ which are quite far away from each other.

- (c) Based on what we learned in class, what would you expect the results to look like? Do the results you got match your expectations? In your answer address the following issues:
- Which sample size should get a smaller training error? What about test error? Do the results match your expectations?
 - What should be the trend in the *training error* as a function of λ (decreasing/increasing/other)? Why? Do the results (for the small sample size) match your expectations?
 - What should be the trend in the *test error* as a function of λ (decreasing/increasing/other)? Why? Do the results (for the small sample size) match your expectations?

Question 3. Let $x \in \mathbb{R}$, $\mathcal{Y} \in \{0, 1\}$. Given a text string c which describes a mathematical condition, define $h_c : x \rightarrow \mathcal{Y}$ as follows: $h_c(x) := \mathbb{I}[x \text{ satisfies } c]$. For instance, if c is the condition “ x is a natural number”, then $h_c(x) = \mathbb{I}[x \in \mathbb{N}]$. For any $n \in \mathbb{N}$, define the hypothesis class $\mathcal{H}_n \subseteq \mathcal{Y}^x$ as follows:

$$\mathcal{H}_n := \{h_c \mid c \text{ is a condition which can be described using at most } n \text{ characters}\}$$

For instance, the condition c given above is in \mathcal{H}_n for all $n \geq 21$. The text of the conditions can include only ASCII characters (there are 128 ASCII characters).

Suppose that the examples x are temperature values, and the labels y indicate whether a doctor gives medicine to a person with this temperature. Suppose that we know that the rule that the doctors use to decide whether to give the medicine has at most 10 characters and that the doctors follow this rule exactly. We get a random independent sample from the distribution D over $x \times \mathcal{Y}$ over temperatures and doctor decisions. We would like to run an ERM algorithm with the hypothesis class \mathcal{H}_{10} to find a condition with an error of at most 0.1. We have unlimited computational resources.

- Explain why running the ERM algorithm with the hypothesis class \mathcal{H}_n for $n < 10$ or for $n > 10$ may be a worse idea than running it with \mathcal{H}_{10} .
- Use PAC bounds to calculate an upper bound on the size of a sample size that would be sufficient to find, with a probability at least 0.99, a condition with an error at most 0.1 on the distribution, using an ERM algorithm \mathcal{H}_{10} . Give a formal description of the guarantee that you are providing.
- Suppose we run ERM with \mathcal{H}_n on the training sample for some $n < 10$. Calculate an upper bound on the size of a sample size that would be sufficient to find, with a probability at least 0.99, a condition with an excess error at most 0.1 on the distribution, using an ERM algorithm with \mathcal{H}_n . Give a formal description of the guarantee that you are providing.

Question 4. Consider the hypothesis class of homogeneous linear predictors $\mathcal{H}_L^d := \{h_w \mid w \in \mathbb{R}^d\}$ where $h_w(x) = \text{sign}(\langle w, x \rangle)$, $\forall x \in \mathbb{R}^d$. Prove that any linearly independent set of d vectors, $u_1, \dots, u_d \in \mathbb{R}^d$, are shattered by \mathcal{H}_L^d .

Reminder: A set of input vectors is said to be shattered by an hypothesis class \mathcal{H} if it can get all the possible labelings by predictors from \mathcal{H} .

Question 5. For a given training sample $S = ((x_1, y_1), \dots, (x_m, y_m))$, consider the following **modified version** of the soft-SVM optimization problem:

$$\lambda \|\mathbf{w}\|_1^2 + \sum_{i=1}^m \ell_h(w, (x_i, y_i)),$$

where $\ell_h(w, (x_i, y_i)) = \max\{0, 1 - y_i \langle w, x_i \rangle\}$ is the *hinge loss* defined in the lecture. Recall that $\|\mathbf{w}\|_1$ is the ℓ_1 norm defined as $\|\mathbf{w}\|_1 := \sum_{i=1}^d |w(i)|$.

Express this optimization problem as a quadratic program in standard form, as we showed in class.

- (a) Write a quadratic minimization problem with constraints that is equivalent to the problem above, using auxiliary variables similar to the ξ_i in the soft-SVM implementation.
- (b) Write what H, u, A, v in the definition of a Quadratic Program should be set to so as to solve the minimization problem you wrote.

Question 6. In this exercise, we will see that without margin assumptions, the Perceptron algorithm might run for a long time, exponential in the dimension d . We define a special sample $S = ((x_1, y_1), \dots, (x_m, y_m))$, where $m = d$, $x_i \in \mathbb{R}^d$, and $y_i \in \{-1, 1\}$, where $y_i := (-1)^{i+1}$, and $x_i \in \mathbb{R}^d$ is defined by:

$$x_i(j) = \begin{cases} (-1)^i, & j < i \\ (-1)^{i+1}, & j = i \\ 0, & j > i. \end{cases}$$

In the following steps, you will prove that for the sample above, the Perceptron performs a number of updates at least exponential in d :

- (a) Prove that in any iteration t of the Perceptron on the given sample S , and for any $i \leq d$, $|w^{(t+1)}(i)| \leq t$.
Hint: use induction on the Perceptron update $w^{(t+1)} \leftarrow w^{(t)} + y_i x_i$.
- (b) Let $w^{(T)}$ be the separator that the Perceptron outputs. Prove that for every coordinate i , $|w^{(T)}(i)| \geq 2^{i-1}$. Hint: prove this by induction on i , using the fact that the separator defined by $w^{(T)}$ labels correctly all the examples in the given sample S .
- (c) Conclude from the two previous items that the number of updates of the Perceptron until it stops and outputs $w^{(T)}$ is exponential in d .

To help with intuition, you may want to implement the Perceptron algorithm defined in class (no need to submit the code), and check its behavior on the sample S for $d = 1, 2, \dots, 10$.