

INF582: Introduction to Text Mining and NLP

Lab session 3: Supervised Document Classification

Lecture: Prof. Michalis Vazirgiannis
Lab: Antoine Tixier and Kostas Skianis

January 25, 2019

1 Introduction

Text categorization has a variety of real-world applications such as sentiment analysis and opinion mining, email filtering, etc. It is also tightly connected to the field of information retrieval (i.e., search engines).

The most common pipeline for text classification is the vector space representation (a.k.a. “bag-of-words”) followed by TF-IDF term weighting. With this approach, each document d_i from the collection $D = \{d_1 \dots d_m\}$ of size m is viewed as an n -dimensional feature vector, where n is the number of unique terms in all processed documents. The set of unique terms $T = \{t_1 \dots t_n\}$ is called the vocabulary.

A classifier is then trained on the document-term matrix of dimension m by n , and is used for classifying new documents. In this lab session, we will compare how using different weighting schemes to populate the feature vectors and different algorithms to classify documents impacts performance.

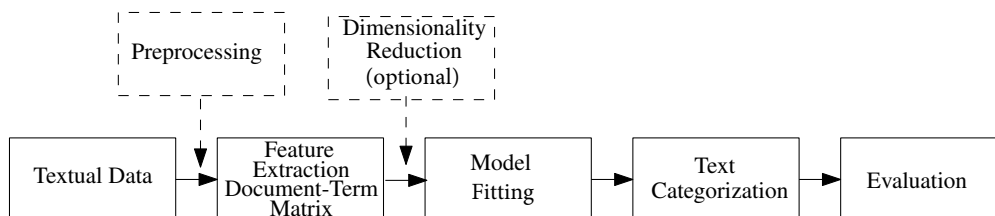


Figure 1: Basic pipeline of the Text Categorization task.

2 Term Frequency-Inverse Document Frequency

As we saw during the first lab, the coordinates of document d in the “bag-of-words” space (made of the $\{t_1 \dots t_n\}$ features) are computed as:

$$\text{weight}(t, d, D) = tf(t, d) \times idf(t, D) \quad (1)$$

Where $tf(t, d)$ is the number of times term t appears in document d , and $idf(t, D) = \ln \left(\frac{m}{1+df(t)} \right)$, with $df(t)$ the number of documents in the collection D that contain t and m the size of D . The intuition behind this scoring function is that frequent words in a document are representative of

that document as long as they are not also very frequent at the corpus level.

Using the raw term frequency $tf(t, d)$ favors long documents, which is not optimal because:

- long documents contain many repetitions of the same terms without necessarily containing more information,
- long documents contain more unique terms, which may artificially increase their similarity with any other document.

In an attempt to fix the aforementioned issues, researchers have normalized the $tf(t, d)$ term by composing (1) a monotonically increasing concave function with (2) a document length penalization function. One of the best known formulas was introduced by Singhal et al. [1996]¹ in the context of information retrieval, and is known as the *pivoted normalization weighting*:

$$tf_p(t, d) = \frac{1 + \ln(1 + \ln(tf(t, d)))}{1 - b + b \times \frac{|d|}{avgdl}} \quad (2)$$

Where $|d|$ is the length of document d (i.e., number of words in d), $avgdl$ is the average document length across the corpus, and $b \in [0, 1]$ is the slope parameter. The use of a strictly increasing concave function at the numerator repeatedly diminishes the gain of seeing an additional occurrence of a term in a document. E.g., it ensures that the gain of increasing tf from 1 to 2 is much larger than the gain of increasing tf from 1000 to 1001. In the denominator, b controls how much documents should be penalized based on their size. For $b = 0$, there is no normalization, while for $b = 1$, it is maximum. A standard default value is $b = 0.2$. The final formula is given by:

$$tf_p idf(t, d, D) = \frac{1 + \ln(1 + \ln(tf(t, d)))}{1 - b + b \times \frac{|d|}{avgdl}} \times \ln\left(\frac{m}{1 + df(t)}\right) \quad (3)$$

3 Term Weighting-Inverse Document Frequency

We capitalize on the graph-of-words representation introduced during the previous lab session to derive a new scoring function, TW-IDF (Rousseau and Vazirgiannis [2013]), where node centrality metrics are used to replace the numerator of $tf_p(t, d)$, and $b = 0.003$.

$$tw_p-idf(t, d) = \frac{tw(t, G_d)}{1 - b + b \times \frac{|d|}{avgdl}} \times idf(t, D)$$

Where G_d is the graph-of-words representation of document d . The absence of concave normalization in the numerator and the very low value of b are justified in the case of graphs-of-words with unweighted edges. Indeed, in that case, edge weights (and thus node metrics) do not increase linearly with the size of the document like raw term frequency. They only increase when new edges are added, which corresponds to new contexts of co-occurrence and depends much less on document size.

Some basic centrality metrics that can intuitively be used to provide the term weights $tw(t, G_d)$ are the:

¹Amit Singhal was the head of Google core search team from 2000 to 2016.

- *Normalized Degree.* The degree is a local metric equal to the number of incident edges of a node. The normalized degree can be computed as follows:

$$\text{degree}_G(v) = \frac{|\mathcal{N}(v)|}{|G| - 1}$$

where $\mathcal{N}(v)$ is the number of connections of node v in graph G and $|G|$ is the number of nodes in G . In its weighted version, the degree of a node is equal to the sum of the weights of its incident edges. It is also possible to consider only the incoming edges of the node, or its outgoing edges, giving respectively the in-degree and out-degree centrality measures.

- *Closeness centrality.* It measures how close a node is to all the other nodes in the graph. It is formally defined as the inverse of the average shortest path distance from the node to any other node in the graph:

$$\text{closeness}_G(v) = \frac{|G| - 1}{\sum_{v' \in V} \text{dist}(v, v')}$$

Where $\text{dist}(v, v')$ is the shortest path distance between nodes v and v' in G . Unlike degree centrality, closeness is a global metric, in that it combines information from all nodes in the graph.

3.1 Term Weighting-Inverse Collection Weighting (TW-ICW)

Instead of multiplying our graph-based component tw_p by the same term $\text{idf}(t, D) = \log \left(\frac{m+1}{\text{df}(t)} \right)$, we can try to compute an IDF at the graph level. The goal is still the same: penalizing the terms that occur frequently at the collection level. To this purpose, we define the collection level graph G_D as the union² of all the graphs in the collection, i.e. $G_D = G_{d_1} \cup G_{d_2} \cup \dots \cup G_{d_m}$. Note that instead of combining graphs constructed independently for each document, we can also obtain G_D by sliding a window over the entire collection while making sure that the window does not overspan documents. The inverse collection weight of a term t is then defined as:

$$\text{icw}(t, G_D) = \log \left(\frac{\max_{v \in G_D} tw(v, G_D) + 1}{tw(t, G_D)} \right)$$

Which gives the final TW-ICW measure as:

$$tw_p - \text{icw}(t, d) = \frac{tw(t, G_d)}{1 - b + b \times \frac{|d|}{\text{avgdl}}} \times \text{icw}(t, G_D)$$

TW-ICW, along with some other variants, were introduced by Skianis et al. [2018].

4 Coding time

We will work with the WebKB dataset. It features academic webpages belonging to four different categories: (1) project, (2) course, (3) faculty, and (4) students, and contains 2,803 documents for training and 1,396 for testing. Documents have already been preprocessed with stopword removal and Porter stemming. The code that you will work with implements the following steps:

²the union of two graphs G_1 and G_2 is defined as the union of their node and edge sets, i.e., $G_1 \cup G_2 = (V_{G_1} \cup V_{G_2}, E_{G_1} \cup E_{G_2})$

- data loading,
- computation of features (TF-IDF, TW-IDF and TW-ICW) for the training set,
- computation of features for the test set. Note that the documents in the test set are represented in the space made of the unique terms in the training set *only* (words in the testing set absent from the training set are disregarded). Similarly, the average document length, the inverse document frequency, etc. should be the ones computed on the training set.
- classifier training/testing. Naive Bayes classifier (McCallum and Nigam [1998]), Multinomial Logistic Regression (Collins et al. [2002]) and linear kernel SVM (Joachims [1998]) are compared,
- get the most/least important features for each class.

Each of the steps above is implemented either entirely or partially in the `document_classification.py` script. The custom functions used can be found in the `library.py` file. As usual, **fill the gaps wherever needed**. You have 3 gaps to fill in `library.py` and 8 in `document_classification.py`. A list of `igraph` methods can be found here: <http://igraph.org/python/doc/igraph.Graph-class.html>.

Questions³: Compare performance for different features and the same classifier, and for the same features but different classifiers. What are the best features/classifiers? Should we invest time on feature extraction or model selection?

References

- Michael Collins, Robert E. Schapire, and Yoram Singer. Logistic regression, adaboost and breiman distances. *Mach. Learn.*, 48(1-3):253–285, September 2002. ISSN 0885-6125. doi: 10.1023/A:1013912006537. URL <http://dx.doi.org/10.1023/A:1013912006537>.
- Thorsten Joachims. Text categorization with support vector machines: Learning with many relevant features. In *Proceedings of the 10th European Conference on Machine Learning*, ECML '98, pages 137–142, London, UK, UK, 1998. Springer-Verlag. ISBN 3-540-64417-2. URL <http://dl.acm.org/citation.cfm?id=645326.649721>.
- Andrew McCallum and Kamal Nigam. A comparison of event models for naive bayes text classification, 1998.
- François Rousseau and Michalis Vazirgiannis. Graph-of-word and tw-idf: New approach to ad hoc ir. In *Proceedings of the 22Nd ACM International Conference on Information & Knowledge Management*, CIKM '13, pages 59–68, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2263-8. doi: 10.1145/2505515.2505671. URL <http://doi.acm.org/10.1145/2505515.2505671>.
- Amit Singhal, Chris Buckley, and Mandar Mitra. Pivoted document length normalization. In *Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '96, pages 21–29, New York, NY, USA, 1996. ACM. ISBN 0-89791-792-8. doi: 10.1145/243199.243206. URL <http://doi.acm.org/10.1145/243199.243206>.

³in-class only. You just have to submit your code

Konstantinos Skianis, Fragkiskos Malliaros, and Michalis Vazirgiannis. Fusing document, collection and label graph-based representations with word embeddings for text classification. In *Proceedings of the Twelfth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-12)*, pages 49–58, 2018.