

COMP527  
Data Mining and Visualisation  
Problem Set 1

Danushka Bollegala

**Question 1**

- A. State the two main types of data mining models. **(2 marks)**
- B. Consider that you measured the height and weight of 100 students for a health survey. For 20 students in your sample you could only measure either their height or weight, but not both values. Assume that we would like to train a binary classifier to predict whether a student is overweight compared to the students in this dataset. Answer the following questions about this experiment.
- (a) State two algorithms that you can use to learn a binary classifier for this purpose. **(2 marks)**
  - (b) What is meant by the *missing-value* problem in data mining? **(3 marks)**
  - (c) State two disadvantages we will encounter if we ignore the 20 instances that we have incomplete measurements for and use the remaining 80 instances to train the classifier. **(4 marks)**
  - (d) The average height of the students in this dataset is 169cm. Provide a reason for and a reason against using the average to fill the missing values. **(4 marks)**
  - (e) Assume that we would like to check whether there is any correlation between the height and the weight of the students in this dataset. How do we check this? **(4 marks)**
  - (f) Given that there is a high correlation between the height and the weight of a student, how can we use this information to overcome the missing-value problem? **(4 marks)**
  - (g) Without having access to a separate test dataset, how can we evaluate the accuracy of our binary classifier? **(2 marks)**

**Question 2** Consider a training dataset consisting of four instances  $(\mathbf{x}_1, 1)$ ,  $(\mathbf{x}_2, 1)$ ,  $(\mathbf{x}_3, -1)$   $(\mathbf{x}_4, -1)$  where  $\mathbf{x}_1 = (1, 1)^\top$ ,  $\mathbf{x}_2 = (-1, 1)^\top$ ,  $\mathbf{x}_3 = (-1, -1)^\top$ , and  $\mathbf{x}_4 = (1, -1)^\top$ . Here,  $\mathbf{x}^\top$  denotes the transpose of vector  $\mathbf{x}$ . We would like to train a binary Perceptron to classify the four instances in this dataset. For this question ignore the bias term  $b$  in the Perceptron and answer the following.

- A. Let us predict an instance  $\mathbf{x}$  to be positive if  $\mathbf{w}^\top \mathbf{x} \geq 0$ , and negative otherwise. Initializing  $\mathbf{w} = (0, 0)^\top$ , show that after observing  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ ,  $\mathbf{x}_3$ , and  $\mathbf{x}_4$  in that order the weight vector will be  $-\mathbf{x}_3 - \mathbf{x}_4$ . **(6 marks)**
- B. If we present the four instances in the reverse order  $(\mathbf{x}_4, -1)$ ,  $(\mathbf{x}_3, -1)$ ,  $(\mathbf{x}_2, 1)$ ,  $(\mathbf{x}_1, 1)$ , to the Perceptron, what would be the final value of weight vector at the end of the first iteration? **(4 marks)**
- C. Normalize each of the four instances  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ ,  $\mathbf{x}_3$ , and  $\mathbf{x}_4$  into unit L2 length. **(4 marks)**
- D. What would be the final weight vector after observing the four instances if you used the L2 normalized training instances instead of the original (unnormalized) instances to train the Perceptron as you did in the part (A) of above? **(4 marks)**
- E. Now, let us re-assign the target labels for this dataset as follows  $(\mathbf{x}_1, 1)$ ,  $(\mathbf{x}_2, -1)$ ,  $(\mathbf{x}_3, 1)$   $(\mathbf{x}_4, -1)$ . Can we use Perceptron algorithm to linearly classify this revised dataset? Justify your answer. **(4 marks)**
- F. Describe a method to learn a binary linear classifier for the revised dataset described in part (E) above. **(3 marks)**