



## **RESIT EXAMINATIONS 2017/18**

### **Data Mining and Visualisation**

**TIME ALLOWED : Two and a Half Hours**

---

#### **INSTRUCTIONS TO CANDIDATES**

Answer **FOUR** questions.

If you attempt to answer more questions than the required number of questions, the marks awarded for the excess questions answered will be discarded (starting with your lowest mark).

**Question 1** Consider a dataset  $\mathcal{D}$  of  $N$  instances, where each instance  $x_i \in \mathcal{D}$  is represented by a three dimensional real-valued vector  $\mathbf{x}_i = (x_{i1}, x_{i2}, x_{i3})^T$ . Moreover, a label  $t_i \in \{-1, 1\}$  is assigned to  $\mathbf{x}_i$ . We would like to learn a binary classifier using  $\mathcal{D}$ . However, for some instances, we do not have  $x_{i3}$  values measured. Answer the following questions.

**A.** Explain what is meant by the *missing value problem* in data mining. **(2 marks)**

*If some features are missing (not measured, unobserved) for some data points in a dataset, then this is called the missing value problem.*

**B.** Compute the  $\ell_2$  norm of  $\mathbf{x}_i$ . **(2 marks)**

$$\|\mathbf{x}_i\|_2 = \sqrt{x_{i1}^2 + x_{i2}^2 + x_{i3}^2}$$

**C.** Write the  $\ell_2$  normalised version of  $\mathbf{x}_i$ . **(2 marks)**

$$\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|_2}$$

**D.** Compute the means  $\mu_1, \mu_2, \mu_3$  and standard deviations  $\sigma_1, \sigma_2, \sigma_3$  for the three features in  $\mathcal{D}$ . **(6 marks)**

$$\mu_1 = \frac{1}{N} \sum_{n=1}^N x_{n1}$$

$$\sigma_1 = \sqrt{\frac{\sum_{n=1}^N (x_{n1} - \mu_1)^2}{N - 1}}$$

$$\mu_2 = \frac{1}{N} \sum_{n=1}^N x_{n2}$$

$$\sigma_2 = \sqrt{\frac{\sum_{n=1}^N (x_{n2} - \mu_2)^2}{N - 1}}$$

$$\mu_3 = \frac{1}{N} \sum_{n=1}^N x_{n3}$$

$$\sigma_3 = \sqrt{\frac{\sum_{n=1}^N (x_{n3} - \mu_3)^2}{N - 1}}$$

**E.** Apply Gaussian scaling on  $\mathbf{x}_i$ . **(2 marks)**

$$\left( \frac{x_{i1} - \mu_1}{\sigma_1}, \frac{x_{i2} - \mu_2}{\sigma_2}, \frac{x_{i3} - \mu_3}{\sigma_3} \right)$$

**F.** Given that  $\mu_3 = 0$  would it be problematic to replace missing values of  $x_{i3}$  to zero? Explain your answer. **(2 marks)**

*Yes, this would be problematic because  $x_{i3} \in \mathbb{R}$ , if we replace missing  $x_{i3}$  values by zero we would not be able to distinguish among the instances for which  $x_{i3}$  was measured but turned out to be zero vs. instances where  $x_{i3}$  is missing.*

- G. As a solution to the missing value problem, we would like to predict  $x_{i3}$  using  $x_{i1}$  and  $x_{i2}$  using the linear relationship  $\hat{x}_{i3} = ax_{i1} + bx_{i2} + c$ , where  $a, b, c \in \mathbb{R}$  are parameters that must be estimated from  $\mathcal{D}$  and  $\hat{x}_{i3}$  is the predicted value for  $x_{i3}$ . Write the squared loss for this prediction problem. **(3 marks)**

$$E(\mathcal{D}) = \sum_{i=1}^N (ax_{i1} + bx_{i2} + c - x_{i3})^2$$

- H. Compute the gradient of the squared loss function w.r.t.  $a, b$  and  $c$ . **(3 marks)**

$$\frac{\partial E}{\partial a} = 2 \sum_{i=1}^N (ax_{i1} + bx_{i2} + c)x_{i1}$$

$$\frac{\partial E}{\partial b} = 2 \sum_{i=1}^N (ax_{i1} + bx_{i2} + c)x_{i2}$$

$$\frac{\partial E}{\partial c} = 2 \sum_{i=1}^N (ax_{i1} + bx_{i2} + c)$$

- I. Write the update rules for  $a, b$  and  $c$  using stochastic gradient descent. **(3 marks)**

$$a^{(k+1)} = a^{(k)} - 2\eta \sum_{i=1}^N (ax_{i1} + bx_{i2} + c)x_{i1}$$

$$b^{(k+1)} = b^{(k)} - 2\eta \sum_{i=1}^N (ax_{i1} + bx_{i2} + c)x_{i2}$$

$$c^{(k+1)} = c^{(k)} - 2\eta \sum_{i=1}^N (ax_{i1} + bx_{i2} + c)$$

**Question 2** We would like to use the Perceptron algorithm to learn a linear classifier  $y = \mathbf{w}^\top \mathbf{x} + b$ , defined by a weight vector  $\mathbf{w} \in \mathbb{R}^d$  and a bias  $b \in \mathbb{R}$  from a training dataset consisting of three instances,  $\{(t_n, \mathbf{x}_n)\}_{n=1}^3$ . Here,  $\mathbf{x}_1 = (0, 0)^\top$ ,  $\mathbf{x}_2 = (1, 1)^\top$  and  $\mathbf{x}_3 = (-1, 1)^\top$ , and the labels are  $t_1 = 1$ ,  $t_2 = -1$  and  $t_3 = 1$ . We predict an instance  $\mathbf{x}$  as positive if  $\mathbf{w}^\top \mathbf{x} + b \geq 0$ , and negative otherwise. The initial values of the weight vector and the bias are set respectively to  $\mathbf{w}^{(0)} = (0, 0)^\top$  and  $b = 0$ . We visit the training instances in the order  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$ . Answer the following questions.

- A.** Plot the dataset in the two-dimensional space. **(2 marks)**  
*The three points form a triangle with  $x_1$  at the origin and  $x_2$  and  $x_3$  mirroring each other on the y-axis.*
- B.** Write the perceptron update rule for a misclassified instance  $(t, \mathbf{x})$ . **(3 marks)**  
 $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + t\mathbf{x}$
- C.** What will be the values of the weight vector and the bias after observing the instance  $\mathbf{x}_1$ . **(3 marks)**  
 $y_1 = \mathbf{w}^{(0)\top} \mathbf{x}_1 + b^{(0)} = 0$ . Therefore, this instance is classified correctly as positive. The weight vector and bias are not updated.  $\mathbf{w}^{(1)} = (0, 0)^\top$ ,  $b^{(1)} = 0$ .
- D.** What will be values of the weight vector and the bias after observing  $\mathbf{x}_2$ . **(4 marks)**  
 $y_2 = 0$ . Therefore,  $\mathbf{x}_2$  will be incorrectly classified as positive. The weight vector and the bias will be update to  $\mathbf{w}^{(2)} = (-1, -1)$ ,  $b^{(2)} = -1$
- E.** What will be the values of the weight vector and the bias after observing  $\mathbf{x}_3$ . **(4 marks)**  
 $y_3 = (-1, 1)^\top (-1, -1) - 1 = -1 < 0$ . Therefore,  $\mathbf{x}_3$  will be incorrectly classified as negative. The updated values will be  $\mathbf{w}^{(3)} = (-2, 2)^\top$ ,  $b^{(3)} = 0$ .
- F.** Is the dataset consisting of  $x_1, x_2, x_3$  linearly separable?. Justify your answer. **(2 marks)**  
*Yes. We have already found a Perceptron with a weight vector and a bias that would correctly classify all three instances in this dataset.*
- G.** Is it the case that a dataset consisting of three points is always linearly separable? If yes, explain your answer. If no, provide a counter example. **(4 marks)**  
*No. For example, if the three datapoints are on a straight line and the middle point has the opposite label than the other two points, then this dataset cannot be linearly separable.*
- H.** Explain a method that you can use to learn a Perceptron from a non-linearly separable dataset. **(3 marks)**  
*Apply a kernel method to project the dataset into a high dimensional feature space and learn a Perceptron in this high dimensional feature space.*

**Question 3** Consider the two sentences  $S_1$  and  $S_2$  given by:

$S_1 =$  I love cake with tea

$S_2 =$  I drink beer with cake

Answer the following questions.

- A.** Represent  $S_1$  and  $S_2$  respectively by feature vectors  $\mathbf{s}_1$  and  $\mathbf{s}_2$ , where elements correspond to the frequency of unigrams. **(4 marks)**

*Let the unigram features be indexed as follows: I=0, love=1, cake=2, with=3, tea=4, drink=5, beer=6. Then we have  $\mathbf{s}_1 = (1, 1, 1, 1, 1, 0, 0)^\top$  and  $\mathbf{s}_2 = (1, 0, 1, 1, 0, 1, 1)^\top$ .*

- B.** Compute the  $\ell_2$  norms of  $\mathbf{s}_1$  and  $\mathbf{s}_2$ . **(4 marks)**

$$\|\mathbf{s}_1\|_2 = \sqrt{5}, \|\mathbf{s}_2\|_2 = \sqrt{5}$$

- C.** Compute the  $\ell_1$  norms of  $\mathbf{s}_1$  and  $\mathbf{s}_2$ . **(4 marks)**

$$\|\mathbf{s}_1\|_1 = 5, \|\mathbf{s}_2\|_1 = 5$$

- D.** Compute the cosine similarity between  $\mathbf{s}_1$  and  $\mathbf{s}_2$ . **(2 marks)**

$$\frac{\mathbf{s}_1^\top \mathbf{s}_2}{\|\mathbf{s}_1\|_2 \|\mathbf{s}_2\|_2} = 3/5$$

- E.** Compute the Manhattan distance between  $\mathbf{s}_1$  and  $\mathbf{s}_2$ . **(2 marks)**

$$|1 - 1| + |1 - 0| + |1 - 1| + |1 - 1| + |1 - 0| + |0 - 1| + |0 - 1| = 4$$

- F.** Assume that for all the unigrams  $u_i$  and bigrams  $u_i u_{i+1}$  that appear in  $S_1$  and  $S_2$  we are given the marginal probabilities respectively  $p(u_i)$  and  $p(u_i u_{i+1})$ . Compute the conditional probability of observing  $u_{i+1}$  given  $u_i$ . **(2 marks)**

$$p(u_{i+1}|u_i) = \frac{p(u_{i+1}, u_i)}{p(u_i)}$$

$$p(u_{i+1}, u_i) = p(u_{i+1} u_i) + p(u_i u_{i+1})$$

$$p(u_{i+1}|u_i) = \frac{p(u_{i+1} u_i) + p(u_i u_{i+1})}{p(u_i)}$$

- G.** Using the Markov assumption, compute the likelihood  $p(S_1)$  and  $p(S_2)$ . **(4 marks)**

$$p(S_1) = p(I)p(\text{love}|I)p(\text{cake}|\text{love})p(\text{with}|\text{love})p(\text{tea}|\text{with})$$

$$p(S_2) = p(I)p(\text{drink}|I)p(\text{beer}|\text{drink})p(\text{with}|\text{drink})p(\text{tea}|\text{with})$$

- H.** Explain how you can use the computation done in part (F) to evaluate whether  $S_2$  is less common than  $S_1$  in English texts written by native speakers. **(3 marks)**

*Compare  $p(S_1)$  and  $p(S_2)$ . If the likelihood of a sentence is small, then it is unlikely to be produced by a native speaker.*

**Question 4** Table 1 shows how four users  $u_1, u_2, u_3, u_4$  purchased four items  $l_1, l_2, l_3, l_4$  in an on-line shopping site over a period of one year. A cell value of 1 indicates that the user corresponding to the row has purchased the item corresponding to the column, and 0 otherwise. Answer the following questions.

	$l_1$	$l_2$	$l_3$	$l_4$
$u_1$	1	0	1	1
$u_2$	1	1	0	0
$u_3$	0	0	1	1
$u_4$	0	1	0	0

Table 1: A table showing four users  $u_1, u_2, u_3, u_4$  who have purchased four items  $l_1, l_2, l_3, l_4$  in an online shopping site over a period of one year.

- A.** Given that the users have been initially clustered into two clusters  $S_1 = \{u_1, u_2\}$  and  $S_2 = \{u_3, u_4\}$ , compute the centroids for the two clusters respectively denoted by  $\mu_1$  and  $\mu_2$ . For this purpose, consider a user is represented by a vector over the items he or she has purchased in the past. **(2 marks)**

$$\mu_1 = (1, 0.5, 0.5, 0.5)^T \text{ and } \mu_2 = (0, 0.5, 0.5, 0.5)^T$$

- B.** Compute Euclidean distances between  $\mu_1$  and each of the four users. **(4 marks)**

$$d(u_1, \mu_1) = \sqrt{0.75}$$

$$d(u_2, \mu_1) = \sqrt{0.75}$$

$$d(u_3, \mu_1) = \sqrt{1.75}$$

$$d(u_4, \mu_1) = \sqrt{0.75}$$

- C.** Compute Euclidean distances between  $\mu_2$  and each of the four users. **(4 marks)**

$$d(u_1, \mu_2) = \sqrt{1.75}$$

$$d(u_2, \mu_2) = \sqrt{1.75}$$

$$d(u_3, \mu_2) = \sqrt{0.75}$$

$$d(u_4, \mu_2) = \sqrt{0.75}$$

- D.** Based on the distances computed in parts (B) and (C), determine the assignment of users to clusters for the next iteration. **(2 marks)**

*Two possible assignments exist.  $S_1 = \{u_1, u_2, u_4\}, S_2 = \{u_3\}$  or  $S_1 = \{u_1, u_2\}, S_2 = \{u_3, u_4\}$*

- E.** Let us denote the probability of a user purchasing an item  $l_j$  when he or she has purchased  $l_i$  by  $p(l_j|l_i)$ . From Table 1, compute  $p(l_1|l_4)$ ,  $p(l_2|l_4)$  and  $p(l_3|l_4)$ . **(3 marks)**

$$p(l_1|l_4) = 0.5, p(l_2|l_4) = 0 \text{ and } p(l_3|l_4) = 1$$

- F. Based on your calculations in part (E), explain what is the best item to recommend to a user who has just purchased  $l_4$ . **(2 marks)**

*$l_3$  because  $p(l_3|l_4) = 1$  and the user is likely to buy  $l_3$  too, given that he/she has already purchased  $l_4$*

- G. Represent the information shown in Table 1 by a bi-partite graph where the users and items are represented by vertices, and an undirected edge is formed between the vertices corresponding to  $u_i$  and  $l_j$  if and only if  $u_i$  has purchased  $l_j$ . **(4 marks)**

- H. Consider a random walker moving along the edges of the graph you created in part (G), where the probability of moving from  $u_i$  to  $l_j$  is given by  $\frac{1}{d(u_i)}$  and the probability of moving from  $l_j$  to  $u_i$  is given by  $\frac{1}{d(l_j)}$ . Here,  $d(x)$  is the out-degree of the vertex  $x$ . Given that the random walker started from  $u_1$ , compute the probability that the random walker will be in  $u_3$  after two time steps. **(4 marks)**

$$p(u_1 \rightarrow l_3)p(l_3 \rightarrow u_3) + p(u_1 \rightarrow l_4)p(l_4 \rightarrow u_3) = 1/3$$

**Question 5** Consider the three points  $x_1 = (0, 1)$ ,  $x_2 = (-1, 0)$  and  $x_3 = (1, 0)$ . We would like to project these three points onto a straight line using principle component analysis. Answer the following questions.

- A.** Compute the total projection error if we project the three points onto the y-axis. **(3 marks)**

$$1+1 = 2$$

- B.** Compute the total projection error if we project the three points onto the x-axis. **(3 marks)**

$$1$$

- C.** Compute the mean  $\bar{x}$  of the three points. **(2 marks)**

$$(0, 1/3)$$

- D.** Compute the covariance matrix for the three points. **(3 marks)**

$$\text{Compute } x_i - \bar{x} \text{ and adding the outer product matrices gives } [[2, 0], [0, 2/3]].$$

- E.** Compute the eigenvalues of the covariance computed in part (D). **(4 marks)**

$$\text{Because the covariance matrix is diagonal we have } \lambda_1 = 2, \lambda_2 = 2/3$$

- F.** Compute the first principle component of the projection. **(3 marks)**

*The eigenvector corresponding to  $\lambda_1$  (the larger eigenvalue) is  $(1, 0)$ . This is the x-axis.*

- G.** Compute the second principle component of the projection. **(3 marks)**

*The eigenvector corresponding to  $\lambda_2$  (the smaller eigenvalue) is  $(0, 1)$ . This is the y-axis.*

- H.** Compute the total variance if we had projected the three points on to the first principle component. **(2 marks)**

$$\frac{(1-0)^2 + (-1-0)^2 + (0-0)^2}{3} = 2/3$$

- I.** Compute the total variance if we had projected the three points on to the second principle component. **(2 marks)**

$$\frac{(1-0.5)^2 + (0-0.5)^2 + (0-0.5)^2}{3} = 0.25$$