

# Text Mining Assignment

Elena Montenegro, Fernando Freire

May 2, 2019

## Contents

<b>1</b>	<b>Modules importation and data loading</b>	<b>2</b>
1.1	Data split . . . . .	4
<b>2</b>	<b>Part I. Construction of an automatic classifier</b>	<b>4</b>
2.1	Pipelines . . . . .	5
2.1.1	Find additional stopwords . . . . .	5
2.1.2	Pipelining methods . . . . .	6
2.2	Main process with prefixed parameters . . . . .	9
2.3	Main process with grid search parameters . . . . .	21
<b>3</b>	<b>Part 2: Construction of a clustering of biology documents</b>	<b>35</b>
3.1	Main process with prefixed parameters . . . . .	36
3.2	Main process with grid search parameters . . . . .	37
3.3	Reference process . . . . .	41

# 1 Modules importation and data loading

## Script 1.0.1 (python)

```
1 import warnings
2 warnings.filterwarnings('ignore')
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 import sys
7 %matplotlib inline
8 from sklearn.feature_extraction.text import CountVectorizer
9 from sklearn.feature_extraction.text import TfidfTransformer
10
11 from sklearn.naive_bayes import MultinomialNB
12 from sklearn.decomposition import TruncatedSVD# SVD = Singular Value Descomposition
13 from sklearn.model_selection import GridSearchCV
14 from sklearn.feature_extraction.text import CountVectorizer
15 from sklearn.feature_extraction.text import TfidfVectorizer
16 from sklearn.preprocessing import StandardScaler, Normalizer, MinMaxScaler, MaxAbsScaler
17 from sklearn.linear_model import LogisticRegression
18 from sklearn.feature_selection import SelectKBest, SelectPercentile, f_classif
19 from sklearn.pipeline import Pipeline
20 from sklearn.model_selection import train_test_split
21 from sklearn import metrics
22 from sklearn.svm import SVC, LinearSVC
23 from sklearn.tree import DecisionTreeClassifier
24 from sklearn.neighbors import KNeighborsClassifier
25 from sklearn import tree
26 from sklearn.feature_extraction import stop_words
27 from sklearn.base import TransformerMixin
28 from sklearn.cluster import KMeans
29 from sklearn.metrics import calinski_harabaz_score, accuracy_score
30 from sklearn.preprocessing import Normalizer, LabelBinarizer, OneHotEncoder
31 from sklearn.metrics import make_scorer
32
33 random_state=0
```

## Script 1.0.2 (python)

```
1 # Data loading
2 #NROWS = sys.maxsize
3 NROWS = 50
4 ## Negative dataset
5 df_neg = pd.read_csv('./practica_clase/PRECISION_MEDICINE/negative_training_abstracts.tsv',
6     ↪ sep='\t',
7     header=None, nrows = NROWS)
8 df_neg.columns = ['Accession number', 'Title', 'Abstract']
9 df_neg['Label'] = '0' #'neg'
10
11 display(df_neg.head())
```

```

12
13 corpus_neg = list(df_neg['Abstract'].values)
14 ### len(corpus_neg) # 4078
15
16 ## Positive
17 df_pos = pd.read_csv('./practica_clase/PRECISION_MEDICINE/positive_training_abstracts.tsv',
18   → sep='\t',
19   header=None, nrows = NROWS)
20
21 df_pos.columns = ['Accession number', 'Title', 'Abstract']
22 df_pos['Label'] = '1' # 'pos'
23 display(df_pos.head())
24
25 # Add corpus
26 df_corpus = df_neg.append(df_pos)
27 display(df_corpus.head())
28
29 # len(corpus) # 8156
30
31 labels = df_corpus['Label']
32 corpus = df_corpus['Abstract']
33 # len(labels) # 8156
34
35 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 Orail 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)

```

1 TEST_SIZE = 0.33
2 X_train, X_test, y_train, y_test = train_test_split(
3     corpus, labels, test_size=TEST_SIZE, 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 these questions:

- 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)

```

1 def get_top_n_words(corpus, n=None):
2     """
3     List the top n words in a vocabulary according to occurrence in a text corpus.
4     """
5     vec = CountVectorizer().fit(corpus)
6     bag_of_words = vec.transform(corpus)
7     sum_words = bag_of_words.sum(axis=0)
8     words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary_.items()]
9     words_freq = sorted(words_freq, key = lambda x: x[1], reverse=True)
10
11     return words_freq[:n]
12
13 def improve_stop_words(X_train, n=50):
14     """
15     """
16     common_words = [i[0] for i in get_top_n_words(X_train, n)]
17     eng_and_custom_stopwords = set(list(stop_words.ENGLISH_STOP_WORDS) + common_words)
18     print(len(eng_and_custom_stopwords))
19     return eng_and_custom_stopwords

```

## 2.1.2 Pipelining methods

Script 2.1.2 (python)

```
1 CLASSIFIERS = ['knn', 'dtree', 'nb', 'lr', 'svc', 'lsvc']
2 CLASSIFIERS_UNSUPERVISED = ['kmeans']
3 REDUCERS = ['svd', 'kbest', 'percentile', None]
4 CV = 4
5 VERBOSE = False
6
7 def create_text_pipeline(reducer='svd', classifier="nb"):
8     """ Create text vectorization pipeline with optional dimensionality reduction """
9     assert reducer in REDUCERS, "ERROR: Reducer %s not supported, only %s" % (reducer,
10     ↪ REDUCERS)
11     assert classifier in CLASSIFIERS + CLASSIFIERS_UNSUPERVISED, \
12         "ERROR: Classifier %s not supported, only %s" % (classifier, CLASSIFIERS +
13     ↪ CLASSIFIERS_UNSUPERVISED)
14     pipeline = [
15         ('vect', TfidfVectorizer()),
16         ('scaler', StandardScaler())
17     ]
18     # Reduce dimensions
19     if reducer == 'svd':
20         pipeline.append(('red_svd', TruncatedSVD()))
21     elif reducer == 'kbest':
22         pipeline.append(('red_kbest', SelectKBest()))
23     elif reducer == 'percentile':
24         pipeline.append(('red_percentile', SelectPercentile()))
25     elif reducer == None:
26         pass
27
28     # Classify
29     if classifier == "nb":
30         if reducer == 'svd':
31             pipeline.append(('mm_scaler', MinMaxScaler()))
32         elif reducer == 'kbest':
33             pipeline.append(('mm_scaler', MaxAbsScaler()))
34         elif reducer == 'percentile':
35             pipeline.append(('mm_scaler', MaxAbsScaler()))
36         elif reducer == None:
37             pass
38         pipeline.append(('clf_' + classifier, MultinomialNB()))
39     elif classifier == "lr":
40         pipeline.append(('clf_' + classifier, LogisticRegression()))
41     elif classifier == "svc":
42         pipeline.append(('clf_' + classifier, SVC()))
43     elif classifier == "lsvc":
44         pipeline.append(('clf_' + classifier, LinearSVC()))
45     elif classifier == "dtree":
46         pipeline.append(('clf_' + classifier, DecisionTreeClassifier()))
47     elif classifier == "knn":
48         pipeline.append(('clf_' + classifier, KNeighborsClassifier()))
49     elif classifier == "kmeans":
```

```

48     pipeline.append(('norm', Normalizer()))
49     pipeline.append(('cluster_kmeans', KMeans()))
50 elif classifier == None:
51     pass
52
53 return Pipeline(pipeline)
54
55 def get_prediction_from_cluster(X, pipeline):
56     """ Transform cluster assignment in y_pred object """
57     def swap_label(label):
58         if label == 1:
59             return '0'
60         elif label == 0:
61             return '1'
62         else:
63             return str(label)
64     labels = pipeline.predict(X_test)
65     labels_predicted = [str(label) for label in labels]
66     predicted = pd.Series(labels_predicted)
67     accuracy = metrics.accuracy_score(y_test, predicted)
68     labels_predicted_reverse = [swap_label(label) for label in labels]
69     predicted_reverse = pd.Series(labels_predicted_reverse)
70     accuracy_reverse = metrics.accuracy_score(y_test, predicted_reverse)
71     if accuracy_reverse > accuracy: predicted = predicted_reverse
72     return predicted
73
74 def get_filtered_params(parameters, pipeline):
75     """ Filter the params that aren't related to steps in the pipeline """
76     filtered_params = {}
77     for param_key in parameters.keys():
78         if param_key.split('__')[0] in pipeline.named_steps.keys():
79             filtered_params[param_key] = parameters[param_key]
80     return filtered_params
81
82 def get_filtered_set(parameters, pipeline):
83     """ Filter the params that aren't related to steps in the pipeline """
84     if type(parameters) == dict:
85         return get_filtered_params(parameters, pipeline)
86     else:
87         filtered_set = []
88         for param_set in parameters:
89             filtered_set.append(get_filtered_params(param_set, pipeline))
90     return filtered_set
91
92 def prediction_metrics(X_train, y_train, X_test, y_test, parameters, results, reducer="svd",
93     ↪ classifier="nb"):
94     """
95     """
96     print("### Reducer: %s Classifier: %s" %(reducer, classifier))
97     pipeline = create_text_pipeline(reducer=reducer, classifier=classifier)
98     pipeline.set_params(*get_filtered_params(parameters, pipeline))
99     if VERBOSE: print("Pipeline", pipeline.named_steps)

```

```

99 pipeline.fit(X_train, y_train)
100 if classifier in CLASSIFIERS_UNSUPERVISED:
101     predicted = get_prediction_from_cluster(X_test, pipeline)
102 else:
103     predicted = pipeline.predict(X_test)
104 print()
105 accuracy = metrics.accuracy_score(y_test, predicted)
106 print("Accuracy", accuracy)
107 clf_rep = metrics.classification_report(y_test, predicted, output_dict=True, digits=2)
108 if VERBOSE: print(clf_rep['micro avg'])
109 if VERBOSE: print(metrics.confusion_matrix(y_test, predicted))
110
111 results.append([reducer, classifier, accuracy] + list(clf_rep['micro avg'].values()))
112
113 def process_classifications(X_train, y_train, X_test, y_test, parameters,
114                             classifiers=CLASSIFIERS, reducers=REDUCERS):
115     """
116     """
117     results = []
118     for classifier in classifiers:
119         for reducer in reducers:
120             prediction_metrics(X_train, y_train, X_test, y_test, parameters, results,
121                               ⇨ reducer, classifier)
122             # Group all results into a dataframe
123             df = pd.DataFrame(results, columns=['reducer', 'classifier', 'accuracy', 'precision',
124                                               ⇨ 'recall', 'f1-score', 'support'])
125             df['classifier'].fillna('None', inplace=True)
126
127             return df
128
129 def prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters_grid, results=[],
130                             reducer="svd", classifier="nb", cv=CV):
131     """
132     """
133     print("### Reducer: %s Classifier: %s" %(reducer, classifier))
134     pipeline = create_text_pipeline(reducer=reducer, classifier=classifier)
135     filtered_params = get_filtered_set(parameters_grid, pipeline)
136     #scoring = {'accuracy': make_scorer(accuracy_score), 'calinski':
137     ⇨ make_scorer(calinski_harabaz_score)}
138     scoring = {'accuracy': make_scorer(accuracy_score)}
139     grid_model = GridSearchCV(pipeline, filtered_params, cv=cv, iid=False, error_score=0,
140                               scoring=None, refit=False)
141     grid_model.fit(X_train, y_train)
142     print()
143     print("Best parameters")
144     for param_name in sorted(grid_model.best_params_.keys()):
145         print("\t%s: %s" % (param_name, grid_model.best_params_[param_name]))
146     pipeline.set_params(**grid_model.best_params_)
147     pipeline.fit(X_train, y_train)
148     if classifier in CLASSIFIERS_UNSUPERVISED:
149         predicted = get_prediction_from_cluster(X_test, pipeline)
150     else:

```



```

148     predicted = pipeline.predict(X_test)
149     print()
150     accuracy = metrics.accuracy_score(y_test, predicted)
151     print("Accuracy", accuracy)
152     clf_rep = metrics.classification_report(y_test, predicted, output_dict=True, digits=2)
153     if VERBOSE: print(clf_rep['micro avg'])
154     if VERBOSE: print(metrics.confusion_matrix(y_test, predicted))
155
156     results.append([reducer, classifier, accuracy] + list(clf_rep['micro avg'].values()))
157
158
159 def process_classifications_grid(X_train, y_train, X_test, y_test, parameters, cv=CV,
160                                classifiers=CLASSIFIERS, reducers=REDUCERS):
161     """
162     """
163     results = []
164     for classifier in classifiers:
165         for reducer in reducers:
166             prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters,
167                                    results, reducer, classifier, cv=cv)
168
169     # Group all results into a dataframe
170     df = pd.DataFrame(results, columns=['reducer', 'classifier', 'accuracy', 'precision',
171                                       'recall', 'f1-score', 'support'])
172     df['classifier'].fillna('None', inplace=True)
173
174     return df

```

## 2.2 Main process with prefixed parameters

### Script 2.2.1 (python)

```

1  VERBOSE = False
2  # First set of parameters
3  param_set_1 = {
4      'vect__norm': None,
5      'vect__smooth_idf': True,
6      'vect__sublinear_tf': True,
7      'vect__max_features': 1000,
8      'vect__min_df': 6,
9      'vect__stop_words': 'english',
10     'vect__strip_accents' : 'unicode',
11     'vect__analyzer' : 'word',
12     'vect__token_pattern': r'\w{1,}',
13     'vect__ngram_range' : (1, 2),
14     'scaler' : None,
15     'red_kbest__k' : 3,
16     'red_percentile__score_func' : f_classif,
17     'red_percentile__percentile' : 10,
18     #'scaler__with_mean' : False,
19     'vect__norm': 'l2',

```

```

20     'red_svd__n_components': 40,
21     'clf_knn__n_neighbors' : 2,
22 }
23
24 # More stop words
25 eng_and_custom_stopwords = improve_stop_words(X_train, 200)
26 param_set_1['vect__stop_words'] = eng_and_custom_stopwords
27
28 data_classifier = process_classifications(X_train, y_train, X_test, y_test, param_set_1,
29                                         reducers=REDUCERS, classifiers=CLASSIFIERS)
30
31 #process_classifications(X_train, y_train, X_test, y_test, param_set_1)

```

## Output

```

462
### Reducer: svd   Classifier: knn

Accuracy 0.6666666666666666
### Reducer: kbest   Classifier: knn

Accuracy 0.7272727272727273
### Reducer: percentile   Classifier: knn

Accuracy 0.6666666666666666
### Reducer: None   Classifier: knn

Accuracy 0.6363636363636364
### Reducer: svd   Classifier: dtree

Accuracy 0.6363636363636364
### Reducer: kbest   Classifier: dtree

Accuracy 0.6060606060606061
### Reducer: percentile   Classifier: dtree

Accuracy 0.696969696969697
### Reducer: None   Classifier: dtree

Accuracy 0.42424242424242425
### Reducer: svd   Classifier: nb

Accuracy 0.8484848484848485
### Reducer: kbest   Classifier: nb

Accuracy 0.6666666666666666
### Reducer: percentile   Classifier: nb

Accuracy 0.7575757575757576
### Reducer: None   Classifier: nb

```

```

Accuracy 0.7878787878787878
### Reducer: svd   Classifier: lr

Accuracy 0.7878787878787878
### Reducer: kbest  Classifier: lr

Accuracy 0.6666666666666666
### Reducer: percentile  Classifier: lr

Accuracy 0.7272727272727273
### Reducer: None   Classifier: lr

Accuracy 0.7878787878787878
### Reducer: svd   Classifier: svc

Accuracy 0.45454545454545453
### Reducer: kbest  Classifier: svc

Accuracy 0.6666666666666666
### Reducer: percentile  Classifier: svc

Accuracy 0.45454545454545453
### Reducer: None   Classifier: svc

Accuracy 0.45454545454545453
### Reducer: svd   Classifier: lsvc

Accuracy 0.8484848484848485
### Reducer: kbest  Classifier: lsvc

Accuracy 0.6666666666666666
### Reducer: percentile  Classifier: lsvc

Accuracy 0.7878787878787878
### Reducer: None   Classifier: lsvc

Accuracy 0.8484848484848485

```

### Script 2.2.2 (python)

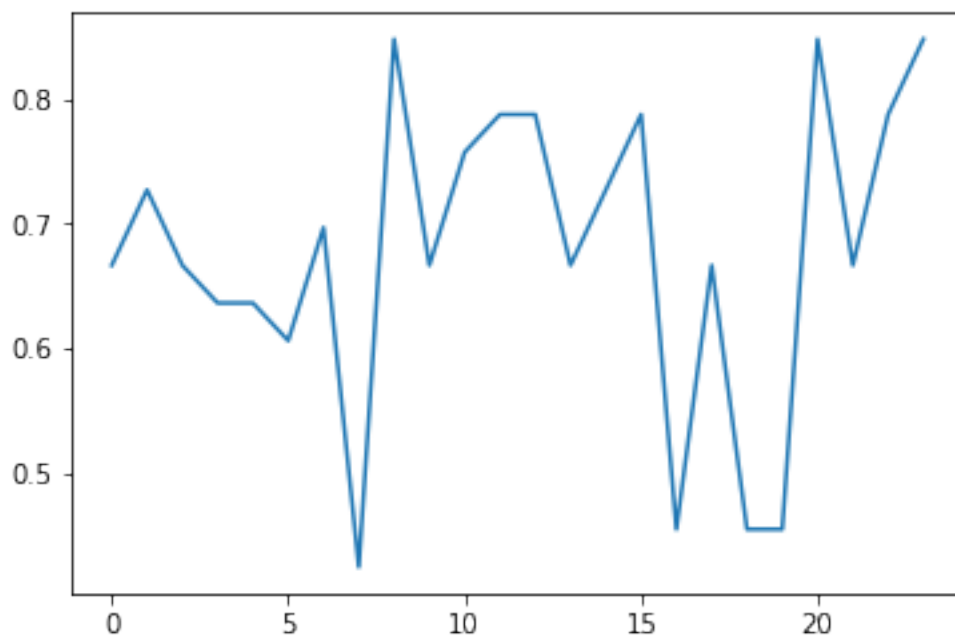
```

1 display(data_classifier)
2 plt.plot(data_classifier['accuracy'].values);

```

	reducer	classifier	accuracy	precision	recall	f1-score	support
0	svd	knn	0.666667	0.666667	0.666667	0.666667	33
1	kbest	knn	0.727273	0.727273	0.727273	0.727273	33
2	percentile	knn	0.666667	0.666667	0.666667	0.666667	33
3	None	knn	0.636364	0.636364	0.636364	0.636364	33
4	svd	dtree	0.636364	0.636364	0.636364	0.636364	33
5	kbest	dtree	0.606061	0.606061	0.606061	0.606061	33
6	percentile	dtree	0.696970	0.696970	0.696970	0.696970	33

7	None	dtree	0.424242	0.424242	0.424242	0.424242	33
8	svd	nb	0.848485	0.848485	0.848485	0.848485	33
9	kbest	nb	0.666667	0.666667	0.666667	0.666667	33
10	percentile	nb	0.757576	0.757576	0.757576	0.757576	33
11	None	nb	0.787879	0.787879	0.787879	0.787879	33
12	svd	lr	0.787879	0.787879	0.787879	0.787879	33
13	kbest	lr	0.666667	0.666667	0.666667	0.666667	33
14	percentile	lr	0.727273	0.727273	0.727273	0.727273	33
15	None	lr	0.787879	0.787879	0.787879	0.787879	33
16	svd	svc	0.454545	0.454545	0.454545	0.454545	33
17	kbest	svc	0.666667	0.666667	0.666667	0.666667	33
18	percentile	svc	0.454545	0.454545	0.454545	0.454545	33
19	None	svc	0.454545	0.454545	0.454545	0.454545	33
20	svd	lsvc	0.848485	0.848485	0.848485	0.848485	33
21	kbest	lsvc	0.666667	0.666667	0.666667	0.666667	33
22	percentile	lsvc	0.787879	0.787879	0.787879	0.787879	33
23	None	lsvc	0.848485	0.848485	0.848485	0.848485	33



## 2.3 Main process with grid search parameters

Script 2.3.1 (python)

```
1 parameters_grid = {
2     'vect__norm': ['l1', 'l2', None],
3     'vect__smooth_idf': [True],
4     'vect__sublinear_tf': [True],
5     'vect__max_features': [900, 1000],
6     'vect__min_df': [1, 5, 6],
7     'vect__max_df': [1., 5., 6],
8     'vect__stop_words': [None, 'english', eng_and_custom_stopwords],
9     'vect__strip_accents' : ['unicode'],
10    'vect__analyzer' : ['word'],
11    'vect__token_pattern': [r'\w{1,}'],
12    'vect__ngram_range' : [(1, 2)],
13    'scaler' : [None],
14    'red_svd__n_components': [2, 30, 40],
15    'clf_knn__n_neighbors' : [2, 5],
16    'red_percentile__score_func' : [f_classif],
17    'red_percentile__percentile' : [10],
18    'red_kbest__k' : [3]
19 }
20
21 eng_and_custom_stopwords = improve_stop_words(X_train, 200)
22 #prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters_grid, reducer='svd',
23 → classifier="knn", cv=2)
24 data_classifier_grid = process_classifications_grid(X_train, y_train, X_test, y_test,
25 → parameters_grid, cv=3)
```

### Output

```
462
### Reducer: svd   Classifier: knn

Best parameters
  clf_knn__n_neighbors: 5
  red_svd__n_components: 2
  scaler: None
  vect__analyzer: 'word'
  vect__max_df: 1.0
  vect__max_features: 900
  vect__min_df: 1
  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,}'

Accuracy 0.9090909090909091
```

```
### Reducer: kbest Classifier: knn
```

```
Best parameters
```

```
clf_knn__n_neighbors: 2
red_kbest__k: 3
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 6
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,}'
```

```
Accuracy 0.9393939393939394
```

```
### Reducer: percentile Classifier: knn
```

```
Best parameters
```

```
clf_knn__n_neighbors: 5
red_percentile__percentile: 10
red_percentile__score_func: <function f_classif at 0x1a1b98f1e0>
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

```
Accuracy 0.8787878787878788
```

```
### Reducer: None Classifier: knn
```

```
Best parameters
```

```
clf_knn__n_neighbors: 5
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

Accuracy 0.9090909090909091

### Reducer: svd Classifier: dtree

Best parameters

```
red_svd__n_components: 2  
scaler: None  
vect__analyzer: 'word'  
vect__max_df: 1.0  
vect__max_features: 900  
vect__min_df: 1  
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,}'
```

Accuracy 0.9393939393939394

### Reducer: kbest Classifier: dtree

Best parameters

```
red_kbest__k: 3  
scaler: None  
vect__analyzer: 'word'  
vect__max_df: 1.0  
vect__max_features: 900  
vect__min_df: 5  
vect__ngram_range: (1, 2)  
vect__norm: 'l1'  
vect__smooth_idf: True  
vect__stop_words: None  
vect__strip_accents: 'unicode'  
vect__sublinear_tf: True  
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

### Reducer: percentile Classifier: dtree

Best parameters

```
red_percentile__percentile: 10  
red_percentile__score_func: <function f_classif at 0x1a1b98f1e0>  
scaler: None  
vect__analyzer: 'word'  
vect__max_df: 6  
vect__max_features: 1000  
vect__min_df: 5  
vect__ngram_range: (1, 2)
```

```
vect__norm: None
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.8787878787878788

### Reducer: None Classifier: dtree

Best parameters

```
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

Accuracy 0.8787878787878788

### Reducer: svd Classifier: nb

Best parameters

```
red_svd__n_components: 30
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 1
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,}'
```

Accuracy 0.9090909090909091

### Reducer: kbest Classifier: nb

Best parameters

```
red_kbest__k: 3
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 6
```



```
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,}'
```

Accuracy 0.5454545454545454

### Reducer: percentile Classifier: nb

Best parameters

```
red_percentile__percentile: 10
red_percentile__score_func: <function f_classif at 0x1a1b98f1e0>
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

Accuracy 0.9090909090909091

### Reducer: None Classifier: nb

Best parameters

```
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 5
vect__ngram_range: (1, 2)
vect__norm: 'l2'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9090909090909091

### Reducer: svd Classifier: lr

Best parameters

```
red_svd__n_components: 30
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
```

```
vect__max_features: 900
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,}'
```

Accuracy 0.9393939393939394

### Reducer: kbest Classifier: lr

Best parameters

```
red_kbest__k: 3
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: None
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

### Reducer: percentile Classifier: lr

Best parameters

```
red_percentile__percentile: 10
red_percentile__score_func: <function f_classif at 0x1a1b98f1e0>
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: 'l2'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

### Reducer: None Classifier: lr

Best parameters

```
scaler: None
```

```
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

Accuracy 0.9393939393939394

### Reducer: svd Classifier: svc

Best parameters

```
red_svd__n_components: 2
scaler: None
vect__analyzer: 'word'
vect__max_df: 5.0
vect__max_features: 900
vect__min_df: 1
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,}'
```

Accuracy 0.9393939393939394

### Reducer: kbest Classifier: svc

Best parameters

```
red_kbest__k: 3
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: None
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

### Reducer: percentile Classifier: svc

Best parameters

```
red_percentile__percentile: 10
red_percentile__score_func: <function f_classif at 0x1a1b98f1e0>
scaler: None
vect__analyzer: 'word'
vect__max_df: 6
vect__max_features: 900
vect__min_df: 5
vect__ngram_range: (1, 2)
vect__norm: None
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.8181818181818182

### Reducer: None Classifier: svc

Best parameters

```
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

Accuracy 0.8787878787878788

### Reducer: svd Classifier: lsvc

Best parameters

```
red_svd__n_components: 30
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
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,}'
```

Accuracy 0.8787878787878788

### Reducer: kbest Classifier: lsvc

Best parameters

```
red_kbest__k: 3
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: 'l2'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.8181818181818182

### Reducer: percentile Classifier: lsvc

Best parameters

```
red_percentile__percentile: 10
red_percentile__score_func: <function f_classif at 0x1a1b98f1e0>
scaler: None
vect__analyzer: 'word'
vect__max_df: 6
vect__max_features: 900
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: 'l2'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9090909090909091

### Reducer: None Classifier: lsvc

Best parameters

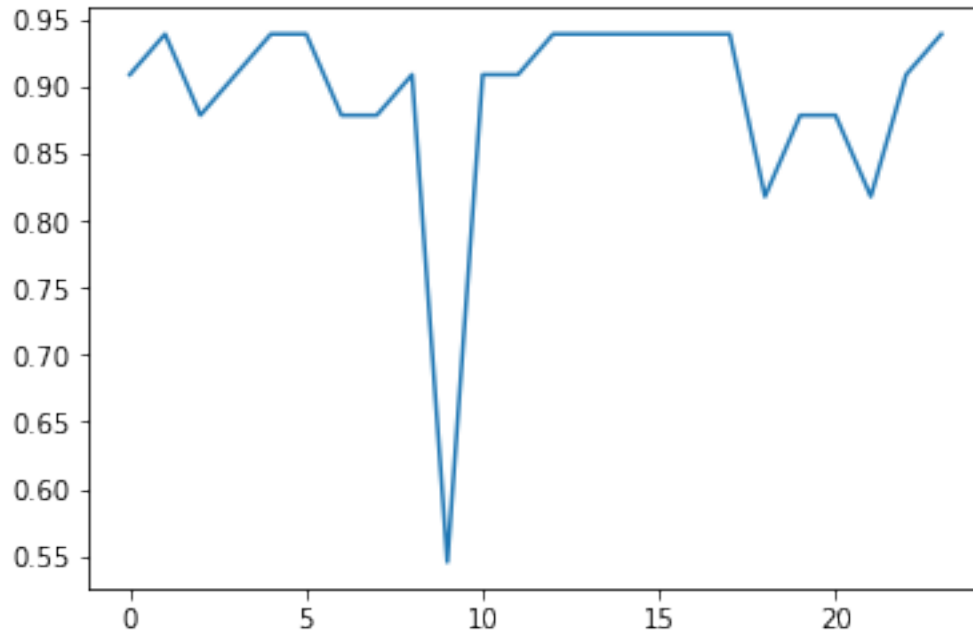
```
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 900
vect__min_df: 5
vect__ngram_range: (1, 2)
vect__norm: 'l2'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

### Script 2.3.2 (python)

```
1 display(data_classifier_grid)
2 plt.plot(data_classifier_grid['accuracy'].values);
```

	reducer	classifier	accuracy	precision	recall	f1-score	support
0	svd	knn	0.909091	0.909091	0.909091	0.909091	33
1	kbest	knn	0.939394	0.939394	0.939394	0.939394	33
2	percentile	knn	0.878788	0.878788	0.878788	0.878788	33
3	None	knn	0.909091	0.909091	0.909091	0.909091	33
4	svd	dtree	0.939394	0.939394	0.939394	0.939394	33
5	kbest	dtree	0.939394	0.939394	0.939394	0.939394	33
6	percentile	dtree	0.878788	0.878788	0.878788	0.878788	33
7	None	dtree	0.878788	0.878788	0.878788	0.878788	33
8	svd	nb	0.909091	0.909091	0.909091	0.909091	33
9	kbest	nb	0.545455	0.545455	0.545455	0.545455	33
10	percentile	nb	0.909091	0.909091	0.909091	0.909091	33
11	None	nb	0.909091	0.909091	0.909091	0.909091	33
12	svd	lr	0.939394	0.939394	0.939394	0.939394	33
13	kbest	lr	0.939394	0.939394	0.939394	0.939394	33
14	percentile	lr	0.939394	0.939394	0.939394	0.939394	33
15	None	lr	0.939394	0.939394	0.939394	0.939394	33
16	svd	svc	0.939394	0.939394	0.939394	0.939394	33
17	kbest	svc	0.939394	0.939394	0.939394	0.939394	33
18	percentile	svc	0.818182	0.818182	0.818182	0.818182	33
19	None	svc	0.878788	0.878788	0.878788	0.878788	33
20	svd	lsvc	0.878788	0.878788	0.878788	0.878788	33
21	kbest	lsvc	0.818182	0.818182	0.818182	0.818182	33
22	percentile	lsvc	0.909091	0.909091	0.909091	0.909091	33
23	None	lsvc	0.939394	0.939394	0.939394	0.939394	33



### 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.

### 3.1 Main process with prefixed parameters

Script 3.1.1 (python)

```

1 param_set_1 = {
2     'vect__smooth_idf': True,
3     'vect__sublinear_tf': True,
4     'vect__max_features': 1000,
5     'vect__min_df': 1,
6     'vect__max_df': 1.,
7     'vect__stop_words': 'english',
8     'vect__strip_accents': 'unicode',
9     'vect__analyzer': 'word',
10    'vect__token_pattern': r'\w{1,}',
11    'vect__ngram_range': (1, 2),
12    #'scaler__with_mean': False,
13    'vect__norm': 'l2',
14    'red_svd__n_components': 100,
15    'clf_knn__n_neighbors': 2,
16    'cluster_kmeans__n_clusters': 2,
17    'red_kbest__k': 3,
18    'red_percentile__score_func': f_classif,
19    'red_percentile__percentile': 10,
20    'scaler': None,
21    'norm': None
22 }
23
24 process_classifications(X_train, y_train, X_test, y_test, param_set_1, reducers=['svd'],
    ↪ classifiers=['kmeans'])

```

#### Output

```

### Reducer: svd   Classifier: kmeans
Pipeline {'vect': TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
    dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
    lowercase=True, max_df=1.0, max_features=1000, min_df=1,
    ngram_range=(1, 2), norm='l2', preprocessor=None, smooth_idf=True,
    stop_words='english', strip_accents='unicode', sublinear_tf=True,
    token_pattern='\\w{1,}', tokenizer=None, use_idf=True,
    vocabulary=None), 'scaler': None, 'red_svd': TruncatedSVD(algorithm='randomized',
    ↪ n_components=100, n_iter=5,

```



```

    random_state=None, tol=0.0), 'norm': None, 'cluster_kmeans': KMeans(algorithm='auto',
    ↪ copy_x=True, init='k-means++', max_iter=300,
    n_clusters=2, n_init=10, n_jobs=None, precompute_distances='auto',
    random_state=None, tol=0.0001, verbose=0)}

```

Accuracy 0.8787878787878788

	precision	recall	f1-score	support
0	1.00	0.73	0.85	15
1	0.82	1.00	0.90	18
micro avg	0.88	0.88	0.88	33
macro avg	0.91	0.87	0.87	33
weighted avg	0.90	0.88	0.88	33

```

[[11  4]
 [ 0 18]]

```

## 3.2 Main process with grid search parameters

### Script 3.2.1 (python)

```

1 eng_and_custom_stopwords = improve_stop_words(X_train, 200)
2 parameters_grid = [
3     {'vect__norm': ['l1', 'l2', None],
4     'vect__smooth_idf': [True],
5     'vect__sublinear_tf': [True],
6     'vect__max_features': [20, 30],
7     'vect__min_df': [1, 5],
8     'vect__max_df': [1., 6],
9     'vect__stop_words': [None, 'english', eng_and_custom_stopwords],
10    'vect__strip_accents': ['unicode'],
11    'vect__analyzer': ['word'],
12    'vect__token_pattern': [r'\w{1,}'],
13    'vect__ngram_range': [(1, 2)],
14    'scaler': [None],
15    'red_svd__n_components': [2, 10, 15],
16    'clf_knn__n_neighbors': [2, 5],
17    'cluster_kmeans__n_clusters': [2],
18    'norm': [None]},
19    # without svd
20    {'vect__norm': ['l1', 'l2', None],
21    'vect__smooth_idf': [True],
22    'vect__sublinear_tf': [True],
23    'vect__max_features': [20, 30],
24    'vect__min_df': [1, 5],
25    'vect__max_df': [1., 6],
26    'vect__stop_words': [None, 'english', eng_and_custom_stopwords],
27    'vect__strip_accents': ['unicode'],
28    'vect__analyzer': ['word'],

```

```

29     'vect__token_pattern': [r'\w{1,}'],
30     'vect__ngram_range' : [(1, 2)],
31     'scaler' : [None],
32     'red_svd__n_components': [2, 10, 15],
33     'clf_knn__n_neighbors' : [2, 5],
34     'cluster_kmeans__n_clusters' : [2],
35     'red__svd' : [None]}
36 ]
37
38 eng_and_custom_stopwords = improve_stop_words(X_train, 200)
39 #prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters_grid,
→   reducer="reducer", classifier="kmeans", cv=CV)
40 process_classifications_grid(X_train, y_train, X_test, y_test, parameters_grid,
→   classifiers=["kmeans"], cv=4)

```

## Output

```

462
462
### Reducer: svd   Classifier: kmeans

Best parameters
  cluster_kmeans__n_clusters: 2
  norm: None
  red_svd__n_components: 2
  scaler: None
  vect__analyzer: 'word'
  vect__max_df: 1.0
  vect__max_features: 30
  vect__min_df: 1
  vect__ngram_range: (1, 2)
  vect__norm: 'l1'
  vect__smooth_idf: True
  vect__stop_words: None
  vect__strip_accents: 'unicode'
  vect__sublinear_tf: True
  vect__token_pattern: '\\w{1,}'

Accuracy 0.7878787878787878
      precision    recall  f1-score   support

     0         1.00      0.53      0.70         15
     1         0.72      1.00      0.84         18

   micro avg       0.79      0.79      0.79         33
   macro avg       0.86      0.77      0.77         33
weighted avg       0.85      0.79      0.77         33

[[ 8  7]
 [ 0 18]]
### Reducer: kbest   Classifier: kmeans

```

Best parameters

```
cluster_kmeans__n_clusters: 2
norm: None
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 30
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: 'l1'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.8181818181818182

	precision	recall	f1-score	support
0	0.91	0.67	0.77	15
1	0.77	0.94	0.85	18
micro avg	0.82	0.82	0.82	33
macro avg	0.84	0.81	0.81	33
weighted avg	0.83	0.82	0.81	33

[[10 5]

[ 1 17]]

### Reducer: percentile Classifier: kmeans

Best parameters

```
cluster_kmeans__n_clusters: 2
norm: None
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 20
vect__min_df: 1
vect__ngram_range: (1, 2)
vect__norm: 'l1'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.696969696969697

	precision	recall	f1-score	support
0	0.86	0.40	0.55	15
1	0.65	0.94	0.77	18

micro avg	0.70	0.70	0.70	33
macro avg	0.76	0.67	0.66	33
weighted avg	0.75	0.70	0.67	33

```
[[ 6  9]
 [ 1 17]]
```

```
### Reducer: None Classifier: kmeans
```

Best parameters

```
cluster_kmeans__n_clusters: 2
norm: None
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 20
vect__min_df: 5
vect__ngram_range: (1, 2)
vect__norm: 'l1'
vect__smooth_idf: True
vect__stop_words: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.7575757575757576

	precision	recall	f1-score	support
0	0.89	0.53	0.67	15
1	0.71	0.94	0.81	18
micro avg	0.76	0.76	0.76	33
macro avg	0.80	0.74	0.74	33
weighted avg	0.79	0.76	0.74	33

```
[[ 8  7]
 [ 1 17]]
```

### 3.3 Reference process

#### Script 3.3.1 (python)

```
1 from sklearn.cluster import KMeans
2 from sklearn.metrics import calinski_harabaz_score
3 from sklearn.preprocessing import Normalizer
4 from sklearn.pipeline import make_pipeline
5 from sklearn.preprocessing import Normalizer
6
7 def get_X_transform(X):
8     vectorizador = TfidfVectorizer(max_df=1., max_features=1000, norm='l2',
```

```

9         min_df=1, stop_words='english',
10         #stop_words=stopwords,
11         #token_pattern=r'(?u)\b[A-Za-z]+\b',
12         #token_pattern=r'(?ui)\b\w*[a-z]+\w*\b',
13         use_idf=True)
14
15     vectorizador = TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
16         encoding='utf-8', input='content',
17         lowercase=True, max_df=1.0, max_features=1000, min_df=1,
18         ngram_range=(1, 2), norm='l2', preprocessor=None, smooth_idf=True,
19         stop_words='english', strip_accents='unicode', sublinear_tf=True,
20         token_pattern='(?u)\b\w*\b', tokenizer=None, use_idf=True,
21         vocabulary=None)
22
23     X = vectorizador.fit_transform(X)
24
25     print(X.shape)
26     n_componentes = 100
27     svd_truncado = TruncatedSVD(n_componentes)
28     normalizador = Normalizer(copy=False)
29
30     lsa = make_pipeline(svd_truncado, normalizador)
31     #lsa = svd_truncado
32
33     X_lsa = lsa.fit_transform(X)
34
35     varianza_explicada = svd_truncado.explained_variance_ratio_.sum()
36     normalizer = Normalizer()
37     X_lsa_norm = normalizer.fit_transform(X_lsa)
38     return X_lsa_norm
39
40 X_km = get_X_transform(X_train)
41 qmetric = calinski_harabaz_score
42
43 Nclusters_max = 15
44 Nrepetitions = 100
45
46 qualities = []
47 inertias = []
48 models = []
49 kini = 1
50 kfin = 20
51 for k in range(kini, kfin+1):
52     print("Evaluando k=%d" % k)
53     km = KMeans(n_clusters=k,
54         init='k-means++', n_init=Nrepetitions,
55         max_iter=500, random_state=2)
56     km.fit(X_km)
57     models.append(km)
58     inertias.append(km.inertia_)
59     if k > 1:
60         qualities.append(qmetric(X_km, km.labels_))

```

```

61         #qualities.append(km.score(X_km))
62     else:
63         qualities.append(0)

```

## Output

```

(67, 1000)
Evaluando k=1
Evaluando k=2
Evaluando k=3
Evaluando k=4
Evaluando k=5
Evaluando k=6
Evaluando k=7
Evaluando k=8
Evaluando k=9
Evaluando k=10
Evaluando k=11
Evaluando k=12
Evaluando k=13
Evaluando k=14
Evaluando k=15
Evaluando k=16
Evaluando k=17
Evaluando k=18
Evaluando k=19
Evaluando k=20

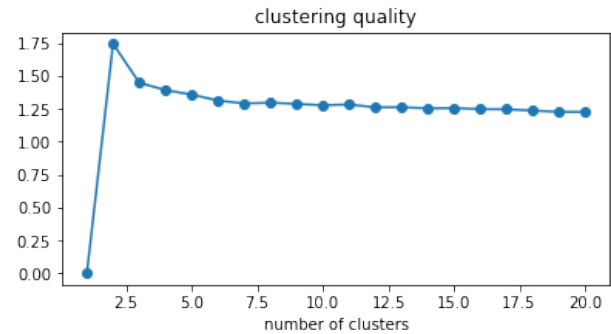
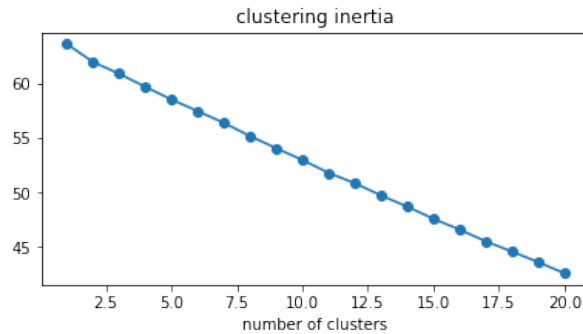
```

## Script 3.3.2 (python)

```

1  fig = plt.figure(figsize=(14,3))
2
3  ax = plt.subplot(1,2,1)
4  plt.plot(range(kini,kfin+1), inertias, marker='o')
5  plt.xlabel('number of clusters')
6  plt.title('clustering inertia')
7
8  ax = plt.subplot(1,2,2)
9  plt.plot(range(kini,kfin+1), qualities, marker='o')
10 plt.xlabel('number of clusters')
11 plt.title('clustering quality')
12 plt.show()
13
14 best = pd.Series(qualities).idxmax() # get index for the best model
15 print("Best number of clusters", best)
16 km = models[best]
17 n_clusters = km.get_params()['n_clusters']
18 clusters = km.labels_
19 print ('Number of clusters of best quality', n_clusters)

```



## Output

Best number of clusters 1  
Number of clusters of best quality 2

## Script 3.3.3 (python)

```
1 # We choose the best option to evaluate the quality of prediction
2 X = X_test
3 y = y_test
4 X_km = get_X_transform(X)
5 labels = km.fit_predict(X_km)
6 #print(labels)
7 # First we try with labels as is
8 labels_predicted = [str(label) for label in labels]
9 predicted = pd.Series(labels_predicted)
10 #print(labels_predicted)
11 print(metrics.classification_report(y, predicted))
12 print(metrics.confusion_matrix(y, predicted))
13
14 # Alternatively we invert the label to match the real labels of each group
15 labels_predicted = [str((label + 1)%2) for label in labels]
16 #print(labels_predicted)
17 predicted = pd.Series(labels_predicted)
18 print(metrics.classification_report(y, predicted))
19 print(metrics.confusion_matrix(y, predicted))
```

## Output

```
(33, 1000)
      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]]

```

	precision	recall	f1-score	support
0	0.10	0.13	0.11	15
1	0.00	0.00	0.00	18
micro avg	0.06	0.06	0.06	33
macro avg	0.05	0.07	0.06	33
weighted avg	0.05	0.06	0.05	33

```

[[ 2 13]
 [18  0]]

```