

Introduction to machine learning

Exercise 3

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: Q1: 17 points. Q2:18 points. Q3:18 points. Q4:12 points. Q5a:10. Q5b,c:5 points each. Q6:15 points.

Question 1. Implement the soft-margin kernel SVM routine described in class, using `cvxopt` quadratic problem solver you used in Question 1 Assignment 2. The algorithm should use the polynomial kernel. You will implement the function `softsvmpoly` in the shell file “`softsvmpoly.py`” which is provided for this exercise in Moodle. function details:

```
def softsvmpoly(l, k, trainX, trainY)
```

The input parameters are:

- `l` - the parameter λ of the soft SVM algorithm.
- `k` - the degree of the polynomial kernel.
- `trainX` - a 2-D matrix of size $m \times d$. Row i in this matrix is a vector with d coordinates that describes 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 from the training sample. You can assume that each label is either -1 or 1 .

The function returns the column vector `alpha` $\in \mathbb{R}^m$, which contains the coefficients found by the algorithm.

Question 2. For this question, use the data file `EX3q2_data.npz` provided on the course website, which contains data points in the domain $\mathcal{X} = \mathbb{R}^2$ and labels in $\{-1, 1\}$, split into a training set and test set.

- We would like to use soft SVM to learn a predictor for this problem. Plot the points in the training set in \mathbb{R}^2 , and color them by their label. Explain why it may be a better idea use kernel soft SVM and not the linear (non-kernel) soft SVM. You can use the function `matplotlib.pyplot.scatter`.
- Run your Polynomial soft SVM code on the training set. Perform 5-fold cross-validation to tune λ and k . Try the values $\lambda \in \{1, 10, 100\}$ and $k \in \{2, 5, 8\}$ — a total of 9 parameter pairs to try. Report the 9 average validation error values for each of the pairs (λ, k) . Report which pair was selected by the cross validation, rerun the training using this pair on the entire training set, and report the test error of the resulting classifier.
- For a general classification problem, give one reason why a polynomial SVM might get a better validation error than linear soft SVM, and one reason why it might get a worse validation error.
- Set $\lambda = 100$ and consider $k \in \{3, 5, 8\}$. For these values, run the polynomial soft SVM on the training set, and plot the resulting predictor in \mathbb{R}^2 as follows: Define a fixed region (roughly the region in which the training data resides), divide it into a fine grid, and color the grid points red or blue, depending on the label predicted by the classifier for each point. You can use the the function `matplotlib.pyplot.imshow` to plot this.

Question 3. Kernel functions. Consider a space of examples $\mathcal{X} = \mathbb{R}^d$. Let $x, x' \in \mathcal{X}$.

- Prove that the following function *cannot be* a kernel function for any feature mapping:

$$K(x, x') := (2x(7) + x(3)) \cdot x'(2).$$

Hint: what property of inner products does this function violate? How can you prove it?

- Prove that the following function *cannot be* a kernel function for any feature mapping:

$$K(x, x') := 5 - (x(1) - x(2))(x'(1) - x'(2)).$$

Hint: consider the case $x = x'$.

- Show that the following function f is a Kernel function, by providing a possible feature mapping Ψ such that $f(x, x') = \langle \Psi(x), \Psi(x') \rangle$:

$$f(x, x') = (x(1)x'(1))^6 + e^{x(3)+x(5)+x'(3)+x'(5)} + 1/(x(1)x(1')) + (x(4)+x(6))(x'(4)+x'(6)).$$

Question 4. Consider a linear combination of $k \in \{1, 2, \dots\}$ **convex** functions $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ for $i \in \{1, \dots, k\}$:

$$g(u) = \sum_{i=1}^k a_i f_i(u)$$

where $a_1, \dots, a_k \in \mathbb{R}$ and $u \in \mathbb{R}^d$.

- (a) Prove that g is not necessarily a convex function for any $a_1, \dots, a_k \in \mathbb{R}$.
- (b) Prove that g is a convex function if $a_1, \dots, a_k \geq 0$.

Question 5. Let $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ where the input domain is $\mathcal{X} = \mathbb{R}^d$ and the labels are from $\mathcal{Y} = \mathbb{R}$. Consider learning of a linear predictor using the following optimization for $\lambda > 0$:

$$\underset{w \in \mathbb{R}^d}{\text{minimize}} \lambda \|w - v\|_2^2 + \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2.$$

where $v \in \mathbb{R}^d$ is a given (fixed) vector.

Organize the training examples in S by defining the input data matrix $\mathbf{X} = [x_1, x_2, \dots, x_d]$ and the

labels vector $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}$.

- (a) Formulate the closed-form expression for the w that solves the optimization in this question. In the closed-form expression you can use only \mathbf{X} , \mathbf{y} , λ , v that were defined above.
- (b) Suppose that Gradient Descent is run on S with a step size η . Calculate the formula for $w^{(t+1)}$ as a function of $w^{(t)}$ and η (where in addition you can use \mathbf{X} , \mathbf{y} , λ , v). Explain the steps of your derivation.
- (c) What would the update step for $w^{(t+1)}$ be in Stochastic Gradient Descent for the same objective?

Question 6. (Do this question after we learn about PCA in class)

In an experiment, several measurements were taken at times $t = 1, 2, \dots, m$. At each time t , the measurements taken were $x_t(1), x_t(2), x_t(3), x_t(4)$. This created a data set $S = \{x_1, \dots, x_m\}$, where x_t is a vector in \mathbb{R}^4 which includes all the measurements from time t . PCA was performed on the data set S to reduce its dimensionality from 4 to 2.

- (a) In one experiment, it turned out that in all times t , $x_t(3) = 3x_t(1) + x_t(2)$, and $x_t(4) = 2x_t(2) - 4x_t(3)$. What will be the error of the PCA in this case? Prove your claim.
- (b) In another experiment, it turned out that at all times t , $x_t(3) = (x_t(1))^2 + (x_t(2))^3$, and $x_t(4) = (x_t(3) - x_t(1))^2$. Show an example of experiment results that satisfy these equations such that the error of the PCA is larger than the error you showed for the experiment in (a). You may choose m as you like.