

# Data Science Lab

## Lab 5 solution

### 1 Introduction

In the following solution, we will go through the pipeline that achieved the baseline result you saw on the leaderboard. Therefore, careful readers will find several things that can be improved and others left apart for further discussions.

#### 1.0.1 Prerequisites

We need to download several contents to complete the lab assignment.

```
[1]: import numpy as np
import nltk
nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to
[nltk_data]   /Users/giuseppe/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

```
[1]: True
```

### 2 Newsgroup clustering

Let's start by loading the dataset.

```
[2]: import os

in_dir = "T-newsgroups"
files = sorted(os.listdir(in_dir), key=int) # note how the part 'key=lambda x:
↳ int(x)' is shortened here
files[:5]
```

```
[2]: ['0', '1', '2', '3', '4']
```

Be careful: `os.listdir` returns filenames in random order. To tackle that, we can just sort them using their integer cast as ordering criteria.

## 2.1 Data preparation step

Here we can make use of the `TfidfVectorizer` class from `scikit-learn`. Starting from the code snippet provided in the laboratory text, we have a custom `LemmaTokenizer` class that can be used in place of the `TfidfVectorizer`'s standard tokenizer. In our tokenizer, we create tokens and lemmatize them through `nlTK` `word_tokenize` and `WordNetLemmatizer` respectively.

```
[3]: from sklearn.feature_extraction.text import TfidfVectorizer
from nltk.tokenize import word_tokenize
from nltk.stem.wordnet import WordNetLemmatizer
from nltk.corpus import stopwords as sw
import re
import string

class LemmaTokenizer(object):
    def __init__(self):
        self.lemmatizer = WordNetLemmatizer()

    def __call__(self, document):
        lemmas = []
        re_digit = re.compile("[0-9]") # regular expression to filter digit
        ↪tokens

        for t in word_tokenize(document):
            t = t.strip()
            lemma = self.lemmatizer.lemmatize(t)

            # remove tokens with only punctuation chars and digits
            if lemma not in string.punctuation and \
                len(lemma) > 3 and \
                len(lemma) < 16 and \
                not re_digit.match(lemma):
                lemmas.append(lemma)

        return lemmas
```

Let's apply now our custom tokenizer to generate the TFIDF representation. We can let the `TfidfVectorizer` read the documents for us, using the `input` parameter to its constructor. We should only care about passing to `fit_transform` the correct path to our files.

```
[4]: stopwords = sw.words('english') + ["'d", "'ll", "'re", "'s", "'ve", 'could',
    ↪'doe', 'ha', 'might', 'must', "n't", 'need', 'sha', 'wa', 'wo', 'would']
tokenizer = LemmaTokenizer()
vectorizer = TfidfVectorizer(input='filename', tokenizer=tokenizer,
    ↪stop_words=stopwords)
```

```
filepath = [os.path.join(in_dir, f) for f in files] # here we need the correct_
↳path
X_tfidf = vectorizer.fit_transform(filepath)
X_tfidf
```

```
[4]: <4000x40684 sparse matrix of type '<class 'numpy.float64''>'
      with 357446 stored elements in Compressed Sparse Row format>
```

Note that we extended the list of stopwords with further non significant words.

### 2.1.1 Principal Component Analysis

As you can see, we are dealing with about 40k different features. These are all the distinct words found in our corpus (i.e. the collection of our news). We can transform and reduce our input feature space through the Principal Component Analysis (PCA) (see Appendix on the laboratory text).

Let's focus on PCA applied by means of the Singular Value Decomposition (SVD). Each component extracted by the decomposition will express a given amount of the variance of our data. In the best case, all the variance is expressed by a low number of new features. As suggested, the `TruncatedSVD` class let us decide how many components we want to retain in the new space. Additionally, it provides you the variance those components encode with some inner attributes.

```
[5]: from sklearn.decomposition import TruncatedSVD
```

Let's try with a low number of features, e.g. 30, and use them to cluster our points. Please remember that the SVD algebraic method transforms the initial space. Hence, you will not end up with the 30 words the best separate your data. Instead, you will have new components (i.e. projections of your points) in a new 30-dimensional space.

The `fit_transform` method performs the decomposition and maps the input data to the new vector space.

```
[6]: svd = TruncatedSVD(n_components=30, random_state=42)
      X_svd = svd.fit_transform(X_tfidf)

      print(f"Total variance explained: {np.sum(svd.explained_variance_ratio_):.2f}")
```

```
Total variance explained: 0.07
```

The 7% is quite low for real applications. We will investigate it in Section 2.4.

## 2.2 Clustering and wordclouds

Now that we have an initial vector representation, we can use it to cluster our documents. We will focus on the K-means algorithm with K starting from 2 and increasing. Next, we can choose the K value that better helps us to spot topics in the final clusters. Doing so, we are impersonating a domain expert that analyzes and evaluates the final subdivision. Remember that, in our context, one clear topic per cluster is the desired outcome.

One simple approach to describe the clusters would exploit wordcloud images. These “clouds” include the most important words in collection of texts for a certain metric. The more a word is significant the bigger it appears in the picture. Fortunately, we already have a “significance” metric: the TFIDF weight of each word. Let’s exploit its very meaning and choose it as the “importance” metric for the wordcloud.

```
[7]: from sklearn.cluster import KMeans
      from wordcloud import WordCloud
      import matplotlib.pyplot as plt
      %matplotlib inline
```

Let’s define a function that clusters our data with K-means and generates one wordcloud per cluster. As you can see in the following code, we will use only the top\_count words with the highest TFIDF value.

```
[8]: def generate_wordclouds(X, X_tfidf, k, word_positions):
      """Cluster X with K-means with the specified k, and generate one wordcloud
      ↪per cluster.

      :input X: numpy.array or numpy.matrix to cluster
      :input X_tfidf: sparse matrix with TFIDF values
      :input k: the k to be used in K-means
      :input word_positions: dictionary with pairs as word index (in vocabulary)
      ↪-> word
      :return cluster_ids: set with the clusters ids
      """

      model = KMeans(n_clusters=k, random_state=42, n_jobs=-1)
      y_pred = model.fit_predict(X)
      cluster_ids = set(y_pred)
      top_count = 100

      for cluster_id in cluster_ids:

          # compute the total tfidf for each term in the cluster
          tfidf = X_tfidf[y_pred == cluster_id]
          tfidf_sum = np.sum(tfidf, axis=0) # numpy.matrix
          tfidf_sum = np.asarray(tfidf_sum).reshape(-1) # numpy.array of shape
          ↪(1, X.shape[1])
          top_indices = tfidf_sum.argsort()[-top_count:]

          term_weights = {word_positions[idx]: tfidf_sum[idx] for idx in
          ↪top_indices}
          wc = WordCloud(width=1200, height=800, background_color="white")
          wordcloud = wc.generate_from_frequencies(term_weights)

          fig, ax = plt.subplots(figsize=(10, 6), dpi=100)
          ax.imshow(wordcloud, interpolation='bilinear')
```









## Cluster 2



It looks like a new topic has popped out. Cluster 0 is now related to space stuff (e.g. space, launch, NASA, shuttle). Also, Cluster 2 now seems to highlight words related to law, government, and justice in general. Note that you cannot expect any fixed order in clusters' Ids, hence Cluster \* is likely going to change "identity" from different runs.

There are at least two reasons that should lead you to increase K. On one hand, we lost the topic identified in Cluster 0 with K=2 (baseball/sport) and that we assumed to find in our final result. On the other hand, Cluster 1 has still no well-defined identity.

### 2.2.3 K = 4

```
[11]: cluster_ids = generate_wordclouds(X_svd, X_tfidf, 4, word_positions)
```































specifically, its sparse data structures to reduce the memory penalty. You can read more about it [here](#).

```
[16]: min_support = 0.3
dist_words = sorted(v for k, v in word_positions.items()) # distinct words in
↳the vocabulary
for cluster_id in cluster_ids:
    print(f"FP-Growth results on Cluster {cluster_id} with min support
↳{min_support}")

    tfidf = X_tfidf[y_pred == cluster_id]
    tfidf[tfidf > 0] = 1 # encoded as binary "presence/absence" representation
↳as required by mlxtend
    df = pd.DataFrame.sparse.from_spmatrix(tfidf, columns=dist_words) # df is a
↳pandas sparse dataframe

    fset = fpgrowth(df, min_support=min_support, use_colnames=True).
↳sort_values(by='support', ascending=False)
    print(fset, '\n')
```

FP-Growth results on Cluster 0 with min support 0.3

	support	itemsets
1	0.615234	(space)
2	0.599609	(writes)
3	0.546875	(article)
10	0.494141	(writes, article)
4	0.378906	(like)
5	0.349609	(also)
9	0.322266	(writes, space)
8	0.320312	(much)
7	0.306641	(year)
0	0.302734	(system)
6	0.300781	(think)

FP-Growth results on Cluster 1 with min support 0.3

	support	itemsets
0	0.672504	(writes)
1	0.604203	(game)
2	0.597198	(article)
16	0.563923	(writes, article)
12	0.502627	(year)
3	0.432574	(last)
4	0.430823	(team)
5	0.390543	(think)
15	0.383538	(game, writes)
6	0.378284	(like)
13	0.376532	(good)
7	0.367776	(time)

8	0.366025	(baseball)
20	0.355517	(year, writes)
17	0.350263	(game, article)
9	0.343257	(first)
10	0.339755	(player)
18	0.327496	(game, writes, article)
19	0.316988	(last, year)
21	0.309982	(year, article)
14	0.304729	(know)
11	0.304729	(well)

FP-Growth results on Cluster 2 with min support 0.3

	support	itemsets
0	0.728129	(writes)
1	0.682369	(article)
16	0.609690	(writes, article)
2	0.557201	(people)
14	0.418573	(right)
7	0.414536	(think)
17	0.410498	(people, writes)
18	0.380888	(people, article)
8	0.379542	(like)
11	0.376851	(government)
9	0.360700	(time)
13	0.345895	(make)
10	0.340511	(thing)
3	0.336474	(even)
19	0.332436	(people, writes, article)
4	0.325707	(know)
5	0.317631	(believe)
20	0.316285	(think, writes)
15	0.312248	(also)
22	0.310902	(writes, right)
6	0.309556	(reference)
21	0.305518	(think, article)
12	0.301480	(well)

FP-Growth results on Cluster 3 with min support 0.3

	support	itemsets
0	0.518399	(writes)
1	0.469641	(article)
2	0.418583	(writes, article)

Here are reported word sets with support at least above 30% sorted by support. As you can see, compared to a wordcloud visualization, cluster topics are harder to identify. Nonetheless, with our prior knowledge coming from wordclouds, we can spot some key words like: \* “space” in Cluster 0 \* “team”, “baseball” and “game” in Cluster 1 \* “government”, “people” in Cluster 2

## 2.4 Further comments

Every previous result was generated with a transformed dataset by means of the Truncated SVD. However, we skipped over the discussion on the number of components. Actually, our new representation expresses only 7% of the explained variance, which is an extremely low number in almost every real-world case.

With this in mind, we can enhance the truncated representation and then explore the results with the best configuration so far ( $K = 4$ ). The choice of the number of components which span the new vector space is typically a matter of the amount of explained variance desired and varies from application to application. To sufficiently distantiate from previous results, let's see how clustering behaves when we retain 80% of the explained variance.

Let's see if 2000 components are enough to our goal.

```
[17]: svd = TruncatedSVD(n_components=2000, random_state=42)
      X_svd = svd.fit_transform(X_tfidf)
```

We can now compute the cumulative sum of the explained variance ratio, per component, and check whether we cross the 80% or not.

```
[18]: cum_variance = np.cumsum(svd.explained_variance_ratio_)
      idx = np.argmax(cum_variance > .8)
      idx
```

```
[18]: 1622
```

Hence, we should transform our data into a new 1622-dimensional space to retain at least 80% of the explained variance. Let's do it and inspect the wordclouds obtained from it.

```
[19]: svd = TruncatedSVD(n_components=1622, random_state=42)
      X_svd = svd.fit_transform(X_tfidf)
      _ = generate_wordclouds(X_svd, X_tfidf, 4, word_positions)
```







