

---

---

## PARAMETERIZATION OF COMMENTS FROM PERUVIAN FACEBOOK AND TWITTER: LEXICAL RESOURCES AND ALGORITHM

**Angels Catena, Mikhail Alexandrov, Roque López**

**Abstract:** *Millions of Facebook and Twitter users send their comments all over the world about products and services, political and economical events, etc. (almost 3 billions each day) The principal problem of opinion mining such information is text parameterization, and in the paper we describe our experience in solution of this problem with Peruvian Facebook and Twitter. We use enriched vocabularies of Spanish SentiStrength and propose a simple algorithm for evaluation of sentiment contribution. The work was completed in the framework of the project WAYRA (Telefonica-Peru). The results proved to be promised: opinion analysis of parameterized texts showed the accuracy about 75% with elementary classifier*

**Keywords:** *opinion mining, SentiStrength, Facebook and Twitter*

**ACM Classification Keywords:** *1.2.7 Natural Language Processing*

---

### Introduction

Social Networks became an important source of information for very different applications. Around 400 million tweets [CNET, <http>] and more than 2.7 billion comments on Facebook [CNBC, <http>] are generated every day. People want to share their opinions about products and services, economical situation and political events and this feedback proves to be very useful for those who offer these products or who are responsible for these events.

The first NLP tools for Opinion Mining dealt with so-called 'ordinary' documents. Speaking 'ordinary' we mean not very short documents with usual normative lexis. As an example we can mention here the well-known program SO CAL. This program was applied for processing messages on forums, where the quality of various products and services were discussed. It provided the average accuracy 80% for binary classification of opinions [Taboada, 2006]. The other example is the program Opininer-UAB. We used this program for processing analytical texts in electronic publications concerning economical crisis in Europe. Opininer-UAB showed the accuracy 75%-80% for binary classification and about 60% for more detail classification with four categories [Catena, 2010].

But we meet the absolutely other situation when we have to work with comments from Facebook and Twitter. They are super-short (usually 2-5 phrases) and they use non-normative lexis (abbreviations, slang, emotional punctuation, etc.). For this reason it is difficult to reach the accuracy we pointed above. The SentiStrength program is one of the best tools we know for Opinion Mining texts from Facebook and Twitter [Thelwall, 2010; Thelwall, 2012]. By the moment there are versions of SentiStrength in English, German, French, Italian, Spanish, and some other languages. The usual Spanish SentiStrength version provides the accuracy 65% and its last modification especially oriented on the Peruvian Facebook and Twitter provides the accuracy 74% [Lopez, 2012]

This paper is the continuation of the previous work. It was completed in the framework of the ambitious project WAYRA by the company Telefonica-Peru [WAYRA-a, <http>; WAYRA-b, <http>]. In this work we concern only the process of text parameterization, which have the decisive role in Opinion Mining texts from social networks. The program is written on Python. The document collection includes 200 comments.



2) **Text**. It is a comment presented in the lineal form without the division on phrases. Each cell contains one word or one sign after correction

N no n m g e est pro m m ( ;- ;- v a L  
 o o e usta e ducto uy alo ( ale a  
 P olv  
 ena idarlo

Having eliminated the repeated elements we have a new list of words and signs

n m g ste e ducto pro uy m alo m ( ;- ale v a l ena p olv  
 o e usta ( ale a ena idarlo

Having eliminated stop-words we have the other list. But stop-words being the part of fixed word combinations are not considered. See here: vale **la pena** (= it is worth).

no sta gu ducto pro uy m alo m ; ale v a l ena p olv  
 -(

3) **Mark-1**. It is the set of denominations, which reflect the type of elements. We use the following list of denominations:

- P is a word from the Vocabulary-P
- C is a word combination from the Vocabulary-C
- R is an emphatic construction from the Vocabulary-R
- N is a negation from the Vocabulary-N
- S is a sign from the Vocabulary-S
- Z means that an element has no any weight

Therefore for our example we have:

N P Z R P S C Z

4) **Mark-2**. This array is directly related with the array **Mark-1**. Namely, here each cell contains the marks and weight of a correspondent element from the **Mark-1**.

1 1 0 1 - - 1 0  
 1 1

These data are taken from the vocabularies.

### 3.2 Processing

On this stage we use two arrays: **Indicator** and **Weight**

5) **Indicator**. It is a set of indicators, which show each moment whether the correspondent element of text has been evaluated or no. The program evaluates a given text from left to right revealing patterns according the list of rules (we mention them later). The elements are marked with '1' step by step.

6) **Weight**. It is the array of points, each pattern of a comment contributes to the total assessment. Here the array **Mark-2** is used. The array **Weight** changes its contents having revealed each pattern. The points are written to the first cell this pattern occupies

Initially we have:

#### Mark-1

N P Z R P S C Z

**Mark-2**

1	1	0	1	-	-	1	0
			1	1			

**Indicator**

0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---

**Weight**

0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---

After the first step:

**Mark-1**

N	P	Z	R	P	S	C	Z
---	---	---	---	---	---	---	---

**Mark-2**

1	1	0	1	-	-	1	0
			1	1			

**Indicator**

1	1	1	0	0	0	0	0
---	---	---	---	---	---	---	---

**Weight**

	0	0	0	0	0	0	0
-1							

After the second step:

**Mark-1**

N	P	Z	R	P	S	C	Z
---	---	---	---	---	---	---	---

**Mark-2**

1	1	0	1	-	-	1	0
			1	1			

**Indicator**

1	1	1	1	1	0	0	0
---	---	---	---	---	---	---	---

**Weight**

	0	0	-	0	0	0	0
-1		2					

After the third step:

**Mark-1**

N	P	Z	R	P	S	C	Z
---	---	---	---	---	---	---	---

**Mark-2**

1	1	0	1	-	-	1	0
			1	1			

**Indicator**

1	1	1	1	1	1	0	0
---	---	---	---	---	---	---	---

**Weight**

	-	0	0	-	0	-	0	0
1			2			1		

After the fourth step:

**Mark-1**

N	P	Z	R	P	S	C	Z
---	---	---	---	---	---	---	---

**Mark-2**

1	1	0	1	-	-	1	0
			1	1			

**Indicator**

1	1	1	1	1	1	1	1
---	---	---	---	---	---	---	---

**Weight**

	-	0	0	-	0	-	1	0
1			2		1			

Then all positive and negative scores are summed. The Table 1 shows the final result of text processing

*Table 1. Results of comment processing*

<b>Positive scores</b>	<b>Number of positive sentiments</b>	<b>Negative scores</b>	<b>Number of negative sentiments</b>
1	1	-4	3

In this example we indirectly used some rules for processing sequences NPZ, RPZ, etc. Here the symbols 'N', 'P', 'Z' mean the types of elements in a phrase under consideration (see section 3.1). This moment the program includes 6 different rules for processing phrases.

## Experiments

### 4.1 Elementary classifier

To demonstrate the effectiveness of proposed algorithm we used 200 comments from the Peruvian Facebook and Twitter. All comments were evaluated by 3 experts using the scale: positive, neutral, negative. Then these comments were parameterized: the part of texts (150) was used for constructing classifier, and the other part (50) for testing.

The elementary threshold-based classifier is formed by the following way

- All parameterized comments are presented in the scale  
 $r = ( PosScore + NegScore ) / ( PosScore + | NegScore | )$ . Here: *PosScore* and *NegScore* are contributions of positive and negative sentiments respectively. Obviously,  $|r| \leq 1$
- Histogram of the learning set (150 comments in our case) is constructed on the basis of variable *r*
- Thresholds for decision-making are manually adjusted to provide the maximum accuracy

Of course, we do not pretend here to have any sophisticated classifier. The modern approaches for opinion mining are known and these approaches have been already described in detail in the publications [Pang, 2008; Taboada, 2011]. The only we want is to show that our algorithm of parameterization allows to obtain good results even with the simplest classifier.

### 4.2 Results of classification

The Table 2 shows some results with the classification of testing set (50 comments)

*Table 2. Results of classification*

<b>Categories and the number of objects</b>	<b>Threshold ds</b>	<b>Rules</b>	<b>Accuracy</b>
Positive ( 42% )	Up =	$I \geq \text{Up}$	Positive
Neutral ( 39% )	0,05	$I <$	Un 72%
Negative ( 19% )	Un = -	Negative	
	0,05	The other ones are Neutral	
Positive ( 42% )	Up = 0,2	$I \geq \text{Up}$	Positive
Undefined -	Un = -0,2	$I <$	Un 77%
Negative ( 19% )		Negative	
		The other ones are Undefined	

The neutral category is not considered

Notes.

- 1) We have to introduce the so-called undefined category to improve the quality of results: sometimes is better to say "I do not know" instead of any erroneous answer
- 2) We suppose the results will be essentially better if to use any more advanced classifier instead of the elementary one

---

### Conclusions

---

The main results of the paper are:

- We proposed method (linguistic resources and algorithm) for parameterization of comments from Facebook and Twitter.
- Experiments with real Peruvian Facebook and Twitter showed the promised results
- The method with modifications can be used for processing comments on other languages

In the future we suppose:

- to extend the list of grammatical rules
- to test the binary scale for sentiment classification

---

### Bibliography

---

[Catena, 2010] A. Catena, M. Alexandrov, N. Ponomareva. Opinion Analysis of Publications on Economics with Limited Vocabulary of Sentiments the Intern. Journal on Social Media MMM: Monitoring, Measurement, and Mining, Brno, Czech Rep., N\_1, 2010, Publ. House 'Konvoj' , pp. 20-31

[CNBC, http] [http://www.cnbc.com/id/45582325/The\\_World\\_s\\_Most\\_Liked\\_Brands](http://www.cnbc.com/id/45582325/The_World_s_Most_Liked_Brands)

[CNET, http] [http://news.cnet.com/8301-1023\\_3-57448388-93/twitter-hits-400-million-tweets-per-day-mostly-mobile/](http://news.cnet.com/8301-1023_3-57448388-93/twitter-hits-400-million-tweets-per-day-mostly-mobile/)

- [Kaurova, 2010] O. Kaurova, M. Alexandrov, N. Ponomareva. The Study of sentiment word granularity for Opinion Analysis (a comparison with Maite Taboada works). Intern. Journal on Social Media MMM: Monitoring, Measurement, and Mining, Brno, Czech Rep., N\_1, 2010, Publ. House 'Konvoj' , pp.45-57
- [Lopez, 2012] R. Lopez, J. Tejada, M. Thelwall. Spanish Sentistrength as a tool for opinion mining Peruvian Facebook and Twitter. (this Proceedings)
- [Pang, 2008] B. Pang, L. Lee. Opinion mining and sentiment analysis. Foundations and trends in Information Retrieval, 1 (1-2), 1-135 (2008)
- [Taboada, 2006] M. Taboada, C. Anthony, K. Voll. Opinion mining and sentiment analysis. Foundations and trends in Information Retrieval, 1(1-2), 1-135 (2008). Creating semantic orientation dictionaries. *Proc. of 5th Intern. Conf. on Language Resources and Evaluation (LREC)*. Italy, 2006, p. 427-432.
- [Taboada, 2011] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, M. Stede. Lexicon-based methods for sentiment analysis. *Computational Linguistics*, 37(2), 267-307 (2011).
- [Thelwall, 2010] M. Thelwall, K. Buckley, G. Paltoglou, D. Cai. Sentiment strength detection in short informal text. *Journal of the American Society for Information Science and Technology*. 61, 2544–2558 (2010).
- [Thelwall, 2012] M. Thelwall, K. Buckley, G. Paltoglou. Sentiment strength detection for the social web. *Journal of the American Society for Inform. Science and Technology*. 63, 163–173 (2012).
- [WAYRA-a, http] WAYRA: <http://innovar.org/?p=772>
- [WAYRA-b, http] WAYRA: <http://www.youtube.com/watch?v=Kc1ho6GLQJU>

---

## Authors' Information

---



**Angels Catena** – Professor, Department of French and Romance Philology, fLexSem Research Group, Autonomous University of Barcelona, 08193 Bellaterra (Barcelona), Spain;

e-mail: [Angels.Catena@gmail.com](mailto:Angels.Catena@gmail.com)

Major Fields of Scientific Research: sentiment analysis, French lexis and semantics sentiment analysis,



**Mikhail Alexandrov** – Professor of the Russian Presidential Academy of national economy and public administration; Prosp. Vernadskogo 82, bld. 1, Moscow, 119571, Russia; research fellow of the fLexSem Research Group, Autonomous University of Barcelona, 08193 Bellaterra (Barcelona), Spain;

e-mail: [MAlexandrov@mail.ru](mailto:MAlexandrov@mail.ru)

Major Fields of Scientific Research: data mining, text mining, mathematical modelling



**Roque López** – Student of System Engineering at San Agustín National University, calle Santa Catalina N° 117 Arequipa, Peru;

e-mail: [ropezc27@gmail.com](mailto:ropezc27@gmail.com)

Major Fields of Scientific Research: natural language processing, text mining, social network analysis