# Text Mining Assignment

Elena Montenegro, Fernando Freire

May 1, 2019

## Contents

# 1 Modules importation and data loading

### Script 1.0.1 (python)

```python
import warnings
warnings.filterwarnings('ignore')
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

%matplotlib inline
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer

from sklearn.naive_bayes import MultinomialNB
from sklearn.decomposition import TruncatedSVD# SVD = Singular Value Descomposition
from sklearn.model_selection import GridSearchCV
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import StandardScaler, Normalizer, MinMaxScaler, MaxAbsScaler
from sklearn.linear_model import LogisticRegression
from sklearn.feature_selection import SelectKBest
from sklearn.pipeline import Pipeline
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.svm import SVC, LinearSVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn import tree
from sklearn.feature_extraction import stop_words

random_state=0
```

### Script 1.0.2 (python)

```python
# Data loading
NROWS = 50
## Negative dataset
df_neg = pd.read_csv('./practica_clase/PRECISION_MEDICINE/negative_training_abstracts.tsv',
    sep='\t',
                    header=None, nrows = NROWS)

df_neg.columns = ['Accession number', 'Title', 'Abstract']
df_neg['Label'] = '0' #'neg'

display(df_neg.head())

corpus_neg = list(df_neg['Abstract'].values)
### len(corpus_neg) # 4078

## Positive
```

```python
df_pos = pd.read_csv('./practica_clase/PRECISION_MEDICINE/positive_training_abstracts.tsv',
→   sep='\t',
                    header=None, nrows = NROWS)

df_pos.columns = ['Accession number', 'Title', 'Abstract']
df_pos['Label'] = '1' # 'pos'
display(df_pos.head())

# Add corpus
df_corpus = df_neg.append(df_pos)
display(df_corpus.head())

# len(corpus) # 8156

labels = df_corpus['Label']
corpus = df_corpus['Abstract']
# len(labels) # 8156

print(len(corpus), len(labels))
```

|   | Accession number | Title | \ |
|---|---|---|---|
| 0 | 29606186 | Can reactivity and regulation in infancy predi... | |
| 1 | 29471205 | Fabrication of bioinspired, self-cleaning supe... | |
| 2 | 29175165 | Functional properties of chickpea protein isol... | |
| 3 | 29098524 | Mechanical dyssynchrony alters left ventricula... | |
| 4 | 27507285 | Reducing the width of confidence intervals for... | |

|   | Abstract | Label |
|---|---|---|
| 0 | A need to identify early infant markers of lat... | 0 |
| 1 | The mechanical properties, corrosion-resistanc... | 0 |
| 2 | In the present study, the effect of Refractanc... | 0 |
| 3 | The impact of left bundle branch block (LBBB) ... | 0 |
| 4 | In the last decade, it has been shown that an ... | 0 |

|   | Accession number | Title | \ |
|---|---|---|---|
| 0 | 27829177 | A naturally occurring variant of HPV-16 E7 exe... | |
| 1 | 27806271 | Functional Analysis of Orai1 Concatemers Suppo... | |
| 2 | 27796307 | KAT2A/KAT2B-targeted acetylome reveals a role ... | |
| 3 | 27795438 | The Cellular DNA Helicase ChlR1 Regulates Chro... | |
| 4 | 27794539 | Human R1441C LRRK2 regulates the synaptic vesi... | |

|   | Abstract | Label |
|---|---|---|
| 0 | Human Papillomavirus E6 and E7 play critical r... | 1 |
| 1 | Store-operated Ca(2+) entry occurs through the... | 1 |
| 2 | Lysine acetylation is a widespread post-transl... | 1 |
| 3 | In papillomavirus infections, the viral genome... | 1 |
| 4 | Mutations in leucine-rich repeat kinase 2 (LRR... | 1 |

```
  Accession number                                           Title  \
0         29606186  Can reactivity and regulation in infancy predi...
1         29471205  Fabrication of bioinspired, self-cleaning supe...
2         29175165  Functional properties of chickpea protein isol...
3         29098524  Mechanical dyssynchrony alters left ventricula...
4         27507285  Reducing the width of confidence intervals for...

                                          Abstract Label
0  A need to identify early infant markers of lat...     0
1  The mechanical properties, corrosion-resistanc...     0
2  In the present study, the effect of Refractanc...     0
3  The impact of left bundle branch block (LBBB) ...     0
4  In the last decade, it has been shown that an ...     0
```

**Output**

```
100 100
```

## 1.1 Data split

Script 1.1.1 (python)

```python
TEST_SIZE = 0.33
X_train, X_test, y_train, y_test = train_test_split(
    corpus, labels, test_size=0.33, random_state=random_state)
```

# 2 Part I. Construction of an automatic classifier

The following parameters can be adjusted in order to try to maximize the quality of the classifier:

- In function TfidfVectorizer:

  - Parameters that affect the vocabulary quality:
    * List of stopwords (one of the options is setting it to None)
    * maxfeatures
    * max_df, min_df
  - Norm (none, 'l1' or 'l2')

- In Latent Semantic Analysis (LSA):

  - n_components
  - not performing LSA

- Classifier model:

  - You can use strategies included in some of the notebooks we used

* Logistic Regression,
* Naïve Bayes,
* decision trees,
* SVC
* or others you learnt from the Machine Learning course (k-nn, neural networks, etc.)

The goal is not to check all possible combinations of these parameters but respond to thesequestions:

- Which tips can you give about constructing an automatic text classifier? What do you recommend to do? What do you recommend not to do?
- What is the best classifier you have obtained?

Your responses to these questions should be illustrated with tables and/or figures and/or screen captures.

## 2.1 Pipelines

### 2.1.1 Find additional stopwords

Script 2.1.1 (python)

```python
def get_top_n_words(corpus, n=None):
    """
    List the top n words in a vocabulary according to occurrence in a text corpus.
    """
    vec = CountVectorizer().fit(corpus)
    bag_of_words = vec.transform(corpus)
    sum_words = bag_of_words.sum(axis=0)
    words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary_.items()]
    words_freq =sorted(words_freq, key = lambda x: x[1], reverse=True)

    return words_freq[:n]

def improve_stop_words(X_train, n=50):
    """
    """
    common_words = [i[0] for i in get_top_n_words(X_train, n)]
    eng_and_custom_stopwords = set(list(stop_words.ENGLISH_STOP_WORDS) + common_words)
    print(len(eng_and_custom_stopwords))
    return eng_and_custom_stopwords
```

### 2.1.2 Pipelining methods

Script 2.1.2 (python)

```python
CLASSIFIERS = ['knn', 'dtree', 'nb', 'lr', 'svc', 'lsvc']
REDUCERS = ['svd', None]
CV = 4

def create_text_pipeline(reducer='svd', classifier="nb"):
    """ Create text vectorization pipeline with optional dimensionality reduction"""
```

```python
        assert reducer in REDUCERS, "ERROR: Reducer %s not supported, only %s" % (reducer,
            REDUCERS)
        assert classifier in CLASSIFIERS, "ERROR: Classifier %s not supported, only %s" %
            (classifier, CLASSIFIERS)
        num_comp = 100
        pipeline = [
            ('vect', TfidfVectorizer()),
            ('scaler', StandardScaler())
        ]

        # Reduce dimensions
        if reducer == 'svd':
            pipeline.append(('red_svd', TruncatedSVD()))
            pipeline.append(('norm', MinMaxScaler()))
        elif reducer == 'kbest':
            pipeline.append(('red_kbest', SelectKBest(k=num_comp)))
            pipeline.append(('norm', MinMaxScaler()))
        elif reducer == 'percentile':
            pipeline.append(('red_percentile', SelectPercentile(f_classif, percentile=num_comp)))
            pipeline.append(('norm', MinMaxScaler()))
        elif reducer == None:
            #pipeline.append(('normalizer', MaxAbsScaler()))
            pass

        # Classify
        if classifier == "nb":
            pipeline.append(('clf_' + classifier, MultinomialNB()))
        elif classifier == "lr":
            pipeline.append(('clf_' + classifier, LogisticRegression()))
        elif classifier == "svc":
            pipeline.append(('clf_' + classifier, SVC()))
        elif classifier == "lsvc":
            pipeline.append(('clf_' + classifier, LinearSVC()))
        elif classifier == "dtree":
            pipeline.append(('clf_' + classifier, DecisionTreeClassifier()))
        elif classifier == "knn":
            pipeline.append(('clf_' + classifier, KNeighborsClassifier()))
            KNeighborsClassifier()
        elif classifier == None:
            pass

        #print("Pipeline", pipeline)
        return Pipeline(pipeline)

def prediction_metrics(X_train, y_train, X_test, y_test, parameters, reducer="svd",
    classifier="nb"):
        """
        """
        print("### Reducer: %s    Classifier: %s" %(reducer, classifier))
        pipeline = create_text_pipeline(reducer=reducer, classifier=classifier)
        # Filter params to only the params related with the pipeline steps
        filtered_params = {}
```

6

```python
56     for param_key in parameters.keys():
57         if param_key.split('__')[0] in pipeline.named_steps.keys():
58             filtered_params[param_key] = parameters[param_key]
59     pipeline.set_params(**filtered_params)
60     pipeline.fit(X_train, y_train)
61     predicted = pipeline.predict(X_test)
62     print(metrics.classification_report(y_test, predicted))
63     print(metrics.confusion_matrix(y_test, predicted))
64
65 def process_classifications(X_train, y_train, X_test, y_test, parameters,
66                             classifiers=CLASSIFIERS, reducers=REDUCERS):
67     """
68     """
69     for classifier in classifiers:
70         for reducer in reducers:
71             prediction_metrics(X_train, y_train, X_test, y_test, parameters, reducer,
                ↪  classifier)
72
73 def prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters_grid,
   ↪  reducer="svd", classifier="nb", cv=CV):
74     """
75     """
76     print("### Reducer: %s   Classifier: %s" %(reducer, classifier))
77     pipeline = create_text_pipeline(reducer=reducer, classifier=classifier)
78     # Filter params to only the params related with the pipeline steps
79     filtered_params = {}
80     for param_key in parameters_grid.keys():
81         if param_key.split('__')[0] in pipeline.named_steps.keys():
82             filtered_params[param_key] = parameters_grid[param_key]
83     grid_model = GridSearchCV(pipeline, filtered_params, cv=cv, iid=False)
84     grid_model.fit(X_train, y_train)
85     for param_name in sorted(filtered_params.keys()):
86         print("\t%s: %r" % (param_name, grid_model.best_params_[param_name]))
87     pipeline.set_params(**grid_model.best_params_)
88     pipeline.fit(X_train, y_train)
89     predicted = pipeline.predict(X_test)
90     print(metrics.classification_report(y_test, predicted))
91     print(metrics.confusion_matrix(y_test, predicted))
92
93 def process_classifications_grid(X_train, y_train, X_test, y_test, parameters, cv=CV,
94                             classifiers=CLASSIFIERS, reducers=REDUCERS):
95     """
96     """
97     for classifier in classifiers:
98         for reducer in reducers:
99             prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters,
                ↪  reducer, classifier, cv=cv)
```

## 2.2 Main process with prefixed parameters

### Script 2.2.1 (python)

```python
# First set of parameters
param_set_1 = {
    'vect__norm': None,
    'vect__smooth_idf': True,
    'vect__sublinear_tf': True,
    'vect__max_features': 50,
    'vect__min_df': 6,
    'vect__stop_words': 'english',
    'vect__strip_accents' : 'unicode',
    'vect__analyzer' : 'word',
    'vect__token_pattern': r'\w{1,}',
    'vect__ngram_range' : (1, 2),
    'scaler__with_mean' : False,
    'vect__norm': 'l2',
    'red_svd__n_components': 40,
    'clf_knn__n_neighbors' : 2
}


# More stop words
#eng_and_custom_stopwords = improve_stop_words(X_train, 200)
#param_set_1['vect__stop_words'] = eng_and_custom_stopwords

#process_classifications(X_train, y_train, X_test, y_test, param_set_1, reducers=['svd'],
    classifiers=['nb'])

process_classifications(X_train, y_train, X_test, y_test, param_set_1)
```

### Output

```
### Reducer: svd   Classifier: knn
             precision    recall  f1-score   support

          0       0.67      0.93      0.78        15
          1       0.92      0.61      0.73        18

  micro avg       0.76      0.76      0.76        33
  macro avg       0.79      0.77      0.76        33
weighted avg       0.80      0.76      0.75        33

[[14  1]
 [ 7 11]]
### Reducer: None   Classifier: knn
             precision    recall  f1-score   support

          0       0.79      1.00      0.88        15
          1       1.00      0.78      0.88        18

  micro avg       0.88      0.88      0.88        33
```

```
      macro avg       0.89      0.89      0.88        33
   weighted avg       0.90      0.88      0.88        33


[[15  0]
 [ 4 14]]
### Reducer: svd   Classifier: dtree
               precision    recall  f1-score   support

            0       0.92      0.80      0.86        15
            1       0.85      0.94      0.89        18

    micro avg       0.88      0.88      0.88        33
    macro avg       0.89      0.87      0.88        33
 weighted avg       0.88      0.88      0.88        33


[[12  3]
 [ 1 17]]
### Reducer: None   Classifier: dtree
               precision    recall  f1-score   support

            0       1.00      0.87      0.93        15
            1       0.90      1.00      0.95        18

    micro avg       0.94      0.94      0.94        33
    macro avg       0.95      0.93      0.94        33
 weighted avg       0.95      0.94      0.94        33


[[13  2]
 [ 0 18]]
### Reducer: svd   Classifier: nb
               precision    recall  f1-score   support

            0       0.82      0.93      0.87        15
            1       0.94      0.83      0.88        18

    micro avg       0.88      0.88      0.88        33
    macro avg       0.88      0.88      0.88        33
 weighted avg       0.89      0.88      0.88        33


[[14  1]
 [ 3 15]]
### Reducer: None   Classifier: nb
               precision    recall  f1-score   support

            0       0.93      0.87      0.90        15
            1       0.89      0.94      0.92        18

    micro avg       0.91      0.91      0.91        33
    macro avg       0.91      0.91      0.91        33
 weighted avg       0.91      0.91      0.91        33


[[13  2]
```

```
 [ 1 17]]
### Reducer: svd    Classifier: lr
            precision    recall  f1-score   support

          0       0.81      0.87      0.84        15
          1       0.88      0.83      0.86        18

   micro avg       0.85      0.85      0.85        33
   macro avg       0.85      0.85      0.85        33
weighted avg       0.85      0.85      0.85        33

[[13  2]
 [ 3 15]]
### Reducer: None    Classifier: lr
            precision    recall  f1-score   support

          0       0.87      0.87      0.87        15
          1       0.89      0.89      0.89        18

   micro avg       0.88      0.88      0.88        33
   macro avg       0.88      0.88      0.88        33
weighted avg       0.88      0.88      0.88        33

[[13  2]
 [ 2 16]]
### Reducer: svd    Classifier: svc
            precision    recall  f1-score   support

          0       0.45      1.00      0.62        15
          1       0.00      0.00      0.00        18

   micro avg       0.45      0.45      0.45        33
   macro avg       0.23      0.50      0.31        33
weighted avg       0.21      0.45      0.28        33

[[15  0]
 [18  0]]
### Reducer: None    Classifier: svc
            precision    recall  f1-score   support

          0       0.87      0.87      0.87        15
          1       0.89      0.89      0.89        18

   micro avg       0.88      0.88      0.88        33
   macro avg       0.88      0.88      0.88        33
weighted avg       0.88      0.88      0.88        33

[[13  2]
 [ 2 16]]
### Reducer: svd    Classifier: lsvc
            precision    recall  f1-score   support
```

```
             0       0.81        0.87       0.84          15
             1       0.88        0.83       0.86          18

    micro avg        0.85        0.85       0.85          33
    macro avg        0.85        0.85       0.85          33
 weighted avg        0.85        0.85       0.85          33

 [[13  2]
  [ 3 15]]
 ### Reducer: None   Classifier: lsvc
             precision    recall  f1-score   support

             0       0.81        0.87       0.84          15
             1       0.88        0.83       0.86          18

    micro avg        0.85        0.85       0.85          33
    macro avg        0.85        0.85       0.85          33
 weighted avg        0.85        0.85       0.85          33

 [[13  2]
  [ 3 15]]
```

## 2.3 Main process with grid search parameters

Script 2.3.1 (python)

```python
1
```

Script 2.3.2 (python)

```python
parameters_grid = {
    'vect__min_df': [5, 6],
    #'vect__max_df': [10, 11],
    'vect__stop_words': (None, 'english', eng_and_custom_stopwords),
    'vect__max_features': [50],
    #'vect__smooth_idf': [True, False],
    'vect__norm': ['l1', 'l2', None]
    #'red_svd__n_components': (50, 100, 200, None),
    #'red_svd__n_components': (10, 20, 30, None),
    #'clf_nb__alpha': (1e-1, 1e-2, 1e-3)
}

parameters_grid = {
    'vect__norm': ['l1', 'l2', None],
    'vect__smooth_idf': [True],
    'vect__sublinear_tf': [True],
    'vect__max_features': [30, 50],
    'vect__min_df': [5,6],
    #'vect__max_df': [7,8],
    'vect__stop_words': [None, 'english', eng_and_custom_stopwords],
```

```python
21      'vect__strip_accents' : ['unicode'],
22      'vect__analyzer' : ['word'],
23      'vect__token_pattern': [r'\w{1,}'],
24      'vect__ngram_range' : [(1, 2)],
25      'scaler__with_mean' : [False],
26      'red_svd__n_components': [2,3],
27      'clf_knn__n_neighbors' : [2, 5]
28 }
29
30 eng_and_custom_stopwords = improve_stop_words(X_train, 200)
31 #prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters_grid, reducer='svd',
   ↪  classifier="knn", cv=2)
32 process_classifications_grid(X_train, y_train, X_test, y_test, parameters_grid, cv=2)
```

## Output

```
### Reducer: svd    Classifier: knn
        clf_knn__n_neighbors: 5
        red_svd__n_components: 2
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 6
        vect__ngram_range: (1, 2)
        vect__norm: None
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

           0       0.88      0.93      0.90        15
           1       0.94      0.89      0.91        18

   micro avg       0.91      0.91      0.91        33
   macro avg       0.91      0.91      0.91        33
weighted avg       0.91      0.91      0.91        33

[[14  1]
 [ 2 16]]
### Reducer: None    Classifier: knn
        clf_knn__n_neighbors: 5
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 50
        vect__min_df: 5
        vect__ngram_range: (1, 2)
        vect__norm: 'l2'
        vect__smooth_idf: True
        vect__stop_words: 'english'
```

```
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
                precision    recall  f1-score   support

             0       0.88      0.93      0.90        15
             1       0.94      0.89      0.91        18

    micro avg       0.91      0.91      0.91        33
    macro avg       0.91      0.91      0.91        33
 weighted avg       0.91      0.91      0.91        33

[[14  1]
 [ 2 16]]
```

### Reducer: svd    Classifier: dtree

```
        red_svd__n_components: 2
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 6
        vect__ngram_range: (1, 2)
        vect__norm: None
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
                precision    recall  f1-score   support

             0       0.93      0.87      0.90        15
             1       0.89      0.94      0.92        18

    micro avg       0.91      0.91      0.91        33
    macro avg       0.91      0.91      0.91        33
 weighted avg       0.91      0.91      0.91        33

[[13  2]
 [ 1 17]]
```

### Reducer: None    Classifier: dtree

```
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 5
        vect__ngram_range: (1, 2)
        vect__norm: 'l1'
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
                precision    recall  f1-score   support
```

```
              0         0.92      0.80      0.86          15
              1         0.85      0.94      0.89          18

   micro avg         0.88      0.88      0.88          33
   macro avg         0.89      0.87      0.88          33
weighted avg         0.88      0.88      0.88          33

[[12  3]
 [ 1 17]]
### Reducer: svd   Classifier: nb
        red_svd__n_components: 2
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 5
        vect__ngram_range: (1, 2)
        vect__norm: 'l2'
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

              0         0.82      0.93      0.87          15
              1         0.94      0.83      0.88          18

   micro avg         0.88      0.88      0.88          33
   macro avg         0.88      0.88      0.88          33
weighted avg         0.89      0.88      0.88          33

[[14  1]
 [ 3 15]]
### Reducer: None   Classifier: nb
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 50
        vect__min_df: 6
        vect__ngram_range: (1, 2)
        vect__norm: 'l2'
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

              0         0.93      0.87      0.90          15
              1         0.89      0.94      0.92          18

   micro avg         0.91      0.91      0.91          33
   macro avg         0.91      0.91      0.91          33
```

```
weighted avg       0.91      0.91      0.91          33

[[13  2]
 [ 1 17]]
### Reducer: svd    Classifier: lr
        red_svd__n_components: 2
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 50
        vect__min_df: 6
        vect__ngram_range: (1, 2)
        vect__norm: 'l2'
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

           0       0.93      0.93      0.93          15
           1       0.94      0.94      0.94          18

   micro avg       0.94      0.94      0.94          33
   macro avg       0.94      0.94      0.94          33
weighted avg       0.94      0.94      0.94          33

[[14  1]
 [ 1 17]]
### Reducer: None    Classifier: lr
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 5
        vect__ngram_range: (1, 2)
        vect__norm: None
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

           0       0.93      0.87      0.90          15
           1       0.89      0.94      0.92          18

   micro avg       0.91      0.91      0.91          33
   macro avg       0.91      0.91      0.91          33
weighted avg       0.91      0.91      0.91          33

[[13  2]
 [ 1 17]]
### Reducer: svd    Classifier: svc
```

```
        red_svd__n_components: 3
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 5
        vect__ngram_range: (1, 2)
        vect__norm: 'l2'
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

           0       0.88      0.93      0.90        15
           1       0.94      0.89      0.91        18

   micro avg       0.91      0.91      0.91        33
   macro avg       0.91      0.91      0.91        33
weighted avg       0.91      0.91      0.91        33

[[14  1]
 [ 2 16]]
### Reducer: None    Classifier: svc
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 5
        vect__ngram_range: (1, 2)
        vect__norm: None
        vect__smooth_idf: True
        vect__stop_words: 'english'
        vect__strip_accents: 'unicode'
        vect__sublinear_tf: True
        vect__token_pattern: '\\w{1,}'
              precision    recall  f1-score   support

           0       0.94      1.00      0.97        15
           1       1.00      0.94      0.97        18

   micro avg       0.97      0.97      0.97        33
   macro avg       0.97      0.97      0.97        33
weighted avg       0.97      0.97      0.97        33

[[15  0]
 [ 1 17]]
### Reducer: svd    Classifier: lsvc
        red_svd__n_components: 3
        scaler__with_mean: False
        vect__analyzer: 'word'
        vect__max_features: 30
        vect__min_df: 6
```

```
          vect__ngram_range: (1, 2)
          vect__norm: None
          vect__smooth_idf: True
          vect__stop_words: 'english'
          vect__strip_accents: 'unicode'
          vect__sublinear_tf: True
          vect__token_pattern: '\\w{1,}'
                  precision    recall  f1-score   support

             0       0.88      0.93      0.90        15
             1       0.94      0.89      0.91        18

    micro avg       0.91      0.91      0.91        33
    macro avg       0.91      0.91      0.91        33
 weighted avg       0.91      0.91      0.91        33

[[14  1]
 [ 2 16]]
### Reducer: None    Classifier: lsvc
          scaler__with_mean: False
          vect__analyzer: 'word'
          vect__max_features: 30
          vect__min_df: 5
          vect__ngram_range: (1, 2)
          vect__norm: None
          vect__smooth_idf: True
          vect__stop_words: 'english'
          vect__strip_accents: 'unicode'
          vect__sublinear_tf: True
          vect__token_pattern: '\\w{1,}'
                  precision    recall  f1-score   support

             0       0.93      0.87      0.90        15
             1       0.89      0.94      0.92        18

    micro avg       0.91      0.91      0.91        33
    macro avg       0.91      0.91      0.91        33
 weighted avg       0.91      0.91      0.91        33

[[13  2]
 [ 1 17]]
```

# 3 Part 2: Construction of a clustering of biology documents

We already know the class information in our dataset (positive and negative) but we will test if an automatic clustering system discovers automatically these classes ("labels"). The objective is to learn strategies that will be very useful when we have to cluster unlabeled documents. Therefore, we "hide" this information (the real class) to the clustering algorithm.

The objective in this section is to check what are the parameters that maximize clustering's quality. The parameters to be taken into account are:

- In function TfidfVectorizer:

    - Vocabulary (larger or smaller)
    - Norm (none, 'l1' or 'l2')

- In Latent Semantic Analysis (LSA):

    - n_components
    - o not performing LSA

- Normalize the data/not normalize it with "Normalizer" (included in the notebook).

The questions to be responded in this part are:

- Which tips can you give about constructing a text clustering with k-means? What do you recommend to do? What do you recommend not to do?

- What is the best clustering you have obtained? The quality of the cluster is the degree of correspondence between real class and assigned cluster. For example:

    - If there are 2 clusters and cluster 0 contains all examples of positive class and cluster 1 contains all examples of negative class, the clustering is perfect.
    - If there are 2 clusters and cluster 1 contains all examples of positive class and cluster 0 contains all examples of negative class, the clustering is also perfect.
    - If there are 2 clusters and cluster 0 contains 50% of examples of positive class and 50% of examples of negative class, and statistics in cluster 1 are similar, the clustering quality is the worst possible.