

# 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 df_pos.columns = ['Accession number', 'Title', 'Abstract']
21 df_pos['Label'] = '1' # 'pos'
22 display(df_pos.head())
23
24 # Add corpus
25 df_corpus = df_neg.append(df_pos)
26 display(df_corpus.head())
27
28 # len(corpus) # 8156
29
30 labels = df_corpus['Label']
31 corpus = df_corpus['Abstract']
32 # len(labels) # 8156
33
34 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
6 def create_text_pipeline(reducer='svd', classifier="nb"):
7     """ Create text vectorization pipeline with optional dimensionality reduction """
8     assert reducer in REDUCERS, "ERROR: Reducer %s not supported, only %s" % (reducer,
9     ↪ REDUCERS)
10    assert classifier in CLASSIFIERS + CLASSIFIERS_UNSUPERVISED, \
11        "ERROR: Classifier %s not supported, only %s" % (classifier, CLASSIFIERS +
12        ↪ CLASSIFIERS_UNSUPERVISED)
13    pipeline = [
14        ('vect', TfidfVectorizer()),
15        ('scaler', StandardScaler())
16    ]
17    # Reduce dimensions
18    if reducer == 'svd':
19        pipeline.append(('red_svd', TruncatedSVD()))
20    elif reducer == 'kbest':
21        pipeline.append(('red_kbest', SelectKBest()))
22    elif reducer == 'percentile':
23        pipeline.append(('red_percentile', SelectPercentile()))
24    elif reducer == None:
25        pass
26
27    # Classify
28    if classifier == "nb":
29        if reducer == 'svd':
30            pipeline.append(('mm_scaler', MinMaxScaler()))
31        elif reducer == 'kbest':
32            pipeline.append(('mm_scaler', MaxAbsScaler()))
33        elif reducer == 'percentile':
34            pipeline.append(('mm_scaler', MaxAbsScaler()))
35        elif reducer == None:
36            pass
37        pipeline.append(('clf_' + classifier, MultinomialNB()))
38    elif classifier == "lr":
39        pipeline.append(('clf_' + classifier, LogisticRegression()))
40    elif classifier == "svc":
41        pipeline.append(('clf_' + classifier, SVC()))
42    elif classifier == "lsvc":
43        pipeline.append(('clf_' + classifier, LinearSVC()))
44    elif classifier == "dtree":
45        pipeline.append(('clf_' + classifier, DecisionTreeClassifier()))
46    elif classifier == "knn":
47        pipeline.append(('clf_' + classifier, KNeighborsClassifier()))
48    elif classifier == "kmeans":
49        pipeline.append(('norm', Normalizer()))
```

```

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

```

```

99     if classifier in CLASSIFIERS_UNSUPERVISED:
100         predicted = get_prediction_from_cluster(X_test, pipeline)
101     else:
102         predicted = pipeline.predict(X_test)
103     print()
104     print("Accuracy", metrics.accuracy_score(y_test, predicted))
105     print(metrics.classification_report(y_test, predicted))
106     print(metrics.confusion_matrix(y_test, predicted))
107
108 def process_classifications(X_train, y_train, X_test, y_test, parameters,
109                             classifiers=CLASSIFIERS, reducers=REDUCERS):
110     """
111     """
112     for classifier in classifiers:
113         for reducer in reducers:
114             prediction_metrics(X_train, y_train, X_test, y_test, parameters, reducer,
115                               ↪ classifier)
116
117 def prediction_metrics_grid(X_train, y_train, X_test, y_test, parameters_grid,
118 ↪ reducer="svd", classifier="nb", cv=CV):
119     """
120     """
121     print("### Reducer: %s Classifier: %s" %(reducer, classifier))
122     pipeline = create_text_pipeline(reducer=reducer, classifier=classifier)
123     filtered_params = get_filtered_set(parameters_grid, pipeline)
124     #scoring = {'accuracy': make_scorer(accuracy_score), 'calinski':
125     ↪ make_scorer(calinski_harabaz_score)}
126     scoring = {'accuracy': make_scorer(accuracy_score)}
127     grid_model = GridSearchCV(pipeline, filtered_params, cv=cv, iid=False, error_score=0,
128                               scoring=None, refit=False)
129     grid_model.fit(X_train, y_train)
130     print()
131     print("Best parameters")
132     for param_name in sorted(grid_model.best_params_.keys()):
133         print("\t%s: %r" % (param_name, grid_model.best_params_[param_name]))
134     pipeline.set_params(**grid_model.best_params_)
135     pipeline.fit(X_train, y_train)
136     if classifier in CLASSIFIERS_UNSUPERVISED:
137         predicted = get_prediction_from_cluster(X_test, pipeline)
138     else:
139         predicted = pipeline.predict(X_test)
140     print()
141     print("Accuracy", metrics.accuracy_score(y_test, predicted))
142     print(metrics.classification_report(y_test, predicted))
143     print(metrics.confusion_matrix(y_test, predicted))
144
145 def process_classifications_grid(X_train, y_train, X_test, y_test, parameters, cv=CV,
146 ↪ classifiers=CLASSIFIERS, reducers=REDUCERS):
147     """
148     """
149     for classifier in classifiers:

```



```

148         for reducer in reducers:
149             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)

```

1  # First set of parameters
2  param_set_1 = {
3      'vect__norm': None,
4      'vect__smooth_idf': True,
5      'vect__sublinear_tf': True,
6      'vect__max_features': 1000,
7      'vect__min_df': 6,
8      'vect__stop_words': 'english',
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_kbest__k': 3,
15     'red_percentile__score_func': f_classif,
16     'red_percentile__percentile': 10,
17     #'scaler__with_mean': False,
18     'vect__norm': 'l2',
19     'red_svd__n_components': 40,
20     'clf_knn__n_neighbors': 2,
21 }
22
23 # More stop words
24 #eng_and_custom_stopwords = improve_stop_words(X_train, 200)
25 #param_set_1['vect__stop_words'] = eng_and_custom_stopwords
26
27 process_classifications(X_train, y_train, X_test, y_test, param_set_1, reducers=REDUCERS,
28 ⇨ classifiers=CLASSIFIERS)
29
29 #process_classifications(X_train, y_train, X_test, y_test, param_set_1)

```

### Output

```

### Reducer: svd   Classifier: knn
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=6,
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=40, n_iter=5,

```

```

random_state=None, tol=0.0), 'clf_knn': KNeighborsClassifier(algorithm='auto',
↳ leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=None, n_neighbors=2, p=2,
weights='uniform'))}

```

Accuracy 0.9393939393939394

	precision	recall	f1-score	support
0	0.88	1.00	0.94	15
1	1.00	0.89	0.94	18
micro avg	0.94	0.94	0.94	33
macro avg	0.94	0.94	0.94	33
weighted avg	0.95	0.94	0.94	33

```

[[15  0]
 [ 2 16]]

```

### Reducer: kbest Classifier: knn

```

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=6,
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_kbest': SelectKBest(k=3, score_func=<function
↳ f_classif at 0x1a1b98f1e0>), 'clf_knn': KNeighborsClassifier(algorithm='auto',
↳ leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=None, n_neighbors=2, p=2,
weights='uniform'))}

```

Accuracy 0.9090909090909091

	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: percentile Classifier: knn

```

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=6,
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_percentile': SelectPercentile(percentile=10,
score_func=<function f_classif at 0x1a1b98f1e0>), 'clf_knn':
↳ KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',

```

```
metric_params=None, n_jobs=None, n_neighbors=2, p=2,
weights='uniform'))}
```

Accuracy 0.696969696969697

	precision	recall	f1-score	support
0	0.61	0.93	0.74	15
1	0.90	0.50	0.64	18
micro avg	0.70	0.70	0.70	33
macro avg	0.75	0.72	0.69	33
weighted avg	0.77	0.70	0.69	33

```
[[14  1]
 [ 9  9]]
```

### Reducer: None Classifier: knn

```
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=6,
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, 'clf_knn': KNeighborsClassifier(algorithm='auto',
↪ leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=None, n_neighbors=2, p=2,
weights='uniform'))}
```

Accuracy 0.9393939393939394

	precision	recall	f1-score	support
0	0.88	1.00	0.94	15
1	1.00	0.89	0.94	18
micro avg	0.94	0.94	0.94	33
macro avg	0.94	0.94	0.94	33
weighted avg	0.95	0.94	0.94	33

```
[[15  0]
 [ 2 16]]
```

### Reducer: svd Classifier: dtree

```
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=6,
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=40, n_iter=5,
random_state=None, tol=0.0), 'clf_dtree': DecisionTreeClassifier(class_weight=None,
↪ criterion='gini', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
```

```

min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best'))}

```

Accuracy 0.8787878787878788

	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: kbest Classifier: dtree

```

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=6,
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_kbest': SelectKBest(k=3, score_func=<function
↪ f_classif at 0x1a1b98f1e0>), 'clf_dtrees':
↪ DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best'))}

```

Accuracy 0.9696969696969697

	precision	recall	f1-score	support
0	1.00	0.93	0.97	15
1	0.95	1.00	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

```

[[14  1]
 [ 0 18]]

```

### Reducer: percentile Classifier: dtree

```

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=6,
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_percentile': SelectPercentile(percentile=10,

```

```
score_func=<function f_classif at 0x1a1b98f1e0>), 'clf_dtrees':
↳ DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best'))
```

Accuracy 0.9393939393939394

	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: None Classifier: dtrees

```
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=6,
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, 'clf_dtrees':
↳ DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best'))
```

Accuracy 0.9393939393939394

	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

```
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=6,
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=40, n_iter=5,
random_state=None, tol=0.0), 'mm_scaler': MinMaxScaler(copy=True, feature_range=(0,
↪ 1)), 'clf_nb': MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)}

```

Accuracy 0.9090909090909091

	precision	recall	f1-score	support
0	0.83	1.00	0.91	15
1	1.00	0.83	0.91	18
micro avg	0.91	0.91	0.91	33
macro avg	0.92	0.92	0.91	33
weighted avg	0.92	0.91	0.91	33

```

[[15 0]
 [ 3 15]]

```

### Reducer: kbest Classifier: nb

```

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=6,
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_kbest': SelectKBest(k=3, score_func=<function
↪ f_classif at 0x1a1b98f1e0>), 'mm_scaler': MaxAbsScaler(copy=True), 'clf_nb':
↪ MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)}

```

Accuracy 0.5454545454545454

	precision	recall	f1-score	support
0	0.50	1.00	0.67	15
1	1.00	0.17	0.29	18
micro avg	0.55	0.55	0.55	33
macro avg	0.75	0.58	0.48	33
weighted avg	0.77	0.55	0.46	33

```

[[15 0]
 [15 3]]

```

### Reducer: percentile Classifier: nb

```

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=6,
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_percentile': SelectPercentile(percentile=10,

```

```
score_func=<function f_classif at 0x1a1b98f1e0>), 'mm_scaler':
↳ MaxAbsScaler(copy=True), 'clf_nb': MultinomialNB(alpha=1.0, class_prior=None,
↳ fit_prior=True)}
```

Accuracy 0.9090909090909091

	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: nb

```
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=6,
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, 'clf_nb': MultinomialNB(alpha=1.0,
↳ class_prior=None, fit_prior=True)}
```

Accuracy 0.9090909090909091

	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

```
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=6,
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=40, n_iter=5,
random_state=None, tol=0.0), 'clf_lr': LogisticRegression(C=1.0, class_weight=None,
↳ dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='warn',
n_jobs=None, penalty='l2', random_state=None, solver='warn',
tol=0.0001, verbose=0, warm_start=False)}
```

Accuracy 0.9393939393939394

	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: kbest Classifier: lr

```
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=6,
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_kbest': SelectKBest(k=3, score_func=<function
↪ f_classif at 0x1a1b98f1e0>), 'clf_lr': LogisticRegression(C=1.0,
↪ class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='warn',
n_jobs=None, penalty='l2', random_state=None, solver='warn',
tol=0.0001, verbose=0, warm_start=False))}
```

Accuracy 0.7575757575757576

	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: percentile Classifier: lr

```
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=6,
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_percentile': SelectPercentile(percentile=10,
score_func=<function f_classif at 0x1a1b98f1e0>), 'clf_lr':
↪ LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='warn',
n_jobs=None, penalty='l2', random_state=None, solver='warn',
tol=0.0001, verbose=0, warm_start=False))}
```



Accuracy 0.8484848484848485

	precision	recall	f1-score	support
0	0.78	0.93	0.85	15
1	0.93	0.78	0.85	18
micro avg	0.85	0.85	0.85	33
macro avg	0.86	0.86	0.85	33
weighted avg	0.86	0.85	0.85	33

```
[[14  1]
 [ 4 14]]
```

### Reducer: None Classifier: lr

```
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=6,
    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, 'clf_lr': LogisticRegression(C=1.0,
    ↪ class_weight=None, dual=False, fit_intercept=True,
    intercept_scaling=1, max_iter=100, multi_class='warn',
    n_jobs=None, penalty='l2', random_state=None, solver='warn',
    tol=0.0001, verbose=0, warm_start=False))}
```

Accuracy 0.9393939393939394

	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: svd Classifier: svc

```
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=6,
    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=40, n_iter=5,
    random_state=None, tol=0.0), 'clf_svc': SVC(C=1.0, cache_size=200, class_weight=None,
    ↪ coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False))}
```

Accuracy 0.45454545454545453

	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: kbest Classifier: svc

```
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=6,
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_kbest': SelectKBest(k=3, score_func=<function
↪ f_classif at 0x1a1b98f1e0>), 'clf_svc': SVC(C=1.0, cache_size=200,
↪ class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
kernel='rbf', max_iter=-1, probability=False, random_state=None,
shrinking=True, tol=0.001, verbose=False)}
```

Accuracy 0.45454545454545453

	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: percentile Classifier: svc

```
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=6,
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_percentile': SelectPercentile(percentile=10,
score_func=<function f_classif at 0x1a1b98f1e0>), 'clf_svc': SVC(C=1.0,
↪ cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
kernel='rbf', max_iter=-1, probability=False, random_state=None,
shrinking=True, tol=0.001, verbose=False)}
```

Accuracy 0.45454545454545453

	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

```
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=6,
    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, 'clf_svc': SVC(C=1.0, cache_size=200,
    ↪ class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False)}
```

Accuracy 0.45454545454545453

	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: svd Classifier: lsvc

```
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=6,
    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=40, n_iter=5,
    random_state=None, tol=0.0), 'clf_lsvc': LinearSVC(C=1.0, class_weight=None,
    ↪ dual=True, fit_intercept=True,
    intercept_scaling=1, loss='squared_hinge', max_iter=1000,
    multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
    verbose=0)}
```

Accuracy 0.9393939393939394

	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: kbest Classifier: lsvc

```
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=6,
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_kbest': SelectKBest(k=3, score_func=<function
↪ f_classif at 0x1a1b98f1e0>), 'clf_lsvc': LinearSVC(C=1.0, class_weight=None,
↪ dual=True, fit_intercept=True,
intercept_scaling=1, loss='squared_hinge', max_iter=1000,
multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
verbose=0)}
```

Accuracy 0.8181818181818182

	precision	recall	f1-score	support
0	0.74	0.93	0.82	15
1	0.93	0.72	0.81	18
micro avg	0.82	0.82	0.82	33
macro avg	0.83	0.83	0.82	33
weighted avg	0.84	0.82	0.82	33

```
[[14  1]
 [ 5 13]]
```

### Reducer: percentile Classifier: lsvc

```
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=6,
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_percentile': SelectPercentile(percentile=10,
score_func=<function f_classif at 0x1a1b98f1e0>), 'clf_lsvc': LinearSVC(C=1.0,
↪ class_weight=None, dual=True, fit_intercept=True,
intercept_scaling=1, loss='squared_hinge', max_iter=1000,
multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
verbose=0)}
```

Accuracy 0.8181818181818182

	precision	recall	f1-score	support
0	0.76	0.87	0.81	15
1	0.88	0.78	0.82	18
micro avg	0.82	0.82	0.82	33
macro avg	0.82	0.82	0.82	33
weighted avg	0.82	0.82	0.82	33

[[13 2]

[ 4 14]]

### Reducer: None Classifier: lsvc

```
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=6,
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, 'clf_lsvc': LinearSVC(C=1.0, class_weight=None,
↪ dual=True, fit_intercept=True,
intercept_scaling=1, loss='squared_hinge', max_iter=1000,
multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
verbose=0)}
```

Accuracy 0.9090909090909091

	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]]

## 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 process_classifications_grid(X_train, y_train, X_test, y_test, parameters_grid, cv=2)

```

## Output

462

### Reducer: svd Classifier: knn

Best parameters

```

    clf_knn__n_neighbors: 2
    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: None
    vect__strip_accents: 'unicode'
    vect__sublinear_tf: True
    vect__token_pattern: '\\w{1,}'

```

Accuracy 0.8787878787878788

	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: kbest Classifier: knn

Best parameters

```
clf_knn__n_neighbors: 5
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: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

	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: 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: 1000
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 0.8181818181818182

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.74	0.93	0.82	15
1	0.93	0.72	0.81	18
micro avg	0.82	0.82	0.82	33
macro avg	0.83	0.83	0.82	33
weighted avg	0.84	0.82	0.82	33

```
[[14 1]
 [ 5 13]]
```

### Reducer: None Classifier: knn

Best parameters

```
clf_knn__n_neighbors: 2
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 1000
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 0.8484848484848485

	precision	recall	f1-score	support
0	0.78	0.93	0.85	15
1	0.93	0.78	0.85	18
micro avg	0.85	0.85	0.85	33
macro avg	0.86	0.86	0.85	33
weighted avg	0.86	0.85	0.85	33

```
[[14 1]
 [ 4 14]]
```

### Reducer: svd Classifier: dtree

Best parameters

```
red_svd__n_components: 2
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 1000
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

	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: 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: 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 0.9090909090909091

	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: 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: 1.0
```

```

vect__max_features: 1000
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

	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: None Classifier: dtree

Best parameters

```

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 0.9090909090909091

	precision	recall	f1-score	support
0	1.00	0.80	0.89	15
1	0.86	1.00	0.92	18
micro avg	0.91	0.91	0.91	33
macro avg	0.93	0.90	0.91	33
weighted avg	0.92	0.91	0.91	33

```

[[12  3]
 [ 0 18]]

```

### Reducer: svd Classifier: nb

Best parameters

```
red_svd__n_components: 40
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 0.9090909090909091

	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: 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: 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 0.5151515151515151

	precision	recall	f1-score	support
0	0.48	1.00	0.65	15
1	1.00	0.11	0.20	18
micro avg	0.52	0.52	0.52	33
macro avg	0.74	0.56	0.43	33

weighted avg	0.77	0.52	0.41	33
--------------	------	------	------	----

```
[[15 0]
 [16 2]]
```

### 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: 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 0.8787878787878788

	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: nb

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: '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

	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

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: None
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'

```

Accuracy 0.9090909090909091

	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: 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

	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: 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: 1000  
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

	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
```

Best parameters

```
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

	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: svd Classifier: svc

Best parameters

```

red_svd__n_components: 2
scaler: None
vect__analyzer: 'word'
vect__max_df: 1.0
vect__max_features: 1000
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

	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: 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

	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: 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: 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: 'english'
vect__strip_accents: 'unicode'
vect__sublinear_tf: True
vect__token_pattern: '\\w{1,}'
```

Accuracy 0.9393939393939394

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

Accuracy 0.9696969696969697

	precision	recall	f1-score	support
0	1.00	0.93	0.97	15
1	0.95	1.00	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

```
[[14  1]
 [ 0 18]]
```

```
### Reducer: svd Classifier: lsvc
```

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

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

```

Accuracy 0.8484848484848485

	precision	recall	f1-score	support
0	0.78	0.93	0.85	15
1	0.93	0.78	0.85	18
micro avg	0.85	0.85	0.85	33
macro avg	0.86	0.86	0.85	33
weighted avg	0.86	0.85	0.85	33

[[14 1]

[ 4 14]]

### 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: 'english'
    vect__strip_accents: 'unicode'
    vect__sublinear_tf: True
    vect__token_pattern: '\\w{1,}'

Accuracy 0.9090909090909091
      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: 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: '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
      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.

### 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,
40 → reducer="reducer", classifier="kmeans", cv=CV)
41 process_classifications_grid(X_train, y_train, X_test, y_test, parameters_grid,
42 → 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*\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

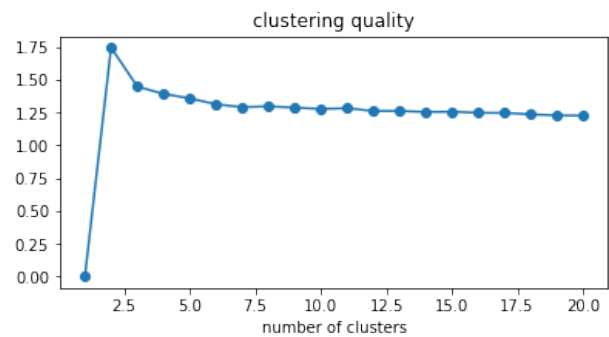
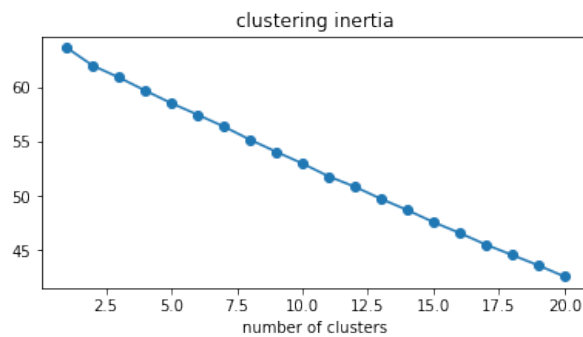
```

(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]]

```