

# MACHINE LEARNING

## FIRST SET OF PROBLEMS

Solve

- at least 4 of the 6 subproblems of Problem 1
- at least 2 of the 4 subproblems of Problem 2

You can earn bonus points if you solve more subproblems than the minimum required.

You have a free choice of programming language.

Work in parties of 3 students each.

Deadline: Tuesday, 11 December 2018.

Printed reports required.

### Problem 1

Consider the generalized linear regression problem defined by the following model:

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_5 x^5 + \eta \quad (1)$$

where  $\eta$  corresponds to white Gaussian noise and the components of the weight vector assume the following values:

$$\theta_0 = 0.2, \theta_1 = -1, \theta_2 = 0.9, \theta_3 = 0.7, \theta_5 = -0.2. \quad (2)$$

In every case below, we consider  $N$  equidistant points  $x_1, x_2, \dots, x_n$  in the interval  $[0,2]$  and use them to create samples for our training set:

$$y_n = \theta_0 + \theta_1 x_n + \theta_2 x_n^2 + \theta_3 x_n^3 + \theta_5 x_n^5 + \eta_n, \quad n = 1, 2, \dots, N \quad (3)$$

where  $\eta_n$  are i.i.d. noise samples originating from a Gaussian distribution with mean 0 and variance  $\sigma_\eta^2$ .

- 1) Using  $N = 20$ ,  $\sigma_\eta^2 = 0.1$  and the structure of the correct model (5th degree polynomial with the coefficient of the 4th power equal to zero), apply the Least Squares method to estimate the parameter vector. Calculate the Mean Square Error of  $y$  over the training set and over a test set comprising of 1000 points randomly selected in the interval  $[0,2]$ .
- 2) For  $N = 20$  and  $\sigma_\eta^2 = 0.1$  apply regression using the Least Squares method and a 2nd degree polynomial. Perform 100 experiments using different noise samples for each experiment. For each point of the training set, calculate the mean and variance of  $y$  over the 100 experiments and plot these quantities on the  $(x,y)$  plane along with the curve obtained by the true model.  
Repeat using a 10th degree polynomial. Compare your results obtained for the 2 different cases (2<sup>nd</sup> versus 10<sup>th</sup> degree polynomial) making special reference to the bias-variance dilemma.

- 3) Repeat experiment (1) above, implementing the Ridge Regression method with various values of  $\lambda$  (instead of the Least Squares Method). Report whether you have observed an improvement of the Mean Square Error for some of the values of  $\lambda$ .
- 4) We encode our prior knowledge for the unknown parameter vector via a Gaussian distribution  $G(\theta)$  with mean  $\theta_0$  equal to the true parameter vector in equation (1) and covariance matrix  $\Sigma_\theta = \sigma_\theta^2 I$ ,  $\sigma_\theta^2 = 0.1$ . Use the structure of the true model and perform full Bayesian Inference in order to evaluate  $y$  for 20 randomly selected test set points belonging to the interval  $[0,2]$  and for two different values of  $\sigma_\eta^2$  (0.05 and 0.15). Plot your estimates and their errors on the  $(x,y)$  plane.
- 5) Repeat experiment (4) using the following mean vector for  $G(\theta)$ :  $\theta_0 = [-10.54, 0.465, 0.0087, -0.093, -0.004]^T$

With  $\sigma_\eta^2 = 0.05$ , perform the experiment four times, using two different values for  $\sigma_\theta^2$  (0.1 and 2) and two different values for  $N$  (20 and 500). Comment on your results.

- 6) Try to recover the true variance of the noise using the Expectation-Maximization method. Construct a training set with  $N = 500$  and  $\sigma_\eta^2 = 0.05$ . Initialize the algorithm with  $\alpha = \sigma_\theta^{-2} = 1$ ,  $\beta = \sigma_\eta^{-2} = 1$ . After convergence, estimate the  $y$ 's and their errors over a test set of 20 points randomly selected in the interval  $[0,2]$ . Plot these quantities on the  $(x,y)$  plane, along with the true model curve.

## Problem 2

- 1) Program and implement a k nearest neighbours classifier (k-NN). Use this classifier to solve the following problems:
  - i. IRIS PLANT DATABASE (Classification of three different kinds of iris plants).
  - ii. PIMA INDIANS DIABETES DATABASE (Classification of pregnant Indians of the Pima tribe according to whether they suffer from diabetes or not).
 The relevant data can be found in the file UCIdata-exercise1.rar.  
 Report on the percentage of correct classification as a function of the number of nearest neighbours. Use cross-validation to obtain the results.
- 2) For the second problem, implement a Bayes classifier assuming that the pdfs of the two classes follow normal distributions. Use reasonable assumptions for the covariance matrices, based on the data. Use cross-validation to evaluate the performance of the classifier and compare it to the performance of the k-NN classifier.
- 3) Implement a naive Bayes classifier assuming Gaussian pdfs for the two classes. Evaluate its performance and compare it to the performance of the classifiers in 1) and 2).
- 4) Implement the perceptron algorithm and use it to perform classification on the IRIS PLANT DATABASE data as follows: Examine whether the data of each class are linearly separable from the data of the combined remaining classes (e.g. if the Iris Setosa data are linearly separable from the combined Iris Versicolor and Iris Virginica data).