

A Case Study of Spanish Text Transformations for Twitter Sentiment Analysis

Eric S. Tellez^{1,3} Sabino Miranda-Jiménez^{1,3} Mario Graff^{1,3}
 Daniela Moctezuma^{1,2} Oscar S. Siodia² Elio A. Villaseñor¹

¹INFOTEC Centro de Investigación e Innovación en Tecnologías de la Información y Comunicación, Circuito Tecnopolo Sur No 112, Fracc. Tecnopolo Pocitos II, Aguascalientes 20313, México

²CentroGEO Centro de Investigación en Ciencias de Información Geoespacial, Circuito Tecnopolo Norte No. 117, Col. Tecnopolo Pocitos II, C.P., Aguascalientes, Ags 20313 México

³CONACyT Consejo Nacional de Ciencia y Tecnología, Dirección de Cátedras, Insurgentes Sur 1582, Crédito Constructor, Ciudad de México 03940 México

This work is published in Expert Systems with Applications
<https://doi.org/10.1016/j.eswa.2017.03.071>

Abstract

Sentiment analysis is a text mining task that determines the polarity of a given text, i.e., its positiveness or negativeness. Recently, it has received a lot of attention given the interest in opinion mining in micro-blogging platforms. These new forms of textual expressions present new challenges to analyze text given the use of slang, orthographic and grammatical errors, among others. Along with these challenges, a practical sentiment classifier should be able to handle efficiently large workloads.

The aim of this research is to identify which text transformations (lemmatization, stemming, entity removal, among others), tokenizers (e.g., words n -grams), and tokens weighting schemes impact the most the accuracy of a classifier (Support Vector Machine) trained on two Spanish corpus. The methodology used is to exhaustively analyze all the combinations of the text transformations and their respective parameters to find out which characteristics the best performing classifiers have in common. Furthermore, among the different text transformations studied, we introduce a novel approach based on the combination of word based n -grams and character based q -grams. The results show that this novel combination of words and characters produces a classifier that outperforms the traditional word based combination by 11.17% and 5.62% on the INEGI and TASS'15 dataset, respectively.

1 Introduction

In recent years, the production of textual documents in social media has increased exponentially; for instance, up to April 2016, Twitter has 320 million active users, and Facebook has 1,590 million users.¹ In social media, people share comments about many disparate topics, i.e., events, persons, and organizations, among others. These facts have had the result of seeing social media as a gold mine of human opinions, and consequently, there is an increased interest in doing research and business activities around opinion mining and sentiment analysis fields.

Automatic sentiment analysis of texts is one of the most important tasks in text mining, where the goal is to determine whether a particular document has either a positive, negative or neutral

¹<http://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/>

opinion². Determining whether a text document has a positive or negative opinion is becoming an essential tool for both public and private companies, [Liu, 2015, Peng et al., 2008]. Given that it is a useful tool to know *what people think* about anything; so, it represents a major support for decision-making processes (for any level of government, marketing, etc.) [Pang and Lee, 2008].

Sentiment analysis has been traditionally tackled as a classification task where two major problems need to be faced. Firstly, one needs to transform the text into a suitable representation, this is known as text modeling. Secondly, one needs to decide which classification algorithm to use; one of the most widely used is Support Vector Machines (SVM). This contribution focus on the former problem, i.e., we are interested in improving the classification by finding a suitable text representation.

Specifically, the contribution of this research is twofold. Firstly, we parametrize our text transformations with different techniques such as: lemmatization, stemming, and entity removal, just to mention a few (Table 3 contains all the transformations explored). This parametrization is used to exhaustively evaluate the entire configurations space to know those transformations that produce the best SVM classifier on two sentiment analysis corpus written in Spanish. Counterintuitively, we found that the complexity of techniques used in the pre-processing step is not correlated with the final performance of the classifier, e.g., a classifier using lemmatization, which is one of the pre-processing techniques having the greatest complexity, might not be one of the systems having the highest performance.

Secondly, we propose a novel approach based on the combination of word based n -grams and character based q -grams. This novel combination of words and characters produces a classifier that outperforms the traditional word based combination by 11.17% and 5.62% on the INEGI and TASS'15 dataset, respectively. Hereafter, we will use n -words to refer to word n -grams, and q -grams to character q -grams just to make a clear distinction between these techniques.

This manuscript is organized as follows. Section 1 introduces the paper and the problem being tackled. Section 2 deals with literature review. The text transformations are described in Section 3, meanwhile the parameters settings and definition of the problem are presented on Section 4. Section 5 describes our experimental results. Finally, Section 6 and Section 7 present the discussion and conclusions of our results along with possible directions for future work.

2 Related Work

The sentiment analysis task has been widely studying due to the interest to know the people's opinions and feelings about something, particularly in social media. This task is commonly tackled in two different ways. The first one involves the use of static resources that summarize the sense or semantic of the task; these knowledge databases contain mostly affective lexicons. These lexicons are created by experts, in psychology or by automated processes, that perform the selection of features (words) along with a corpus of labeled text as done in [Ghiassi et al., 2013]. Consequently, the task is solved by trying to detect how the affective features are used in a text, and how these features can be used to predict the polarity of a given text.

The alternative approach states the task as a text classification problem. This includes several distinguished parts like the pre-processing step, the selection of the vectorization and weighting schemes, and also the classifier algorithm. So, the problem consists of selecting the correct techniques in each step to create the sentiment classifier. Under this approach, the idea is to process the text in a way that the classifier can take advantage of the features to solve the problem. Our contribution focus in this later approach; we describe the best way to pre-process, tokenize, and vectorize the text, based on a fixed set of text-transformation functions. For simplicity, we fix our classifier to be Support Vector Machines (SVM). SVM is a classifier that excels in high dimensional datasets as is the case of text classification, [Joachims, 1998]. This section reviews the related literature.

²Albeit, there are other variations considering intermediate levels for sentiments, e.g. more *positive* or less *positive*

There are several works in the sentiment analysis literature which use several representations; such as dictionaries [Alam et al., 2016], [Khan et al., 2016]; text content and social relations between users [Wu et al., 2016]; relationships between meanings of a word in a corpus [Razavi et al., 2014]; co-occurrence patterns of words [Saif et al., 2016], among others.

Focusing on the n -grams technique, a method that considers the local context of the word sequence and the semantic of the whole sentence is proposed in [Cui et al., 2015]. The local context is generated via the “bag-of- n -words” method, and the sentence’s sentiment is determined based on the individual contribution of each word. The word embedding is learned from a large monolingual corpus through a deep neural network, and the n -words features are obtained from the word embedding in combination with a sliding window procedure.

A hybrid approach that uses n -gram analysis for feature extraction together with a dynamic artificial neural network for sentiment analysis is proposed in [Ghiassi et al., 2013]. Here, a dataset over 10,000,000 of tweets, related to Justin Bieber topic, was used. As a result, a Twitter-specific lexicon with a reduced feature set was obtained.

The work presented in [Han et al., 2013] proposes an approach for sentiment analysis which combines an SVM classifier and a wide range of features like bag-of-word (1-words, 2-words) and part-of-speech (POS) features, etc., as well as votes derived from character n -words language models to achieve the final result. The authors concluded that lexical features (1-words, 2-words) produce the better contributions.

In [Tripathy et al., 2016] different classifiers and representations were applied to determine the sentiment in movie reviews, taken from internet blogs. The classifiers tested were Naive Bayes, maximum entropy, stochastic gradient descent, and SVM. These algorithms use n -words, for n in $\{1, 2, 3\}$ and all the combinations. Here, the results show that the value of n increases the classification accuracy decreases, i.e., using 1-words and 2-words the result achieved is better than using 3-words, 4-words, and 5-words.

Regarding the use of q -grams; in [Aisopos et al., 2011] a method that captures textual patterns is introduced. This method creates a graph, whose nodes correspond to q -grams of a document and their edges denoted the average distance between them. A comparative analysis on data from Twitter is performed between three representation models: term vector model, q -grams, and q -grams graphs. The authors found that vector models are faster, but q -grams (especially 4-grams) perform better in terms of classification quality.

With the purpose to attend sentiment analysis in Spanish tweets, a number of works has been presented in the literature, e.g. several sizes of n -grams and some polarity lexicons combined with a Support Vector Machine (SVM) was used in [Almeida, 2015]. Another approach which uses polarity lexicons with a number of features related to n -words, part-of-speech tag, hashtags, emoticon and lexicon resources is described in [Araque et al., 2015].

Features related to lexicons and syntactic structures are commonly used, for example, [Alvarez-López et al., 2015], [Cámara et al., 2015], [de la Vega and Vea-Murguía, 2015], [Borda and Saladich, 2015],[Deas et al., 2015]. In the other hand, features related to word vectorization, e.g. Word2Vec and Doc2Vec, are also used in several works, such as [Díaz-Galiano and Montejo-Ráez, 2015, Valverde et al., 2015].

Following with the Spanish language, in the most recent TASS (Taller de Análisis de Sentimientos ’16) competition, was presented some works still using polarity dictionaries and vectorization approach; such is the case of [Casasola Murillo and Marín Raventós, 2016], where an adaptation of Turney dictionary [Turney, 2002] over 5 millions of Spanish tweets was generated. Furthermore, [Casasola Murillo and Marín Raventós, 2016] in the step of vectorization uses n -grams and skip-grams in combination with this polarity dictionary. [Quirós et al., 2016] proposes the use of word embedding with SVM classifier. Despite the explosion of words using word embeddings, the classical word vectorization is still in use, [citementejo2016participacion](#).

A new approach is using ensembles or a combination of several techniques and classifiers, e.g. the work presented in [Cerón-Guzmán and de Cali, 2016] proposes an ensemble built on the combination of systems with the lowest correlation between them. [y Ferran Pla, 2016] presents another ensemble method where the Tweetmotif’s tokenizer, [O’Connor et al., 2010], is used in conjunction with Freeling [Padró and Stanilovsky, 2012]. These tools create a vector space that is



Figure 1: Generic treatment of input text to obtain the input vectors for the classifier.

the input for an SVM classifier.

It can be seen that one of the objectives of the related work is to optimize the number of n -words or q -grams (almost tackled as independent approaches), to increase performance; clearly, there is not a consensus. This lack of agreement motivates us to perform an extensive experimental analysis of the effect of the parameters (including n and q values), and so, we determined the best parameters on the Twitter databases employed.

3 Text Representation

Natural Language Processing (NLP) is a broad and complex area of knowledge having many ways to represent an input text [Giannakopoulos et al., 2012, Sammut and Webb, 2011]. In this research, we select the widely used vector representation of a text given its simplicity and powerful representation. Figure 1 depicts the procedure used to transform a text input into a vector. There are three main blocks: the first one transforms the text into another text representation, then the text is tokenized, and, finally, the vector is calculated using a weighting scheme. The resulting vectors are the input of the classifier.

In the following subsections, we described the text transformation techniques used which have a counterpart in many languages, the proper implementation of them rely heavily on the targeted language, in our case study the Spanish language. The interested reader looking for solutions in a particular language is encouraged to follow the relevant linguistic literature for its objective language, in addition to the general literature in NLP [Jurafsky and Martin, 2009, Bird et al., 2009, Sammut and Webb, 2011].

3.1 Text Transformation Pipeline

One of the contributions of this manuscript is to measure the effects that each different text transformation has on the performance of a classifier. This subsection describes the text transformations explored whereas the particular parameters of these transformations can be seen in Table 3.

3.1.1 TFIDF (tfidf)

In the vector representation, each word, in the collection, is associated with a coordinate in a high dimensional space. The numeric value of each coordinate is sometimes called the *weight* of the word. Here, $\text{tf} \times \text{idf}$ (Term Frequency-Inverse Document Frequency) [Baeza-Yates and Ribeiro-Neto, 2011] is used as bootstrapping weighting procedure. More precisely, let $D = \{D_1, D_2, \dots, D_N\}$ be the set of all documents in the corpus, and f_w^i be the frequency of the word w in document D_i . tf_w^i is defined as the normalized frequency of w in D_i

$$\text{tf}_w^i = \frac{f_w^i}{\max_{u \in D_i} \{f_u^i\}}.$$

In some way, tf describes the importance of w , locally in D_i . On the other hand, idf gives a global measure of the importance of w ;

$$\text{idf}_w = \log \frac{N}{|\{D_i \mid f_w^i > 0\}|}.$$

The final product, $\text{tf} \times \text{idf}$, tries to find a balance between the local and the global importance of a term. It is common to use variants of tf and idf instead of the original ones, depending in the application domain [Sammut and Webb, 2011]. Let v_i be the vector of D_i , a weighted matrix TFIDF of the collection D is created by concatenating all individual vectors, in some consistent order. Using this representation, a number of machine learning methods can be applied; however, the plain transformation of text to TFIDF poses some problems. On one hand, all documents will contain common terms having a small semantic content such as articles and determiners, among others. These terms are known as *stopwords*. The bad effects of stopwords are controlled by TFIDF, but most of them can be directly removed since they are fixed for a given language. On the other hand, after removing stopwords, TFIDF will produce a very high dimensional vector space, $O(N)$ in Twitter, since new terms are commonly introduced (e.g. misspellings, URLs, hashtags). This will rapidly yield to the *Curse of Dimensionality*, which makes hard to learn from examples since any two random vectors will be orthogonal with high probability. From a more practical point of view, a high dimensional representation will also impose huge memory requirements, at the point of being impossible to train a typical implementation of a machine learning algorithm (not being designed to use sparse vectors).

3.1.2 Stopwords (*del-sw*)

In many languages, like Spanish, there is a set of extremely common words such as determiners or conjunctions (*the* or *and*) which help to build sentences but do not carry any meaning themselves. These words are known as *Stopwords*, and they are removed from the text before any attempt to classify them. A stop list is built using the most frequent terms from a huge document collection. We used the Spanish stop list included in NLTK Python package [Bird et al., 2009].

3.1.3 Spelling

Twitter messages are full of slang, misspelling, typographical and grammatical errors among others; however, in this study, we focus only on the following transformations:

Punctuation (*del-punc*). This parameter considers the use of symbols such as question mark, period, exclamation point, commas, among other spelling marks.

Diacritic (*del-diac*). The Spanish language is full of diacritic symbols, and its wrong usage is one of the main sources of orthographic errors in informal texts. Thus, this parameter considers the use or absence of diacritical marks.

Symbol reduction (*del-d1*, *del-d2*). Usually, twitter messages use repeated characters to emphasize parts of the word to attract user's attention. This aspect makes the vocabulary explodes. Thus, we applied two strategies to deal with these phenomena: the first one replaces the repeated symbols by one occurrence of the symbol, and the second one replaces the repeated symbols by two occurrences to keep the word emphasize at the minimal level.

Case sensitive (*lc*). This parameter considers letters to be normalized in lowercase or to keep the original text. The aim is to cut the words that are the same in uppercase and lowercase.

3.1.4 Stemming (*stem*)

Stemming is a heuristic process in Information Retrieval field that chops off the end of words and often includes the removal of derivational affixes. This technique uses the morphology of the language coded in a set of rules; to find out word stems and reduce the vocabulary collapsing derivationally related words. In our study, we use the Snowball Stemmer for the Spanish language implemented in NLTK package [Bird et al., 2009].

3.1.5 Lemmatization (*lem*)

Lemmatization process is a complex task from Natural Language Processing that determines the lemma of a group of word forms, i.e., the dictionary form of a word. For example, the words *went* and *goes* are the verb conjugations of the verb *go*; and these words are grouped under the same lemma *go*. To apply this process, we use Freeling tool [Padró and Stanilovsky, 2012] as Spanish lemmatizer. All texts are prepared by the *Error correction* process before applying lemmatization to obtain the best results of part-of-speech identification.

Error correction Freeling is a tool for text analysis, but the assumption is that text is well-written. However, language used in Twitter is very informal, with slang, misspellings, new words, creative spelling, URLs, specific abbreviations, hashtags (which are especial words for tagging in Twitter messages), and emoticons (which are short strings and symbols that express different emotions). These problems are treated to prepare and standardize tweets for the lemmatization stage to get the best results. All words in each tweet are checked to be a valid Spanish word or are reduced according to the rules for Spanish word formation.

In general, words or tokens with invalid duplicated vowels or consonants are reduced to valid or standard Spanish words, e.g., (*ruiidoooo* → ruido (noise); *jajajaaa* → jaja; *jijijji* → jaja). We used an approach based on Spanish dictionary, a statistical model for common double letters, and heuristic rules for common interjections. In general, the duplicated vowels or consonants are removed from the target word; the resulting word is looked up in a Spanish dictionary (approximately 550,000 entries) to be validated, it is included in Freeling. For words that are not in the dictionary are reduced at least with valid rules for Spanish word formation. Also, colloquial words and abbreviations are transformed using a regular expression based on a dictionary of those sort of words, figure 2 illustrates some rules. The text on the left side of the arrow is replaced by the text of the right side. Twitter tags such as user names, hashtags (topics), URLs, and emoticons are handled as special tags in our representation to keep the structure of the sentence.

<i>tqm</i> → <i>te quiero mucho</i> (I love you so much), <i>compu</i> → <i>computadora</i> (computer).
--

Figure 2: Expansion of colloquial words and abbreviations.

In Figure 3, we can see the lemmatized text after applying Freeling. As we mentioned, the text is prepared with the Error correction step (see Figure 3(a)) then Freeling is applied to normalize words. Figure 3(b) shows Freeling’s output where each token has the original word followed by the slash symbol and its lexical information. The lexical information can be read as follows; for instance, token *orgulloso/AQOMS0* (proud) stands for adjective as part of speech (AQ), masculine gender (M), and singular number (S); the token *querer/VMIP1S0* (to want) stands for lemmatized main verb as part of speech (VM), indicative mood (I), present time (P), singular form of the first person (1S); *positive_tag/NTE0000* stands for noun tag as part of speech, and so on.

Lexical information is used to identify entities, stopwords, content words among others, it depends on the settings of the other parameters. The words are filtered based on heuristic rules that take into account the lexical information shown in Fig. 3(b). Finally, lexical information is removed in order to get the lemmatized text depicted on Figure 3(c).

3.1.6 Negation (*neg*)

Spanish negation markers might change the polarity of the message. Thus, we attached the negation clue to the nearest word, similar to the approaches used in [Sidorov et al., 2013]. A set of rules was designed for common Spanish negation structures that involve negation markers, namely, *no* (not), *nunca*, *jamás* (never), and *sin* (without). The rules are processed in order, and, when one of them matches, the remaining rules are discarded. We have two sorts of rules; it depends on the input text. If the text is not parsed by Freeling, a few rules (regular expressions) are applied to negate the nearest word to the negation marker using only the information on the

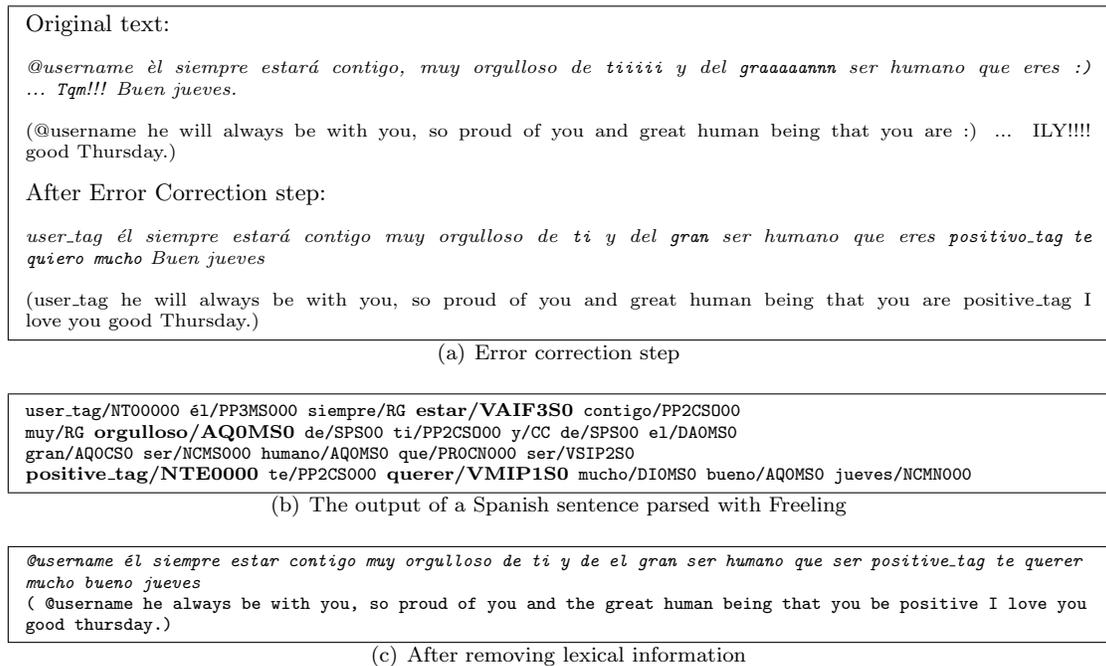


Figure 3: A step-by-step lemmatization of a tweet.

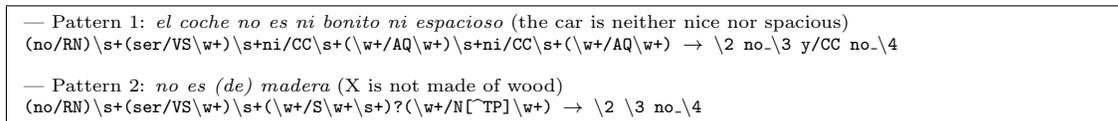


Figure 4: An example of negation rules

text, e.g., avoiding pronouns and articles. The second approach uses a set of fine-grained rules to take advantage of the lexical information, approximately 50 rules were designed considering the negation markers. The negation marker is attached to the closest word to the marker.

In the box below, Pattern 1 and Pattern 2 are examples of negation rules (regular expressions). A rule consists of two parts: the left side of the arrow represents the text to be matched, and the right side of the arrow is the structure to be replaced. All rules are based on a linguistic motivation taking into account lexical information. The set of negation rules are available³.

For example, in the sentence *El coche no es ni bonito ni espacioso* (The car is neither nice nor spacious), the negation marker *no* is attached to its two adjectives *no_bonito* (not nice) and *no_espacioso* (not spacious), as it is showed in Pattern 1, the negation marker is attached to group 3 (\3) and group 4 (\4) that stand for adjective position because of the coordinating conjunction *ni*. The number of group is identified by parenthesis in the rule from left to right. Negation markers are attached to content words (nouns, verbs, adjectives, and some adverbs), e.g., ‘*no seguir*’ (*do not follow*) is replaced by ‘*no_seguir*’, ‘*no es bueno*’ (*it is not good*) is replaced by ‘*es no_bueno*’, ‘*sin comida*’ (*without food*) is replaced by ‘*no_comida*’. Figure 4 exemplifies a pair of these negation rules.

3.1.7 Emoticon (*emo*)

In the case of emotions, we classify more than 500 popular emoticons, including text emoticons, and the whole set of emoticons (close to 1600) defined by [Unicode, 2016] into three classes: positive,

³<http://ws.ingeotec.mx/~sadit/>

:)	:D	:P	→	positive
:(:-(:'(→	negative
:-	U_U	-.-	→	neutral
emoticon without polarity			→	unicode-text

Table 1: An excerpt of the mapping table from Emoticons to its polarity words.

negative or neutral, which are replaced by a polarity word or definition associated to the emoticon according to the Unicode standard. The emoticons considered as positive are replaced by the word *positive*, negative emoticons are replaced by the word *negative*, neutral emotions are replaced by the word *neutral*. Emoticons that do not have a polarity, or are ambiguous, are replaced by the associated Unicode text. Table 1 shows an excerpt of the dictionary that maps emoticons to their corresponding polarity class.

3.1.8 Entity (*del-ent*)

We consider entities to be proper names, hashtags, urls or nicknames. However, nicknames (see *usr* parameter, Table 3) is a particular feature in Twitter messages; thus, user names is another parameter to see the effect on the classification system. User names, urls and numbers (see *url*, *num*) parameters, Table 3) could be grouped under an especial generic name. Entities such as user names and hashtags are identified directly by its corresponding especial symbol @ and #, and proper names are identified using Freeling, the lexical information used to identify a proper name is “NP0000”.

3.1.9 Word-based n-grams (*n-words*)

N-words are widely used in many NLP tasks, and they have also been used in sentiment analysis by [Sidorov et al., 2013, Cui et al., 2015]. N-words are word sequences. To compute the n-words, the text is tokenized and n-word are calculated from tokens. NLTK Tokenizer is used to identified word tokens. For example, let $T = \text{"the lights and shadows of your future"}$, its 1-words (unigrams) are each word alone, and its 2-words (bigrams) set are the sequences of two words, the set (W_2^T) , and so on. For example, let $W_2^T = \{\text{the lights, lights and, and shadows, shadows of, of your, your future}\}$, then, given a text of m words, we obtain a set with at most $m - n + 1$ elements. Generally, n-words are used up to 2 or 3-words because it is uncommon to find good matches of word sequences greater than three or four words [Jurafsky and Martin, 2009].

3.1.10 Character-based q-grams (*q-grams*)

In addition to the traditionally n-words representation, we represent the resulting text as q -grams. A q -grams is an agnostic language transformation that consists in representing a document by all its substring of length q . For example, let $T = \text{abra.cadabra}$, its 3-grams set are

$$Q_3^T = \{\text{abr, bra, ra-, a.c, _ca, aca, cad, ada, dab}\},$$

so, given text of size m characters, we obtain a set with at most $m - q + 1$ elements. Notice that this transformation handle white-spaces as part of the text. Since there will be q -grams connecting words, in some sense, applying q -grams to the entire text can capture part of the syntactic information in the sentence. The rationale of q -grams is to tackle misspelled sentences from the approximate pattern matching perspective [Navarro and Raffinot, 2002], where it is used for efficient searching of text with some degree of error.

A more elaborated example shows why the q -gram transformation is more robust to variations of the text. Let $T = \text{I.like.vanilla}$ and $T' = \text{I.lik3.vanila}$, clearly, both texts are different and a plain algorithm will simply associate a low similarity between both texts. However, after

```

original text:
pésiiiimo auto :( @autoX fallan frenos y sistema de entretenimiento; no lo compren
after text transformation:
pesim aut :( _user fal fren y sistem de entreten ; lo no_compr
computed 1-word:
{pesim, aut, :(, _user, fal, fren, y, sistem, de, entreten, ;, lo, no_compr }

```

(a) An example of configuration for INEGI benchmark for word n-grams

```

original text:
pésiiiimo auto :( @autoX fallan frenos y sistema de entretenimiento; no lo compren
after text transformation:
pesiiiimo auto _negativo _user fallan frenos y sistema de entretenimiento ; lo no_compren
computed 4-grams:
{ _pes, pesi, esii, siii, iiii, iim, iimo, imo_, mo_a, o_au, _aut, auto, uto_, to_, o_n, _ne, _neg, nega, egat, gati, ativ,
tivo, ivo_, vo_, o_u, _us, _use, user, ser_, er_f, r_fa, _fal, fall, alla, llan, lan_, an_f, n_fr, _fre, fren, reno, enos, nos_,
os_y, s_y_, _y_s, y_si, _sis, sist, iste, stem, tema, ema_, ma_d, a_de, _de_, de_e, e_en, _ent, entr, ntre, tret, rete, eten,
teni, enim, nimi, imie, mien, ient, ento, nto_, to_-, o_-;_, _;l, ;_lo, _lo_, lo_n, o_no, _no_, no_c, o_co, _com, comp, ompr,
mpre, pren, ren_ }

```

(b) An example of configuration for INEGI benchmark for q-grams (i.e., 4-grams)

Figure 5: Examples of text representation.

extracting its 3-grams, the resulting objects are more similar:

$$Q_3^T = \{I_l, _li, lik, ike, ke_, e_v, _va, van, ani, nil, ill, lla\}$$

$$Q_3^{T'} = \{I_l, _li, lik, ik3, k3_, 3_v, _va, van, ani, nil, ila\}$$

Just to fix ideas, let these two sets to be compared using the Jaccard’s coefficient as similarity, i.e.

$$\frac{|Q_3^T \cap Q_3^{T'}|}{|Q_3^T \cup Q_3^{T'}|} = 0.448.$$

These sets are more similar than the ones resulting from the original text split as words

$$\frac{|\{I, like, vanilla\} \cap \{I, lik3, vanila\}|}{|\{I, like, vanilla\} \cup \{I, lik3, vanila\}|} = 0.2$$

The assumption is that a machine learning algorithm knowing how to classify T will do a better job classifying T' using q -grams than a plain representation. This fact is used to create a robust method against misspelled words and other deliberated modifications to the text.

3.2 Examples of Text Transformation Stage

In order to illustrate the text transformation pipeline, we show the examples in Figure 5(a) and Figure 5(b). In Figure 5(a) we can see the resulting text representation for a configuration for words on INEGI bechmark, i.e., the parameters used to transform the original text into the new representation are stemming (*stem*), reduced repeated symbols up to one symbol (*del-d1*), the removal of diacritic (*del-diac*), and coarsening users (*usr*), and negations (*neg*). The final text representation is based on 1-words.

The other example, Figure 5(b), is a configuration for character 4-gram representation on the same benchmark using the following parameters: the removal of diacritic (*del-diac*), coarsening emoticons (*emo*), coarsening users (*usr*), changing words into lowercase (*lc*), negations (*neg*), and TFIDF is used to weight the tokens, it has no text representation. The final representation is based on character 4-grams, and the underscore symbol is used as space character to separate words and it is part of the token in which it appears.

4 Benchmarks and Parameter Settings

At this point, we are in the position to analyze the performance of described text representations on sentiment analysis benchmarks. In particular, we test our representations in the task of

Table 2: Datasets details from each competition tested in this work

benchmark		classes				total
name	part	positive	neutral	negative	none	
INEGI	train	2,908	986	1,110	409	5,413
	test	26,911	8,868	9,571	3,361	48,711
						54,124
TASS'15	train	2,884	670	2,182	1,482	7,218
	test	22,233	1,305	15,844	21,416	60,798
						68016

determining the global polarity —four polarity levels: positive, neutral, negative, and none (no sentiment)— of each tweet in two benchmarks.

Table 2 describes our benchmarks. The INEGI benchmark consists on tweets geo-referenced to Mexico; the data was collected and labeled between 2014 and 2015 by the Mexican Institute of Statistics and Geography (INEGI). The INEGI’s tweets come from the general population without any filtering beyond its geographic location. INEGI benchmark has a total of 54,124 tweets (in the Spanish language). The tagging process of INEGI dataset was conducted through a web application (called pioanalysis⁴, it was designed by the personnel of the Institute). Each tweet was displayed and human tagged it as positive, neutral, negative or none. After this procedure, every tweet was tagged by several humans, the label with major consensus was assigned as a final tagged. We discard tweets being on tie.

On the other hand, our second benchmark is the one used in TASS’15 workshop (Taller de Análisis de Sentimientos en la SEPLN) [Román et al., 2015]. Here, the whole corpus contains over 68,000 tweets, written in Spanish, related to well-known personalities and celebrities of several topics such as politics, economy, communication, mass media, and culture. These tweets were acquired between November 2011 and March 2012. The whole corpus was split into a training set (about 10%) and test set (remaining 90%). Each tweet was tagged with its global polarity (positive, negative or neutral) or no polarity at all (four classes in total). The tagging process was done in a semi-automatically way where a baseline machine learning algorithm classifies them, and then all the tagged tweets are manually checked by human experts; for more details of this database construction see [Román et al., 2015].

We partitioned INEGI in 10% for training and 90% for testing, following the setup of TASS’15; this large test-set pursues the generality of the method. Hereafter, we name the test set as the gold-standard, and we interchange both names as synonyms. The accuracy is the major score in both benchmarks, again because TASS’15 uses this score as its measure. We also report the macro-F1 score to help to understand the performance on heavily unbalanced datasets, see 2.

In general, both benchmarks are full of errors, and these errors vary from simple mistakes to deliberate modification of words and syntactic rules. However, it is worth to mention that INEGI is a collection of an open domain, and moreover, it comes from the general public; then we can see the frequency of misspellings and grammatical errors as a major difference between INEGI and TASS’15.

Figure 6 shows the size of the vocabulary as the number of words in the collection increases. The Heaps’ law, [Baeza-Yates and Ribeiro-Neto, 2011], states that the growth of the vocabulary follows $O(n^\alpha)$ for $0 < \alpha < 1$, for a document of size n . The figure illustrates the growth rate of our both benchmarks, along with a well-written set of documents, i.e., classic Books of the Spanish literature from the Gutenberg project [Gutenberg, 2016]. The Books collections curve is below than any of our collections; its growth factor is clearly smaller. The precise values of α for each collection are $\alpha_{\text{TASS'15}} = 0.718$, $\alpha_{\text{INEGI}} = 0.756$, and $\alpha_{\text{Books}} = 0.607$, these values were determined

⁴<http://cienciadedatos.inegi.org.mx/pioanalysis/#/login>

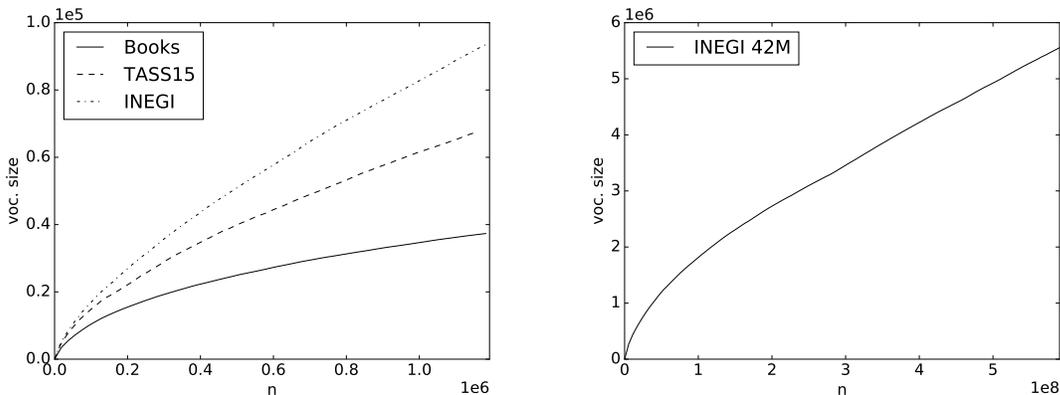


Figure 6: On the left, the growth of the vocabulary in our benchmarks and a collection of books from the Gutenberg project. On the right, the vocabulary growth in 42 million tweets.

with a regression over the formulae.⁵ There is a significant difference between the three collections, and it corresponds to the high amount of errors in TASS’15, and, the higher one in INEGI.

4.1 Parameters of the text transformations

As described in Section 3 the different text transformation methods explored in this research. Table 3 complements this description by listing the different values these transformations have. From the table, it can be observed that most parameters are either the use or absence of the particular transformation with the exceptions n -words and q -grams.

Based on the different values of the parameters, we can count the number of different text transformation which is $7 \times 2^{15} = 229,369$ configurations (the constant 7 corresponds to the number of tokenizers). Evaluating all these setups, for each benchmark, is computationally expensive. Also, we perform the same exhaustive in the test set to compare the achieved result and the best possible under our approach. Along with these experiments, we also evaluate a number of experiments to prove and compare a series of improvements. In the end, we evaluated close to one million configurations. For instance, using an Intel(R) Xeon(R) CPU E5-2630 v2 @ 2.60GHz workstation, we need ~ 12 minutes in average for a single configuration, running on a single core. Therefore, it needs roughly 24 years of computing time. Nonetheless, we used a small cluster to compute all configurations in some weeks. Notice that the time of determining the part-of-the-speech, needed by parameters *stem* and *lem*, is not reported since it was executed only once for all texts and loaded from a cache whenever is needed. The lemmatization step needs close to 56 minutes to transform the INEGI dataset in the same hardware.

5 Experimental Analysis

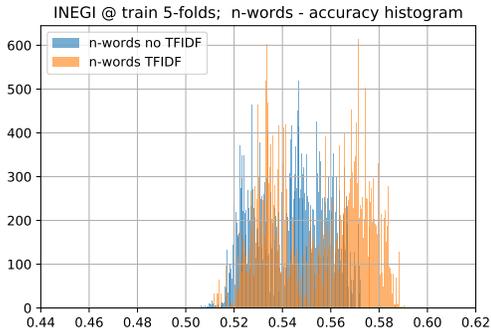
This section is devoted to describe and analyze the performance of the configuration space, provide the sufficient experimental evidence to prove that q -gram tokenizers are better than n -words, at least under the sentiment analysis domain in Spanish. Furthermore, we also provide the experimental analysis for the combination of tokenizers, which improves the whole performance without moving too far from our text classifier structure.

We use both training and test datasets in our experiments. The performance on the training set is computed using 5-fold cross validation, and the performance on test set is computed directly on the gold-standard. As previously described, training and test are disjoint sets, see Table 2 for details of our benchmarks. As mentioned, the classifier was fixed to be SVM; we use the

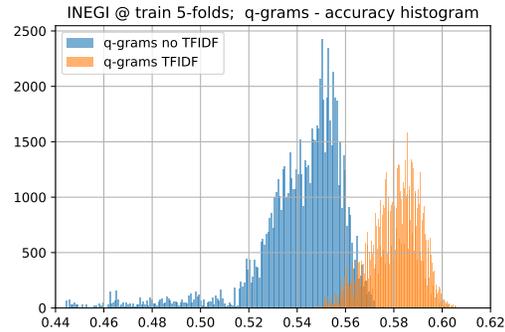
⁵The tweets were slightly normalized removing all URLs and standardizing all characters to lowercase.

Table 3: Parameter list and a brief description of their functionality

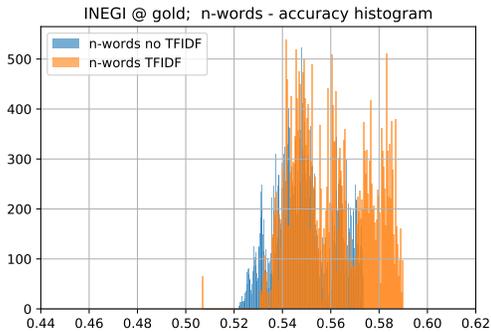
weighting schemes / removing common words		
name	values	description
tfidf	yes, no	After the text is represented as a bag of words, it determines if the vectors are weighted using the TFIDF scheme. If it is <i>no</i> then the term frequency in the text is used as weight.
del-sw	yes, no	Determines if the stopwords are removed. It is related to TFIDF in the sense that a proper weighting scheme assigns a low weight for common words.
morphological reductions		
name	values	description
lem	yes, no	Determines if words sharing a common root are replaced by its root.
stem	yes, no	Determines if words are stemmed.
transformations based on removing or replacing substrings		
name	values	description
del-punc	yes, no	The punctuation symbols are removed if <i>del-punc</i> is <i>yes</i> , they are left untouched otherwise.
del-ent	yes, no	Determines if entities are removed in order to generalize the content of the text.
del-d1	yes, no	If it is enabled then the sequences of repeated symbols are replaced by a single occurrence of the symbol.
del-d2	yes, no	If it is enabled then the repeated sequences of two symbols are replaced by a single occurrence of the sequence.
del-diac	yes, no	Determines if diacritic symbols, e.g., accent symbols, should be removed from the text.
coarsening transformations		
name	values	description
emo	yes, no	Emoticons are replaced by its expressed emotion if it is enabled.
num	yes, no	Determines if numeric words are replaced by a common identifier.
url	yes, no	Determines if URLs are left untouched or replaced by a unique url identifier.
usr	yes, no	Determines if users mentions are replaced by a unique user identifier.
lc	yes, no	Letters are normalized to be lowercase if it is enabled
handling negation words		
name	values	description
neg	yes, no	Determines if negation operators in the text are normalized and directly connected with the modified object.
tokenizing the transformation		
name	values	description
n-words	{1, 2}	Determines the number of words used to describe a token.
q-grams	{3, 4, 5, 6, 7}	Determines the length in characters of the q -grams (q).



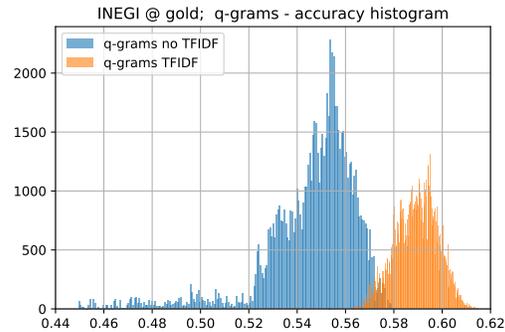
(a) Performance for n -words in training subset



(b) Performance for q -grams in training subset



(c) Performance for n -words in gold standard



(d) Performance for q -grams in gold standard

Figure 7: Accuracy’s histogram, by tokenizer’s class, for the INEGI benchmark. The performance on the training set was computed with 5-folds. We select to divide each figure to show the effect of TFIDF, which it is essential for q -grams’s performance.

implementation from the Scikit-learn project [Pedregosa et al., 2011] using a linear kernel. We use the default parameters of the library; no additional tuning was performed in this sense.

5.1 A Performance Comparison of n -words and q -grams

Figure 7 shows the histogram of accuracies for our configuration-space in both training and test partitions. Figures 7(a) and 7(c) show the performance of configurations with n -words as tokenizer (unigrams and bigrams), for training and test datasets respectively. It is possible to see that the form is preserved, and also that TFIDF configurations can perform slightly better than those using only the term frequency. However, the accuracy range being shared by both kinds of configurations is large.

In contrast, Figure 7(b) shows the performance of configurations with q -grams as tokenizers. Here, the improvement of the TFIDF class is more significant than those configurations not using TFIDF; also, the performance achieved by the q -grams with TFIDF is consistently better than the performance of the all n -word configurations in our space. This is also valid for the test dataset, see Figure 7(d).

Figure 8 shows the performance of INEGI on configurations using q -grams as tokenizers. On the left, Figures 8(a) and 8(c) show the performance of configurations without TFIDF. In train, the best performance is close to 0.57, and less than 0.58 in the test set. The best performing tokenizer is 7-grams. When TFIDF is allowed, Tables 8(b) and 8(d), the best performances are achieved, in both training and test, close to 0.61 in the training set and higher in the gold-standard. The best

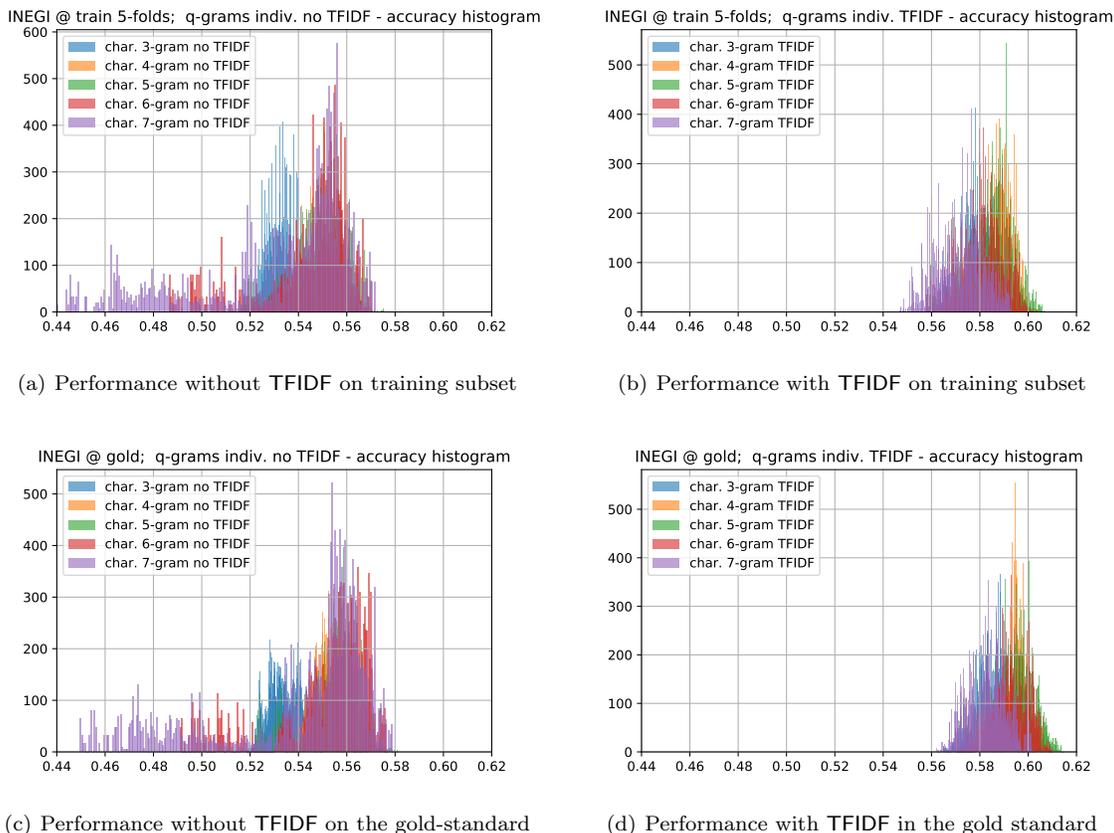


Figure 8: Accuracy’s histogram for q-gram configurations in the INEGI benchmark. As before, the performance on the training set was computed with 5-folds.

configurations are those with 5-grams and 6-grams. The 5-grams is consistently better, it achieves accuracy values of 0.6065 and 0.6148 for training and test sets, respectively.

5.1.1 Performance on the TASS’15 benchmark

The performance on TASS’15 is similar to that found in the INEGI benchmark; however, TASS’15 shows a higher sparsity of the accuracy along the range on n -words, ranging from 0.35 to close than 0.61. In the training set, the best performances are achieved using TFIDF.

The best configurations are those using q -grams, as depicted in Figure 9(b) and 9(d), where accuracy values achieve close to 0.63 in both training and test sets. In contrast to INEGI and the training set of TASS’15, the best performing q -gram tokenizer has no TFIDF, however the configurations with TFIDF are tightly concentrated which means that is more easy to pick a good configuration under a random selection, or by the insight of an expert.

Figure 10 shows a finer analysis of the performance of q -grams tokenizers in TASS’15. We can observe that 5-grams appear as the best in the training set and in the gold-standard with TFIDF, but the best performing configuration uses 6-grams tokenizers and no TFIDF; please note that TFIDF has the best accuracy on the training set, so we have not way to know this behaviour without testing all possible configurations in the gold-standard. Also, the difference between the best TFIDF and the best no-TFIDF configurations is of around 0.005; that is quite small to discard the current bias that suggest to use TFIDF configurations.

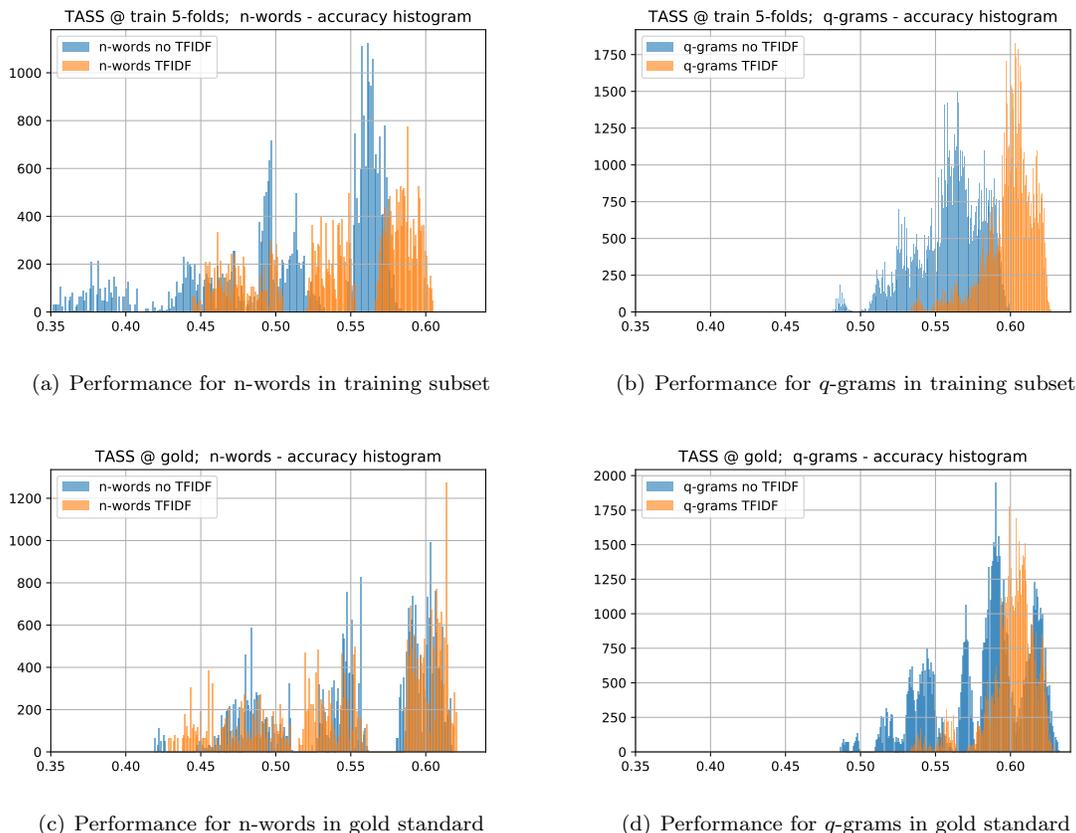


Figure 9: Accuracy’s histogram, by tokenizer’s class, for the TASS benchmark. The performance on the training set was computed with 5-folds. We select to divide each figure to show the effect of TFIDF, which it is essential for q -grams’s performance.

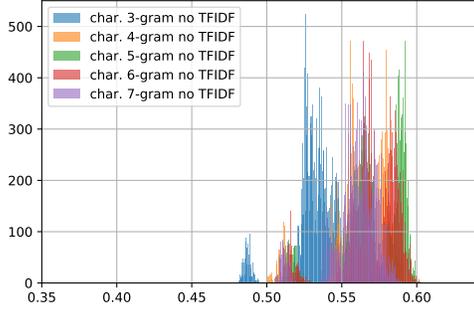
5.2 Top- k Analysis

This section focus on the structural analysis of the best k configurations (based on the accuracy score) of our previous results. We call this technique top- k analysis, and it describes the configurations with the empirical probability of a parameter to be enabled among the best k configurations. The score values are defined as the minimum among the set. The main idea is to discover patterns on the composition of best performing configurations. As we double k at each row, then k and $2k$ share k configurations which produces a smoothly convergence to 0.5 for each probability. At the best of our knowledge, this kind of analysis has never been used in the literature.

All tables in this subsection are induced by the accuracy score (i.e., best k as measured with accuracy). Also, we display the macro-F1 score as a secondary measure of performance that can help to describe the behaviour of unbalanced multi-class datasets. We omit to show the tokenizer probabilities in favor of Figures 8 and 10; please remind that almost all top configurations use q -grams.

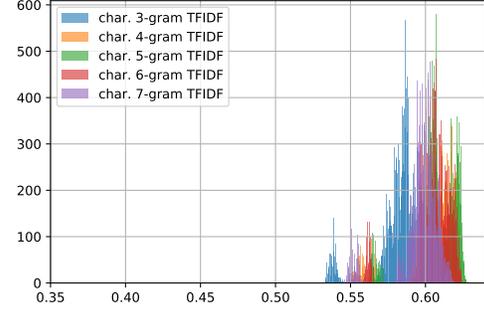
Table 4 shows the composition of INEGI’s best configurations in both training and test sets. As previously shown, almost all best setups enable TFIDF, and properly handle emoticons and users. The parameters *del-sw*, *lem*, *del-d1*, *del-d2*, *num*, and *url*, are almost deactivated in both training and test sets. The rest of the parameters (*stem*, *del-diac*, *del-ent*, and *neg*) do not remain between training and test sets. However, the later set of parameters are disabled in the gold-standard best configurations, excepting for *neg*. Such fact supports the idea that faster configurations also can produce excellent performances. Please notice that lemmatization (*lem*) and stemming (*stem*) are

TASS @ train 5-folds; q-grams indiv. no TFIDF - accuracy histogram



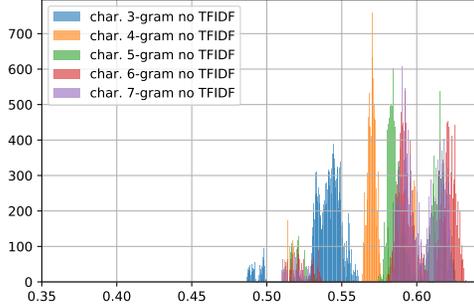
(a) Performance without TFIDF in training subset

TASS @ train 5-folds; q-grams indiv. TFIDF - accuracy histogram



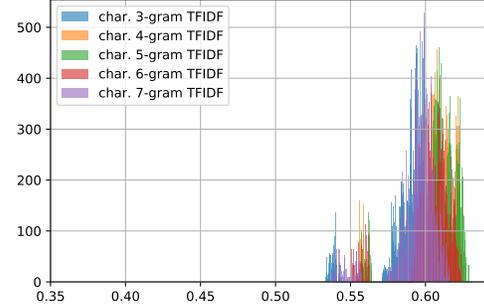
(b) Performance with TFIDF in training subset

TASS @ gold; q-grams indiv. no TFIDF - accuracy histogram



(c) Performance without TFIDF in the gold standard

TASS @ gold; q-grams indiv. TFIDF - accuracy histogram



(d) Performance with TFIDF in the gold standard

Figure 10: Accuracy's histogram for q-gram configurations in the TASS benchmark. As before, the performance on the training set was computed with 5-folds.

Table 4: Analysis of the k best configurations for the INEGI benchmark in both training and test datasets.

k	accuracy	macro-F1	tfidf	del-sw	lem	stem	del-d1	del-d2	del-punc	del-diac	del-ent	emo	num	url	usr	lc	neg
1	0.6065	0.4524	1.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	1.00	1.00	0.00	0.00	1.00	0.00	1.00
2	0.6065	0.4524	1.00	0.00	0.00	1.00	0.00	0.00	0.50	0.00	1.00	1.00	0.00	0.00	1.00	0.00	0.50
4	0.6065	0.4524	1.00	0.00	0.00	1.00	0.00	0.00	0.50	0.00	1.00	1.00	0.00	0.00	1.00	0.00	0.50
8	0.6059	0.4511	1.00	0.00	0.00	1.00	0.00	0.00	0.50	0.00	1.00	1.00	0.50	0.00	1.00	0.00	0.50
16	0.6058	0.4568	1.00	0.19	0.00	1.00	0.19	0.00	0.44	0.13	0.81	1.00	0.38	0.19	1.00	0.19	0.5625
32	0.6052	0.4507	1.00	0.31	0.00	0.69	0.25	0.00	0.47	0.19	0.56	1.00	0.31	0.38	1.00	0.44	0.5312
64	0.6047	0.4516	1.00	0.22	0.00	0.78	0.44	0.00	0.50	0.33	0.66	1.00	0.38	0.19	1.00	0.53	0.5156
128	0.6037	0.4643	1.00	0.20	0.00	0.77	0.45	0.03	0.50	0.31	0.53	1.00	0.42	0.28	1.00	0.58	0.4922
256	0.6024	0.4489	1.00	0.14	0.00	0.77	0.36	0.09	0.50	0.40	0.51	1.00	0.44	0.43	1.00	0.62	0.5078
512	0.6008	0.4315	1.00	0.17	0.00	0.73	0.42	0.17	0.50	0.43	0.41	1.00	0.41	0.48	0.99	0.62	0.5098

a) Performance on the training dataset (5-folds)

k	accuracy	macro-F1	tfidf	del-sw	lem	stem	del-d1	del-d2	del-punc	del-diac	del-ent	emo	num	url	usr	lc	neg
1	0.6148	0.4442	1.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	1.00	0.00	0.00	1.00	1.00	1.00
2	0.6148	0.4442	1.00	0.00	0.00	0.00	0.00	0.00	0.50	1.00	0.00	1.00	0.00	0.00	1.00	1.00	1.00
4	0.6136	0.4405	1.00	0.00	0.00	0.00	0.00	0.00	0.50	1.00	0.00	1.00	0.50	0.00	1.00	1.00	1.00
8	0.6135	0.4545	1.00	0.00	0.00	0.25	0.00	0.00	0.62	0.75	0.00	1.00	0.50	0.00	1.00	1.00	0.88
16	0.6134	0.4546	1.00	0.00	0.00	0.38	0.00	0.00	0.50	0.88	0.00	1.00	0.62	0.38	1.00	1.00	0.81
32	0.6130	0.4528	1.00	0.00	0.00	0.50	0.06	0.00	0.50	0.94	0.00	1.00	0.44	0.44	1.00	1.00	0.62
64	0.6119	0.4403	1.00	0.12	0.00	0.44	0.19	0.00	0.50	0.72	0.00	1.00	0.44	0.41	1.00	1.00	0.62
128	0.6112	0.4547	1.00	0.30	0.00	0.48	0.27	0.00	0.50	0.61	0.00	1.00	0.50	0.52	1.00	0.98	0.61
256	0.6099	0.4379	1.00	0.35	0.00	0.46	0.37	0.00	0.50	0.50	0.05	1.00	0.46	0.48	1.00	0.92	0.53
512	0.6083	0.4479	1.00	0.27	0.00	0.52	0.27	0.05	0.50	0.51	0.13	1.00	0.45	0.48	1.00	0.75	0.55

b) Performance on the gold-standard dataset

also disabled, which are the linguistic operations with higher computational costs in our pipeline of text transformations.

Table 5: Analysis of the k best configurations (top- k) for the TASS’15 benchmark in both training and test datasets.

k	accuracy	macro-F1	tfidf	del-sw	lem	stem	del-d1	del-d2	del-punc	del-diac	del-ent	emo	num	url	usr	lc	neg
1	0.6286	0.4951	1.00	1.00	0.00	0.00	1.00	0.00	1.00	1.00	0.00	0.00	1.00	1.00	1.00	1.00	1.00
2	0.6286	0.4951	1.00	1.00	0.00	0.00	1.00	0.00	0.50	1.00	0.00	0.00	1.00	1.00	1.00	1.00	1.00
4	0.6281	0.4947	1.00	0.75	0.00	0.00	0.75	0.00	0.50	1.00	0.00	0.50	0.75	1.00	1.00	0.75	1.00
8	0.6279	0.4895	1.00	0.50	0.00	0.00	0.50	0.00	0.50	1.00	0.00	0.50	0.50	1.00	1.00	0.50	1.00
16	0.6270	0.4864	1.00	0.38	0.00	0.00	0.38	0.06	0.44	0.88	0.00	0.50	0.50	0.88	1.00	0.50	1.00
32	0.6265	0.4884	1.00	0.25	0.00	0.00	0.34	0.06	0.47	0.69	0.00	0.38	0.59	0.62	1.00	0.62	1.00
64	0.6258	0.4852	1.00	0.20	0.00	0.00	0.42	0.22	0.50	0.62	0.00	0.48	0.56	0.48	0.94	0.61	0.88
128	0.6254	0.4862	1.00	0.20	0.00	0.00	0.46	0.27	0.48	0.67	0.00	0.56	0.59	0.38	0.81	0.68	0.77
256	0.6247	0.4846	1.00	0.21	0.00	0.12	0.38	0.32	0.50	0.66	0.02	0.47	0.60	0.42	0.77	0.73	0.69
512	0.6240	0.4848	1.00	0.14	0.00	0.24	0.38	0.36	0.50	0.65	0.02	0.47	0.61	0.42	0.77	0.67	0.63

a) Performance in the training dataset (5-folds)

k	accuracy	macro-F1	tfidf	del-sw	lem	stem	del-d1	del-d2	del-punc	del-diac	del-ent	emo	num	url	usr	lc	neg
1	0.6330	0.5101	0.00	1.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	1.00	1.00	1.00
2	0.6330	0.5101	0.00	1.00	0.00	0.00	1.00	0.00	0.50	0.00	0.00	0.00	0.00	1.00	1.00	1.00	1.00
4	0.6326	0.5099	0.00	1.00	0.00	0.00	1.00	0.00	0.50	0.00	0.00	0.50	0.00	1.00	1.00	1.00	1.00
8	0.6317	0.5104	0.00	1.00	0.00	0.00	1.00	0.00	0.50	0.00	0.00	0.25	0.00	0.50	1.00	1.00	0.75
16	0.6315	0.5082	0.00	1.00	0.00	0.00	1.00	0.00	0.50	0.38	0.00	0.38	0.00	0.38	1.00	1.00	0.88
32	0.6315	0.5069	0.00	0.69	0.00	0.12	0.78	0.16	0.47	0.38	0.00	0.56	0.00	0.62	0.69	1.00	0.81
64	0.6311	0.5071	0.00	0.62	0.00	0.12	0.83	0.19	0.48	0.39	0.00	0.66	0.12	0.56	0.62	1.00	0.81
128	0.6302	0.5061	0.00	0.56	0.00	0.25	0.80	0.25	0.48	0.43	0.00	0.72	0.22	0.62	0.56	1.00	0.77
256	0.6296	0.5054	0.00	0.38	0.00	0.27	0.69	0.34	0.50	0.47	0.00	0.73	0.23	0.62	0.38	0.94	0.65
512	0.6286	0.5048	0.06	0.39	0.00	0.34	0.68	0.42	0.50	0.46	0.02	0.75	0.38	0.57	0.36	0.80	0.66

b) Performance in the gold-standard dataset

Table 5 shows the top- k analysis for TASS’15. Again, TFIDF is a common ingredient of the majority of the better configurations in the training set; however, the best ones deactivate this parameter to use only the frequency of the term; reflected in a minimum improvement. The transformations that remain active in both training and set are *del-sw*, *del-d1*, *url*, *usr*, *lc*, and *neg*. The deactivated ones in both sets are *lem*, *stem*, *del-d2*, and *del-ent*, and *emo*. The rest of the parameters that change between training and test sets are *tfidf*, *del-diac*, and *num*. Note that as k grows, *del-punc* and *emo*, are close to be random choices. It is counterintuitive to see the *emo* parameter outside the top- k items, the same happens for the *del-ent* parameter. The *emo* parameter is used to map emoticons and emojis to sentiments, and *del-ent* is an heuristic designed to generalize the sentiment expression in the text (see Table 3). This behaviour remember us that, in the end, everything depends on the particular distribution of the dataset. In general, it is clear that there is no a rule-of-thumb to compute the best configuration. Therefore, a probabilistic approach, as it is the output of top- k analysis, is useful to reduce the cost of the exploration of the configuration space.

5.3 Improving the Performance with Combination of Tokenizers

In previous experiments, we performed an exhaustive evaluation of the configuration space; then, to improve over our results we need to modify the configuration space. Instead of adding more complex text transformations, we decide to use more than one tokenizer per configuration. More detailed, there exists 127 possible combinations of tokenizers, that is, the powerset of

$$\{2\text{-words, 1-words, 3-grams, 4-grams, 5-grams, 6-grams, 7-grams}\},$$

minus the empty set. For this experiment, we only applied the expansion of tokenizers to the best configurations found in the previous experiments, since performing an exhaustive analysis of the new configuration space becomes unfeasible. The hypothesis is that the previous best configurations will be also compose some of the best configurations in the new space, this is a fair assumption that never get worst under an exhaustive analysis.

Figure 11(a) shows the performance of 4064 configurations that correspond to all combinations of tokenizers over the top-32 configurations on the training set, see Table 4. The performance in

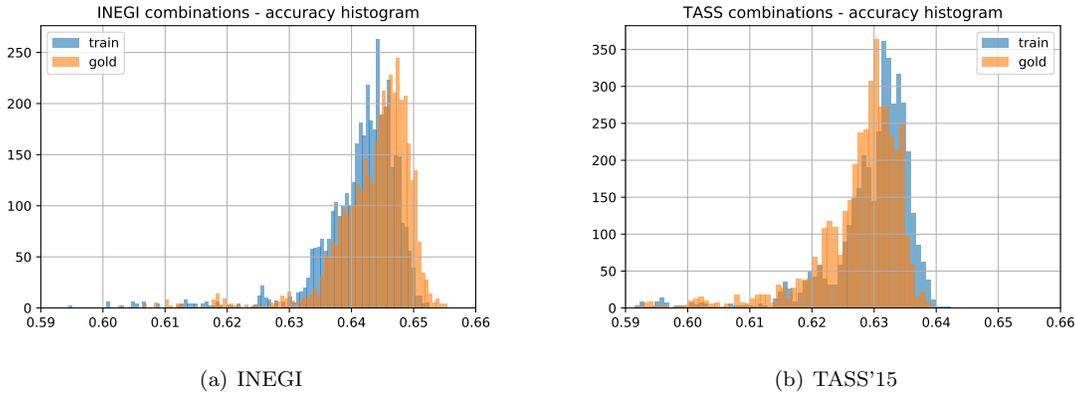


Figure 11: Accuracy’s histogram for combination of tokenizers.

Table 6: Analysis of the top- k combinations of tokenizers for both INEGI and TASS’15 benchmarks. We consider both n -words and q -grams.

INEGI										TASS 15									
k	accuracy	macro-F1	n=2	n=1	q=3	q=4	q=5	q=6	q=7	k	accuracy	macro-F1	n=2	n=1	q=3	q=4	q=5	q=6	q=7
1	0.6553	0.5287	1.00	1.00	1.00	1.00	0.00	0.00	1.00	1	0.6391	0.4997	0.00	1.00	1.00	1.00	0.00	1.00	0.00
2	0.6550	0.5270	1.00	1.00	1.00	0.50	0.00	0.00	1.00	2	0.6391	0.4995	0.00	1.00	1.00	1.00	0.00	1.00	0.00
4	0.6549	0.5281	0.50	1.00	1.00	0.75	0.00	0.00	1.00	4	0.6391	0.4997	0.00	1.00	1.00	1.00	0.00	1.00	0.00
8	0.6542	0.5268	0.63	1.00	1.00	0.62	0.00	0.00	1.00	8	0.6383	0.5020	0.00	1.00	1.00	0.75	0.50	0.75	0.00
16	0.6538	0.5263	0.75	0.94	1.00	0.75	0.00	0.00	1.00	16	0.6380	0.4966	0.25	1.00	1.00	0.75	0.38	0.63	0.13
32	0.6527	0.5241	0.66	0.84	1.00	0.59	0.00	0.06	0.88	32	0.6373	0.4972	0.18	1.00	1.00	0.63	0.50	0.69	0.19
64	0.6519	0.5235	0.56	0.77	1.00	0.52	0.09	0.06	0.84	64	0.6363	0.4940	0.30	1.00	0.94	0.75	0.55	0.58	0.17
128	0.6510	0.5258	0.65	0.61	0.99	0.55	0.18	0.25	0.78	128	0.6356	0.4937	0.32	0.97	0.94	0.77	0.53	0.66	0.26
256	0.6502	0.5205	0.61	0.64	0.97	0.58	0.22	0.31	0.79	256	0.6347	0.4927	0.39	0.96	0.88	0.81	0.56	0.66	0.35
512	0.6492	0.5172	0.62	0.66	0.96	0.55	0.30	0.40	0.74	512	0.6338	0.4954	0.42	0.92	0.86	0.73	0.57	0.56	0.39

both training and test sets is pretty similar, and significantly better than that achieved with a single tokenizer (Table 4). In Table 6 and Figure 11 we can see a significant improvement with respect to single tokenizers. The top- k analysis for the test set is listed in Table 6. In this table we focus on describe the composition of the tokenizers, instead of the text transformations. The analysis shows that 1-words, 2-words, 3-grams, 4-grams, and 7-grams are commonly present on the best configurations.

We found that TASS’15 also improves its performance under the combination of tokenizers, as Figure 11(b) illustrates. In this case, the performance in the gold standard does not surpasses the performance on the training set, as is the case of INEGI, but it is pretty close. Table 6 shows the composition of the configurations, here we can observe that best performances use 1-words, and 3-grams, 4-grams and 6-grams. It is interesting to note that 2-words are not used for the top-8 configurations, in contrast to the best configurations for INEGI.

As mentioned, any datasets will need to adjust the configuration and search for the best combination in the training set, and then, apply to their particular gold-standard. This is a costly procedure, but it is possible to reduce the search space to a sample lead by the probability models of the top- k analysis. The presented top- k analysis are particularly useful for sentiment analysis in Spanish, other languages may present different models but they are beyond the scope of this manuscript.

It is worth to mention that the best performance is high dependent of the particular dataset; however, based on Tables 4 and 5, it is interesting to note that simpler configurations are among the best performing ones when q -grams are used as tokenizers. This allows to create a model that reduces the computational cost and even improves the performance of the top-1 of both, INEGI and TASS’15, datasets with a single tokenizer. We create a configuration created by activate *tfidf*, *emo*, *num*, *usr*, and *lc*; and deactivate *del-sw*, *lem*, *stem*, *del-d1*, *del-d2*, *del-punc*, *del-diac*, *del-ent*, and *neg*. All the activated parameters are relatively simple to implement, even without

Table 7: Top- k analysis of a configuration handcrafted to reduce the computational cost.

INEGI										TASS 15										
k	accuracy	macro-F1	n=2	n=1	q=3	q=4	q=5	q=6	q=7	k	accuracy	macro-F1	n=2	n=1	q=3	q=4	q=5	q=6	q=7	
1	0.6546	0.5279	1.00	1.00	1.00	1.00	0.00	0.00	1.00	1	0.6364	0.4971	0.00	1.00	1.00	1.00	0.00	1.00	0.00	0.00
2	0.6538	0.5268	1.00	1.00	1.00	0.50	0.00	0.00	1.00	2	0.6357	0.4943	0.00	1.00	1.00	1.00	0.50	0.50	0.50	0.50
4	0.6525	0.5266	0.50	1.00	1.00	0.50	0.00	0.00	1.00	4	0.6350	0.4920	0.00	1.00	1.00	0.75	0.25	0.50	0.50	0.50
8	0.6519	0.5257	0.63	0.75	1.00	0.50	0.25	0.00	1.00	8	0.6343	0.4948	0.25	1.00	1.00	0.75	0.50	0.63	0.25	0.25
16	0.6513	0.5237	0.69	0.75	0.94	0.56	0.25	0.31	0.88	16	0.6336	0.4943	0.38	0.94	0.94	0.69	0.63	0.63	0.25	0.25
32	0.6503	0.5270	0.59	0.66	0.94	0.56	0.31	0.41	0.75	32	0.6319	0.4890	0.44	0.84	0.81	0.78	0.59	0.53	0.44	0.44
64	0.6478	0.5225	0.55	0.61	0.78	0.61	0.47	0.59	0.67	64	0.6296	0.4842	0.47	0.69	0.73	0.70	0.59	0.55	0.47	0.47
96	0.6435	0.5250	0.55	0.55	0.61	0.60	0.54	0.54	0.57	96	0.6252	0.4895	0.49	0.58	0.60	0.63	0.57	0.53	0.50	0.50
120	0.6412	0.5128	0.54	0.54	0.55	0.54	0.55	0.54	0.55	120	0.6207	0.4748	0.50	0.54	0.55	0.56	0.56	0.54	0.51	0.51
127	0.5736	0.3946	0.50	0.50	0.50	0.50	0.50	0.50	0.50	127	0.5471	0.4154	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50

external libraries. Note that leaving out stemming and lemmatization dramatically reduces many times evaluation time.

Table 7 shows the performance on the test set. The best configuration based on single tokenizer is 0.6148 and 0.6330 for INEGI and TASS’15, respectively; the best performance for combination of tokenizers is 0.6553 and 0.6391, in the same order. For our handcrafted configuration we reach an accuracy of 0.6546 for INEGI, and 0.6364 for TASS’15. This is very competitive if we take into account that the model selection is reduced to evaluate 127 configurations, and also, each evaluation is pretty fast among other alternatives.

The performances of this simple configurations are pretty close to the best possible ones with our scheme, that is, the gold-standard performance shown in Tables 4 and 5 while it can be easily implemented and optimized.

5.4 Performance Comparison on the TASS’15 Challenge

In the end, a sentiment classifier is a tool that helps to discover the opinion of a crowd of people, the effectiveness is crucial. So, there exists many researchers interested in the field, and for instance, TASS’15 ([Román et al., 2015]) is a forum that gathers many practitioners and researchers for the Spanish version of the problem. As described in §2, the problem is commonly tackled with the use of affective dictionaries, distant supervision methods to increase the knowledge database, word-embedding techniques, complex linguistic tools like lemmatizers, deep learning based classifiers, among other sophisticated techniques. Beyond the use of the SVM, there is no complex procedure that limits the adoption of our approach only to expert users.

However, the question is, how good our approach is as compared with both the state-of-the-art and the state-of-the-technique? We use the TASS’15 benchmark to answer this question. Section 2 reviews several of the best papers in the workshop. Figure 12 shows the official scores of TASS’15 participants, the best scores achieve 0.72 and the worst ones are below 0.43. The gross of the participants are between 0.59 and 0.61; there lies the best sentiment classifier based on n -words (0.6051). The best configuration that uses q -grams, as a single tokenizer, surpasses that range, i.e., 0.6330. The classifiers based on the combination of tokenizers produce a slightly better performances, and our configuration handcrafted for speed is not too distant from these performances, as figure shows.

The magnitude of the improvement is tightly linked to the dataset; for instance, as compared with the best n -words sentiment classifier, the performance of INEGI is improved in 11.17% after applying the combination of tokenizers. In the case of TASS’15, the improvement is of 5.62%, smaller but significant in any case. It is important to take into account this effect in the design of new sentiment classifiers.

6 Discussion

In this study, we covered many traditional techniques used to prepare text representations for sentiment analysis. The majority of them are too simple to be aware of their complexities. However, it is important to know its contribution to the solution of the task being tackled, as we

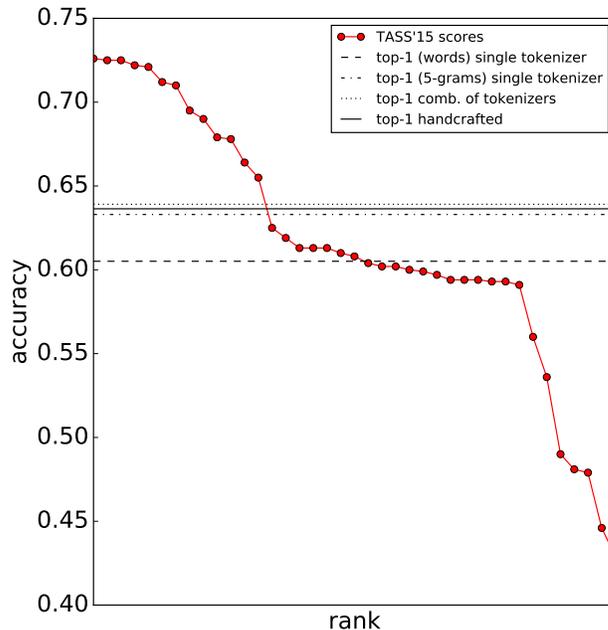


Figure 12: Comparison of our sentiment classifiers with the final scores of TASS'15.

showed, sometimes applying some technique is counterproductive. Therefore, the transformation pipeline should be carefully prepared. Other techniques, like lemmatization and stemming, are too complex to be implemented each time they are needed; therefore, a mature implementation should be used. However, as our experimental results support, for the sentiment analysis task in Spanish, there is no need to use these complex linguistic techniques if our approach, based on the combination of tokenizers, is used.

More detailed, a lemmatizer is tightly linked to the language being processed, we use Freeling by [Padró and Stanilovsky, 2012] for Spanish, and it is designed to work on mostly well-written text. The stemming procedure is another sophisticated tool, in our case, we used the Snowball for Spanish, available in NLTK package by [Bird et al., 2009]. Since it is based mostly on the removal of suffixes, then it is more robust to errors than a lemmatizer. Both techniques are computationally expensive, and both are not used by best-performing configurations; therefore, they should not be applied when the text is full of errors. This is the case of Twitter, the source of our data.

From the perspective of practitioners, the simpler approach is to find the best tokenizer's combination as applied to a set of simple setups; this gives us 127 combinations if our {2-word, 1-word, 3-gram, 4-gram, 5-gram, 6-gram, 7-gram} set is used. Supported by the patterns found in our top- k analysis, the combinations should have at least three tokenizers, and 1-words and 3-grams can always be selected. So, if the complexity of the model selection is an issue, only $\binom{5}{3} + \binom{5}{4} + \binom{5}{5} = 16$ combinations are needed.

7 Conclusions

We were able to improve the performance of our sentiment classifiers significantly. Our approach is simple; given a good *initial configuration*, we can enhance its performance using a set of tokenizers that include both n -words and q -grams. We exhaustively prove the superiority of q -grams over n -words, at least for our case of study (sentiment analysis in the Spanish language). At first

glance, large q -grams ($q = 5, 6, \text{ or } 7$) are quasi-words; however, the q -grams are sliding windows over the entire text, meaning that many times they cover the connection between two words or even three words. In relatively large words, the suffixes and prefixes are captured, when q is small, affixes and word's root are also captured. Nonetheless, this process creates many noisy substrings, and that is the reason behind our best configurations almost always use TFIDF, which weights the tokens to reduce this effect. It is necessary to produce a better process to filter out tokens that not contribute beyond creating larger vectors.

However, a naïve implementation of the multiple tokenizers will multiply the necessary memory, i.e., actually it increases the memory needs by a factor of q for q -grams. This can be a problem on very large collections. Further research is needed to solve this issue.

The *initial configuration* can be a little tricky. In this study, we provide several top- k analysis; the tables produced can be seen as probabilistic models to create good performing classifiers. These models should be valid at least for Spanish. In practice, this means that we need to evaluate the performance of a few dozens of configurations to select the best performing one among them. In a modern multicore computing architecture, this means a relatively fast procedure.

Finally, we conjectured that our approach would generalize to different languages because it works using a few language-specific techniques. However, this claim should be supported by experimental evidence. Also, we provide a list of simple rules to find a sentiment classifier based on our findings; nonetheless, the best setup is dependent of the dataset, the classes, and many others task-dependent properties. In this paper, our approach consists in performing an exhaustive evaluation of the parameter's space and then expand the search using a combination of tokenizers. We will require a faster algorithm to find good setups on large configuration's spaces that work on different languages. Finally, we want to make evident that we used SVM as classifier because of its popularity in the community, this paper mainly focuses on the treatment of the text regardless, so the proper selection and tuning of the classifier is left as future work.

Acknowledgements

The authors would like to thank the anonymous reviewers for their valuable comments and suggestions to improve the quality of this manuscript. We want to thank the *Instituto Nacional de Estadística y Geografía* (INEGI) of México for granting access to the labeled-benchmark of INEGI and its dataset of geolocated tweets, especially to Gerardo Leyva, Juan Muñoz, Alfredo Bustos, Silvia Fraustro, and Abel Coronado. We also would like to thank Julio Villena-Roman for kindly give us access to the gold-standards of TASS'15.

References

References

- [Aisopos et al., 2011] Aisopos, F., Papadakis, G., and Varvarigou, T. (2011). Sentiment analysis of social media content using n-gram graphs. In *Proceedings of the 3rd ACM SIGMM International Workshop on Social Media, WSM '11*, pages 9–14, New York, NY, USA. ACM.
- [Alam et al., 2016] Alam, M. H., Ryu, W.-J., and Lee, S. (2016). Joint multi-grain topic sentiment: Modeling semantic aspects for online reviews. *Information Sciences*, 339:206–223.
- [Almeida, 2015] Almeida, A. (2015). Deustotech internet at tass 2015: Sentiment analysis and polarity classification in spanish tweets. *Comité organizador*, page 23.
- [Alvarez-López et al., 2015] Alvarez-López, T., Juncal-Martínez, J., Gavilanes, M. F., Costa-Montenegro, E., González-Castano, F. J., Cerezo-Costas, H., and Celix-Salgado, D. (2015). Gti-gradient at tass 2015: A hybrid approach for sentiment analysis in twitter. In *TASS SE-PLN*, pages 35–40.

- [Araque et al., 2015] Araque, O., Corcuera, I., Román, C., Iglesias, C. A., and Sánchez-Rada, J. F. (2015). Aspect based sentiment analysis of spanish tweets. In *TASS SEPLN*, pages 29–34.
- [Baeza-Yates and Ribeiro-Neto, 2011] Baeza-Yates, R. A. and Ribeiro-Neto, B. A. (2011). *Modern Information Retrieval*. Addison-Wesley, 2nd edition.
- [Bird et al., 2009] Bird, S., Klein, E., and Loper, E. (2009). *Natural Language Processing with Python*. O’Reilly Media.
- [Borda and Saladich, 2015] Borda, I. M. and Saladich, J. C. (2015). Bittenpotato: Tweet sentiment analysis by combining multiple classifiers. *Comite organizador*, pages 71–72.
- [Cámara et al., 2015] Cámara, E. M., Cumbreras, M. Á. G., Martín-Valdivia, M. T., and López, L. A. U. (2015). Sinai-emma: Vectores de palabras para el análisis de opiniones en twitter. In *TASS SEPLN*, pages 41–46.
- [Casasola Murillo and Marín Raventós, 2016] Casasola Murillo, E. and Marín Raventós, G. (2016). Evaluación de modelos de representación de texto con vectores de dimensión reducida para análisis de sentimiento. In *Proceedings of TASS 2016: Workshop on Sentiment Analysis at SEPLN*, volume 1702, pages 23–28. CEUR Workshop Proceedings.
- [Cerón-Guzmán and de Cali, 2016] Cerón-Guzmán, J. A. and de Cali, S. (2016). Jacerong at tass 2016: An ensemble classifier for sentiment analysis of spanish tweets at global level. In *Proceedings of TASS 2016: Workshop on Sentiment Analysis at SEPLN*, page 35.
- [Cui et al., 2015] Cui, Z., Shi, X., and Chen, Y. (2015). Sentiment analysis via integrating distributed representations of variable-length word sequence. *Neurocomputing*.
- [de la Vega and Veá-Murguía, 2015] de la Vega, M. and Veá-Murguía, J. (2015). Ensemble algorithm with syntactical tree features to improve the opinion analysis. *Comité organizador*, pages 53–54.
- [Deas et al., 2015] Deas, M. S., Biran, O., McKeown, K., and Rosenthal, S. (2015). Spanish twitter messages polarized through the lens of an english system. In *TASS SEPLN, CEUR Workshop Proceedings*, pages 81–86.
- [Díaz-Galiano and Montejo-Ráez, 2015] Díaz-Galiano, M. and Montejo-Ráez, A. (2015). Participación de sinai dw2vec en tass 2015. In *Proceedings of the Sentiment Analysis Workshop at SEPLN (TASS2015)*, pages 59–64.
- [Ghiassi et al., 2013] Ghiassi, M., Skinner, J., and Zimbra, D. (2013). Twitter brand sentiment analysis: A hybrid system using n-gram analysis and dynamic artificial neural network. *Expert Systems with Applications*, 40(16):6266 – 6282.
- [Giannakopoulos et al., 2012] Giannakopoulos, G., Mavridi, P., Paliouras, G., Papadakis, G., and Tserpes, K. (2012). Representation models for text classification: A comparative analysis over three web document types. In *Proceedings of the 2Nd International Conference on Web Intelligence, Mining and Semantics, WIMS ’12*, pages 13:1–13:12, New York, NY, USA. ACM.
- [Gutenberg, 2016] Gutenberg (2016). Gutenberg project. In <https://www.gutenberg.org/>.
- [Han et al., 2013] Han, Q., Guo, J., and Schuetze, H. (2013). Codex: Combining an SVM classifier and character n-gram language models for sentiment analysis on twitter text. In *Proceedings of the 7th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2013, Atlanta, Georgia, USA*, pages 520–524.
- [Joachims, 1998] Joachims, T. (1998). Text categorization with support vector machines: Learning with many relevant features. In *European conference on machine learning*, pages 137–142. Springer.

- [Jurafsky and Martin, 2009] Jurafsky, D. and Martin, J. H. (2009). *Speech and Language Processing (2Nd Edition)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- [Khan et al., 2016] Khan, F. H., Qamar, U., and Bashir, S. (2016). Sentimi: Introducing point-wise mutual information with sentiwordnet to improve sentiment polarity detection. *Applied Soft Computing*, 39:140–153.
- [Liu, 2015] Liu, B. (2015). *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. Cambridge University Press. ISBN: 1-107-01789-0. 381 pages.
- [Navarro and Raffinot, 2002] Navarro, G. and Raffinot, M. (2002). *Flexible Pattern Matching in Strings – Practical on-line search algorithms for texts and biological sequences*. Cambridge University Press. ISBN 0-521-81307-7. 280 pages.
- [O’Connor et al., 2010] O’Connor, B., Krieger, M., and Ahn, D. (2010). Tweetmotif: Exploratory search and topic summarization for twitter. In *ICWSM*, pages 384–385.
- [Padró and Stanilovsky, 2012] Padró, L. and Stanilovsky, E. (2012). Freeling 3.0: Towards wider multilinguality. In *LREC2012*.
- [Padró and Stanilovsky, 2012] Padró, L. and Stanilovsky, E. (2012). Freeling 3.0: Towards wider multilinguality. In *Proceedings of the Language Resources and Evaluation Conference (LREC 2012)*, Istanbul, Turkey. ELRA.
- [Pang and Lee, 2008] Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2):1–135.
- [Pedregosa et al., 2011] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- [Peng et al., 2008] Peng, T., Zuo, W., and He, F. (2008). Svm based adaptive learning method for text classification from positive and unlabeled documents. *Knowledge and Information Systems*, 16(3):281–301.
- [Quirós et al., 2016] Quirós, A., Segura-Bedmar, I., and Martinez, P. (2016). Labda at the 2016 tass challenge task: using word embeddings for the sentiment analysis task. In *Proceedings of TASS 2016: Workshop on Sentiment Analysis at SEPLN*, page 29.
- [Razavi et al., 2014] Razavi, A. H., Matwin, S., Koninck, J. D., and Amini, R. R. (2014). Dream sentiment analysis using second order soft co-occurrences (SOSCO) and time course representations. *J. Intell. Inf. Syst.*, 42(3):393–413.
- [Román et al., 2015] Román, J. V., Morera, J. G., Ángel García Cumberras, M., Cámara, E. M., Valdivia, M. T. M., and López, L. A. U. (2015). Overview of tass 2015. *CEUR Workshop Proceedings*, 1397:13–21.
- [Saif et al., 2016] Saif, H., He, Y., Fernandez, M., and Alani, H. (2016). Contextual semantics for sentiment analysis of twitter. *Information Processing & Management*, 52(1):5–19. Emotion and Sentiment in Social and Expressive Media.
- [Sammut and Webb, 2011] Sammut, C. and Webb, G. I. (2011). *Encyclopedia of Machine Learning*. Springer Publishing Company, Incorporated, 1st edition.
- [Sidorov et al., 2013] Sidorov, G., Miranda-Jiménez, S., Viveros-Jiménez, F., Gelbukh, A., Castro-Sánchez, N., Velásquez, F., Díaz-Rangel, I., Suárez-Guerra, S., Treviño, A., and Gordon, J. (2013). Empirical study of machine learning based approach for opinion mining in tweets. In *Proceedings of the 11th Mexican International Conference on Advances in Artificial Intelligence - Volume Part I, MICAI’12*, pages 1–14, Berlin, Heidelberg. Springer-Verlag.

- [Tripathy et al., 2016] Tripathy, A., Agrawal, A., and Rath, S. K. (2016). Classification of sentiment reviews using n-gram machine learning approach. *Expert Syst. Appl.*, 57:117–126.
- [Turney, 2002] Turney, P. D. (2002). Thumbs up or thumbs down?: Semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, pages 417–424, Stroudsburg, PA, USA. Association for Computational Linguistics.
- [Unicode, 2016] Unicode (2016). Unicode emoji chart.
- [Valverde et al., 2015] Valverde, J., Tejada, J., and Cuadros, E. (2015). Comparing supervised learning methods for classifying spanish tweets. In *TASS SEPLN, CEUR Workshop Proceedings*, volume 1397, pages 87–92.
- [Wu et al., 2016] Wu, F., Huang, Y., and Song, Y. (2016). Structured microblog sentiment classification via social context regularization. *Neurocomputing*, 175:599–609.
- [y Ferran Pla, 2016] y Ferran Pla, L.-F. H. (2016). Elirf-upv en tass 2016: Análisis de sentimientos en twitter. In *Proceedings of TASS 2016: Workshop on Sentiment Analysis at SEPLN*, pages 47–51.