

CMPT 459 Fall 2020
Data Mining
Martin Ester
TAs: Madana Krishnan Vadakandara Krishnan
and Rhea Rodrigues

Assignment 2

Total Marks: 100

Assignment 2.1 (55 marks)

We want to cluster categorical data, i.e. data that have categorical attribute domains. The task of this assignment is to develop the equivalent of the k -means algorithm for categorical data, called the k -modes algorithm. We assume the following distance function (Hamming distance) for pairs of categorical objects:

$$dist(x, y) = \sum_{i=1}^d \delta(x_i, y_i) \text{ with } \delta(x_i, y_i) = \begin{cases} 0 & \text{if } x_i = y_i \\ 1 & \text{else} \end{cases}$$

Instead of the means, we choose the mode m as the representative of a cluster C . In every attribute, the mode takes the attribute value that is most frequent within the cluster. More formally, the mode m of cluster C is defined as follows:

$$m = (m_1, \dots, m_d) \text{ such that } \forall j, 1 \leq j \leq d, \forall a \in A_j : |\{x \in C \mid x_j = m_j\}| \geq |\{x \in C \mid x_j = a\}|$$

a) (30 marks)

Prove that m is the object minimizing the cluster cost

$$TD(C, m) = \sum_{p \in C} dist(p, m)$$

Hint: first formulate the intuition of the proof, then formalize it. The proof can be performed by contradiction.

Assignment 2.1

- a) A mode of $m = \{m_1, m_2, m_3, \dots, m_n\}$ is a vector/cluster $C = \{p_1, p_2, p_3, \dots, p_n\}$ that minimises

$$TD(C, m) = \sum_{p \in C} \text{dist}(p, m)$$

$$\Rightarrow TD(C, m) = \sum_{\substack{i=1, \\ p \in C}}^n \text{dist}(p, m)$$

here C is not necessarily an element of m .

Let $n_{c_{k,j}}$ be the number of objects having the k^{th} category $c_{k,j}$ in attribute A_j and $\text{fr}(A_j = c_{k,j} | m) = \frac{n_{c_{k,j}}}{n}$ the relative frequency of category $c_{k,j}$ in m , and n is the total number of objects in m . (Here, A_j refers to each attribute A_j describes a domain of values associated with a defined semantic & data type)

For the dissimilarity measure:

$$\text{dist}(x, y) = \sum_{i=1}^n \delta(x_i, y_i) \text{ with } \delta(x_i, y_i) = \begin{cases} 0 & \text{if } x_i = y_i \\ 1 & \text{else} \end{cases}$$

$$\Rightarrow \sum_{\substack{i=1, \\ p \in C}}^n \text{dist}(p, m_i) = \sum_{i=1}^n \delta(n_i, p_i)$$

$$\begin{aligned} \Rightarrow \sum_{\substack{i=1 \\ p \in C}}^n \text{dist}(p, m_i) &= \sum_{i=1}^n \sum_{j=1}^l \delta(n_{i,j}, p_j) \\ &= \sum_{j=1}^l \left(\sum_{i=1}^n \delta(n_{i,j}, p_j) \right) \end{aligned}$$

$$= \sum_{i=1}^d n \left(1 - \frac{n_{p_j}}{n} \right)$$

$$= \sum_{i=1}^d n \left(1 - \text{fr}(A_j = p_j | m) \right)$$

because $n(1 - \text{fr}(A_j = p_j | m)) \geq 0$ for $1 \leq j \leq n$,
 $\sum_{\substack{i=1 \\ p \in C}}^n \text{dist}(p, m_i)$ is minimised iff every $n(1 - \text{fr}(A_j = p_j | m))$

is minimal. Thus $\text{fr}(A_j = p_j | m)$ must be maximal.

This defines a way to find C from a given m as it allows k-means paradigm to be used to cluster categorical data. This implies that mode of a dataset X is not unique.

b) (25 marks)

Provide a pseudo-code for the k -modes algorithm with the following header:

$k\text{-modes}(\text{dataset}, \text{numberOfClusters})$

To initialize the k cluster representatives, take a random sample of k objects from the dataset.

```

1  K-modes clustering - Categorical data
2
3  Code outline:
4  Algorithm: k-modes(dataset, numOfClusters)
5  Input: Data, k
6  Output: Cj where 1 <= j <= k
7
8  Step 1: Randomly select the K initial cluster centers such that Cj, j = 1,2,...,K
9  Step 2: Find the matching dissimilarity between the each K initial cluster modes and each data objects using the minizing cost function TD(C,m)
10 Step 3: Evaluate the fitness using the dissimilarity distance function (Hamming distance)
11 Step 4: Find the minimum mode values in each data object i.e. finding the objects nearest to the initial cluster modes.
12 Step 5: Assign the data objects to the nearest cluster modes.
13 Step 6: Update the modes by applying the frequency based method on newly formed clusters
14 Step 7: Recalculate the similarity between the objects and the updated modes
15 Step 8: Repeat step 4 - 5 until no changes in the cluster ship of the objects

```

```

17 Pseudocode:
18
19 Key:
20 k = numOfClusters
21 X = data
22
23 Algorithm: get_distance(x, c)
24 Input: numOfClusters, sizeOfDataset
25 Output: distance
26
27 dist <- 0
28 IF x is not equal to c do
29   dist += 1
30
31
32 Algorithm: kmodes(dataset, numOfClusters)
33 Input: X, k
34 Output: Cj where 1 ≤ j ≤ k
35
36 C1 <- empty array
37 Randomly select k cluster centres from X with n objects
38 C <- random k clusters
39 FOR i from 0 to range of n do
40   minDist <- initialize with some large value
41   FOR j from 0 to k do
42     Calculate the distance between ith data point and jth mode vector using dissimilarity distance function and assign that data point to appropriate cluster whose cluster mode vector is closer
43     to it and update mode vector of corresponding cluster and also find the
44     distribution of mode categories between clusters
45
46     dist <- get_distance(C[row index, features], X[row index, features])
47     IF dist < minDist do
48       minDist <- dist
49       clustNum <- j
50       X[i,d] <- clustNum
51       C[j,d] <- clustNum
52     ENDIF
53   ENDFOR
54   result <- indexes where X[:d] == j
55   mode_info <- array of modal values
56   C[j] <- reshape modal array
57 ENDFOR
58
59 WHILE c1 is not equal to c do
60   c1 <- c
61   FOR i from 0 to range of n do
62     minDist <- initialize with some large value
63     FOR j from 0 to k do
64       Calculate the distance between ith data point and jth mode vector using dissimilarity distance function (6) and assign that data point to appropriate cluster whose cluster mode vector
65       is closer to it and update mode vector of corresponding cluster and also find the distribution of mode categories between clusters
66
67       dist <- get_distance(C[row index, features], X[row index, features])
68       IF dist < minDist do
69         minDist <- dist
70         clustNum <- j
71         X[i,d] <- clustNum
72         C[j,d] <- clustNum
73       ENDIF
74     ENDFOR
75     result <- indexes where X[:d] == j
76     mode_info <- array of modal values
77     C[j] <- reshape modal array
78   ENDFOR
79 ENDFOR
80 IF old_modes is equal to new_modes
81   then break the loop
82 ENDIF
83

```

Assignment 2.2 (45 marks)

Implement the K-modes algorithm in python according to your pseudo-code. ____

Your program will accept two inputs: the dataset and numberOfClusters.

As output, write a new csv file in which every row consists of the example and the corresponding cluster number which was assigned to it by your algorithm.

You will be using the Mushroom Data set which you can download [her e](#). Read the [Data Set Description](#) found on the website to understand the data. Note that the first column is the class label, followed by 22 columns of categorical attributes.

You must handle missing values in the dataset. Mention how you handled the missing values in your PDF Report (see submission details below).

Ans 2.2:

To handle missing data:

1. First we check for null or NaN values in the data. On the first glance it doesn't show anything, so after reading the dataset description we find that only column 11 has missing values in form of '?'
2. To start imputation, we first replace all '?' with NaN attributes so we can identify them easily.
3. Create a new column and replace 1 if the category is NaN else 0. This column is an importance column to the imputed category as with this we can keep track of what values were imputed.
4. We replace NaN attributes with mode of the column as this is the easiest way to impute the data

With this method we capture the importance of missingness in the data. Although if the missing data is too large, there might be a potential to misunderstand the data. This method comes with curse of dimensionality; if 10 columns have missing data, we'll have to create 10 additional columns to capture them hence if the number with missing data increase, then the number of columns increases to capture the missing data.

Follow these instructions for your implementation:

- Prepare the data
 - Handling the missing values (5 marks)
 - Brief description of the above in the report (5 marks)
- Initialization of centroids (5 marks)
- Assignment of observations to a cluster (5 marks)
- Compute new centroids (10 marks)
- Terminating criteria for iterations (10 marks)
- Saving results to a file (5 marks)

You can use libraries including math, numpy, scipy, random, etc. Do not use the scikit-learn libraries. Implement your own logic to compute the distance between observations. You will be marked on the correctness of your implementation.

Submission details

1. Submit a PDF Report for the solutions to Assignment 2.1. In your report, for Assignment 2.2, briefly explain how you decided to handle the missing data and why you handled it in this way.
2. For Assignment 2.2 submit a python program that reads the data set as input (assume the data is in the same directory as the program).

Pseudocode if you are not able to read from the photo!

Pseudocode:

key:

k = numOfClusters

X = data

Algorithm: get_distance(x, c)

Input: numOfClusters, sizeOfDataset

Output: distance

dist ← 0

IF x is not equal to c do

dist += 1

Algorithm: kmodes(dataset, numOfClusters)

Input: X, k

Output: C_j where 1 ≤ j ≤ k

C1 ← empty array

Randomly select k cluster centres from X with n objects

C ← random k clusters

FOR i from 0 to range of n do

minDist ← initialize with some large value

FOR j from 0 to k do

Calculate the distance between ith data point and jth mode vector using dissimilarity distance function and assign that data point to appropriate cluster whose cluster mode vector is closer to it and update mode vector of corresponding cluster and also find the distribution of mode categories between clusters

dist ← get_distance(C[row index, features], X[row index, features])

IF dist < minDist do

minDist ← dist

clustNum ← j

X[i,d] ← clustNum

C[j,d] ← clustNum

ENDIF

ENDFOR

FOR j from 0 to k do

```

        result <- indexes where X[:d] == j
        mode_info <- array of modal values
        C[j] <- reshape modal array
    ENDFOR
ENDFOR

WHILE c1 is not equal to c do
c1 <- c
FOR i from 0 to range of n do
minDist <- initialize with some large value
    FOR j from 0 to k do
        Calculate the distance between ith data point and jth mode vector using dissimilarity
        distance function (6) and assign that data point to appropriate cluster whose cluster mode vector is
        closer to it and update mode vector of corresponding cluster and also find the distribution of mode
        categories between clusters

        dist <- get_distance(C[row index, features], X[row index, features])
    IF dist < minDist do
        minDist <- dist
        clustNum <- j
        X[i,d] <- clustNum
        C[j,d] <- clustNum
    ENDIF
    ENDFOR
    FOR j from 0 to k do
        result <- indexes where X[:d] == j
        mode_info <- array of modal values
        C[j] <- reshape modal array
    ENDFOR
ENDFOR
IF old_modes is equal to new_modes
    then break the loop
ENDIF
ENDWHILE

```