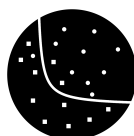


PROGETTO PER SISTEMI INTELLIGENTI

Prof. Nunzio Alberto Borghese
6 CFU

Luca Cappelletti

Project Documentation
Year 2017/2018



Magistrale Informatica
Università di Milano
Italy
10 novembre 2018

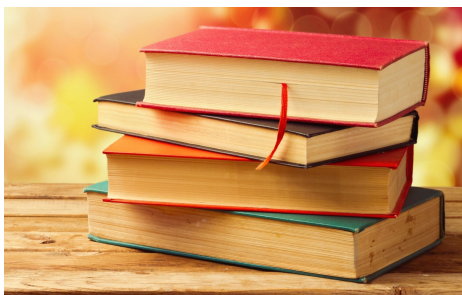
Indice

1	Introduction	3
1.1	Why is the Zipf interesting?	3
1.2	What is attempted in this project?	3
2	The Zipf	4
3	The classifier	5
3.1	Training the classifier	5
3.1.1	Convert documents to vector representation	5
3.1.2	Choosing representative points	5
3.1.3	Completing the classifier	6
3.2	Classifying a document	6
4	Results	7
4.1	Reproducing the results	7
4.2	What are the diagrams representing	7
4.2.1	Precision varying with neighbors number	7
4.2.2	ROC Curves	7
4.2.3	Confusion matrices	8
4.2.4	Truncated SVD reduction	8
4.2.5	Most defining word clouds	8
4.2.6	Representatives points usage	8
4.3	Authors	9
4.3.1	Precision varying with neighbors number	9
4.3.2	ROC Curves	9
4.3.3	Confusion matrices	9
4.3.4	Truncated SVD reduction	10
4.3.5	Most defining word clouds	10
4.3.6	Representatives points usage	10
4.4	Literary periods	11
4.4.1	Precision varying with neighbors number	11
4.4.2	ROC Curves	11
4.4.3	Confusion matrices	11
4.4.4	Truncated SVD reduction	12
4.4.5	Most defining word clouds	12
4.4.6	Representatives points usage	12
4.5	Recipes websites	13
4.5.1	Precision varying with neighbors number	13
4.5.2	ROC Curves	13
4.5.3	Confusion matrices	13
4.5.4	Truncated SVD reduction	14
4.5.5	Most defining word clouds	14
4.5.6	Representatives points usage	14
4.6	Newspaper websites	15
4.6.1	Precision varying with neighbors number	15
4.6.2	ROC Curves	15
4.6.3	Confusion matrices	15
4.6.4	Truncated SVD reduction	16
4.6.5	Most defining word clouds	16

4.6.6	Representatives points usage	16
4.7	Recipes websites or non recipes websites	17
4.7.1	Precision varying with neighbors number	17
4.7.2	ROC Curves	17
4.7.3	Confusion matrices	18
4.7.4	Truncated SVD reduction	18
4.7.5	Most defining word clouds	19
4.7.6	Representatives points usage	19
4.8	Nutritional values or non nutritional values	20
4.8.1	Precision varying with neighbors number	20
4.8.2	ROC Curves	20
4.8.3	Confusion matrices	21
4.8.4	Truncated SVD reduction	21
4.8.5	Most defining word clouds	22
4.8.6	Representatives points usage	22
5	Conclusions	23

1.1 Why is the Zipf interesting?

The Zipf law is an empirical law that refers to the fact that in nature, many data can be approximated by a Zipfian distribution, for example texts, some images¹, even sounds in spoken languages². It is therefore of interest to identify ways to exploit this relatively simple way to convert documents into representative vectors in problems such as classifications.



(a) Books follow a zipf distribution



(b) Image formats such as jpeg follows a zipf distribution



(c) Sounds in spoken language follow a zipf distribution

Figura 1.1: Examples of things in nature that follows the zipf distribution

1.2 What is attempted in this project?

In this project was attempted to classify textual documents using the bag of words model, ignoring any semantic structure, to their class using a general clustering approach with algorithms such as KMeans and CURE. No heuristics relative to the nature of the texts was used, exception made for the removal of frequent stop words, such as articles (the, a, etc...).

¹<https://www.dcs.warwick.ac.uk/bmvc2007/proceedings/CD-ROM/papers/paper-288.pdf>

²<https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0033993>

2

The Zipf

The zipf function z converts any document $d \subseteq \mathbb{D}^m$ comprised of elements into a representative vector $\underline{v} \subseteq [0, 1]^n$, $n \leq m$ based on the frequency of said elements in the given document. Taken in consideration the set of elements $d \subseteq \mathbb{D}$ that comprise the document and $d_{\neq} \subseteq d$ the set of distinct elements of the set d , for any given element $d_{\neq_i} \in d_{\neq}$, the value assigned by the zipf function is:

$$z(d_{\neq_i}) : \mathbb{D} \rightarrow \mathbb{R} = \frac{\#\{d_i \in d : d_{\neq_i} = d_i\}}{\#d}$$

3

The classifier

3.1 Training the classifier

3.1.1 Convert documents to vector representation

Given a set of training class-labeled elements T , we convert the documents, treated as *bag of elements*, to enumeration-sorted frequency vectors: the most common words in the language are dropped before elaborating the documents.

Given a document composed of $n = \#d$ elements $d = \{e_1, e_2, \dots, e_n\}$, first we define $d_{\neq} \subseteq d$ as the set of distinct elements in d . Then, we map to every distinct element its normalized cardinality in the set d :

$$z(e_i) = \frac{\#\{e_j \in d : e_j = e_i\}}{n}$$

Then, we proceed to map to a common enumeration the elements, using as full-set of elements the entire training set elements, so that $\forall e_i, e_j \in T, \exists ! i, j \in \mathbb{N} : i = j \Leftrightarrow e_i = e_j$.

3.1.2 Choosing representative points

Given an arbitrary percentage of points $p \in [0, 1]$ and an arbitrary distance in percentage $\alpha \in [0, 1]$, for each class of points $C_j \in T$ we choose using k -Means:

$$k = \lceil \#C_j \cdot p^2 \rceil$$

This way the centroids surely distribute among the different points, following their density. If the points are, in truth, a unique cluster, the centroids will distribute themselves in an uniform fashion.

Then, we select the most distant points in every cluster in a number equal to $\lfloor \#\{p \in C_j : p \in Q_i\} \cdot p \rfloor$. We move each point p of an amount proportional to the distance from the point to its centroid $c_i : \alpha \cdot L^2(p, c_i)$ towards its centroid.

Using a metric to choose k

For the high dimensionality and number of the vectors, iterating multiple times $KMeans$ searching for an optimal k following any given metric has an high time cost. For this reason, an attempt using **PCA** to reduce the dimensionality and predict the number of clusters in high dimensionality space using a density metric was made.

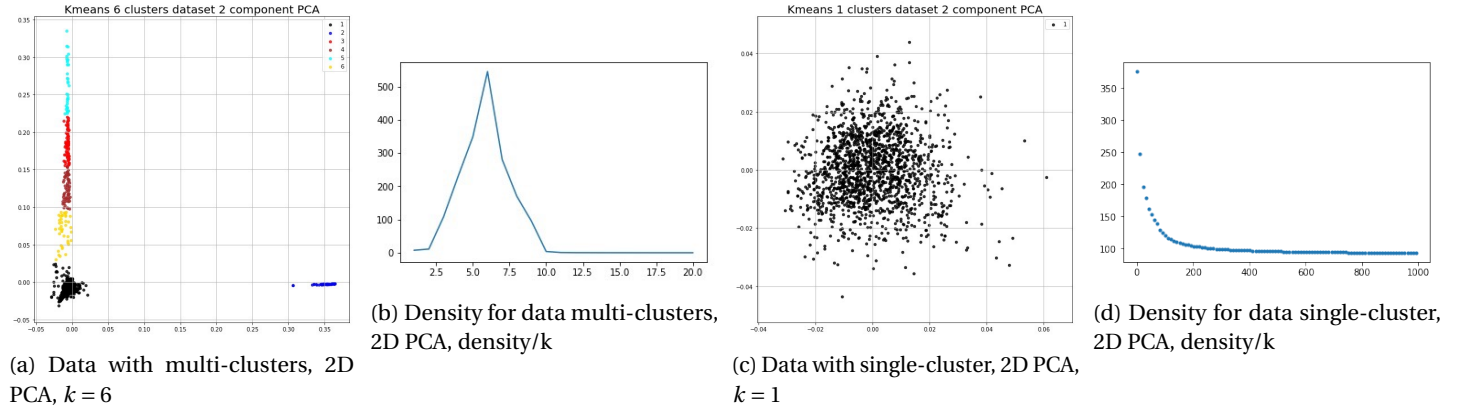
The density was defined as follows:

$$\bar{\rho}_{jk} = \frac{1}{k} \sum_{i=1}^k \rho_{jk_i} = \frac{1}{k} \sum_{i=1}^k \left(\frac{\#\{\underline{v} \in C_{r_j} : \underline{v} \in Q_i\}}{\#C_{r_j}} \right)^k \cdot \frac{1}{r_{Q_i}^2} \quad r_{Q_i}^2 = \begin{cases} \frac{1}{n} \sum_{h=1}^n (\underline{c}_i - \underline{p}_h)^2 & \text{if } n \neq 0 \\ 1 & \text{else} \end{cases}$$

Where r_{Q_i} is the approximated radius of the cluster Q_{k_i} , using the farthest n frontier points p_f .

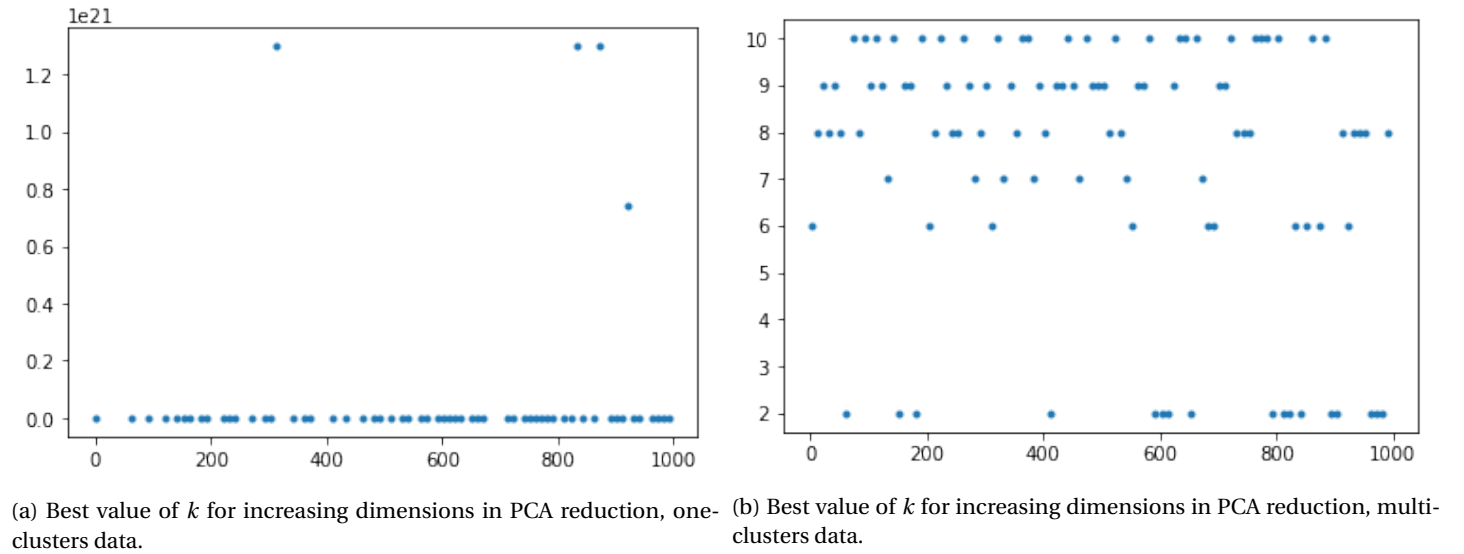
Does the prediction hold? While the metric on two dimensions reduction seemed to work, when iterated on increasingly larger number of dimensions, with the exception of strongly clustered data, it did not much better than a random selection.

As follows, with 2D PCA density metrics the heuristic seems to be successful to identify the number of clusters k necessary to describe the given classes.



The prediction on any given PCA reduction however is not valid for different dimension numbers, on multi-clusters classes.

Any number of clusters k high enough ($k > 5$) is no more precise than a density metric. With increasing number of clusters, as shown below, it becomes increasingly unreliable.



3.1.3 Completing the classifier

The classifier model is now finished, comprised of every class-labeled representative point.

3.2 Classifying a document

To classify a given a document d we proceed as follows:

1. Convert the document d to a zipf representative vector, using the common enumeration: $\underline{v} = z(d)$.
2. The document is then classified with the same label as the closest representative point in the classifier model.

4

Results

4.1 Reproducing the results

Test were run on multiple datasets, all containing non-structured textual documents. All tests were run via the [test.py](#) file, available in the github repository, with the following parameters:

Seed The random uniform generator used to split files, run Kmeans etc... were initialized to 1242.

Training percentage Of every dataset, 70% of the documents were used for training the classifier. The remaining 30% were used for testing. The split of the training and testing datasets was done selecting documents following an uniform random distribution initialized with the seed above.

Kmeans clusters k For every test, 10 clusters were used.

Kmeans iterations For every test 300 iterations were run before assuming convergence. Better results may be obtained with more iterations.

CURE representatives percentage For every test 10% of the training points were used as representative points.

CURE distance percentage For every test chosen representative points were moved 20% towards their respective centroids.

Stopwords Two stopwords list were used: an italian one and an english one. These are available on the repository of the project for reproducibility reason, but are easily available online.

The complete dataset is available at the following url:

<https://mega.nz/#!XCI3iCiZ!EpBtedozqBivPcWYi8oAj6jx38LAJnUaNDfjvgzJi3I>.

The already trained classifiers are available at the following url:

<https://mega.nz/#!yLAR0CSI!4g6jNfTeTpoEbss3P4jmuJyvTvv0Wmtw1sZQydQS1Qs>.

4.2 What are the diagrams representing

4.2.1 Precision varying with neighbors number

This graph shows the precision on growing number of neighbors, meaning that the classification is weighted on the given n number of neighbors and their distances from the considered point. On the proposed datasets, a low number of neighbors yielded best results. The other graphs shown are relative to the best results obtained exploring the neighbor space, but all results are available in the repository documentation. *Not all results were clearly possible to be included in this document as its size would increase a hundred fold.*

4.2.2 ROC Curves

The ROC curve, or receiver operating characteristic curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. For each class in every dataset a ROC curve is drawn to compare each class with every other one.

4.2.3 Confusion matrices

A confusion matrix, also known as an error matrix, is a specific table layout that allows visualization of the performance of an algorithm. Each row of the Matrix represents the instances in a predicted class while each column represents the instances in an actual class.

4.2.4 Truncated SVD reduction

A simple dimensionality reduction to visualize high dimensional datasets in 2D.

4.2.5 Most defining word clouds

The most important coordinates in the clusters used to classify texts for a given class.

4.2.6 Representatives points usage

An enumeration of counters of usages of clusters to classify texts for a given class.

4.3 Authors

The “authors” dataset contains 315 texts from 3 authors: D. H. Lawrence, Mark Twain and Oscar Wilde. Each author has 105 documents.

4.3.1 Precision varying with neighbors number

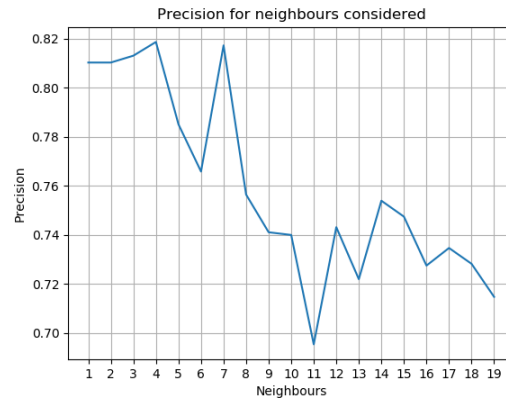


Figure 4.1: Precision scores

4.3.2 ROC Curves

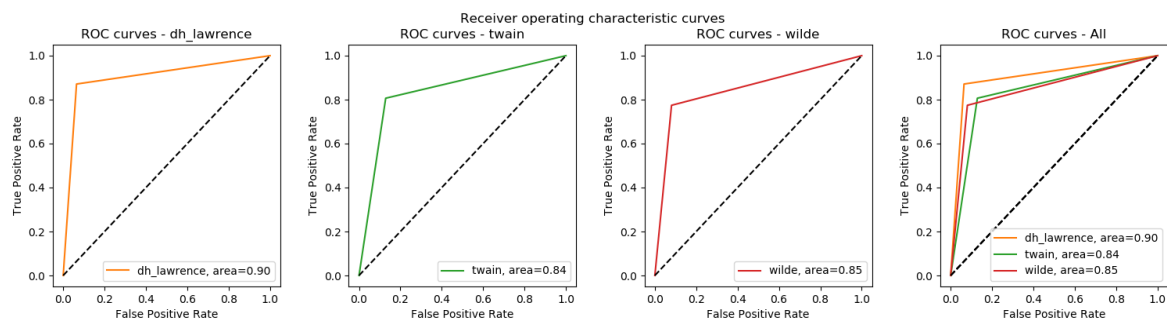


Figure 4.2: ROC Curves

4.3.3 Confusion matrices

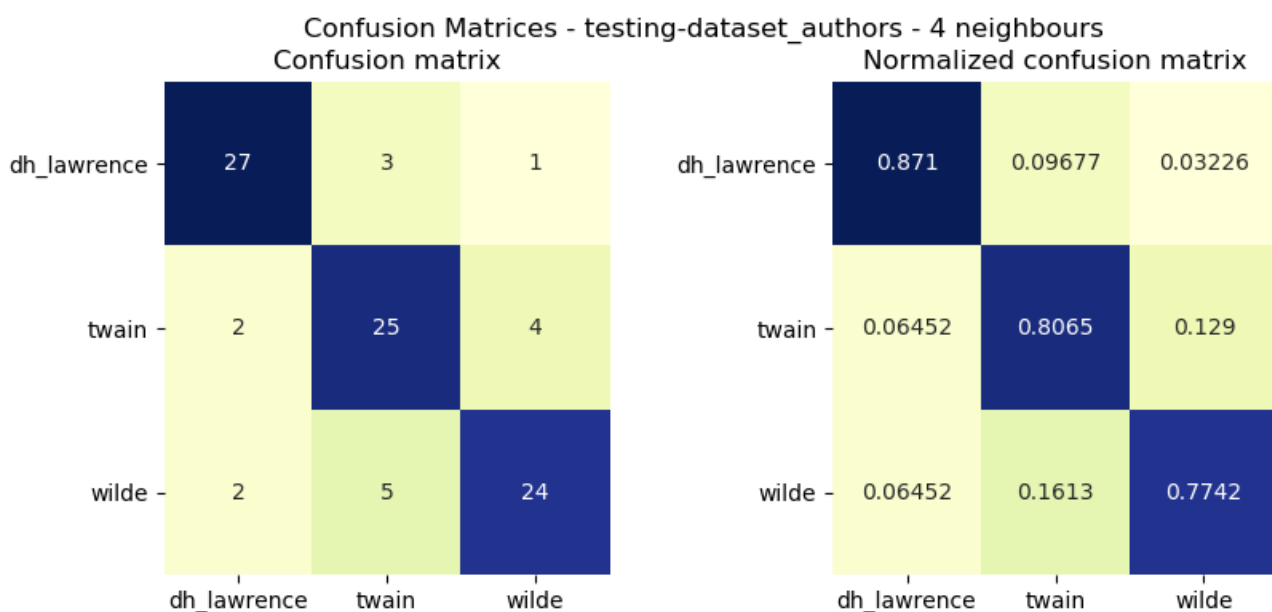


Figure 4.3: Confusion matrices for dataset “authors”

4.3.4 Truncated SVD reduction

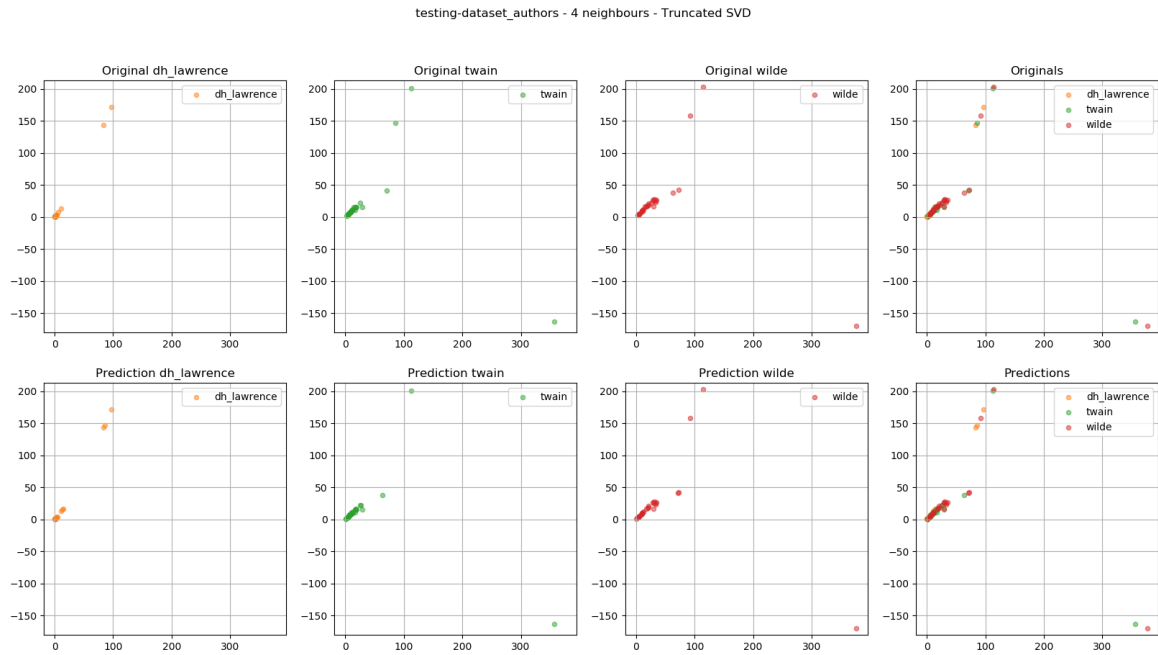


Figura 4.4: Dimensionality reduction using truncated SVD in dataset “authors”

4.3.5 Most defining word clouds



Figura 4.5: Word clouds

4.3.6 Representatives points usage

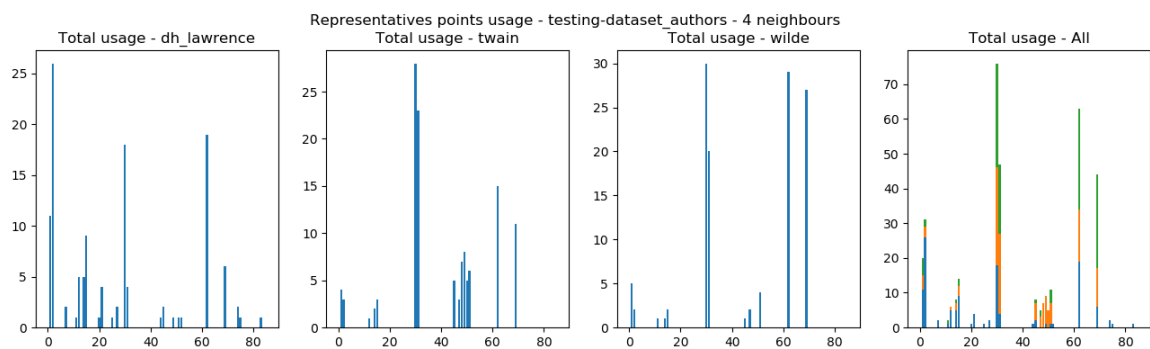


Figura 4.6: Representatives points usage

4.4 Literary periods

The “literary periods” dataset contains 1060 texts from 4 literary periods: Modernism, Realism, Romanticism and Naturalism. Each period has 265 documents.

4.4.1 Precision varying with neighbors number

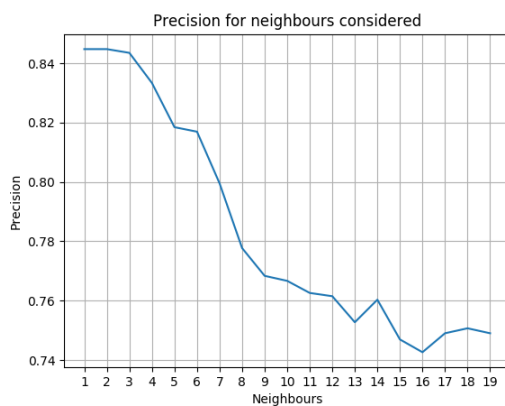


Figure 4.7: Precision scores

4.4.2 ROC Curves

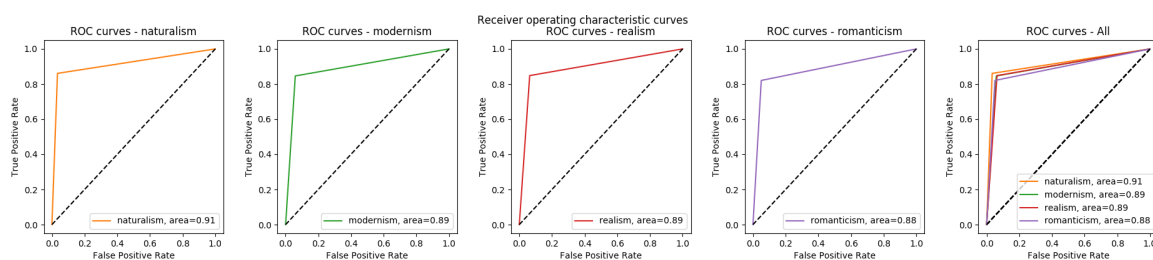


Figure 4.8: ROC Curves

4.4.3 Confusion matrices

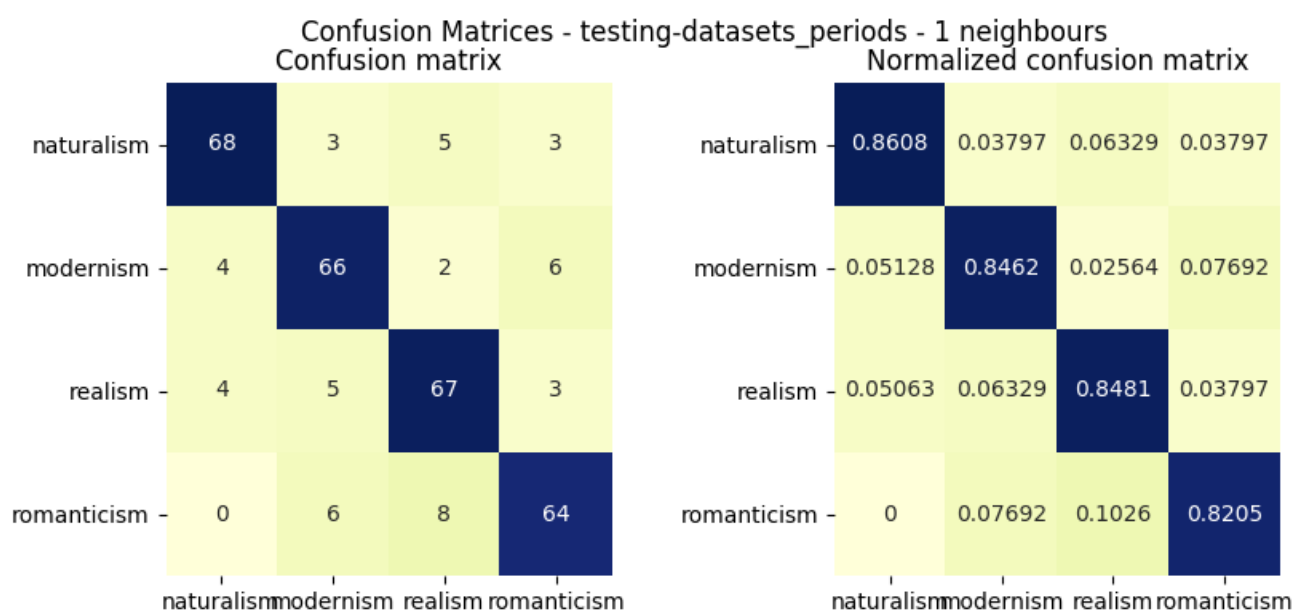


Figure 4.9: Confusion matrices for dataset “periods”

4.4.4 Truncated SVD reduction

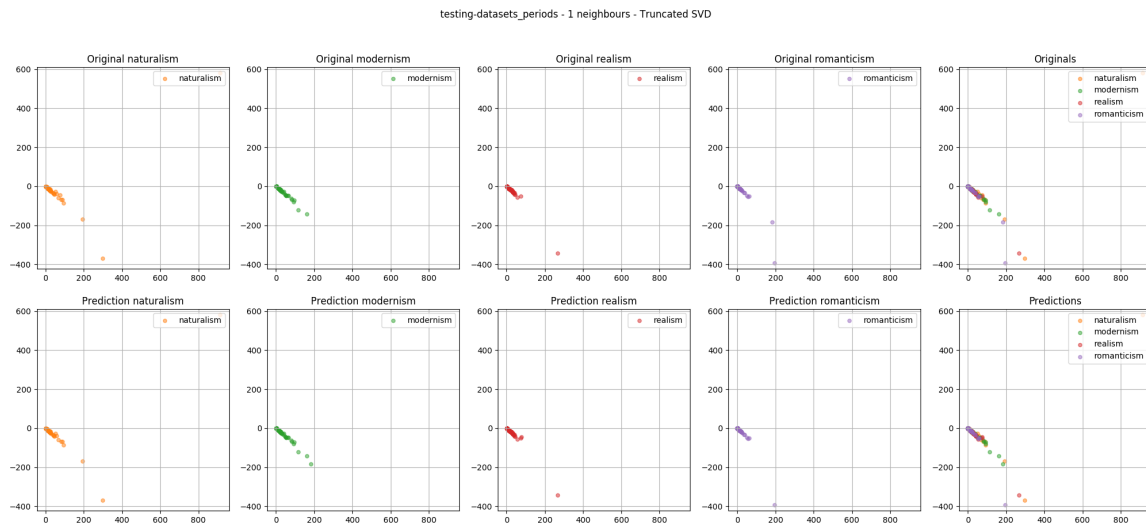


Figura 4.10: Dimensionality reduction using truncated SVD in dataset “periods”

4.4.5 Most defining word clouds

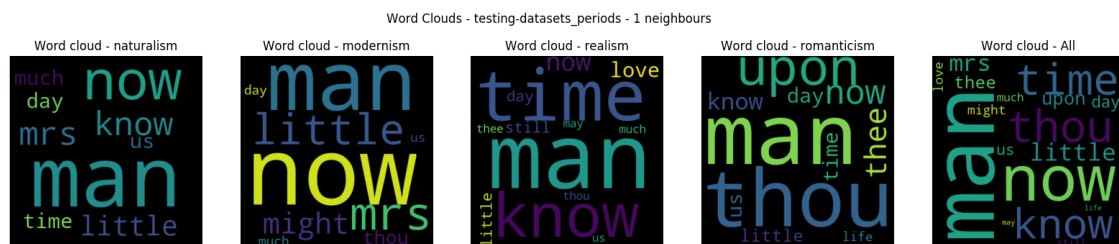


Figura 4.11: Word clouds

4.4.6 Representatives points usage

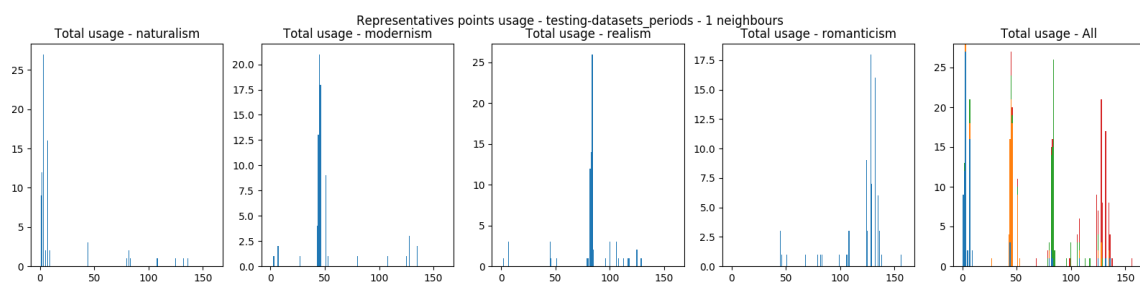


Figura 4.12: Representatives points usage

4.5 Recipes websites

The “recipes websites” dataset contains 12768 texts from 4 Italian websites dedicated to recipes: Misya, Salepepe, Zafferano and Lacucinaitaliana. Every website in the dataset has 3192 texts.

4.5.1 Precision varying with neighbors number

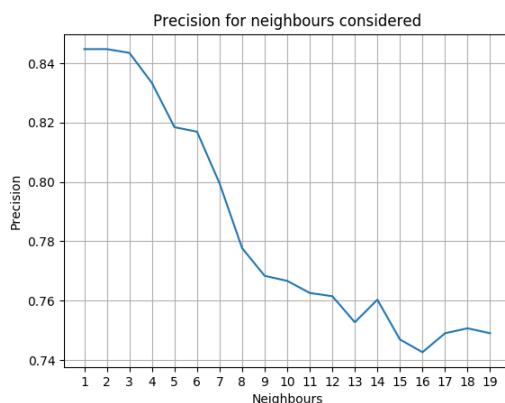


Figura 4.13: Precision scores

4.5.2 ROC Curves

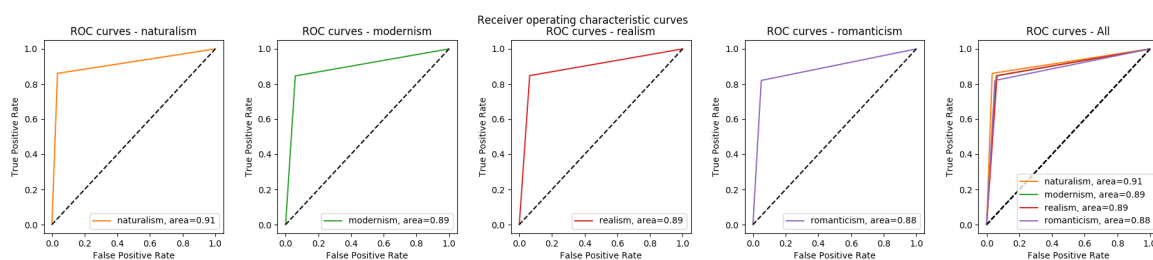


Figura 4.14: ROC CURves

4.5.3 Confusion matrices

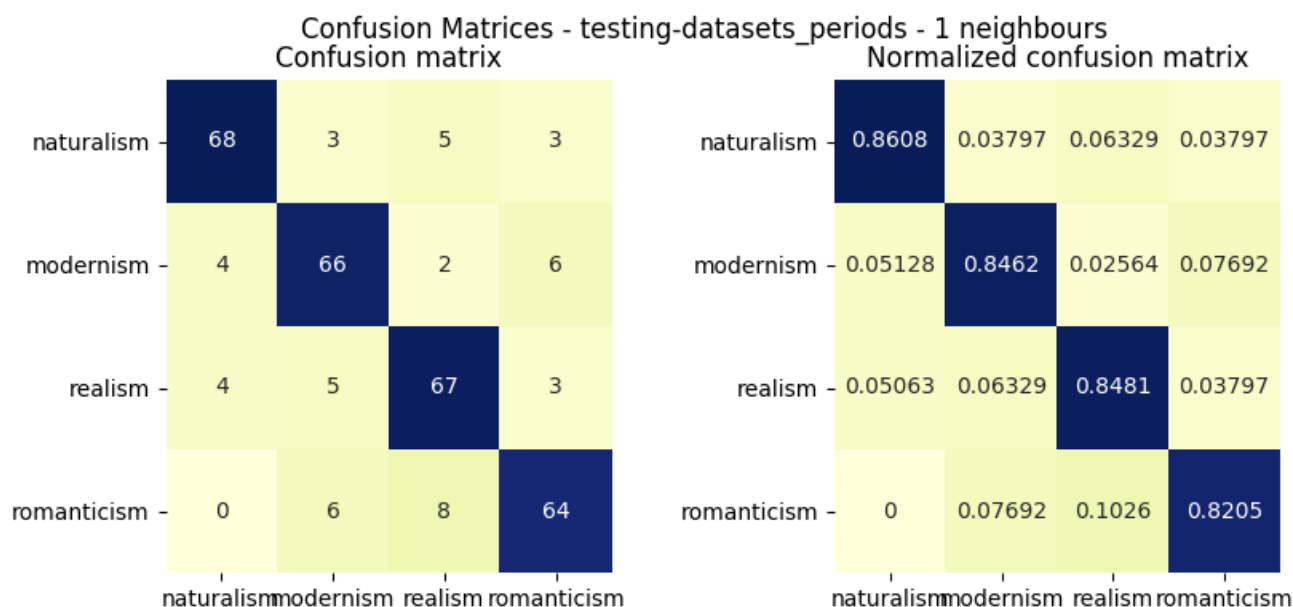


Figura 4.15: Confusion matrices for dataset “periods”

4.5.4 Truncated SVD reduction

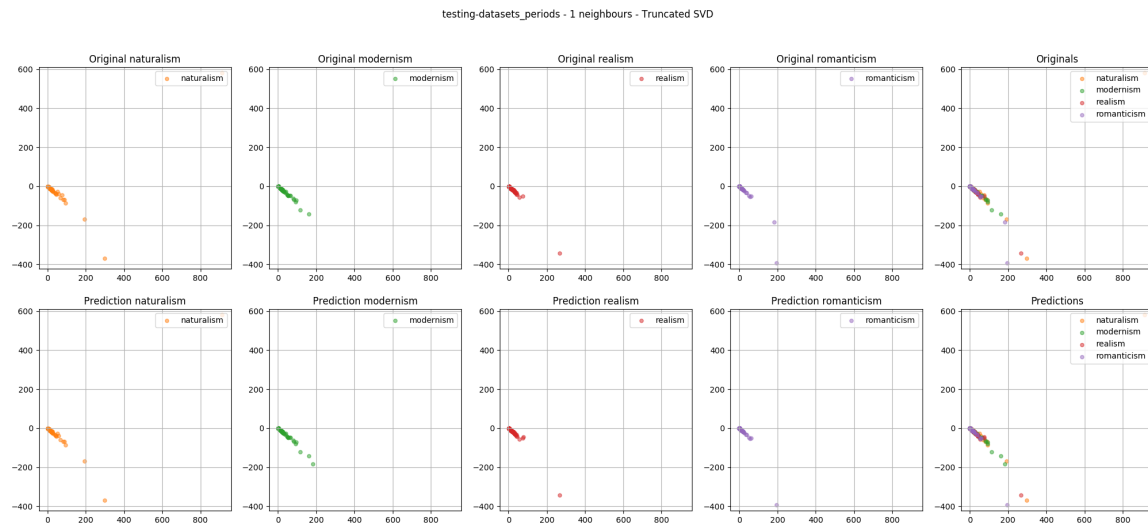


Figura 4.16: Dimensionality reduction using truncated SVD in dataset “periods”

4.5.5 Most defining word clouds

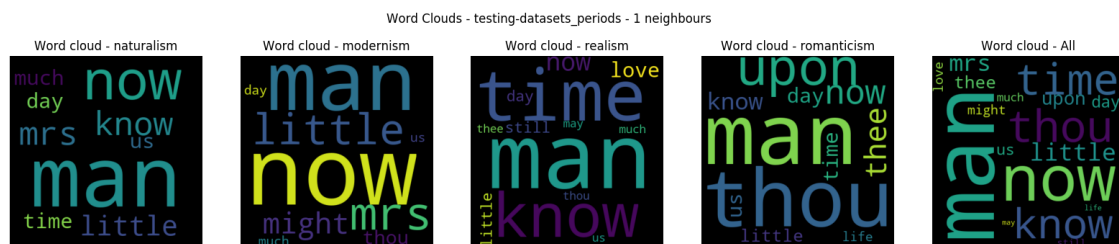


Figura 4.17: Word clouds

4.5.6 Representatives points usage

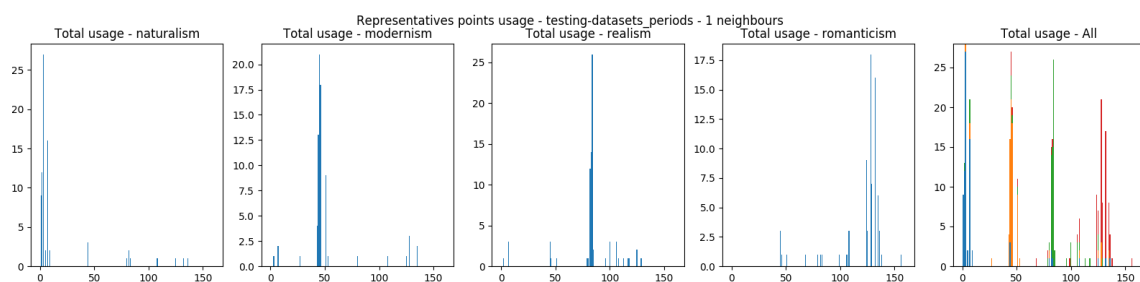


Figura 4.18: Representatives points usage

4.6 Newspaper websites

The “newspaper websites” dataset contains 28000 texts from 4 italian news websites: Repubblica.it, Moto.it, Scienze.it and Gazzetta.it. Every website in the dataset has 7000 texts.

4.6.1 Precision varying with neighbors number

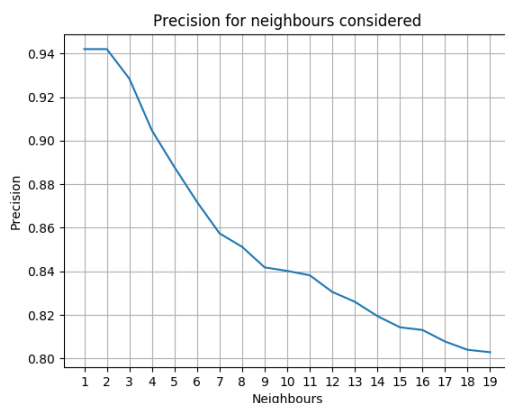


Figura 4.19: Precision scores

4.6.2 ROC Curves

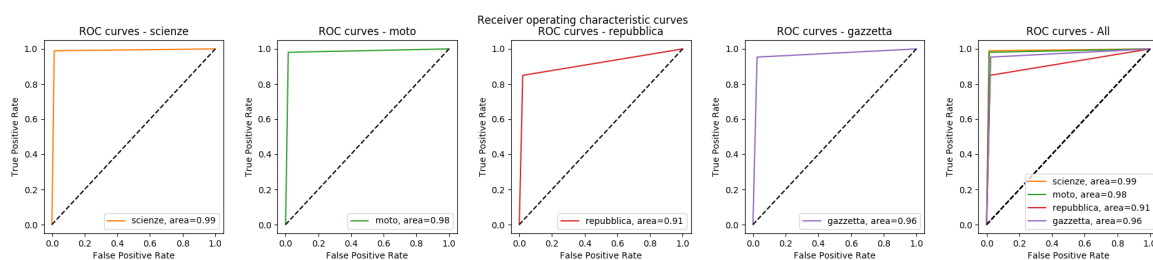


Figura 4.20: ROC CURves

4.6.3 Confusion matrices

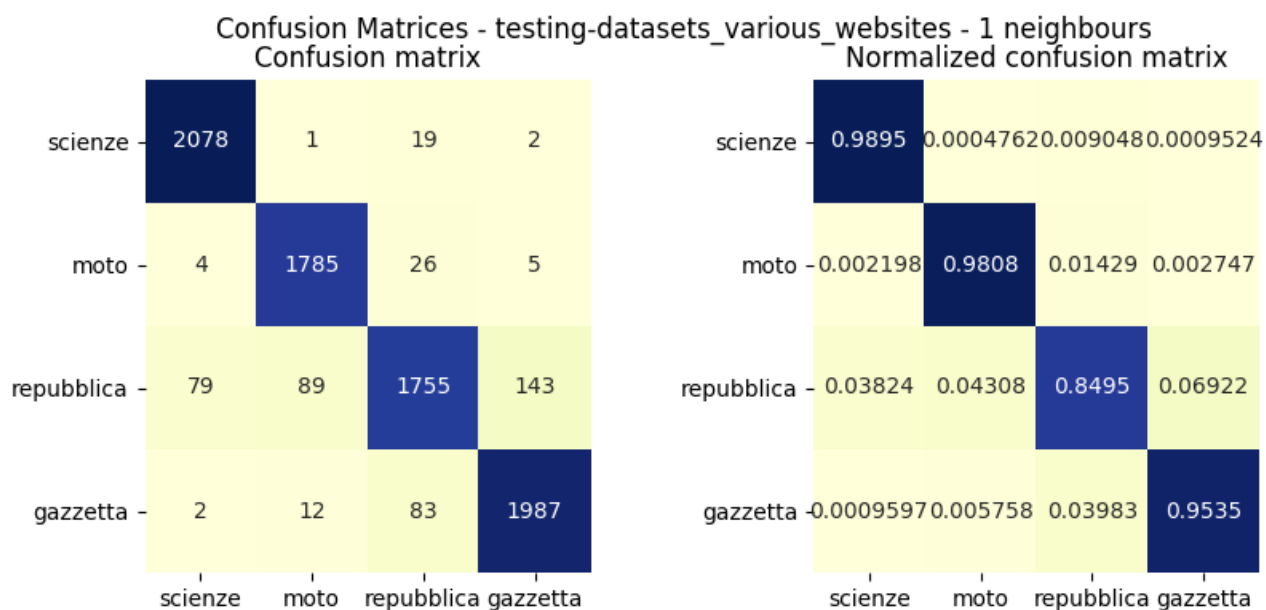


Figura 4.21: Confusion matrices for dataset “newspaper websites”

4.6.4 Truncated SVD reduction

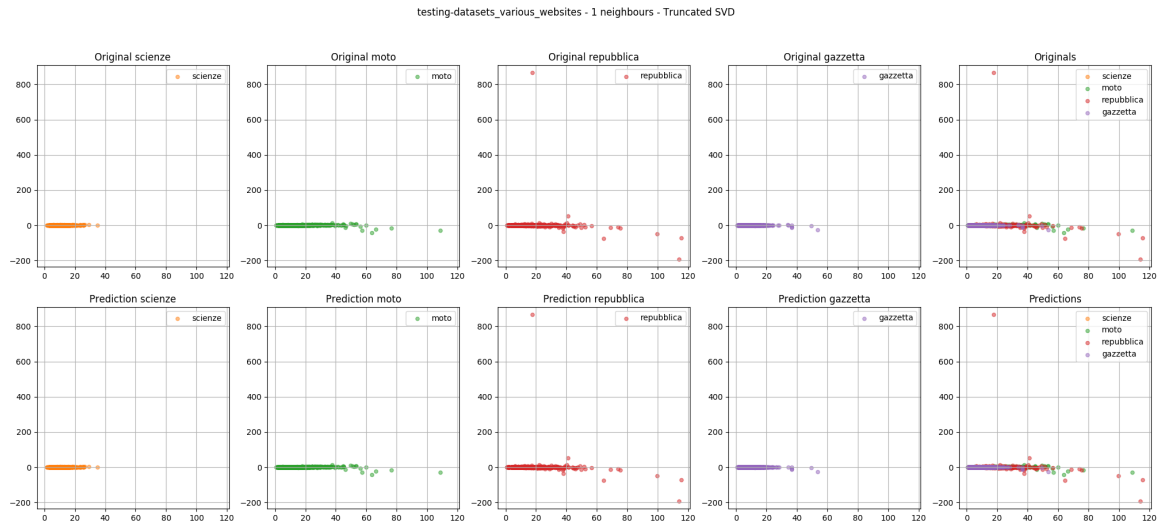


Figura 4.22: Dimensionality reduction using truncated SVD in dataset “newspaper websites”

4.6.5 Most defining word clouds

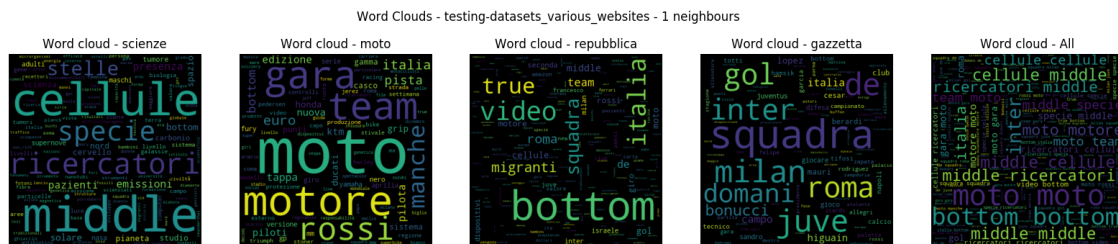


Figura 4.23: Word clouds

4.6.6 Representatives points usage

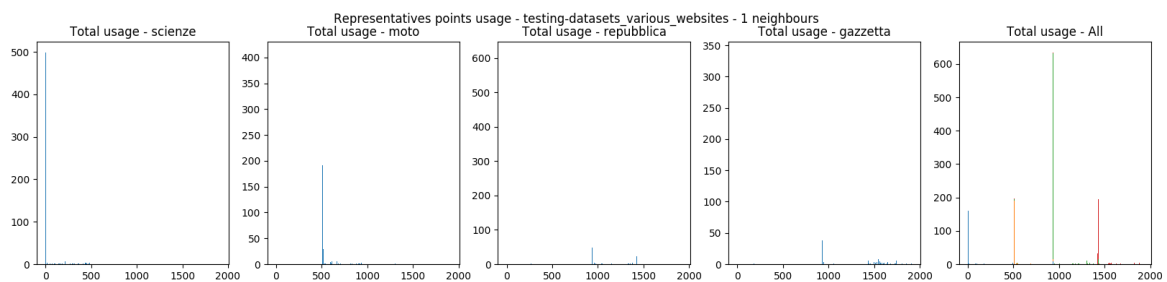


Figura 4.24: Representatives points usage

4.7 Recipes websites or non recipes websites

This dataset contains 36000 texts, with two classes: recipes and non-recipes. Each class has 18000 texts.

4.7.1 Precision varying with neighbors number

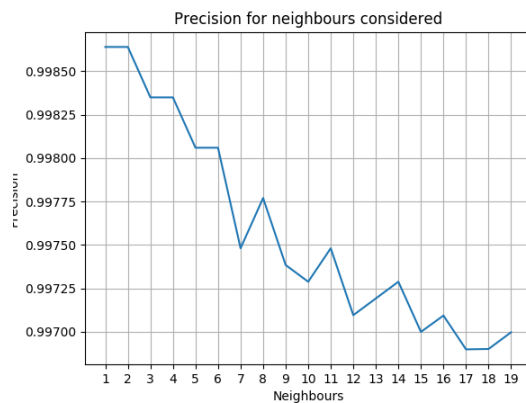


Figura 4.25: Precision scores

4.7.2 ROC Curves

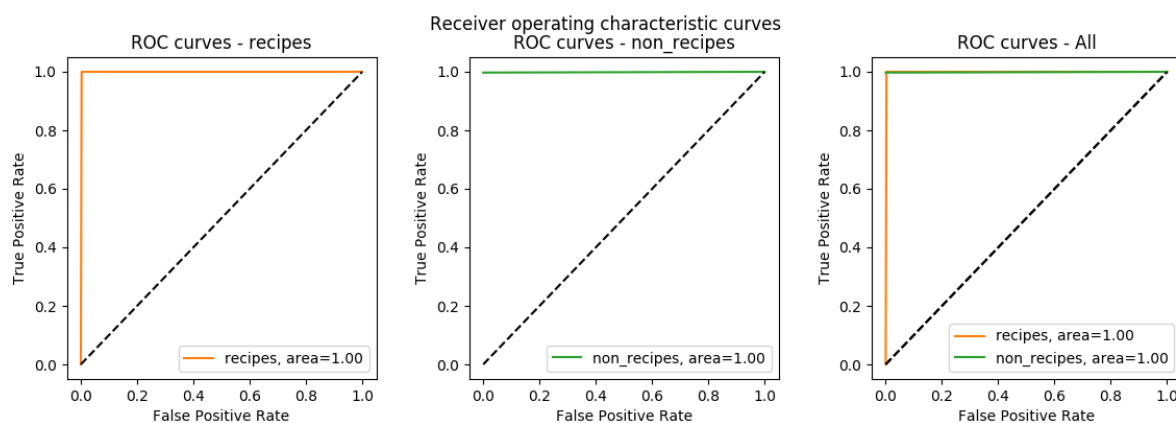


Figura 4.26: ROC CURves

4.7.3 Confusion matrices

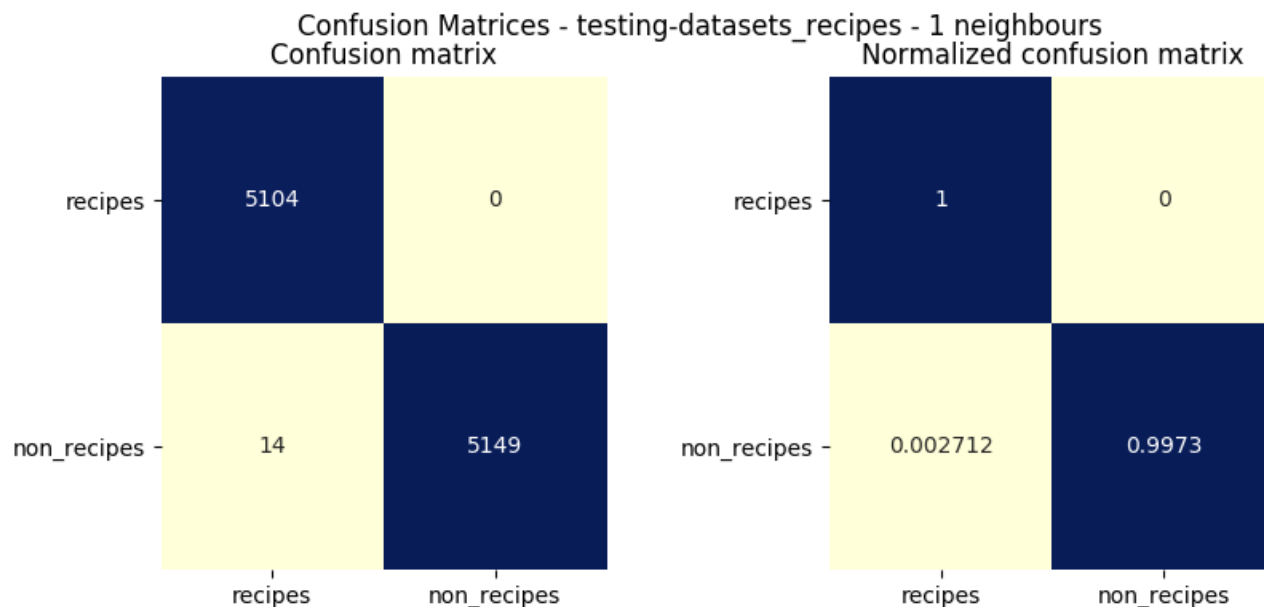


Figura 4.27: Confusion matrices for dataset “recipes”

4.7.4 Truncated SVD reduction

testing-datasets_recipes - 1 neighbours - Truncated SVD

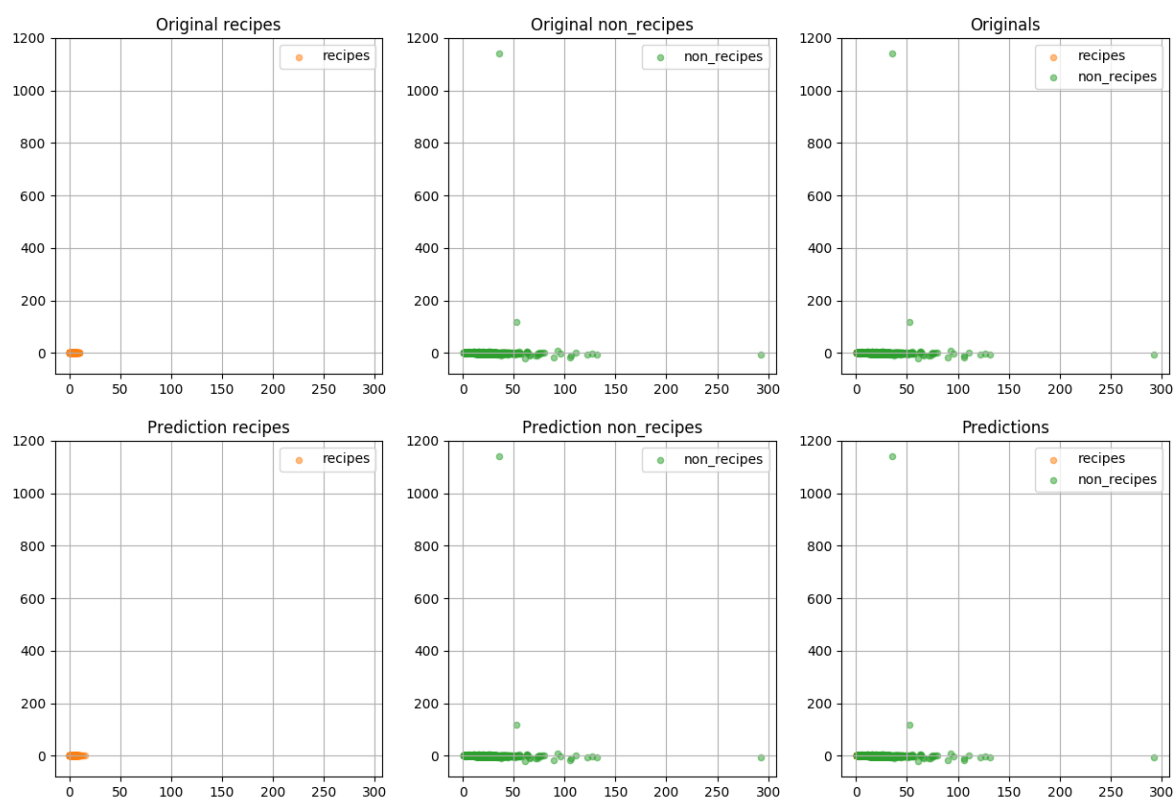


Figura 4.28: Dimensionality reduction using truncated SVD in dataset “recipes”

4.7.5 Most defining word clouds

Word Clouds - testing-datasets_recipes - 1 neighbours

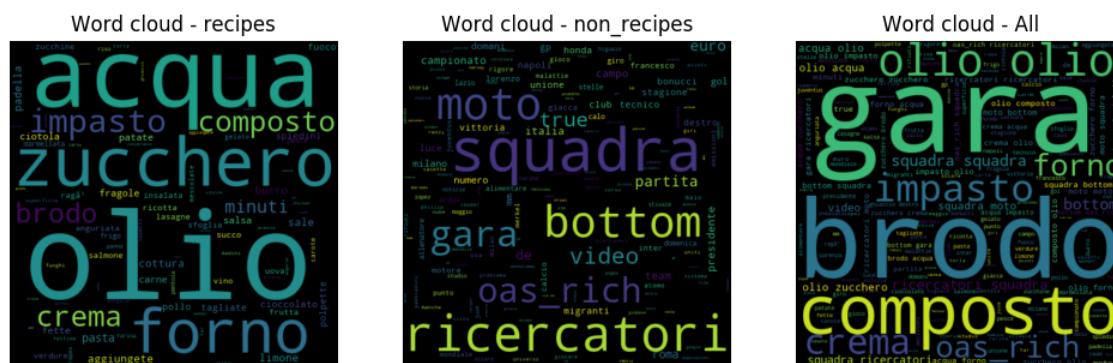


Figura 4.29: Word clouds

4.7.6 Representatives points usage

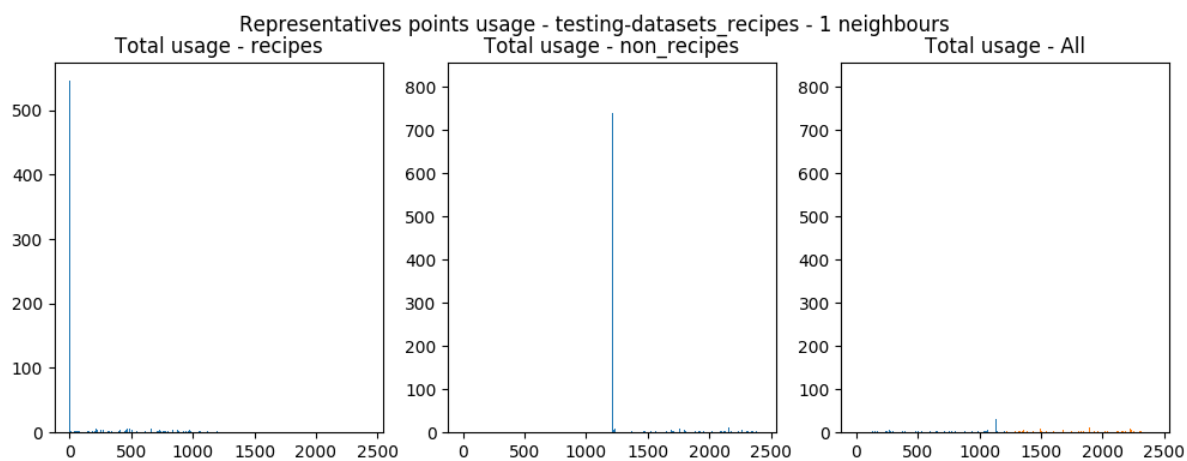


Figura 4.30: Representatives points usage

4.8 Nutritional values or non nutritional values

This dataset contains 51029 texts, with two classes: nutritional values and non-nutritional values. The first one has 5050 and the second one about ten times as many: 45979.

4.8.1 Precision varying with neighbors number

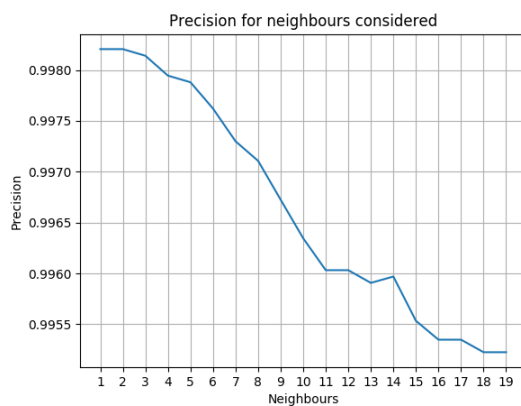


Figura 4.31: Precision scores

4.8.2 ROC Curves

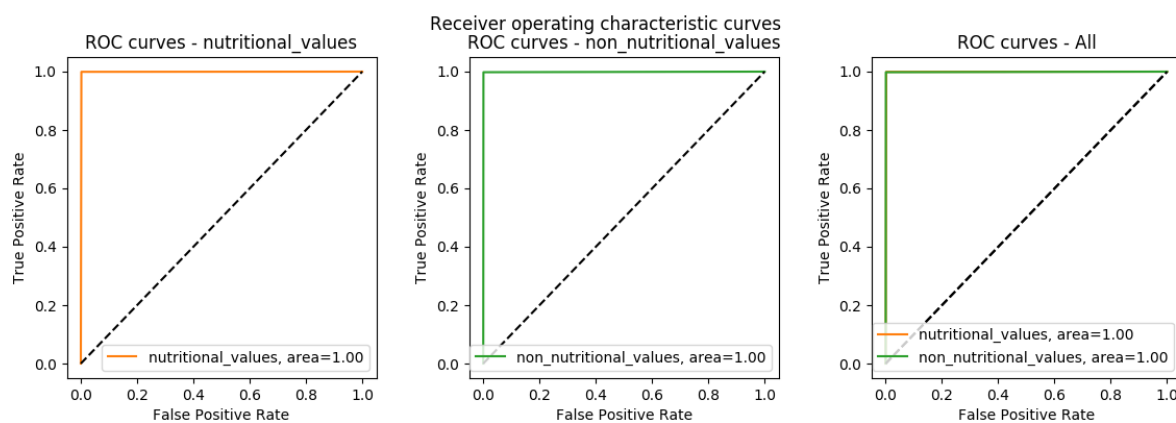


Figura 4.32: ROC CURves

4.8.3 Confusion matrices

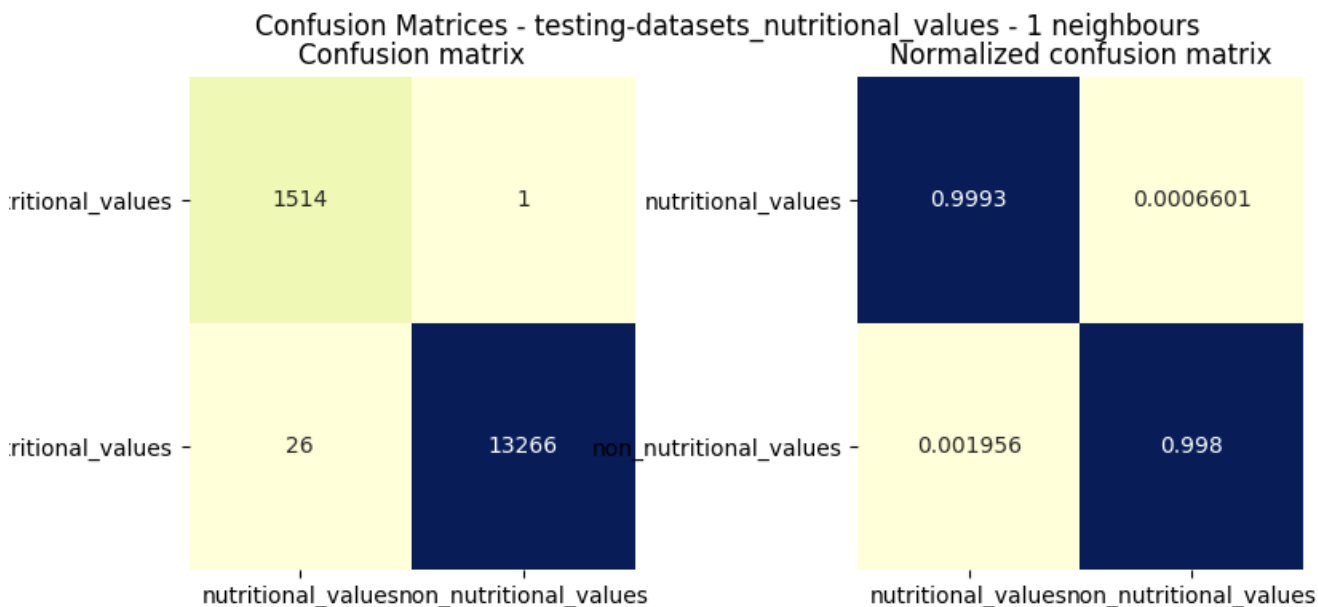


Figura 4.33: Confusion matrices for dataset “nutritional values”

4.8.4 Truncated SVD reduction

testing-datasets_nutritional_values - 1 neighbours - Truncated SVD

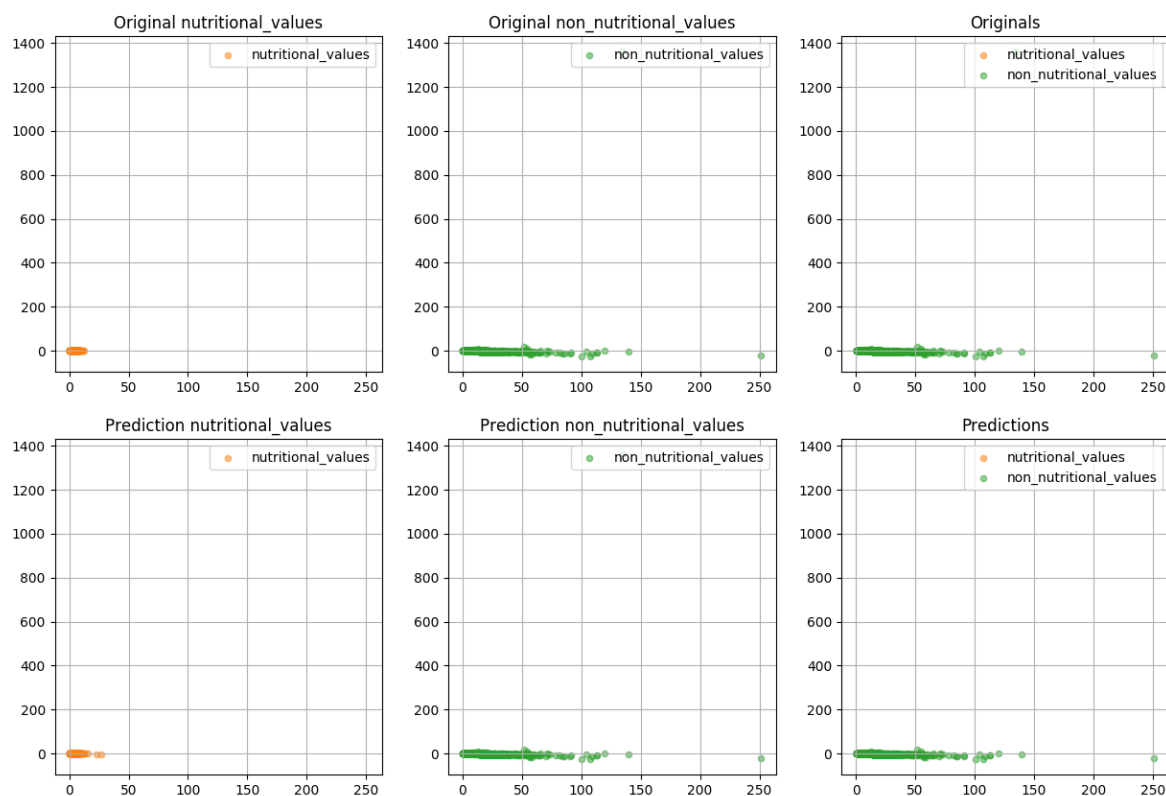


Figura 4.34: Dimensionality reduction using truncated SVD in dataset “nutritional values”

5

Conclusions

This approach yields arguably good classifiers with both big and small datasets, with equal or unbalanced classes.