



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 $\mathbf{x}_i = (x_{i1}, x_{i2}, x_{i3})^\top$. 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)**
- B. Compute the ℓ_2 norm of \mathbf{x}_i . **(2 marks)**
- C. Write the ℓ_2 normalised version of \mathbf{x}_i . **(2 marks)**
- D. Compute the means μ_1, μ_2, μ_3 and standard deviations $\sigma_1, \sigma_2, \sigma_3$ for the three features in \mathcal{D} . **(6 marks)**
- E. Write the result of the Gaussian scaling for \mathbf{x}_i . **(2 marks)**
- F. Given that $\mu_3 = 0$ would it be problematic to replace missing values of x_{i3} by zero? Explain your answer. **(2 marks)**
- G. As a solution to the missing value problem, we would like to predict x_{i3} using x_{i1} and x_{i2} assuming 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)**
- H. Compute the gradient of the squared loss function w.r.t. a, b and c . **(3 marks)**
- I. Write the update rules for a, b and c using stochastic gradient descent. **(3 marks)**

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)**
- B. Write the perceptron update rule for a misclassified instance (t, \mathbf{x}) . **(3 marks)**
- C. What will be the values of the weight vector and the bias after observing the instance \mathbf{x}_1 ? **(3 marks)**
- D. What will be values of the weight vector and the bias after observing \mathbf{x}_2 ? **(4 marks)**
- E. What will be the values of the weight vector and the bias after observing \mathbf{x}_3 ? **(4 marks)**
- F. Is the dataset consisting of $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$ linearly separable?. Justify your answer. **(2 marks)**
- 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)**
- H. Explain a method that you can use to learn a Perceptron from a non-linearly separable dataset. **(3 marks)**

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)**
- B.** Compute the ℓ_2 norms of \mathbf{s}_1 and \mathbf{s}_2 . **(4 marks)**
- C.** Compute the ℓ_1 norms of \mathbf{s}_1 and \mathbf{s}_2 . **(4 marks)**
- D.** Compute the cosine similarity between \mathbf{s}_1 and \mathbf{s}_2 . **(2 marks)**
- E.** Compute the Manhattan distance between \mathbf{s}_1 and \mathbf{s}_2 . **(2 marks)**
- 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 $p(u_i u_{i+1})$. Express the conditional probability of observing u_{i+1} given u_i in terms of $p(u_{i+1} u_i)$, $p(u_i u_{i+1})$ and $p(u_i)$. **(2 marks)**
- G.** Using the Markov assumption, compute the likelihood $p(S_1)$ and $p(S_2)$. **(4 marks)**
- 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)**

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 μ_1 and μ_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)**
- B. Compute Euclidean distances between μ_1 and each of the four users. **(4 marks)**
- C. Compute Euclidean distances between μ_2 and each of the four users. **(4 marks)**
- D. Based on the distances computed in parts (B) and (C), determine the assignment of users to clusters for the next iteration. **(2 marks)**
- 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)**
- 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)**
- 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 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)**

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)**
- B.** Compute the total projection error if we project the three points onto the x -axis. **(3 marks)**
- C.** Compute the mean \bar{x} of the three points. **(2 marks)**
- D.** Compute the covariance matrix for the three points. **(3 marks)**
- E.** Compute the eigenvalues of the covariance computed in part (D). **(4 marks)**
- F.** Compute the first principle component of the projection. **(3 marks)**
- G.** Compute the second principle component of the projection. **(3 marks)**
- H.** Compute the total variance if we had projected the three points on to the first principle component. **(2 marks)**
- I.** Compute the total variance if we had projected the three points on to the second principle component. **(2 marks)**