

# DATA MINING

## TEST 2

### INSTRUCTIONS:

- this test consists of 6 questions
- you may attempt 5 questions.
- maximum marks = 100
- strike out the question you are NOT attempting

**Question 1 (20 marks)**

Learning algorithms can be classified as *supervised* or *unsupervised*.

- a) Describe the difference between *supervised* and *unsupervised* learning. [5]
- b) What is an R dataframe and what is its role in supervised learning. [5]
- c) Discuss the use of training and testing subsets for supervised learning. [5]
- d) Explain why the `iris` dataset is popular for testing supervised learning algorithms. [5]

## Question 1 ...

### Marking schedule

- a) In both cases we try to predict the value of one attribute in a dataset observation using information from the other attributes in the observation. In supervised learning, the attribute we are trying to predict is used to train the learning algorithm whilst in unsupervised learning the target attribute is hidden from us.
- b) A data frame is an enhanced 2d array, (a list of equal length column vectors), with each column representing some attribute and each row representing a measurement of a collection of attributes. both rows and columns can be indexed by names or numbers. various sub-setting routines are available.
- c) In supervised learning, to prevent over-learning, the dataset is split into training and testing subsets. (split proportions vary) the algorithm is trained on one subset and evaluated on the other.
- d) The iris dataset consists of 5 attributes with the target attribute being *Species*. The task is to predict species from the remaining 4 attributes. There are 3 classes of species, one of which is linearly separable from the other two which are not linearly separable. Thus the dataset is useful for testing linear and non-linear classifiers.

**Question 2 (20 marks)**

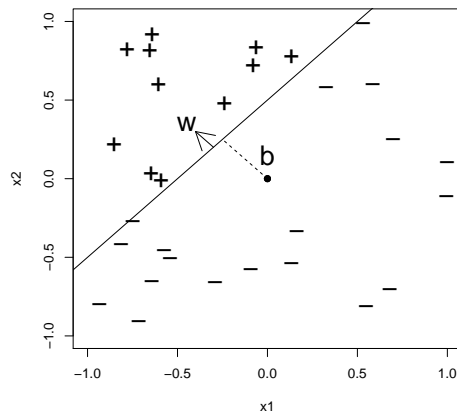
Consider a training set of  $n$  observations,  $\{\vec{x}_i, y_i\}_{i=1\dots n}$ , where the  $\vec{x}_i$  are 2 dimensional feature vectors and  $y_i \in \{-1, +1\}$  is the target attribute that places each observation in one of two possible classes.

In linear classification, we seek to divide the two classes by a linear separator in the feature space.

- a) Explain with the aid of a diagram, the role played by the parameters  $b$  and  $\vec{w}$  for a linear separator. [5]
- b) Prove that the decision boundary generated by  $(\vec{w}, b)$  is identical to the decision boundary generated by  $(s\vec{w}, sb)$  for any scaling parameter  $s$ . [5]
- c) Dr. Magwasa says that the perceptron algorithm is too complicated and instead we should just use the linear classifier  $\vec{w} = \vec{u} - \vec{v}$  where  $\vec{u}$  is the average of the positive examples and  $\vec{v}$  is the average of the negative examples.  
Draw a simple sketch of a small linearly separable dataset that will cause Dr. Magwasa's suggestion to fail. [5]
- d) Write down (but do NOT attempt to prove) the theorem studied in class concerning the convergence of the perceptron algorithm. Make sure you state all conditions under which the theorem holds. [5]

## Question 2 ...

### Marking schedule



a)

In the diagram feature vectors are two-dimensional. We wish to find a line that separates the positive observations from the negative ones. Any line can be specified by a unit normal vector,  $\vec{w}$ , and a signed perpendicular distance from the origin,  $b$ . Classification of a point  $\vec{x}$  is then accomplished according to the value of  $\text{sign}(\vec{x} \cdot \vec{w} - b)$ . In the diagram above, a particular  $\vec{w}$  and  $b$  is given that accomplishes the separation. Note that in the above diagram  $b$  is negative with respect to the direction of  $\vec{w}$ .

b) Classification proceeds as follows:

if  $x \cdot w + b > 0$ , then  $x$  is positive, otherwise  $x$  is negative

now if  $s > 0$ , then

$$x \cdot (sw) + (sb) = s(x \cdot w + b)$$

which is  $> 0$  if  $x \cdot w + b > 0$ .

therefore scaling by  $s$  of  $w$  and  $b$  does not change classification

c) Consider the positive observations at  $(1, 4)$ ,  $(1, -2)$  with mean at  $(1, 1)$  and negative observations at  $(-1, -4)$ ,  $(-1, 2)$  with mean at  $(-1, -1)$ . According to Dr. Magwasa, the normal to the separating line is now  $(1, 1) - (-1, -1) = (2, 2)$  and the best offset is zero which will not separate  $(1, -2)$  from  $(-1, 2)$ .

d) Convergence theorem, due to Novikoff in 1962.

**Assumptions:** Let  $R = \max \|\vec{x}_i\|$  and suppose that the learning task is solvable via a separator that passes through the origin. i.e. there exists some vector  $\vec{w}^*$  of unit length and some  $\delta > 0$  such that  $Y_i(\vec{w}^* \cdot \vec{x}_i) > \delta$  for all  $i$ .

**Theorem:** Under these assumptions, the perceptron algorithm converges after at most  $(\frac{R}{\delta})^2$  updates.

**Question 3 (20 marks)**

- a) Define the *scatter matrix*,  $S$ , for the collection of points  $\{\vec{x}_i\}_{i=1\dots n} \in R^2$  and write down the dimensions of  $S$ . [5]
- b) Explain what is computed by the *action* of the scatter matrix,  $S$ , on a *unit vector*,  $\vec{w}$ . [5]
- c) R.A. Fisher defined a linear discriminant as the vector  $\vec{w}$  that maximizes the ratio:

$$J(\vec{w}) = \frac{\vec{w}^T S_B \vec{w}}{\vec{w}^T S_W \vec{w}}$$

explain in English the roles of the matrices  $S_B$  and  $S_W$  in the above expression and why it is a good idea to maximize this ratio. [5]

- d) What is *Mahalanobis* separation and what role does it play in the construction of a linear discriminant? [5]

### Question 3 ...

#### Marking schedule

- a) In this case the scatter matrix is a  $2 \times 2$  matrix defined by  $S = \sum_i (\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T$ . Note that the  $\vec{x}_i$  are column vectors with mean  $\vec{\mu}$ .
- b)  $\vec{w}^T S \vec{w}$  is the action of  $S$  on  $\vec{w}$  and this computes the the *variance* of the projections of the  $\vec{x}_i - \vec{\mu}$  onto  $\vec{w}$ .
- c)  $S_B$  is the between-class scatter (the scatter of the class means)  
 $S_W$  is the within-class scatter (computed as the sum of the scatter matrices for each class)
- d) Once the optimal separation direction  $\vec{w}$  is found, we must project the feature vectors onto  $\vec{w}$  and then using the *Mahalanobis* bias  $b$  to separate the two classes:

$$b = \frac{\mu_+ \sigma_- + \mu_- \sigma_+}{\sigma_+ + \sigma_-}$$

where the various  $\mu$  and  $\sigma$  are the means and standard deviations of the two classes.

Assuming that the distribution of the projection onto  $\vec{w}$  is double humped in the form of two overlapping normal distributions, this Mahalanobis bias will yield the point of overlap which provides the best offset for the separation plane.

**Question 4 (20 marks)**

The entropy of a sample  $D$  with respect to a target variable of  $k$  possible classes is defined as:

$$H(D) = - \sum_{i=1}^k P(C_i|D) \log_k(P(C_i|D))$$

Consider the following dataset,  $D$ , of training examples:

$X_1$	$X_2$	$C$
T	T	+
T	T	+
T	F	-
F	F	-
F	T	-
F	F	-

- a) Compute  $P(C_+|D)$  [3]
- b) What is the entropy of this dataset with respect to the target attribute  $C$ ? [3]
- c) Write down a general expression for *information gain* in a decision split and then use it to compute the information gain of an  $X_2$  decision split on  $D$ . [14]



**Question 4 ...****Marking schedule**

- a)  $\frac{1}{3}$   
 b)  $H(D) = \frac{1}{3} \log 3 + \frac{2}{3} \log \frac{3}{2}$   
 c) Define the weighted entropy of a decision split as:

$$H(D_L|D_R) = \frac{|D_L|}{|D|} H(D_L) + \frac{|D_R|}{|D|} H(D_R)$$

and then the **information gain** for the split is given by:

$$IG(D, D_L, D_R) = H(D) - H(D_L|D_R)$$

A split on  $X_2$  will have weighted entropy:

$$\frac{1}{2} \left( \frac{2}{3} \log \frac{3}{2} + \frac{1}{3} \log 3 \right) + \frac{1}{2} (0)$$

and so information gain in an  $X_2$  split is:

$$\frac{1}{6} \log 3 + \frac{1}{3} \log \frac{3}{2}$$

**Question 5 (20 marks)**

- a) Contrast k-means clustering with hierarchical clustering outlining similarities and differences. [6]
- b) Construct an example dataset that exemplifies the following quote from wikipedia:  
 k-medoid clustering is more robust to outliers as compared to k-means clustering because it minimizes a sum of pairwise dissimilarities instead of a sum of squared Euclidean distances. [6]
- c) Perform k-medoid ( $k = 2$ ) clustering on the following distance matrix:

D	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$
$X_1$	0	1.91	2.23	3.14	4.25	3.37
$X_2$	1.91	0	2.15	1.82	2.41	2.58
$X_3$	2.23	2.15	0	3.12	3.83	4.64
$X_4$	3.14	1.82	3.12	0	1.9	2.66
$X_5$	4.25	2.41	3.83	1.9	0	3.12
$X_6$	3.37	2.58	4.64	2.66	3.12	0

[8]

## Question 5 ...

### Marking Schedule

- a) Both k-means and hierarchical clustering are unsupervised learning techniques.  
 k-means is a *top down* clustering algorithm where k clusters are assigned at random and then the algorithm iterates between an assignment phase where observations are assigned to nearest cluster mean and a mean re-computation phase. The algorithm suffers from having to know how many cluster to look for from the start although there are methods to estimate optimum number of clusters  
 hierarchical clustering is a *bottom up* approach in which all observations start in their own cluster and then nearest clusters are merged. Distances between the updated set of clusters are computed and the merging process continues. The hierarchy is complete when only one cluster remains.
- b) Suppose you want to cluster on one dimension with  $k = 2$ . One cluster has most of its members around 1000 and the other around -1000; but there is an outlier (or noise) at 100000. It obviously belongs to the cluster around 1000 but k-means will put the center point away from 1000 and towards 100000. This may even make some of the members of the 1000 cluster (say a member with value 500) to be assigned to the -1000 cluster. k-medoid will select one of the members around 1000 as the medoid, it'll probably select one that is bigger than 1000, but it will not select an outlier.
- c) Choose randomly k entities to be the medoids  $m_1, m_2$ .  
 Lets choose  $X_3$  (Lets call this cluster 1) and  $X_5$  (Cluster 2).  
 Assign a given entity to the cluster represented by its closest medoid.  
 Cluster 1 will be made of entities  $(X_1, X_2, X_3)$  - just check your table, these are closer to  $X_3$  than to  $X_5$ , cluster 2 will be  $(X_4, X_5, X_6)$ .  
 Update the medoids. A medoid of a cluster should be the entity with the smallest sum of distances to all other entities within the same cluster.  $X_2$  will be the new medoid for cluster 1, and  $X_4$  for cluster 2. Now what you have to do repeat until convergence.  
 So,  
 Assign each entity to the cluster of the closest medoid, now these are  $X_2$  and  $X_4$ . Cluster one is now made of entities  $(X_1, X_2, X_3, X_6)$ , Cluster 2 will be  $(X_4, X_5)$ . (there was a change in the entities in each cluster, so iterations must continue.  
 The entity with the smallest sum of distances in cluster one is still  $X_2$ , in cluster 2 they are the same, so  $X_4$  stays.  
 Another iteration  
 As there was no change in the medoids, the clusters will stay the same. This means its time to stop the iterations  
 Output: cluster 1 has entities  $(X_1, X_2, X_3, X_6)$ , and cluster 2 has entities  $(X_4, X_5)$ .

**Question 6 (20 marks)**

The R package, `tm`, provides a suite of routines for undertaking text mining tasks.

- a) In the `tm` package text for data mining is commonly stored in a *corpus*. Describe the *corpus* data structure and explain how a corpus is indexed and how it is stored in a directory. [5]
- b) Outline 5 procedures that are commonly used for cleaning a corpus. [5]
- c) What is a *term-document matrix* and what is it used for? [5]
- d) Outline one real life example, where you might want to make use of clustering in a text mining operation. [5]

## Question 6 ...

### Marking schedule

- a) A corpus is a collection of documents. If `myCorpus` is a corpus containing `n` documents then each document can be accessed via integer indexing, `myCorpus[[i]]`, or named indexing `myCorpus[["ovid.2.txt"]]`. If the corpus is written to disc then all the documents are stored in a directory and each document name is used as that documents filename.
- b) Via the `tm_map` function, the `tm` packages provides common text cleaning operations:
  - 1) `tolower`, transform all characters to lower case.
  - 2) `stripWhitespace`, remove white space
  - 3) `removePunctuation`, get rid of punctuation characters
  - 4) `removeWords`, `myStopwords`, get rid of common words that may distort later analysis
  - 5) `gsub`, `pattern="search_string"`, `replacement="repacement_string"`, homegrown cleaning
- c) A term-document matrix represents the relationship between terms (or words) and documents, where each row stands for a term and each column for a document, and an entry is the number of occurrences of the term in the document. Correlations, associations, clustering and classification is usually carried out on the term document matrix and not the original corpus.
- d) Consider each work of Shakespeare as a corpus document and each work of Marlow as a corpus document. Combine the two corpora into one corpus and cluster the documents. Hopefully two clusters appear and any disputed work can be attributed to the author of the nearest cluster.