

# HOMEWORK 3

>>Sean(Xiaoyu) Sun<<  
>>9078202463<<

**Instructions:** Although this is a programming homework, you only need to hand in a pdf answer file. There is no need to submit the latex source or any code. You can choose any programming language, as long as you implement the algorithm from scratch.

Use this latex file as a template to develop your homework. Submit your homework on time as a single pdf file to Canvas. Please check Piazza for updates about the homework.

## 1 A Simplified 1NN Classifier

You are to implement a 1-nearest-neighbor learner for classification. To simplify your work, your program can assume that

- each item has  $d$  continuous features  $\mathbf{x} \in \mathbb{R}^d$
- binary classification and the class label is encoded as  $y \in \{0, 1\}$
- data files are in plaintext with one labeled item per line, separated by whitespace:

$$\begin{array}{cccc} x_{11} & \dots & x_{1d} & y_1 \\ & & \dots & \\ x_{n1} & \dots & x_{nd} & y_n \end{array}$$

Your program should implement a 1NN classifier:

- Use Mahalanobis distance  $d_A$  parametrized by a positive semidefinite (PSD) diagonal matrix  $A$ . For  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ ,

$$d_A(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_A = \sqrt{(\mathbf{x} - \mathbf{x}')^\top A (\mathbf{x} - \mathbf{x}')}.$$

We will specify  $A$  in the questions below. (Hint:  $d$  is dimension while  $d_A$  with a subscript is distance)

- If multiple training points are the equidistant nearest neighbors of a test point, you may use any one of those training points to predict the label.
- You do not have to implement kd-tree.

## 2 Questions

1. (5 pts) What is the mathematical condition on the diagonal elements for a diagonal matrix  $A$  to be PSD?
2. (5 pts) Given a training data set  $D$ , how do we preprocess it to make each feature dimension mean 0 and variance 1? (Hint: give the formula for  $\hat{\mu}_j, \hat{\sigma}_j$  for each dimension  $j$ , and explain how to use them to normalize the data. You may use either the  $\frac{1}{n}$  or  $\frac{1}{n-1}$  version of sample variance. You may assume the sample variances are non-zero.)
3. (5 pts) Let  $\tilde{\mathbf{x}}$  be the preprocessed data. Give the formula for the Euclidean distance between  $\tilde{\mathbf{x}}, \tilde{\mathbf{x}}'$ .
4. (5 pts) Give the equivalent Mahalanobis distance on the original data  $\mathbf{x}, \mathbf{x}'$  by specifying  $A$ . (Hint: you may need  $\hat{\mu}_j, \hat{\sigma}_j$ )
5. (5 pts) Let the diagonal elements of  $A$  be  $a_{11}, \dots, a_{dd}$ . Define a diagonal matrix  $L$  with diagonal  $\sqrt{a_{11}}, \dots, \sqrt{a_{dd}}$ . Define  $\tilde{\mathbf{x}} = L\mathbf{x}$ . Prove that  $d_I(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = d_A(\mathbf{x}, \mathbf{x}')$  where  $I$  is the identity matrix.

6. (5 pts) Geometrically, what does  $Lx$  do to the point  $x$ ? Explain in simple English.
7. (10 pts) Let  $U$  be any orthogonal matrix. Define  $\tilde{x} = ULx$ . (i) Prove that  $d_I(\tilde{x}, \tilde{x}') = d_A(x, x')$  again. (ii) Geometrically, what does  $ULx$  do to the point  $x$ ? Explain in simple English.
8. (20 pts) Use the whole D2z.txt as training set. Use Euclidean distance (i.e.  $A = I$ ). Visualize the predictions of 1NN on a 2D grid  $[-2 : 0.1 : 2]^2$ . That is, you should produce test points whose first feature goes over  $-2, -1.9, -1.8, \dots, 1.9, 2$ , so does the second feature independent of the first feature. You should overlay the training set in the plot, just make sure we can tell which points are training, which are grid.
9. (To normalize, or not to normalize?) Start from D2a.txt. Perform 5-fold cross validation.
  - (a) (5 pts) Do not normalize the data. Report 1NN cross validation error rate for each fold, then the average (that's 6 numbers).
  - (b) (5 pts) Normalize the data. Report 1NN cross validation error rate (again 6 numbers). (Hints: Do not normalize the labels! The relevant quantities should be estimated from the training portion, but applied to both training and validation portions. This should happen 5 times. Also, you would either change  $x$  into  $\tilde{x} = Lx$  but then use Euclidean distance on  $\tilde{x}$ , or do not change  $x$  but use an appropriate  $A$ ; don't mix the two.)
  - (c) (5 pts) Look at D2a.txt, explain the effect of normalization on CV error. Hint: the first 4 features are different than the next 2 features.
10. (Again. 10 pts) Repeat the above question, starting from D2b.txt.
11. (5 pts) What do you learn from Q9 and Q10?
12. (Weka, 10 pts) Repeat Q9 and Q10 with Weka. Convert appropriate data files into ARFF format. Choose classifiers / lazy / IBk. Set  $K = 1$ . Choose 5-fold cross validation. Let us know what else you needed to set. Compare Weka's results to your Q9 and Q10.