



Data Representation, Reduction and Analysis

Juan Jose Soriano Escobar

Redona Brahimetaj

Master in Applied Computer Science and Engineering
Vrije Universiteit Brussels
Brussels Belgium
January 15, 2017

List of Figures

A Introduction

B Cleaning of Tweets

One of the most crucial parts to start with and that has a big influence in the result of our project, is data processing step. A csv file containing 2000 tweets was provided. It was required to preprocess the data to make them good enough for the 'learning' step. Below we will provide more details regarding the way how we cleaned the data.

B.1 Pre-Processing

We cleaned the data by applying several filters to them. We would like to mention that for this part, we have adapted the code that we already did before on Distributed Computing and Storage Architecture project. The first filtering process we did was tokenization. We chopped the data provided into small pieces so that after it would be easier to remove the stop-words like determiners, the coordinating conjunctions and prepositions. Another pre-processing step that we did was trying to keep only the root of the words by applying stemming. In this way the clusters would be correctly implemented since it would be easier to cluster and find similarities between words that are the same. This is a very important part and that has a high influence in the result. Something that we would like to emphasize is the fact that we did not remove the hashtags. We consider that they are a very helpful indicator when it comes to the properly group the similar words together. Below is presented the code that shows the data-cleaning.

Listing 1: Python Cleaning Function

```
1 from nltk.stem import WordNetLemmatizer
2 from nltk.tokenize import TweetTokenizer
3 from nltk.corpus import stopwords
4 tknz = TweetTokenizer() #Tokenization
5 tweets = CSVHelper.load_csv("Tweets_2016London.csv") #Load the provided CSV file
6 clean_tweets = []
7 #specify the format of the URL so we can remove it later
8 url_expression = 'http[s]?[:]?/[?:(?:[a-zA-Z][0-9]|[$-@.&+][!*\(\)\.,]|(?:%[0-9a-fA-F][0-9a-fA-F]))'
9 #remove emojis
10 emoji_pattern = re.compile("[
11     u\"\\U0001F600-\\U0001F64F\" # emoticons
12     u\"\\U0001F300-\\U0001F5FF\" # symbols & pictographs
13     u\"\\U0001F680-\\U0001F6FF\" # transport & map symbols
14     u\"\\U0001F1E0-\\U0001F1FF\"
15     u\"\\U00002600-\\U000027BF\"
16     u\"\\U0001f300-\\U0001f64F\"
17     u\"\\U0001f680-\\U0001f6FF\"
18     u\"\\u2600-\\u27BF\""]+", flags=re.UNICODE)
19 wordnet_lemmatizer = WordNetLemmatizer()
20 for t in tweets:
21     stop = set(stopwords.words('english')) #stop words!
22     fragments = tknz.tokenize(t)
23     clean_fragments = []
24     for f in fragments:
25         if f not in stop: # not included in the stop words
26             f = emoji_pattern.sub(r'', f)
27             f = f.lower() #lowercase fragment
28             f = re.sub(r'[.,"!~_~:|?\\']+', ' ', f, flags=re.MULTILINE) # Special characters
29             f = re.sub(r'\\.\\.\\. ', ' ', f, flags=re.MULTILINE) # 3 dots
30             f = re.sub(url_expression, ' ', f, flags=re.MULTILINE) # links
31             f = re.sub(r'@[a-zA-Z,0-9 ]*', ' ', f, flags=re.MULTILINE) #clean at person references
32             f = re.sub(r'RT @[a-zA-Z]*: ', ' ', f, flags=re.MULTILINE) #Remove retweets
33             f = wordnet_lemmatizer.lemmatize(f)
34             if f:
35                 clean_fragments.append(f) #we append the cleaned tweet to clean_fragment vector
36     clean_tweets.append(" ".join(clean_fragments)) #append the result to clean_tweets
37
38 # Export the cleaned tweets to a new csv file.Pandas Library is used
39 import pandas as pd
40 df = pd.DataFrame(clean_array)
```

```
41 df.to_csv("file_path.csv")
```

```
42
```

C Noise Removal

The data that we have in disposition, are tweets that are retrieved from people that just wrote their opinions or feelings without thinking to much on what they were writing. This means that in this dataset, the tweets contain noise and we will have to remove the noise from tweets. In order to do so, 2 implementations of different algorithms are requested, K-means and DBSCAN. To implement the first algorithm, we will realize the need of having a consensus-matrix so that for each time that we run the same algorithm, we will count the number of times that two data points were clustered together. In this way, at the final table we would realise that if two data points are clustered together many times, it means that they are related a lot to each other and as a result they should be on the same cluster.

Before running the K-means algorithm, we need to represent the text documents as mutually comparable vectors as it was requested. To do so, we used the TF-IDF which ranks the importance of a string in its contextual text corpus. We carefully studied and understood how does TF-IDF works and initially we implemented it by ourselves and tested it in a simple example. We took 3 documents where we tested it. But at the end we decided to use the one that was already implemented by scikit-learn even if ours was correct as well. Below we will present the code that shows the implementation we did for it.

Listing 2: TF-IDF implementation

```
1 def tf(word, blob): # function to count the term frequency.
2     return blob.words.count(word) / len(blob.words)
3
4 def n_containing(word, bloblist): # count the word ocurance in the document list.
5     return sum(1 for blob in bloblist if word in blob)
6
7 def idf(word, bloblist): #inverse document frequency.
8     return math.log(len(bloblist) / (1 + n_containing(word, bloblist)))
9
10 def tfidf(word, blob, bloblist):
11     return tf(word, blob) * idf(word, bloblist)
12
```

After representing the tweets as mutually comparable vectors by using the TF-IDF implementation from scikit-learn, we were requested to use k-means and a consensus matrix for noise removal. As it was already mentioned, for the similarity metric we were allowed to choose either Euclidean distance or the Cosine distance. After some researches, we decided that it is better to use the cosine distance in our case. It was marked as more fast and efficient in comparison with the euclidian-distance and it was especially suggested when we have to deal with text data. Reference: <https://cmry.github.io/notes/euclidean-v-cosine> We understood how the cosine distance works and after that we used the one that is already provided by scikit-learn.

After doing this step, we had to run the K-means algorithm. We used the code that was already provided to us but we did several modifications. We run the algorithm 9 times in order to construct our consensus matrix. For each time we run the algorithm, we will store on the consensus matrix, the number of times these nodes were matched together on the same cluster.

An important additional post-processing step that was requested, was that if a pair (i, j) of tweets did not cluster together more than 10 percent of the total number of runs, the position (i, j) in the consensus matrix is set to 0. So after we had the results of the consensus matrix, we replaced all the values '1' with '0' since we run the algorithm 9 times and approximately the value that we need to change to 0 was 1. The main reason why the consensus matrix is helpful for us is to detect the noise. We initially implemented a threshold value by finding the average of all the entries in the constructed consensus matrix excluding the diagonals. After doing so, we found the average of each row and in case this average was less than the threshold value, we consider this as a noise point. In this way, we would have clusters with more cleaner tweets since we are going to remove all these noise points. On each time that we run the K-Means algorithm, we will remove the points that were marked as noise from the consensus matrix. In order to decide for the number of cluster, we decided to compute the sum of squered error (SSE).

$$SSE = \sum_{i=1}^K d(x, c_i)^2$$

After that we plotted k against the SSE, and we noticed that the error decreases as k gets larger and this

was because when the number of clusters increases, they should be smaller, so distortion is also smaller. After visually checking the plot, we decided that this number started to normalize for $K=NUMBERR$

The code that shows what is explained above is presented in here:

CODE OF NOISE REMOVAL

CODE for the K-means algorithm

Listing 3: K-Means Algorithm

```
1 def kmeans(X, n_clusters):
2     # initialize labels and prev_labels. prev_labels will be compared with labels to check if the
3     # have been reached.
4     prev_labels = np.zeros(X.shape[0])
5     labels = np.zeros(X.shape[0])
6     indices_results = []
7     # init random indices
8     indices = np.random.choice(X.shape[0], n_clusters, replace=False)
9     indices_results.append(indices)
10    # assign centroids using the indices
11    # centroids = X[indices]
12
13    # the iterative algorithm goes here
14    while (True):
15        # calculate the distances to the centroids
16        distances = get_distances(X, indices)
17        print(labels, indices)
18
19        # assign labels
20        labels = assign_labels(distances, indices)
21
22        # stopping condition
23        # if np.array_equal(labels, prev_labels):
24        #     break
25        print(indices_results)
26        print(indices_results.count(indices))
27
28        if indices_results.count(indices) > 2:
29            break
30
31        # calculate new centroids
32        for cluster_indx in range(len(indices)):
33            # members = X[labels == indices[cluster_indx]]
34            members = [i for i, x in enumerate(labels) if x == indices[cluster_indx]]
35            indices[cluster_indx] = new_centroid(X, members, indices)
36            #centroids[cluster_indx,:] = np.mean(members,axis=0)
37        indices_results.append(indices)
38        # keep the labels for next round's usage
39        # prev_labels = np.argmax(distances, axis=1)
40        #prev_labels = labels
41        indices_results.append(indices)
42
43    return labels, indices
44
45
```

D Clustering

E Visualization

F Conclusion