

# A Supervised Term-Weighting Method and its Application to Variable Extraction from Digital Media

Mariano Maisonnave<sup>†</sup>, Fernando Delbianco<sup>‡</sup>, Fernando Tohmé<sup>‡</sup>, and Ana Maguitman<sup>†</sup>

<sup>†</sup> Departamento de Ciencias e Ingeniería de la Computación, Universidad Nacional del Sur,  
Instituto de Ciencias e Ingeniería de la Computación (UNS-CONICET),  
Bahía Blanca, Argentina

<sup>‡</sup> Departamento de Economía, Universidad Nacional del Sur,  
Instituto de Matemática de Bahía Blanca (UNS-CONICET),  
Bahía Blanca, Argentina

**Abstract.** Successful modeling and prediction depend on effective methods for the extraction of domain-relevant variables. This paper proposes a methodology for identifying domain-specific terms. The proposed methodology relies on a collection of documents labeled as relevant or irrelevant to the domain under analysis. Based on the labeled document collection, we propose a supervised technique that weights terms based on their descriptive and discriminating power. Finally, the descriptive and discriminating values are combined into a general measure that, through the use of an adjustable parameter, allows to independently favor different aspects of retrieval such as maximizing precision or recall, or achieving a balance between both of them. The proposed technique is applied to the economic domain and is empirically evaluated through a human-subject experiment involving experts and non-experts in Economy. It is also evaluated as a term-weighting technique for query-term selection showing promising results. We finally illustrate the potential of the proposal as a first step for identifying different types of associations between words.

**Keywords:** Term Weighting, Variable Extraction, Information Retrieval, Query-Term Selection

## 1 Introduction

A great number of machine learning and data science applications require identifying domain- or topic-relevant terms. For instance, automatic query formulation requires selecting good query terms; classification requires extracting good features, and in general, any modeling and prediction task requires mechanisms for variable extraction as an initial step to build useful representations. Also, term weighting is a crucial component of these representations since the importance of a term for a domain or topic can usually be numerically estimated and such weights have an impact on the task to be carried out.

Several term-weighting schemes have been proposed in the literature with varying degree of success. Most of these methods apply an unsupervised approach to determine term importance. This is the case of the widely-used TF-IDF weighting scheme, where

terms are weighted based on local (TF) and global (IDF) term frequencies, but no class-label information is used to compute these weights. This scheme is limited when it comes to identifying terms that are important for general topics or domains because it has the constraints of being document dependent (as it is based on the document and not on the general topic or domain) and label independent (as it is independent of the topic or domain label). Other term weighting methods take a supervised approach to assess the importance of a term in a class. However, term importance is typically taken as a fixed value independent of the task at hand. This represents a limitation because the importance of a term depends on whether the term is needed for query construction, clustering, classification, document summarization, among other tasks. Even for a specific task, such as is the case of query construction, a term may be more or less effective depending on whether the application requires high recall (e.g., looking for all relevant literature about a given topic) or high precision (e.g., looking for a specific piece of information such as a date, place or name). For example, a term that is a useful descriptor for a topic of interest, and therefore useful for attaining high recall, may lack discriminating power, resulting in low precision, unless it is combined with other terms that can discriminate between good and bad results.

This paper proposes a methodology that can be applied to identify domain- or topic-relevant variables from labeled documents. Two forms of relevance are distinguished, namely the relevance of a term as a descriptor, or *descriptive relevance*, and the relevance of a term as a discriminator, or *discriminative relevance*. Guided by this distinction, we propose two weighting schemes that account for these two notions of relevance. These weights are then combined into a parameter-dependent measure to which we refer to as  $FDD_{\beta}$ , accounting for a general notion of relevance. As we will show in the experiments, the  $FDD_{\beta}$  measure offers an advantage over several state-of-the-art term-weighting schemes as its parameter can be adjusted to emphasize different aspects of relevance (i.e., descriptive and discriminative relevance). As a consequence, the  $FDD_{\beta}$  measure has the practical implication of being able to favor either precision or recall, as well as to achieve a balance between both.

The paper is structured as follows. Section 2 briefly describes background concepts and reviews existing term-weighting schemes. Section 3 presents our novel term-weighting scheme, to which we refer to as  $FDD_{\beta}$ . Section 4 describes the data collection used in our analysis and evaluates  $FDD_{\beta}$  through a user study and as a query-term selection mechanisms. Section 5 illustrates the application of the proposal to extract economic relevant variables from digital media. Finally, section 5 presents the conclusions and outlines future research work.

## 2 Background and Related Work

Term weighting has been widely used in text classification and information retrieval. For historical reasons, term-weighting methods in text classification were originally borrowed from the information retrieval area, which traditionally applied unsupervised techniques. These traditional term-weighting schemes were designed to improve both recall and precision in the retrieval task. Based on these considerations, Salton and Buckley (1988) claimed that at least three main factors are required in any term weight-

ing scheme. The first is a local factor that stands for the presence of the term in the document. This factor represents whether the term appears at all, and how many times it does. It represents the idea that frequent terms are semantically close to the content of the document. Such a factor is designed to improve recall. The second factor is a global value associated with each term, which represents how frequent the term is in the document collection, in such a way that frequent terms are penalized. The rationale for using this penalizing factor is that common terms are poor discriminators, and as a consequence, they are not useful to tell apart among different documents containing them. It is known that using this factor helps to achieve higher precision. Note, however, that this might be at the expenses of a drop in recall. Finally, the terms are sometimes corrected by a normalization factor.

The simplest local factor is the binary one, which only measures the presence or absence of the term in the document (with values 1 or 0). Another simple and highly-used factor is *term frequency* (TF), which counts the number of times a term appears in a document. It relies on the assumption that most frequent terms are closely related to the content of the document. Leopold and Kindermann (2002) propose *inverse term frequency* (ITF) as an alternative to the classic TF. The ITF weight is based on Zipf's Law and normalizes the local factor to the interval [0,1]. On the other hand, Debole and Sebastiani (2004) propose another variation for the local factor, with a logarithmic transformation in which the terms that are extremely frequent do not increase at the same rate as in TF. Hassan et al. (2007) present a new local factor using a variant of *TextRank* (Mihalcea and Tarau, 2004) as a scoring function, which recursively increases the importance of a term by determining the degree of connectivity between other terms using co-occurrence as a way to measure connectivity. *TextRank* is based on the renowned *PageRank* algorithm (Page et al., 1999).

A simple global factor can be computed by counting the number of documents in the corpus that contain the term. We refer to this factor as *term global frequency* (TGF). The best known global factor is the *inverse document frequency* (IDF) function (Salton and Buckley, 1988), which relies on the assumption that terms that occur in many documents are not good for discrimination. The TF-IDF formulation is a widely used weighting scheme because it reaches a good balance between the local (TF) and the global (IDF) factor. Tokunaga and Makoto (1994) propose a variant of IDF named *weighted inverse document frequency* (WIDF) that penalizes frequent terms by taking into account the number of times they occur in each document of a collection. A variant of TF-IDF called *modified inverse document frequency* (MIDF) that combines TF and WIDF is proposed in (Deisy et al., 2010). According to the authors, MIDF outperforms TF-IDF in text classification. Also, they remark the ability of MIDF to adapt to dynamic document corpora.

While unsupervised weighting schemes have proved to be useful in many scenarios, these methods do not take full advantage of class information, which is available as part of the training set in a class-labeled collection. The design of term-weighting methods that exploit class information gained increasing attention, giving rise to different forms of supervised term-weighting schemes (Debole and Sebastiani, 2004; Lan et al., 2005; Wang and Zhang, 2013; Deng et al., 2014; Chen et al., 2016; Verberne et al., 2016; Fattah and Sohrab, 2016; Feng et al., 2018). A simple method that uses

class information can be computed by counting the number of documents in a class that contain the term. We use TGF\* to refer to this method. Another supervised weighting scheme is the *inverse class frequency factor* (ICF), which relies on the assumption that a term that occurs in documents from a single class are good discriminants of that class. Conversely, terms that appear in documents from different classes contribute poorly to the identification of the class of the documents. So, this factor penalizes a term proportionally to the number of different classes in which the term appears. Other functions from traditional information theory such as *mutual information* (MI), *chi-squared* ( $\chi^2$ ), *information gain* (IG) and *gain ratio* (GR) can be used as supervised term-weighting scores to capture the idea that the most valuable terms for categorization under a class are those that are distributed most differently in the sets of positive and negative examples of the class. A classic feature scoring function that is commonly used as a global term-weighting factor is the *odds ratio* (OR) (van Rijsbergen et al., 1981). This score is based on the conditional probability of a term occurring given a class. Another supervised technique known as *category relevance factor* (CRF) computes a factor that stands for the discriminating power of a feature to a class (Deng et al., 2002). Some feature selection techniques that were adapted for term weighing are the *Galavotti-Sebastiani-Simi* coefficient (GSS) (Galavotti et al., 2000) and *entropy-based category coverage difference* (ECCD) (Largerone et al., 2011). Liu et al. (2009) propose a probabilistic-based technique (Prob) that involves two ratios directly related to the term's strength in representing a category. These ratios are such that one of them increases if the term appears in a lot of documents of the class (descriptive power), while the other tends to be higher if the term appears only in documents of the class (discriminating power). Another scheme uses a *relevancy frequency factor* (RF) (Lan et al., 2009) that takes into account term distribution across classes. According to this scheme, the higher the concentration of high-frequency terms in the positive category than in the negative one, the greater the contribution to classification. Domeniconi et al. (2015) propose a supervised variant of IDF called *inverse document frequency excluding category* (IDFEC). Similar to IDF, IDFEC penalizes frequent terms, but different from IDF it avoids penalizing those terms occurring in several documents belonging to the same class. Another variant also proposed in (Domeniconi et al., 2015) results from combining IDFEC and RF, resulting in the IDFEC\_B scheme.

Table 1 shows the definitions of the main scores presented above using the following notation (Lan et al., 2005; Domeniconi et al., 2015):

- $A$  denotes the number of documents that belong to class  $c_k$  and contain term  $t_i$ .
- $B$  denotes the number of documents that belong to class  $c_k$  but do not contain the term  $t_i$ .
- $C$  denotes the number of documents that do not belong to class  $c_k$  but contain the term  $t_i$ .
- $D$  denotes the number of documents that do not belong to class  $c_k$  class and do not contain the term  $t_i$ .
- $N$  denotes the total number of documents in the collection (i.e.,  $N = A + B + C + D$ ).

Note that some formulations include the expression  $\max(X, 1)$  to prevent the possibility of undefined values, such as divisions by zero or  $\log(0)$ .

Name	Formulation
TGF	$A + C$
IDF	$\log(N/(A + C))$
TGF*	$A$
MI	$\log((N \times \max(A, 1))/((A + B)(A + C)))$
$\chi^2$	$N((AD - BC)^2/((A + C)(B + D)(A + B)(C + D)))$
OR	$\log((\max(A, 1) \times D)/\max(B \times C, 1))$
IG	$(A/N) \log(\max(A, 1)/(A + C)) - ((A + B)/N) \log((A + B)/N) + (B/N) \log(B/(B + D))$
GR	$IG/(-(A + B)/N \log((A + B)/N) - ((C + D)/N) \log((C + D)/N))$
GSS	$\log(2 + ((A + C + D)/\max(C, 1)))$
Prob	$\log(1 + (A/B)(A/C))$
RF	$\log(2 + (A/\max(C, 1)))$
IDFEC	$\log((C + D)/\max(C, 1))$
TGF-IDFEC	$(A + C)(\log((C + D)/\max(C, 1)))$
TGF*-IDFEC	$A \times (\log((C + D)/\max(C, 1)))$
IDFEC_B	$\log(2 + (A + C + D)/\max(C, 1))$

**Table 1.** Definitions of term weighting schemes.

### 3 A Novel Supervised Term-Weighting Score

Based on the idea that class labels convey useful information for term weighting and on the fact that the importance of a term in a topic or domain depends on the specific objectives at hand (e.g., attaining high recall, high precision or both), we distinguish two relevancy scores. The first score represents the importance of a term to describe the class or topic, and we refer to it as *descriptive relevance* (DESCR). Given a term  $t_i$  and a class  $c_k$  the DESCR score is expressed as:

$$\text{DESCR}(t_i, c_k) = \frac{|d_j : t_i \in d_j \wedge d_j \in c_k|}{|d_j : d_j \in c_k|},$$

which is equivalent to  $A/(A + B)$ , using the notation adopted in the previous section. The descriptive relevance of a term in a class stands for a simple idea: those terms that occur in many documents of a given class are good descriptors of that class. As a consequence, we compute it as the portion of documents in the class that contain the given term.

The second relevancy score represents the importance of a term to discriminate a class or topic, and we call it *discriminative relevance*. For a term  $t_i$  and a class  $c_k$  the DISCR score is expressed as:

$$\text{DISCR}(t_i, c_k) = \frac{|d_j : t_i \in d_j \wedge d_j \in c_k|}{|d_j : t_i \in d_j|},$$

which is equivalent to  $A/(A + C)$ . The discriminative relevance of a term in a class is based on the idea that a term is a good discriminator of a class if it tends to occur only in documents of that class. We compute it as the portion of documents that contain the given term that belong to the class. The DESCR and DISCR scores can be seen as the supervised versions of the semi-supervised techniques proposed in (Maguitman et al., 2004) to compute the descriptive and discriminative power of a term in a topic.

We propose to combine the DESCR and DISCR scores by means of the following general term relevancy formula:

$$\text{FDD}_\beta(t_i, c_k) = (1 + \beta^2) \frac{\text{DISCR}(t_i, c_k) \times \text{DESCR}(t_i, c_k)}{(\beta^2 \times \text{DISCR}(t_i, c_k)) + \text{DESCR}(t_i, c_k)}.$$

The  $\text{FDD}_\beta$  measure is derived from the  $F_\beta$  formula traditionally used in information retrieval to give  $\beta$  times more importance to recall than to precision:

$$F_\beta(t_i, c_k) = (1 + \beta^2) \frac{\text{precision}(t_i, c_k) \times \text{recall}(t_i, c_k)}{(\beta^2 \times \text{precision}(t_i, c_k)) + \text{recall}(t_i, c_k)}.$$

By using a  $\beta$  value higher than 1 in the  $\text{FDD}_\beta$  function we can weight descriptive relevance higher than discriminative relevance (by placing more emphasis on terms that help achieving good recall) while a  $\beta$  smaller than 1 weights descriptive relevance lower than discriminative relevance (by placing more emphasis on terms that help achieving good precision).

We will show next that  $\text{FDD}_\beta$ , can serve the purpose of approximating term relevancy in a topic. This score can be computed for any collection of documents labeled as relevant or irrelevant to the given topic. We will show next that despite the simplicity of the  $\text{FDD}_\beta$  score, it is highly effective both as an estimator of expert assessments of relevance and for guiding the selection of good query terms. In particular, we will show how the tunable parameter  $\beta$  offers a means to favor different objectives in the information retrieval task.

## 4 Evaluation

The goal of this section is to compare the  $\text{FDD}_\beta$  weighting score against other term-weighting schemes. The evaluation comprises a human-subject study and an experiment for assessing the effectiveness of the evaluated techniques in information retrieval. The evaluations were carried out on the economic domain using a labeled collection of news articles and human subjects' relevance assessments, as described next. The labeled collection of news articles and the human subjects' relevance assessments are available for download at <http://ir.cs.uns.edu.ar/datasets>.

### 4.1 Data Collection

The *The Guardian* newspaper (<https://www.theguardian.com/>) was selected as a source to collect a set of digital news. *The Guardian* is a British daily newspaper with an open platform that allows accessing over 1.9 million pieces of content, including full-text news articles. A simple Python script was developed to collect news articles through an API provided by the newspaper. Only news coming from the Politics, World news, Business and Society sections were collected. Although several fields are available for each news article, only the news titles and full body text were used. A simple preprocessing step was carried out to eliminate stopwords and punctuation marks, as well as to transform the text into a sequence of lowercase terms. A total of 1689 news articles

corresponding to January 2013 were manually labeled by two experts in Economy as relevant (537) or irrelevant (1152) to the economic domain. To complete the labeling task, both experts read the news articles and agreed on whether each of them was relevant or not. It is worth mentioning that the manual labeling stage was important due to the fact that news identified by the experts as economic relevant do not exactly correspond to those from the Business section (418 out of 512) but also some of them were in the Politics (39 out of 290), World news (43 out of 650), and Society (37 out of 237) sections. The total number of terms in these news articles is 38511. However, to reduce the dimensionality of the dataset only those terms that occur in at least six news articles were considered, resulting in a set of 10373 terms. The collection of 1689 expert-labeled news articles was used as the training set. Also, a reduced set consisting of 100 expert-labeled news articles (not included in the training set) was used as the validation set.

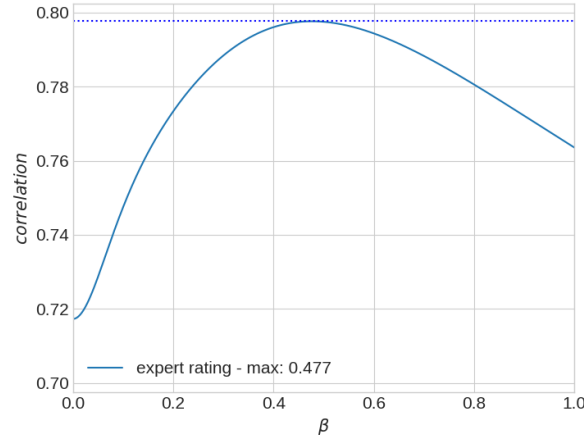
#### 4.2 Validation by User Study

Eight volunteer subjects were recruited for an experiment conducted online. The group of subjects included four volunteers with no background in Economy and four others with a Ph.D. degree in Economy. We refer to the first group as *non-experts* and to the second group as *experts*. A set of 50 terms (10 lists of 5 terms each) and another set of 100 terms (20 lists of 5 terms each) were strategically selected from the 10373 terms of the dataset. The selection was made based on the distribution of term frequency in light of Zipf's law. The goal was to avoid providing low-frequency terms (which are many) more chances of being selected than high-frequency terms (which are a few). To complete an initial parameter-adjusting stage, two of the experts were asked to agree on the economic relevance of each of the words from the 50-term set. The experts were asked to rate these terms with a score ranging from 0 (economic irrelevant) to 5 (economic very relevant). We used these ratings and the labeled collection to learn the best  $\beta$  value for the  $FDD_\beta$  method. As can be seen in figure 1 the highest Pearson correlation between the expert ratings and the  $FDD_\beta$  values was 0.797671, which was achieved for  $\beta = 0.477$ .

To complete the validation stage we asked the eight volunteer subjects to rate the 100 terms using a 0-5 scale, and we computed DESCR, DISCR,  $FDD_{0.477}$  and the 15 weighting schemes listed in table 1 for these terms. In the first place, we tested the level of agreement between pairs of users belonging to the non-expert group and between pairs of users in the expert group. Table 2 presents the means and standard deviations obtained as a result of such analysis. It is possible to observe that there is a high level of agreement in both groups, being this agreement higher in the expert group.

non-expert	experts
$\mu = 0.839475, \sigma = 0.037791$	$\mu = 0.876390, \sigma = 0.009438$

**Table 2.** Means ( $\mu$ ) and standard deviations ( $\sigma$ ) of correlations to test agreement among non-experts and among experts.



**Fig. 1.** Learning the optimal  $\beta$  value (maximum correlation equal to 0.797671 for  $\beta = 0.477$ ).

Table 3 presents results on the level of agreement between the two different groups of users (non-experts and experts), and on the level of agreement between each of these groups (non-experts and experts) and  $FDD_{0.477}$ .

non-experts and experts	non-experts and $FDD_{0.477}$	experts and $FDD_{0.477}$
$\mu = 0.80383, \sigma = 0.053205$	$\mu = 0.685598, \sigma = 0.054969$	$\mu = 0.752352, \sigma = 0.018904$

**Table 3.** Means ( $\mu$ ) and standard deviations ( $\sigma$ ) of correlations computed between non-experts and experts, non-experts and  $FDD_{0.477}$ , and experts and  $FDD_{0.477}$ .

Finally, to compare the effectiveness of the weighting schemes as predictors of subjects' judgments of term relevancy we computed the Pearson correlation coefficients between the averaged ratings assigned by the subjects and those computed by each of the weighting schemes. Table 4 summarizes these correlations. The reported values correspond to the correlations between each of the methods and the different groups of users. In all these cases we observe that  $FDD_{0.477}$  outperforms the other methods, being TGF\*-IDFEC the second most effective one in estimating human subjects' relevance assessments.

### 4.3 Retrieval Effectiveness

In this section, we analyze the performance of  $FDD_{\beta}$  as a mechanism for query-term selection, and we compare it with other state-of-the-art weighting schemes. In the first place, the training set described in section 4.1 was used to select the top-rated terms based on  $FDD_{\beta}$  by assigning different values to the parameter  $\beta$ . Simple queries were generated using the selected terms and then evaluated by means of the classical recall,

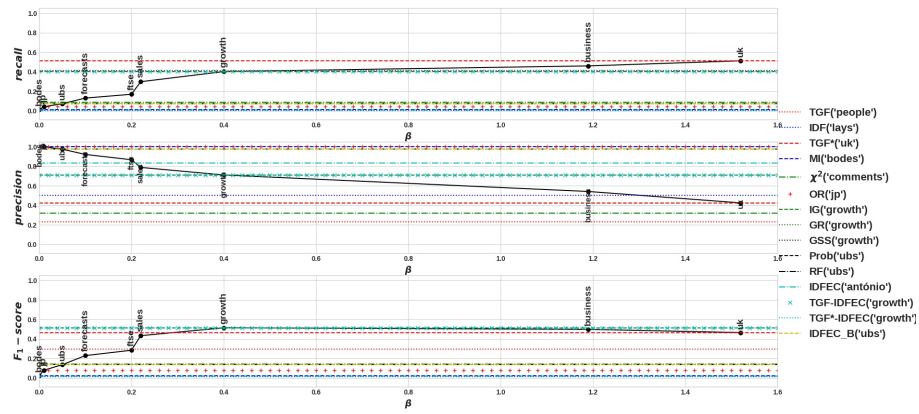


Method	non-expert (averaged)	expert (averaged)	non-expert and expert (averaged)
TGF	0.283553	0.365037	0.332324
IDF	-0.488816	-0.563704	-0.539138
TGF*	0.574110	0.642607	0.623198
MI	0.697053	0.659659	0.694604
$\chi^2$	-0.164537	-0.087771	-0.128992
OR	0.432627	0.306599	0.378188
IG	0.663296	0.705736	0.701123
GR	0.663296	0.705736	0.701123
GSS	0.722761	0.757015	0.757807
Prob	0.654187	0.697007	0.691990
RF	0.472824	0.407394	0.450543
IDFEC	-0.226397	-0.325872	-0.283050
TGF-IDFEC	0.603975	0.676551	0.655882
TGF*-IDFEC	0.721871	0.774026	0.766110
IDFEC.B	-0.221061	-0.320304	-0.277466
DESCR	0.574110	0.642607	0.623198
DISCR	0.662481	0.610804	0.651848
FDD <sub>0.477</sub>	<b>0.735456</b>	<b>0.791969</b>	<b>0.782264</b>

**Table 4.** Correlations between methods and ratings obtained by averaging non-expert, expert and all human subjects' scores.

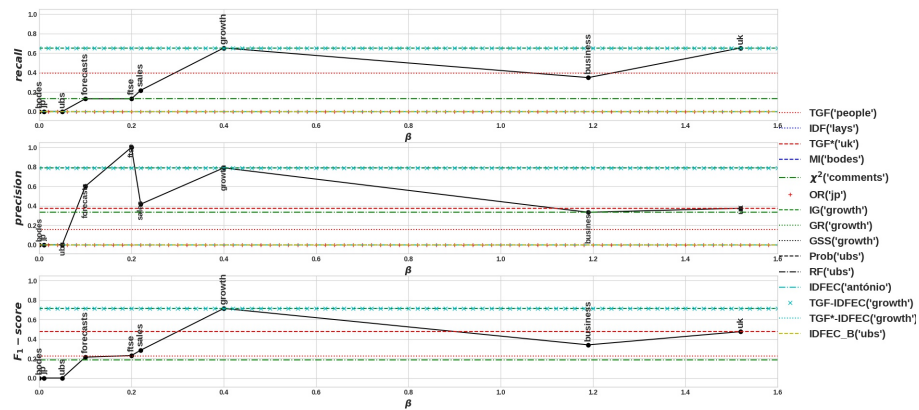
precision and  $F_1$  metrics. The results are shown in figure 2. As expected, the highest recall using the  $FDD_\beta$ -based term selection mechanisms is obtained with larger values of  $\beta$  while the highest precision is obtained for smaller values. Note, for instance, that terms such as `uk` occur often in relevant news articles given the fact that news were collected from a British newspaper. As a result, the term `uk` results in a high-recall query. However, `uk` is not a good discriminator for the Economy domain, resulting in a low-precision query. On the other hand, terms such as `adp`, `jp`, `ubs`, `forecasts` and `ftse` are not good descriptors but tend to occur only in relevant news articles. This means that they are good discriminators, offering a mechanism to ensure high precision, although usually at the expense of low recall. Other terms, such as `sales`, `growth` and `business` achieve a balance between descriptive and discriminative relevance, resulting in a good  $F_1$  score. The query term with the highest  $F_1$  score is `growth`, which is the term achieving the best  $FDD_\beta$  for a range of  $\beta$  values that begins approximately at 0.4 and ends close to 1.2. Note that this range includes 0.477, which is the value that yields the highest correlation between  $FDD_\beta$  scores and experts' relevance assessments. Based on this preliminary analysis the best  $FDD_\beta$  achieves an  $F_1$  score as high as the one obtained with the two most effective state-of-the-art weighting schemes (TGF-IDFEC and TGF\*-IDFEC). The top-rated term according to the three weighting schemes is `growth`. It is interesting to note that for small  $\beta$  values  $FDD_\beta$  outperforms these two methods in terms of precision while for large  $\beta$  values  $FDD_\beta$  outperforms these two methods in terms of recall.

The validation set described in section 4.1 was used to determine if the best queries identified using the training set were effective on a different set. The resulting recall, precision and  $F_1$  metrics computed on the validation set are shown in figure 3. Given that the validation set was small, some of the most discriminating terms identified during the training stage (`adp` and `ubs`) were absent from the validation set, resulting in an empty answer set when used as query terms. However, those terms with a good balance between descriptive and discriminative relevance (`sales`, `growth` and `business`)



**Fig. 2.** Effectiveness on the training set of queries generated based on term weighting schemes. The black solid curve corresponds to the effectiveness of query terms selected using  $FDD_{\beta}$  on the training with different  $\beta$  values.

achieve the highest  $F_1$  scores when used as query terms on the validation set. This preliminary analysis indicates that the proposed method does not overfit the training data.



**Fig. 3.** Effectiveness on the validation set of queries generated based on term weighting schemes. The black solid curve corresponds to the effectiveness of query terms selected using  $FDD_{\beta}$  on the training set with different  $\beta$  values

As can be seen in the reported results,  $FDD_{\beta}$  performs consistently well, not only as an estimator of human subjects’ relevance assessments but also as a method for guiding the selection of good query terms. In this section we illustrate the application of the proposed method on the economic domain as a mechanism to extract variables from digital media with the ultimate goal of building models of prediction, explanation and description. Figure 4 shows a word-cloud visualization with the top-ranked terms based on the training data using  $FDD_{0.477}$  as weighting scheme. To avoid overcharging the figure only terms with  $FDD_{0.477} > 0.7$  are shown. This visualization can support the construction of knowledge models, by helping in the process of choosing relevant variables, which is typically the initial step in any modeling task.



A subsequent modeling step would be to identify different types of dependency relations between these variables. For instance, some *causal relations* that can be recognized are investment-growth-gdp, spending-market-recovery and sales-companies-investment-gdp. Other types of relations, such as *close associations* are illustrated by credit-debt-banks, recession-decline, trading-stock-ftse and debt-bank-investors-trading-recession. A possible *simultaneity relation* is given by *market-prices*. It is also interesting to note

that *christmas*, one of the words selected by the method, may capture *seasonality in a casual series*. Automatically identifying these types of relations is a challenging problem that we plan to address as part of our ongoing research work. In particular, we plan to investigate into the problem of finding causal relations with the purpose of automatically building different types of networks, such as Bayesian networks (Pearl, 2014).

## 6 Conclusions and Future Work

In this paper we presented a methodology for identifying domain-specific terms. As part of the proposed methodology we defined a novel supervised term-weighting scheme called  $FDD_{\beta}$ , which is based on the notions of descriptive and discriminative relevance. Preliminary evaluations show that  $FDD_{\beta}$  achieves good performance as an estimator of human subjects' relevance judgments and as a mechanism for selecting good query terms. Also, it offers the flexibility of adapting to different goals, such as achieving high recall, high precision, or a balance between both. This flexibility represents an important advantage over the analyzed state-of-the-art weighting schemes.

The proposed technique was evaluated on the economic domain with promising results and we anticipate that it will also achieve good performance on other domains. Also, we plan to test  $FDD_{\beta}$  on specific topics (as is the case of the topic of a news article), as opposed to general domains (as is the case of Economy). Another important future task will be to validate  $FDD_{\beta}$  on larger data sets, such as those available as part of the TREC collection ([https://trec.nist.gov/data/test\\_coll.html](https://trec.nist.gov/data/test_coll.html)). The proposed weighting scheme will also be evaluated on classification tasks, which will open new challenges.

## Acknowledgment

This work was supported by CONICET (PIP 11220120100487), MinCyT (PICT 2014-0624) and Universidad Nacional del Sur (PGI-UNS 24/N029).

## Bibliography

- Chen, K., Zhang, Z., Long, J., and Zhang, H. (2016). Turning from tf-idf to tf-igm for term weighting in text classification. *Expert Systems with Applications*, 66:245 – 260.
- Debole, F. and Sebastiani, F. (2004). Supervised term weighting for automated text categorization. In *Text mining and its applications*, pages 81–97. Springer.
- Deisy, C., Gowri, M., Baskar, S., Kalaiarasi, S., and Ramraj, N. (2010). A novel term weighting scheme midf for text categorization. *Journal of Engineering Science and Technology*, 5(1):94–107.
- Deng, Z.-H., Luo, K.-H., and Yu, H.-L. (2014). A study of supervised term weighting scheme for sentiment analysis. *Expert Systems with Applications*, 41(7):3506–3513.
- Deng, Z.-H., Tang, S.-W., Yang, D.-Q., Zhang, M., Wu, X.-B., and Yang, M. (2002). A linear text classification algorithm based on category relevance factors. In *International Conference on Asian Digital Libraries*, pages 88–98. Springer.
- Domeniconi, G., Moro, G., Pasolini, R., and Sartori, C. (2015). A study on term weighting for text categorization: A novel supervised variant of tf. idf. In *DATA*, pages 26–37.
- Fattah, M. and Sohrab, M. (2016). Combined term weighting scheme using ffnn, ga, mr, sum, and average for text classification. *International Journal of Scientific and Engineering Research*, 7(8):2031–2040.
- Feng, G., Li, S., Sun, T., and Zhang, B. (2018). A probabilistic model derived term weighting scheme for text classification. *Pattern Recognition Letters*, 110:23 – 29.
- Galavotti, L., Sebastiani, F., and Simi, M. (2000). Experiments on the use of feature selection and negative evidence in automated text categorization. In *International Conference on Theory and Practice of Digital Libraries*, pages 59–68. Springer.
- Hassan, S., Mihalcea, R., and Banea, C. (2007). Random walk term weighting for improved text classification. *International Journal of Semantic Computing*, 1(04):421–439.
- Lan, M., Sung, S.-Y., Low, H.-B., and Tan, C.-L. (2005). A comparative study on term weighting schemes for text categorization. In *Neural Networks, 2005. IJCNN'05. Proceedings. 2005 IEEE International Joint Conference on*, volume 1, pages 546–551. IEEE.
- Lan, M., Tan, C. L., Su, J., and Lu, Y. (2009). Supervised and traditional term weighting methods for automatic text categorization. *IEEE transactions on pattern analysis and machine intelligence*, 31(4):721–735.
- Largerion, C., Moulin, C., and Géry, M. (2011). Entropy based feature selection for text categorization. In *Proceedings of the 2011 ACM Symposium on Applied Computing*, pages 924–928. ACM.
- Leopold, E. and Kindermann, J. (2002). Text categorization with support vector machines. how to represent texts in input space? *Machine Learning*, 46(1-3):423–444.
- Liu, Y., Loh, H. T., and Sun, A. (2009). Imbalanced text classification: A term weighting approach. *Expert systems with Applications*, 36(1):690–701.

- Maguitman, A., Leake, D., Reichherzer, T., and Menczer, F. (2004). Dynamic extraction topic descriptors and discriminators: towards automatic context-based topic search. In *Proceedings of the thirteenth ACM international conference on Information and knowledge management*, pages 463–472. ACM.
- Mihalcea, R. and Tarau, P. (2004). Textrank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*.
- Page, L., Brin, S., Motwani, R., and Winograd, T. (1999). The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab.
- Pearl, J. (2014). *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier.
- Salton, G. and Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513–523.
- Tokunaga, T. and Makoto, I. (1994). Text categorization based on weighted inverse document frequency. In *Special Interest Groups and Information Process Society of Japan (SIG-IPSI)*. Citeseer.
- van Rijsbergen, C., Harper, D., and Porter, M. (1981). The selection of good search terms. *Information Processing & Management*, 17(2):77 – 91.
- Verberne, S., Sappelli, M., Hiemstra, D., and Kraaij, W. (2016). Evaluation and analysis of term scoring methods for term extraction. *Information Retrieval Journal*, 19(5):510–545.
- Wang, D. and Zhang, H. (2013). Inverse-category-frequency based supervised term weighting schemes for text categorization. *Journal of Information Science and Engineering*, 29(2):209–225. cited By 28.