# Improving Subjectivity Detection for Spanish Texts using Subjectivity Word Sense Disambiguation based on Knowledge

Marco Antonio Sobrevilla Cabezudo, Nora La Serna Palomino, Rolando Maguiña Perez

Facultad de Ingeniería de Sistemas e Informática
Universidad Nacional Mayor de San Marcos
Lima, Perú
msobrevillac@gmail.com, {nlasernap, rmaguinap}@unmsm.edu.pe

*Abstract*—In this paper, we present a Sentence-level Subjectivity Detection method for Spanish using Subjectivity Word Sense Disambiguation (SWSD) based on Knowledge. We use a classic method of Word Sense Disambiguation, using the Spanish WordNet included in Mutlilingual Central Repository 3.0 and the WordNet-Pr as Knowledge base. Because of the alignment between the WordNet and the SentiWordNet, we use this latter as semantic resource annotated with polarity values to determine when a word expresses subjectivity and objectivity, defining subjectivity levels using a fuzzy clustering algorithm previously. Due to the few resources focused on Sentiment Analysis for Spanish, the Semcor corpus was used for analyzing the attributes to be used. Finally, a Rule-based classifier was created to detect subjective sentences. This method was executed over a Spanish corpus, created in this work. The results show that our approach contributes positively to Subjectivity Detection task, despite of using resources created for English.

*Keywords—Sentiment Analysis, Subjectivity Detection, Subjectivity Word Sense Disambiguation*

## I. Introduction

Sentiment Analysis aims at analyzing the opinions about a product or entity [1]. It can be divided into 3 analysis levels: (1) Document-level, classifying a text (which contains one or more sentences) into neutral, positive or negative; (2) Sentence-level, classifying sentences into positive and negative; and (3) Feature-level, analyzing and classifying opinions about the features of a determined entity. The Sentence-level Sentiment Analysis task can be divided into 2 sub-tasks: the first, subjectivity detection and, the second, the polarity analysis.

Subjectivity detection is the task responsible for determining if a sentence expresses an opinion, emotion, evaluation, etc. [2]. This task has demonstrated to be useful to other Natural Language Processing task as Question-Answering systems [3], Opinion Summarization systems [4] and Information Retrieval systems [5].

The importance of the Subjectivity Detection task is mentioned in [6], arguing that the problem of distinguishing between subjective and objective sentences has demonstrated to be more difficult than the next problem, i.e., polarity classification, because a sentence could not be subjective (not expressing an opinion) and be classified as a positive or negative sentence. This makes us to think that improvements in

Subjectivity Detection could benefit to Polarity Classification task positively.

For English, there are a lot of studies about Subjectivity Detection using different techniques and approaches [1]. Classic methods of Subjectivity Detection rely in subjectivity lexicons for determining if a sentence is subjective or objective. This focus shows some problems. For example, two sentences are presented below:

- *Ese niño es un **dolor de cabeza**.* (That child is a headache).

- *Uno de los síntomas del resfriado es el **dolor de cabeza**.* (One of cold symptoms is headache)

If a subjective lexicon contains the expression "dolor de cabeza", a classic method could determine that the first sentence is subjective and the second too, however, the second sentence is objective. This problem occurs because subjective lexicons associate subjectivity to words, instead of sense words.

To solve this problem, some sense-focused sentimental lexicons have been created. Examples of these are SentiWordNet [7], WN-Affect [8] and Micro-WordNet Opinion [9].

SentiWordNet is a lexical resource which contains all of senses included in the WordNet-Pr [10] annotated with polarity values (positive, negative and objective). WordNet-Pr is the most used sense repository in Word Sense Disambiguation task. This repository contains nouns, verbs, adjectives and adverbs organized into a synonym set, called synsets, which represents a sense of a word.

For Spanish, there are few studies about Subjectivity Detection task and some resources focused on Sentiment Analysis but these are not focused on sense words. An example of lexicon is proposed in [11]. This lexicon contains 2 036 words marked with probability to be associated at least one basic emotion, like joy, anger, fear, sadness, surprise, and disgust. Recently, Word Sense Disambiguation methods have been used to help Subjectivity Detection, carrying improvements in this task [12][13][14].

In this work, a Rule-based Subjectivity Detection method is proposed and a Subjectivity Word Sense Disambiguation (SWSD) method based on Knowledge is used to support the

Subjectivity Detection task for Spanish sentences. The SWSD method uses a Word Sense Disambiguation algorithm based on graphs to obtain the senses for each content word in a sentence. This method uses the Multilingual Central Repository 3.0 (MCR3) [15], which includes a WordNet for Spanish, as sense repository, and the WordNet-Pr 3.0 as knowledge base. Then, a mapping from MCR3 to SentiWordNet is performed to obtain the subjectivity levels for each word. Due to few of corpora for Spanish, the Semcor corpus is used to extract information about attributes and parameters used in the proposed method. Finally, a Rule-based method is used to classify Spanish sentences into subjective and objective over a Spanish corpus created in this work. Additionally, an experiment over a Spanish corpus composed by informal text was performed.

Some results of this work are the followings: the Word Sense Disambiguation task contributes positively to Subjectivity Detection and the use of resources developed for other languages (in this case, for English) can be useful to develop methods in Spanish or other languages. Other contribution of this work is the creation of a subjectivity corpus in Spanish.

The paper is organized as follows. Related works are presented in Section II. Section III describes the creation of FilmAffinity corpus. The pre-processing of resources in this work is described in Section IV. Section V is dedicated to describing our proposal, whereas Section VI contains the descriptions and results analysis of the experiments.. Section VII presents an experiment performed in an informal corpus. Finally, conclusions and final remarks are presented in Section VIII.

## II. RELATED WORKS

A study that examines the effects of adjectives on subjectivity detection is presented in [16]. In that work, a lexicon of adjectives was created using log-linear function in a corpus. This lexicon contains a set of semantically oriented adjectives and a set of gradable adjectives. The method classifies a sentence as subjective if the sentence contains one of adjectives in the lexicon as least. The results obtained were that the method obtained a high precision (0.70) and the authors concluded that adjectives and its variations are good indicators of subjectivity.

Two bootstrapping methods are adopted in [17] for creating a list of subjective nouns. Then, the authors used the Naïve Bayes algorithm with two configurations for detecting subjectivity: (1) using set of attributes used in the literature and (2) using the same set of attributes and adding the list of subjective nouns. The results showed that incorporating the list of subjective nouns, the method obtained a better performance.

A method that evaluates the improvements into a subjective classifier and an objective classifier using syntactic patterns is proposed in [18]. In that method, firstly, the Subjectivity Detection method proposed in [17] was executed over an unannotated corpus to obtain a set of subjective and objective sentences. Then, the algorithm proposed in [17] was used to extract syntactic patterns from the annotated corpus. Finally, two Rule-based classifiers were created (one for subjectivity detection and other for objectivity detection). These classifiers used a subjectivity list and the obtained syntactic patterns.

The results showed that the use of syntactic patterns carried improvements into the methods.

A method that uses Fuzzy Sets was proposed in [19]. In this method, the log-linear function and the Fuzzy Set Theory were used to extract subjective words of a corpus and classify these into highly subjective word, subjective word and lowly subjective word. Then, a Rule-based classifier was executed using this subjective lexicon. The results were that the use of Subjectivity Level showed a better performance.

A study about the relation between subjectivity and word senses is introduced in [12]. The study demonstrates that the word senses can be annotated with subjectivity and that subjective knowledge can be improving the performance of Word Sense Disambiguation methods.

The Subjective Word Sense Disambiguation (SWSD) is introduced in [13], this task consists in automatically determining which words in a piece of text are being used in subjective senses and which are being used in objective senses.

An unsupervised Subjective Word Sense Disambiguation method for English is presented in [14]. This method uses a clustering-based Word Sense Disambiguation method to obtain the senses of the words included in a sentence. Then, the subjectivity levels are obtained from the word senses, using the SentiWordNet. Finally, a Rule-based classifier was created. This classifier identifies a sentence as subjective if the sentence contains one highly subjective word or two subjective words. The authors demonstrated that the use of disambiguation methods is useful to improve the currently methods.

## III. FILMAFFINITY CORPUS

Due to the few resources related to Subjectivity Detection in Spanish, the creation of a Spanish corpus which contains subjective and objective sentences was performed in this work.

The methodology proposed in [20] was used for the creation of this corpus, called FilmAffinity corpus. This methodology consists in: (1) extracting sentences from a website, (2) changing all sentences to lowercase and remove sentences or snippets that have less than 10 tokens, and (3) grouping by type of sentence (subjective or objective).

For this work, a sentences set was extracted from Spanish version of the FilmAffinity website[1] [2]. In this process, the movie summaries were used to extract objective sentences and the user reviews to extract subjective sentences. Then, the sentences set (subjective and objective sentences) was processed using the step 2 of the proposed methodology, and finally, a set of 2500 objective sentences and 2500 subjective sentences was selected.

Two sentences (1 subjective and 1 objective, respectively) extracted from FilmAffinity corpus are shown below.

- *Obra maestra del cine de la provocación, áspera y visceral.* (Cinematic masterpiece of provocation, rough and visceral.)

---

[1] Available in http://www.filmaffinity.com/es/main.html
[2] FilmAffinity is a website which shows information about movies in different languages.

- *Un policía desanimado se dispone a resolver el asesinato de un colega de la policía que había sido su mejor amigo.* (A dejected police are solving the murder of a police colleague who had been his best friend.)

## IV. PRE-PROCESSING OF RESOURCES

A previous step to the application of our proposal was the pre-processing of the resources and the creation of attributes to be used. Thus, the SentiWordNet was processed to handle subjectivity and objectivity scores, whereas Semcor corpus supported to obtain the parameters for our proposal.

### A. Processing the SentiWordNet

For the processing of the SentiWordNet, firstly, we defined a set of subjectivity levels. In this case, 4 sets (Highly Subjective, Subjective, Lowly Subjective and Objective) were proposed.

Considering that the SentiWordNet contains all of senses included in the WordNet-Pr annotated with polarity values (positive, negative and objective), the first step was to convert these polarity values into subjectivity values. For this, the sum of positive and negative values for a sense was identified as subjective value and the last value, as objective value.

After defining the subjectivity levels and getting the subjectivity values of all synsets, we associated every synset to subjectivity levels. The last level (objective level) was easy to get because it occurs when the subjective score is zero. The first three levels were difficult to distinguish, for this reason, the fuzzy c-means algorithm was proposed to identify the subjective sets. The fuzzy c-means is a clustering algorithm that does not identify the total belonging for elements to every set, instead this, defines membership degrees for every set. The parameters used in the Fuzzy c-means were: 2 for fuzzification parameter (used in the literature frequently) value and the minimum error was 0.0009.

In order to obtain the subjectivity level for every synset, the Principle of Maximum Membership was applied, i.e., the subjectivity level that showed the maximum membership degree was selected. Some sense examples are shown in Table I. The first of this examples shows that the maximum membership degree belongs to the High Subjectivity level (HS: 0.50), therefore, the synset 01586752-{good} is classified as HS.

TABLE I: MEMBRESHIP DEGREES FOR SUBJECTIVITY LEVELS SET, LOWLY SUBJECTIVE (LS), SUBJECTIVE (MS) AND HIGHLY SUBJECTIVE (HS) OF SYNSETS AND SUBJECTIVITY LEVEL SELECTED

| Synset | LS | MS | HS | Set |
|---|---|---|---|---|
| 01586752-{good} | 0.21 | 0.29 | 0.50 | HS |
| 00801125-{war} | 0.60 | 0.29 | 0.11 | LS |

After the process, the synets in the SentiWordNet were distributed as below (shown in Table II):

TABLE II: SENSES HIGHLY SUBJECTIVE (HS), SUBJECTIVE (MS), LOWLY SUBJECTIVE (LS) AND OBJECTIVE (O) IN SENTIWORDNET

| HS | MS | LS | O |
|---|---|---|---|
| 13783 | 8297 | 7015 | 88564 |

### B. Analyzing the Semcor Corpus

Before analyzing the corpus, the attributes identification was necessary. Thus, the attributes were defined as the union of open morphosyntactic classes (noun or N, verb or V, adjective or A, and adverb or R) with the subjectivity levels (Highly Subjective, Subjective and Lowly Subjective), creating 12 attributes.

Due to the few resources for Spanish, the Semcor corpus was used to extract the information related to attributes to be used into the classifier. The reason to use this corpus was that Semcor is a sense-annotated corpus, thus, this enabled the building of Subjectivity Word Sense Disambiguation methods.

The way to analyze this corpus is similar to the used in [14] and is described as follow: firstly, OpinionFinder tool [21] was applied over a subset of the Semcor to subjectively annotate the sense annotation presented in this corpus. OpinionFinder was selected because it obtains better results on precision (91.7%) in the Subjectivity Detection task. Secondly, the content words included in a sentence were grouped into the defined attributes, and then, a set of all sentences included in the subset of Semcor was analyzed by a feature selection method. The Gain Ratio algorithm [22] from the Weka tool[3] was used to obtain the subjective weight of all attributes.

The results of the feature selection method are presented in the Table III. All attributes presented different weights, the Highly Subjective verbs and the Highly Subjective adjectives presented the highest contribution to identify a subjective sentence, and the Lowly Subjective verbs and the Highly Subjective adverbs presented the lowest contribution.

TABLE III: SUBJECTIVE WEIGHTS FOR ALL ATTRIBUTES

| Attribute | Subjective Weight |
|---|---|
| $A_{LS}$ | 0.0566 |
| $A_{MS}$ | 0.0885 |
| $A_{HS}$ | 0.1499 |
| $V_{LS}$ | 0.0357 |
| $V_{MS}$ | 0.0680 |
| $V_{HS}$ | 0.0730 |
| $N_{LS}$ | 0.0509 |
| $N_{MS}$ | 0.0816 |
| $N_{HS}$ | 0.1244 |
| $R_{LS}$ | 0.0606 |
| $R_{MS}$ | 0.0640 |
| $R_{HS}$ | 0.0488 |

## V. OUR PROPOSAL

The proposed method uses a Word Sense Disambiguation method for helping the Subjectivity Detection task. This approach is used because, as mentioned in [14], the subjectivity is
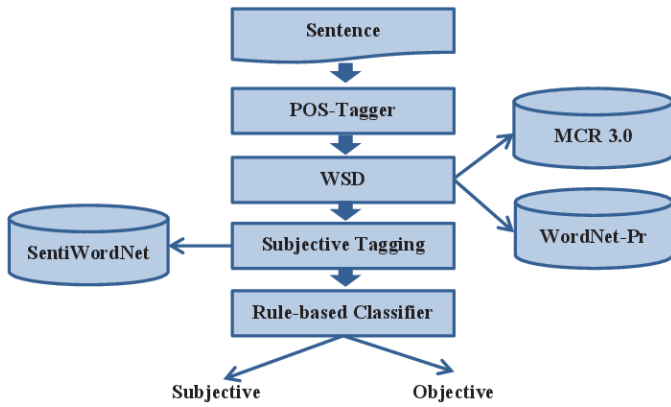
---

[3]Available in http://www.cs.waikato.ac.nz/ml/weka/

Fig. 1: Architecture of the Subjective Classifier



Fig. 2: Graph used by method proposed in [24]



Fig. 3: Method to get the synsets of the WordNet-Pr

more related the sense words than words. The Figure 1 shows the architecture of the Subjective Classifier.

This method works as follow: Firstly, the sentence is POS-Tagged, using TreeTagger tool [23] with the Spanish model. Then, a Knowledge-based Word Sense Disambiguation algorithm is applied to disambiguate all content words (nouns, verbs, adjectives and adverbs), using the Spanish WordNet, including in the Multilingual Central Repository 3.0, as sense repository. Then, when all content words are disambiguated, subjectivity levels are obtained by mapping between the Spanish WordNet and the processed SentiWordNet. Finally, subjectivity levels of all words are grouped into the attributes and these are introduced into the subjectivity classifier, which determines if the sentence is subjective or objective.

*A. Word Sense Disambiguation for Spanish*

The proposed Subjective Word Sense Disambiguation method (SWSD) follows a fine-grained approach. This method uses a Word Sense Disambiguation (WSD) method and a mapping to the processed SentiWordNet, separately.

The used WSD method was similar to the proposed in [24][4]. This method disambiguates all content words included in a sentence using whole WordNet-Pr graph with gloss-tag as knowledge resource, and the PageRank algorithm [25] to rank the senses associated to words to be disambiguated. This method works as below: firstly, this method obtains the synsets of all content words. Secondly, the method executes the PageRank algorithm, considering the probability mass into the mentioned synsets. Finally, the method selects the synset with the highest score in the graph. An example of this method for the sentence *"La película fue interesante"* is presented in the Figure 2.

Due to the target language was Spanish and the Knowledge Base used for the method execution, i.e., the WordNet-Pr, was created for English, a previous step was necessary to obtain the English synsets. The Figure 3 shows this step: for each word included in a sentence was obtained all senses (represented by synsets) contained in the Spanish WordNet, including the MCR 3.0. Finally, these senses were introduced into the WSD algorithm.
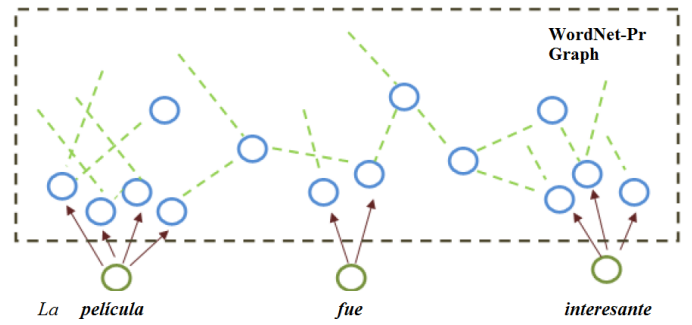
[4]Available in http://ixa2.si.ehu.es/ukb/

As the proposed method for Subjectivity Detection needs the subjectivity level for each word, processed SentiWordNet in Subsection A of Section IV was used for mapping the senses obtained by the WSD algorithm for each word (in a sentence) to subjectivity levels.

*B. Subjectivity Sentence Classifier*

We use a Rule-based classifier to classify sentences into subjective or objective. This method is similar to that proposed in [18]. In this method, every word of the sentence is disambiguated using the algorithm proposed in [24] and then a weight is assigned depending on the attribute to which it belongs. If the sum of all weights is greater than a threshold, the sentence is classified as subjective. Equation 1 is used for classifying a sentence:

$$SD(f) = \begin{cases} subjective & \text{if } \sum n_{i=1} weight(w_i) \geq \lambda \\ objective & \text{in other case} \end{cases} \quad (1)$$

Where $w_i$ is the sense word used in the sentence $f$ and:

$$weight(w_i) = \begin{cases} 0.0357 & \text{if } w_i \text{ is } V_{LS} \\ 0.0680 & \text{if } w_i \text{ is } V_{MS} \\ 0.0730 & \text{if } w_i \text{ is } V_{HS} \\ 0.0509 & \text{if } w_i \text{ is } N_{LS} \\ 0.0816 & \text{if } w_i \text{ is } N_{MS} \\ 0.1244 & \text{if } w_i \text{ is } N_{HS} \\ 0.0566 & \text{if } w_i \text{ is } A_{LS} \\ 0.0885 & \text{if } w_i \text{ is } A_{MS} \\ 0.1499 & \text{if } w_i \text{ is } A_{HS} \\ 0.0606 & \text{if } w_i \text{ is } R_{LS} \\ 0.0640 & \text{if } w_i \text{ is } R_{MS} \\ 0.0488 & \text{if } w_i \text{ is } R_{HS} \end{cases} \qquad (2)$$

The best value of $\lambda$ is 0.220; this value was obtained from an experimental analysis in the Semcor corpus processed in Subsection B of the Section IV.

## VI. Results and Discussions

The realized experiments were conducted in order to evaluate 2 hypotheses. Firstly, SWSD based on Knowledge has a positive impact over Subjectivity Detection task. Secondly, to evaluate the performance changes of Subjectivity Detection method between Spanish and English, believing that English resources are useful for methods implemented in other languages.

In our experiments, we use two corpora: the FilmAffinity corpus (created in this work) and, the Movie Review Dataset [20]. The Movie Review Dataset is a movie-domain corpus for English (as our corpus), which contains 5000 movie summaries, annotated as objective sentences, and 5000 user reviews, annotated as subjective sentences. In order to construct a baseline which to evaluate the impact of applying SWSD on the proposed method, the same method, but without applying WSD, was used. In this baseline, the subjectivity score was defined as the mean of subjective scores of all senses for a word. Then, the subjective level was obtained from the applying of subjective score on the fuzzy membership function defined in the Fuzzy C-means, as it was done for senses in SentiWordNet. Other used baseline was the Most Frequent Sense method (MFS). This method uses the first sense of a word as the select sense. The results of the methods are shown in Table IV and Table V (the results in bold are the best).

As it can be seen in Table IV, the best results were obtained by the MFS method. It is an expected result because the MFS method is a WSD method difficult to be outperformed by other Knowledge-based methods. The results of the proposal method outperform the baseline that does not use WSD, generally (when these are evaluated by average F-Measure and accuracy).

One point to be noted is that performance in subjective precision of the proposed method is better than the baseline. This is good for our method because we expected getting better results in Subjectivity Detection, thus, this confirms our hypothesis that WSD carries improvements into Subjectivity Detection methods. In the case of objective precision the opposite occurs. The same way, this occurs with recall measure. The

F-measure obtained for subjective sentences in the baseline is slightly different from that obtained in the proposed method. In the case of F-measure for objective sentences, the proposed method far outperforms the baseline.

In Table V, it can be seen that proposed method outperforms the baseline too. One important result is obtained comparing the results in Table IV and Table V, it can be noted that the results for Spanish sentences are better than English sentences. This verifies the hypothesis that can be used English resources for creating classifiers in Spanish, or maybe, in other languages, without negatively impact.

## VII. Experiment in Twitter Corpus

In addition to the study of the FilmAffinity corpus, an experiment over a corpus which contains Twitter comments was performed. This corpus is proposed in [11]. This corpus contains sentences grouped into 4 categories: positive, negative, neutral and informative. The way to adequate this corpus to a Subjective corpus was the following: We considered the positive and negative sentences as subjective sentences and the remainder as objective sentences.

The purpose of this experiment was to evaluate how the change of domain impact the proposal method. The results of this experiment are presented in Table VI.

In an overall view, the baseline method was the best method of the three. One interesting result that can be seen in Table VI is that the precision in our method was the best of all methods. This is satisfactory for us, because our method is focused on the Subjectivity Detection.

In case of the recall, our method was one of the worst (outperforming the MFS method in 0.04%). This resulted in a lower value in F measure. Some of the reasons for the lower values in recall and F-measure were the followings: (1) the lenght of the sentences in Twitter was very small, thus, this could cause that some clues are not detected; (2) expressions like emoticons and hashtags are most used in Twitter and these are not recognized by our method; and (3) the language in Twitter is informal, thus, some words are not indexed in sense repositories like WordNet-Pr.

## VIII. Conclusions and Final Remarks

In this paper, we presented a fine-grained SWSD method based on Knowledge for Subjectivity Detection task for Spanish texts. This method uses the WordNet for Spanish included in MCR 3.0 and the SentiWordNet as sense repositories. The results of our experiments show that Subjectivity Detection using a Fine-Grained SWSD-based approach based on Knowledge outperforms a baseline where the disambiguation is not used; therefore, the WSD may carry improvements into Subjectivity Detection task.

In our approach, a fine-grained method for SWSD was used, obtaining good results. An attractive direction for this work is the modification of this method, changing to a coarse-grained method, i.e., grouping the senses of SentiWordNet by subjective level to which to be belonging. Other direction for this work is the evaluation of the performance of the WSD method. Currently, this is not possible because the FilmAffinity

TABLE IV: RESULTS OBTAINED USING THE BASELINE, THE MOST FREQUENT SENSE METHOD AND THE PROPOSED METHOD (A-S) ON FILMAFFINITY CORPUS

| Method | $P_S$ | $R_S$ | $F_S$ | $P_O$ | $R_O$ | $F_O$ | $F_{AVG}$ | Ac. |
|---|---|---|---|---|---|---|---|---|
| Baseline | 57.51% | **91.24%** | 70.55 | **78.82%** | 32.60% | 46.12 | 58.34 | 61.92% |
| MFS | **65.51%** | 78.32% | **71.34** | 73.05% | **58.76%** | **65.13** | **68.24** | **68.54%** |
| A-S | 64.99% | 76.92% | 70.45 | 71.73% | 58.56% | 64.48 | 67.47 | 67.74% |

TABLE V: RESULTS OBTAINED ON MOVIE REVIEW DATASET

| Method | $P_S$ | $R_S$ | $F_S$ | $P_O$ | $R_O$ | $F_O$ | $F_{AVG}$ | Ac. |
|---|---|---|---|---|---|---|---|---|
| MFS | 55.81% | 67.72% | 61.19 | 58.96% | **46.38%** | **51.92** | 56.56 | 57.05% |
| A-S | **56.47%** | **78.74%** | **65.77** | **64.89%** | 39.30% | 48.95 | **57.36** | **59.02%** |

TABLE VI: RESULTS OBTAINED USING THE TWITTER CORPUS

| Method | $P_S$ | $R_S$ | $F_S$ | $P_O$ | $R_O$ | $F_O$ | $F_{AVG}$ | Ac. |
|---|---|---|---|---|---|---|---|---|
| Baseline | 61.69% | **33.02%** | **43.01** | **54.19%** | 79.45% | 64.43 | **53.72** | **56.20%** |
| MFS | 65.05% | 18.77% | 29.13 | 52.47% | 89.89% | 66.26 | 47.70 | 54.28% |
| A-S | **66.02%** | 18.81% | 29.28 | 52.59% | **90.29%** | **66.47** | 47.90 | 54.51% |

corpus has not sense-annotation but a next step in this work is the sense-annotation of this corpus.

The used attributes in the proposal presented different weights or importance degrees associated to Subjectivity Detection task. This is an interesting point because some studies use all morphosyntactic classes with the same weight, when not all of these are important or have the same importance. The proposed method uses all of attributes set, obtaining good results. Future experiments may be realized using an attributes subset that to be more associated with the subjectivity detection or a combination of attributes.

Despite of use a corpus created for English (Semcor) to extract attribute values and rules for Subjectivity Detection, and the use of a Knowledge Base in English (WordNet-Pr) to obtain the senses for all words, the results of our classifier for Spanish was not affected negatively, with respect to the results of classifier for English. This could give a direction for using these resources for other languages without negatively affecting the performance.

In case of corpora from informal texts, like Twitter, it is necessary a deep analysis to work with emoticons and hashtags, and recognize some words of informal contexts and its senses.

One contribution of this work is the creation of the FilmAffinity corpus. This corpus is available in https://msobrevillac.wordpress.com/corpora/.

## REFERENCES

[1] B. Liu, "Sentiment analysis and subjectivity," in *Handbook of Natural Language Processing, Second Edition. Taylor and Francis Group, Boca*, 2010.

[2] J. Wiebe and E. Riloff, "Creating subjective and objective sentence classifiers from unannotated texts," in *Proceedings of the 6th International Conference on Computational Linguistics and Intelligent Text Processing*, ser. CICLing'05. Berlin, Heidelberg: Springer-Verlag, 2005, pp. 486–497.

[3] E. Lloret, A. Balahur, M. Palomar, and A. Montoyo, "Towards a unified approach for opinion question answering and summarization," in *Proceedings of the 2Nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis*, ser. WASSA '11. Stroudsburg, PA, USA: Association for Computational Linguistics, 2011, pp. 168–174.

[4] G. Murray and G. Carenini, "Summarizing spoken and written conversations," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, ser. EMNLP '08. Stroudsburg, PA, USA: Association for Computational Linguistics, 2008, pp. 773–782.

[5] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Found. Trends Inf. Retr.*, vol. 2, pp. 1–135, 2008.

[6] R. Mihalcea, C. Banea, and J. Wiebe, "Learning multilingual subjective language via cross-lingual projections," in *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, 2007.

[7] S. Baccianella, A. Esuli, and F. Sebastiani, "Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining," in *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, N. C. C. Chair), K. Choukri, B. Maegaard, J. Mariani, J. Odijk, S. Piperidis, M. Rosner, and D. Tapias, Eds. Valletta, Malta: European Language Resources Association (ELRA), 2010.

[8] R. Valitutti, "Wordnet-affect: an affective extension of wordnet," in *In Proceedings of the 4th International Conference on Language Resources and Evaluation*, 2004, pp. 1083–1086.

[9] S. Cerini, V. Compagnoni, A. Demontis, M. Formentelli, and G. Gandini, *Language resources and linguistic theory: Typology, second language acquisition, English linguistics.* Milano, IT: Franco Angeli Editore, 2007, ch. Micro-WNOp: A gold standard for the evaluation of automatically compiled lexical resources for opinion mining.

[10] G. A. Miller, "Wordnet: A lexical database for english," *Commun. ACM*, vol. 38, no. 11, pp. 39–41, Nov. 1995.

[11] G. Sidorov, S. Miranda-Jiménez, F. Viveros-Jiménez, A. Gelbukh, N. Castro-Sánchez, F. Velásquez, I. Díaz-Rangel, S. Suárez-Guerra, A. Treviño, and J. Gordon, "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*, ser. MICAI'12. Berlin, Heidelberg: Springer-Verlag, 2013, pp. 1–14.

[12] J. Wiebe and R. Mihalcea, "Word sense and subjectivity," in *Proceed-*

*ings of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, ser. ACL-44. Stroudsburg, PA, USA: Association for Computational Linguistics, 2006, pp. 1065–1072.

[13] C. Akkaya, J. Wiebe, and R. Mihalcea, "Subjectivity word sense disambiguation," in *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*. Singapore: Association for Computational Linguistics, 2009, pp. 190–199.

[14] R. Ortega, A. Fonseca, Y. Gutiérrez, and A. Montoyo, "Improving subjectivity detection using unsupervised subjectivity word sense disambiguation," *Procesamiento del Lenguaje Natural*, vol. 51, no. 0, pp. 179–186, 2013.

[15] A. Gonzalez-Agirre, E. Laparra, and G. Rigau, "Multilingual central repository version 3.0," in *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, N. C. C. Chair), K. Choukri, T. Declerck, M. U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk, and S. Piperidis, Eds. Istanbul, Turkey: European Language Resources Association (ELRA), 2012.

[16] V. Hatzivassiloglou and J. M. Wiebe, "Effects of adjective orientation and gradability on sentence subjectivity," in *Proceedings of the 18th Conference on Computational Linguistics*, ser. COLING '00. Stroudsburg, PA, USA: Association for Computational Linguistics, 2000, pp. 299–305.

[17] E. Riloff, J. Wiebe, and T. Wilson, "Learning subjective nouns using extraction pattern bootstrapping," in *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, ser. CONLL '03. Stroudsburg, PA, USA: Association for Computational Linguistics, 2003, pp. 25–32.

[18] J. Wiebe and E. Riloff, "Creating subjective and objective sentence classifiers from unannotated texts," in *Proceedings of the 6th International Conference on Computational Linguistics and Intelligent Text Processing*, ser. CICLing'05. Berlin, Heidelberg: Springer-Verlag, 2005, pp. 486–497.

[19] X. Wang and G. Fu, "Learning lexical subjectivity strength for chinese opinionated sentence identification." in *CICLing (1)*, ser. Lecture Notes in Computer Science, A. F. Gelbukh, Ed., vol. 7181. Springer, 2012, pp. 580–590.

[20] B. Pang and L. Lee, "A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts," in *Proceedings of the 42Nd Annual Meeting on Association for Computational Linguistics*, ser. ACL '04. Stroudsburg, PA, USA: Association for Computational Linguistics, 2004.

[21] T. Wilson, J. Wiebe, and P. Hoffmann, "Recognizing contextual polarity in phrase-level sentiment analysis," in *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing*, ser. HLT '05. Stroudsburg, PA, USA: Association for Computational Linguistics, 2005, pp. 347–354.

[22] J. R. Quinlan, *C4.5: Programs for Machine Learning*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1993.

[23] H. Schmid, "Probabilistic part-of-speech tagging using decision trees," 1994.

[24] E. Agirre and A. Soroa, "Personalizing pagerank for word sense disambiguation," in *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics*, ser. EACL '09. Stroudsburg, PA, USA: Association for Computational Linguistics, 2009, pp. 33–41.

[25] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine," in *Proceedings of the Seventh International Conference on World Wide Web 7*, ser. WWW7. Amsterdam, The Netherlands: Elsevier Science Publishers B. V., 1998, pp. 107–117.