# Exploratory Analysis of a New Corpus for Political Alignment Identification of Argentinian Journalists

Viviana Mercado[1], Andrea Villagra[1], and Marcelo Luis Errecalde[1,2]

[1] Laboratorio de Tecnologías Emergentes - Unidad Académica Caleta Olivia
Universidad Nacional de la Patagonia Austral
{vmercado, avillagra}@uaco.unpa.edu.ar
[2] Laboratorio de Investigación y Desarrollo en Inteligencia Computacional
Universidad Nacional de San Luis
merrecalde@unsl.edu.ar

**Abstract.** Political alignment identification is an author profiling task that aims at identifying political bias/orientation in people' writings. As usual in this kind of field, a key aspect is to have available adequate data sets so that the data mining and machine learning approaches can obtain reliable and informative results. This article takes a step in this direction by introducing a new corpus for the study of political alignment in documents of Argentinian journalists. The study also includes several kinds of analysis of documents of pro-government and opposition journalists such as sentiment analysis, topic modelling and the analysis of psycholinguistic indicators obtained from the *Linguistic Inquiry and Word Count* (LIWC) system. From the experimental results, interesting patterns could be observed such as the topics both types of journalists write about, how the sentiment polarities are distributed and how the writings of pro-government and opposition journalists differ in the distinct LIWC categories.

**Keywords:** Text Mining, Exploratory Data Analysis (EDA), Author Profiling, Journalist Political Alignment, Sentiment Analysis, Topic Modelling, LIWC

## 1 Introduction

*Political alignment identification* (PAI) in a text or document is a form of *author profiling* (AP), one of the main tasks of *authorship analysis* (AA) together with authorship attribution/determination, plagiarism detection and style inconsistency detection. PAI, the same as other AP tasks like the detection of depressed people or with different personality traits, paedophiles and suicides is a challenging task within the automatic analysis of texts since it involves, in general, the use of representations of texts that capture stylistic and content aspects of their authors. In this context, a particular area within the PAI is that which is oriented to the study of political orientation in texts written by journalists, and which we will refer to now as *journalistic texts*. We will consider as journalistic texts that information that a journalist publishes in various media such as a personal blog, an article written in a mass media such as a newspaper or the content expressed in a book of his authorship.

The PAI has been applied to texts generated by regular users of social media such as Twitter [1, 2] although more recently it has been done with the documents produced by journalists [3]. In [4], political speech in Twitter has been analyzed with LIWC during the 2008' German electoral campaign. The same tool, LIWC, was used to determine

the psychological state and personality of the candidates for the presidency and vice presidency of the US in the 2004 campaign [5] and the language used by the New York's mayor, R. Giuliani [6] throughout his term. Regarding texts in Spanish language, in [7] the linguistic style of the candidates of the main political parties in the Spanish general elections of 2008 and 2011 is analyzed. On the other hand, the Spanish dictionary of LIWC was applied to analyze the political speech and tweets of the candidates in the elections of Galicia in 2012 [8].

The previous approaches are related to our work but, as far as we know, there are no PAI studies of journalistic texts in Spanish. In this work, we will make a first approach to the PAI in journalistic texts in Spanish, in particular, of texts generated by Argentinian journalists. The task, in this case, will be to group all the documents of "pro-government" journalists on one side and "opponents" on the other one. In that way, it will allow in the future to visualize it as a binary classification ("pro-government" versus "opponent") problem. Our objectives, in the long run, will be to answer the following research questions:

1. What are the appropriate forms of document representation for this task?
2. What are the most effective learning algorithms to use with those representations?
3. What is the impact of the dimensionality reduction approaches in the representations of the documents?
4. How related are the results obtained with similar studies with journalistic texts written in other languages?

In the present article, a first step is taken to achieve these objectives by introducing a new corpus for the study of political alignment in documents of Argentinian journalists. The study also includes several kinds of analysis of documents of pro-government and opposition journalists such as sentiment analysis, topic modelling and the analysis of psycholinguistic indicators obtained from the LIWC system.

The rest of the article is organized as follows: Section 2 describes the PAI corpus introduced in the present article with statistics and metrics of the whole data set and of each of both involved classes; Section 3 gives some results obtained from a topic modelling and sentiment analysis; Section 4 goes further in analyzing the PAI corpus taking into account psycholinguistic indicators obtained from the LIWC system. Section 5 finishes this article by giving the main conclusions obtained from our study and some future work.

## 2    Corpus Description

Our work was focussed on generating a collection of Argentinian journalistic documents obtained from news blogs, online newspapers, books, etc. It consists of 196 documents belonging to 10 journalists: 5 of them that clearly support the actions of the Argentine government in the period 2012 to 2015 and 5 of them that explicitly express themselves against the government in that period. The data set was split into two groups of documents according to the political orientation of the journalists. Thus, 98 documents belonging to the 5 pro-government journalists were selected for the *gov* (pro-government) class and the 98 remaining documents of the opposition journalists were used to build the *oppo* (opposition) class. In that way, a balanced corpus with 2 classes was obtained.

To select the documents some guidelines were taken into account:

– Texts correspond to Spanish documents written by Argentinian journalists.

- Texts refer to different political aspects related to the Argentinian government in the period 2012-2015, such as government actions, politicians' declarations, corruption cases, treatment of laws, etc.
- All the documents contain " formal text ", that is to say, they do not present common " informal " aspects of content from social media such as abbreviations, slang expressions, typos, hyperlinks, labels, figures, and emoticons.
- From each journalist, between 18 and 20 documents were taken from his/her personal blog, articles in online newspapers or digital books of his/her authorship.
- Each journalist was clearly identified as pro-government or opponent.
- The same proportion of male and female journalists were kept between both categories.

After collecting the documents, they were manually labeled as belonging to the two above-mentioned classes *gov*, and *oppo*. Table 1 shows information about how the documents were finally distributed in both classes and what were the sources (online newspaper, blog or digital books) from they were obtained.

| Class | Newspapers | Blogs | Books | Total |
|-------|-----------|-------|-------|-------|
| *gov* | 50 | 46 | 2 | 98 |
| *oppo* | 60 | 36 | 2 | 98 |

Table 1: Distribution of documents in classes and source.

## 2.1 Corpus statistics

Before proceeding with a more elaborated analysis of the corpus, some basic statistics were obtained to get some insights about the general characteristics of the documents. First of all, the *number of words* per document was analyzed for each document/article of both classes. Table 2 shows the minimum, maximum, mean and standard deviation values for the number of words in the documents of the *gov*, *oppo* classes and the whole corpus (*gov + oppo*).

| Class | Minimum | Maximum | Mean | St. Dev. |
|-------|---------|---------|------|----------|
| *gov* | 139 | 36619 | 1865.18 | 5031.56 |
| *oppo* | 236 | 3423 | 1243.01 | 733.71 |
| *gov + oppo* | 139 | 36619 | 1554.09 | 3608.91 |

Table 2: Number of words in documents: minimum, maximum, mean and standard deviation values per class.

As we can see, although the sizes of the shortest documents (Min) are similar for both classes (139 vs 236), they differ considerably in the longest ones (36619 vs 3423). That can be observed more clearly in Figures 1 and  2 that show the number of words (at the left) and the histogram (at the right) for the documents in the *gov* and the *oppo* classes. There, for instance, Figure 1a shows that there are some documents in the class

*gov* whose sizes exceeds 34000 words. However, as the Figure 1b confirms, most of the documents in the *gov* class do not exceed 5000 words with only a couple of documents (that correspond to books) whose sizes are between 34000 and 36000 words. Figure 2 shows that the *oppo* class has a "smoother" distribution in the number of words of its documents, with sizes that oscillate between 200 and 3500 words approximately.
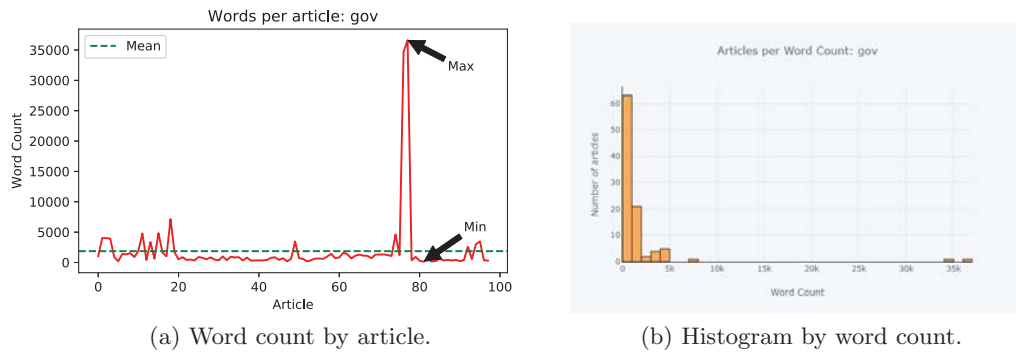


(a) Word count by article.



(b) Histogram by word count.

Fig. 1: Class *gov*: number of words per article



(a) Word count by article.
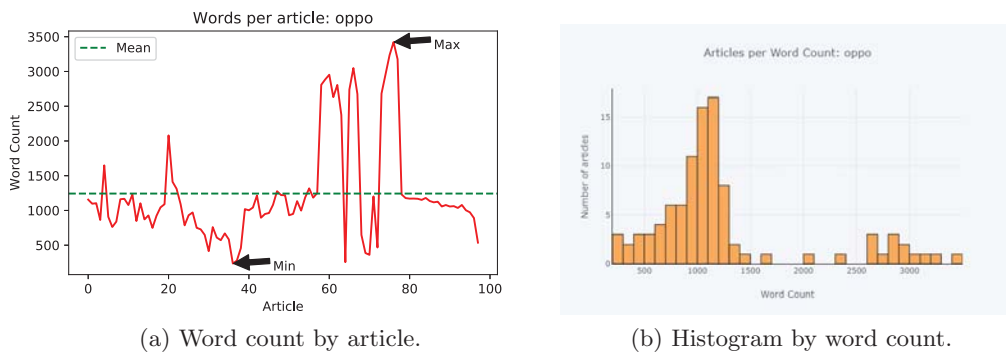


(b) Histogram by word count.

Fig. 2: Class *oppo*: number of words per article

Regarding the total number of words and the size of the vocabulary of the whole corpus and of each class, we will first introduce the following notation: let $D_{\mathcal{C}}$ be the set of documents belonging to our corpus $\mathcal{C}$; let $D_{\mathcal{G}}$ and $D_{\mathcal{O}}$ be the sets of documents belonging to the pro-government and opposition journalists respectively, $D_{\mathcal{G}} \subset D_{\mathcal{C}}$, $D_{\mathcal{O}} \subset D_{\mathcal{C}}$, $D_{\mathcal{C}} = D_{\mathcal{G}} \cup D_{\mathcal{O}}$, $D_{\mathcal{G}} \cap D_{\mathcal{O}} = \emptyset$. Let $|D_{\mathcal{C}}|$ be the total number of words in our corpus $\mathcal{C}$, with similar meanings for $|D_{\mathcal{G}}|$ and $|D_{\mathcal{O}}|$. Besides, let $\mathcal{V}_{\mathcal{C}}$, $\mathcal{V}_{\mathcal{G}}$, and $\mathcal{V}_{\mathcal{O}}$ be the *vocabularies*[3] of $D_{\mathcal{C}}$, $D_{\mathcal{G}}$, and $D_{\mathcal{O}}$, respectively. Table 3 gives some statistics related to these collection of documents.

---

[3] The vocabulary of a collection of documents is the set of *distinct* words that appear in that collection.

| $|D_\mathcal{C}|$ | $|\mathcal{V}_\mathcal{C}|$ | $|D_\mathcal{G}|$ | $|\mathcal{V}_\mathcal{G}|$ | $|\mathcal{V}_\mathcal{G}|/|D_\mathcal{G}|$ | $|D_\mathcal{O}|$ | $|\mathcal{V}_\mathcal{O}|$ | $|\mathcal{V}_\mathcal{O}|/|D_\mathcal{O}|$ | $|\mathcal{V}_\mathcal{G} \cap \mathcal{V}_\mathcal{O}|$ | $|\mathcal{V}_\mathcal{G} \setminus \mathcal{V}_\mathcal{O}|$ | $|\mathcal{V}_\mathcal{O} \setminus \mathcal{V}_\mathcal{G}|$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 280343 | 24323 | 167844 | 18497 | 0.11 | 112499 | 13146 | 0.11 | 7320 | 11177 | 5826 |

Table 3: Statistics on the documents of the whole corpus ($D_\mathcal{C}$), documents of pro-government ($D_\mathcal{G}$) and opposition ($D_\mathcal{O}$) journalists.

One of the first things we can observe from Table 3 is that although when the number of words in the whole corpus is high ($|D_\mathcal{C}| = 280343$), the number of *distinct* words (the size of the *vocabulary*) is relatively small ($|\mathcal{V}_\mathcal{C}| = 24323$). That differs from the size of vocabularies in texts from social media which usually are bigger. A possible cause of this is that writings in social media are usually informal and prone to have abbreviations and typos, increasing in that way the number of distinct words. Another interesting datum from Table 3 is that the vocabulary of pro-government journalists is considerably bigger than the one of opposition journalists ($|\mathcal{V}_\mathcal{G}| = 18497$ $|\mathcal{V}_\mathcal{O}| = 13146$). One cause of this might be that due to the greater number of words in the documents of pro-government journalists ($|D_\mathcal{G}| > |D_\mathcal{O}|$) we will probably have a greater number of distinct words. For this reason, as vocabulary richness estimation it is frequently used the ratio between the size of vocabulary and the number of words in the collection. In that case, we can see that those metrics for pro-government ($|\mathcal{V}_\mathcal{G}|/|D_\mathcal{G}|$) and opposition ($|\mathcal{V}_\mathcal{O}|/|D_\mathcal{O}|$) journalists are the same.

Finally, it is worth to note that the intersection between pro-government and opposition vocabularies ($|\mathcal{V}_\mathcal{G} \cap \mathcal{V}_\mathcal{O}|$) is very small compared to the vocabulary of the whole corpus ($|\mathcal{V}_\mathcal{C}|$). That means that many words are used by a group and not by the other one, and the other way around. This point can be easily observed in Table 3 where the number of words used by pro-government and not by opposition journalists ($|\mathcal{V}_\mathcal{G} \setminus \mathcal{V}_\mathcal{O}|$) is 11177 and $|\mathcal{V}_\mathcal{O} \setminus \mathcal{V}_\mathcal{G}| = 5826$. For instance, some words of $\mathcal{V}_\mathcal{G} \setminus \mathcal{V}_\mathcal{O}$ are "milicos", "globitos", "latinoamericana", "egoísmo", and "ultraderecha" and some words of $\mathcal{V}_\mathcal{O} \setminus \mathcal{V}_\mathcal{G}$ are "monárquica", "tribunera", "dilapidado", "negociaron", "hitlerismo", and "avaricia".

Another analysis that is usually informative is to measure the "relevance" of the terms in the corpus according to some specific metric. For instance, an approach is estimating the importance of a term according to the weight that it would receive in a particular document representation scheme, such as *tf-idf*. This scheme (*tf-idf*) is a weighted model commonly used for information retrieval problems. It is an unsupervised model in the sense that when weighting a term in a document, it does not take into account any information about the class that document belong to; for instance, in our corpus, if we take as terms the word uni-grams, the terms with the highest *tf-idf* value are: "colegio", "años", "comisión", "madre", "dijo", "perón", "día", "plata", "dos", "decía", "chica", "mujeres", "dice", "mamá", "después", "flaco", "casa", "alicia", "néstor", and "cristina". Taking as terms words 2-grams, the terms with the highest *tf-idf* value are: "próximo gobierno", "cinco años", "santa fe", "derechos humanos", "clase media", "muerte néstor", "años después", "día siguiente", "néstor kirchner", "muchas veces", "provincia buenos", "primera vez", "cristina dijo", "procurador general", "gils carbó", "cristina fernández", "buenos aires", "santa cruz", "néstor cristina" and, "río gallegos". Finally, "magdalena ruiz guiñazú", "economía axel kicillof", "ministro economía axel", "asignación universal hijo", "josé pablo feinmann", "joaquín morales solá", "gobernador provincia buenos", "da mucha bronca", "triple crimen general", "derechos humanos cidh", "cristina fernández kirchner", "mil millones dólares", "manuel abal medina", "juan manuel abal", "madres plaza mayo", "comisión interamericana derechos", "aumento mínimo imponible", "ciudad buenos aires", "interamericana derechos humanos",

and "provincia buenos aires" are the terms with the highest $tf\text{-}idf$ values when word 3-grams are used as terms.

A second approach to measure the importance of terms, is considering *supervised* metrics that capture the importance of each term concerning its class/category, such as $\chi^2$ and *information gain*, among others. They are usually used in *feature selection* processes to determine what are the most informative features to be preserved in the document representation. Here, we will calculate $\chi^2$ scores for all the terms consisting in word 2-grams and the top 20 are shown in Figure 3.



Fig. 3: Top 20 word 2-grams according to $\chi^2$ scores.

A third alternative that usually shows interesting information about the importance of features, is to obtain the *coefficients* of a model learned in a training stage of a prediction task with the whole corpus. In that case, those coefficients reflect how the learned model *weights each feature* for that task. In that way, we can look at the largest coefficients, and see which words these correspond to. For instance, Figure 4 shows a bar chart with the 25 largest and 25 smallest coefficients of a logistic regression model, with the bars showing the size of each coefficient. Word 2-grams are used as terms and the negative coefficients on the left belong to terms that according to the model are indicative of pro-government alignment, while the positive coefficients on the right belong to terms that according to the model indicate an article written by an opposition journalist. Most of the terms are quite intuitive, in the sense that reflect aspects very recognizable in texts from both political alignment, such as "grandes medios", "medios dominantes", and "nacional popular" indicating pro-government journalists, while "ruta dinero", "lázaro baez" and, "lavado dinero" indicate opposition documents.

## 3    Topic Modelling and Sentiment Analysis

*Topic modeling* is an umbrella term describing a class of text analysis methods whose task is assigning each document to one or multiple *topics*, usually without supervision. A good example of this is news data, which might be categorized into topics like "politics", "sports", "finance", and so on. Intuitively, a topic is a group of words that appear together frequently. In that context, "topics" obtained by a topic modeling process might
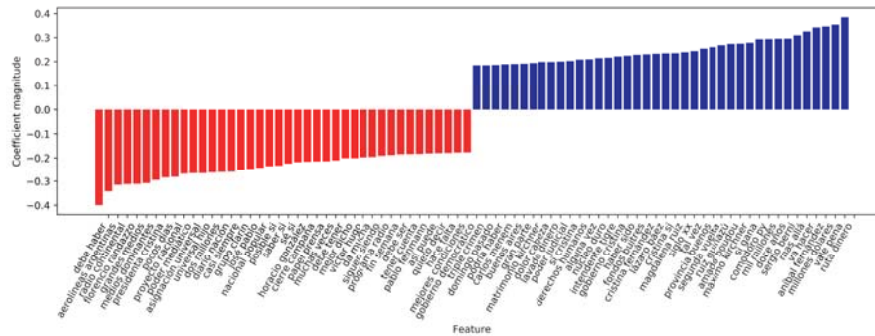
Fig. 4: Largest and smallest coefficients of logistic regression trained on tf-idf features.
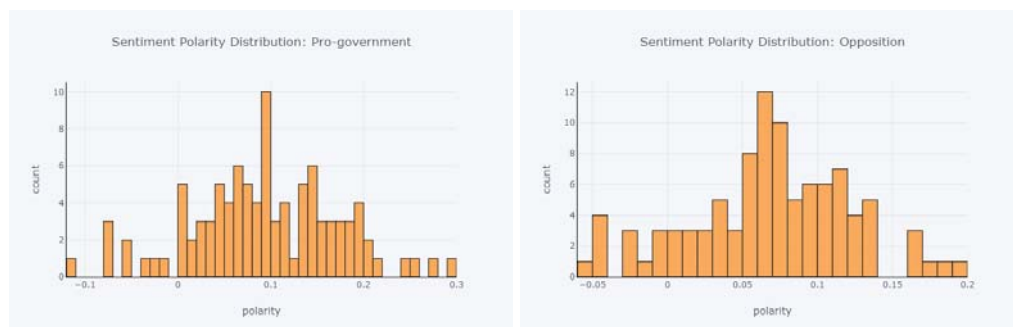
not be what we would normally call a topic in everyday speech. In other words, the obtained groups (topics) might or might not have a semantic meaning clearly identifiable by a person. Often, when people talk about topic modeling, they refer to one particular decomposition method called *Latent Dirichlet Allocation* (LDA). LDA not only tries to find groups of words (the topics) that appear together frequently. It also requires that each document can be understood as a "mixture" of a subset of the topics.

As an example, applying LDA to the pro-government and opposition documents and setting the number of topics to 100, several topics with intuitive meaning are obtained. Table 4 shows the first 20 words of some of those topics, three of the pro-government documents on the left column and three of the opposition documents on the right one. There, it can be observed that pro-government topics have to do with women's rights (***Topic #9***), Argentine debt with vulture funds (***Topic #27***) and social security plans (***Topic #49***), while opposition topics are related to communication media and journalists (***Topic #46***), some events related to what was popularly known as "the cause of the ephedrine" (***Topic #89***), and the relationship between the official Argentine cult and the Pope and some politicians (***Topic #91***).

| Pro-government topics | Opposition topics |
|---|---|
| ***Topic #9:*** voto mujeres décadas femenino incluso ciclo siglo luchas derecho quiera derrota banderas fitzgerald capital feministas siglos diez políticamente evita consigue | ***Topic #46*** radio mitre canal censura lanata tn intento sabemos intervenir oyentes diario adecuación marcelo pánico convertirse puesta usureros colegas clarín clarin |
| ***Topic #27:*** fondos deuda buitre final tema ministro documento kicillof soberana necesidad procesos comunicado litigiosidad previsibilidad reestructuración anexo conformes soberanas deudas líderes | ***Topic #89:*** aníbal desmentida triple kilos crimen efedrina quilmes melnyk clara importación agosto negocio granero morsa pérez nombres indispensable junto lanatta quizá |
| ***Topic #49:*** auh asignación pobreza fondos plan implicó reparación ése cfk octubre previsionales diputados pasaba narváez proyectos dirigencia impulso región decreto corporaciones | ***Topic #91:*** papa francisco iglesia ayuda mirada hombres página uca bergoglio cuervo vaticano guillermo michetti sienten larroque carrió explican quedaron opositores alegría |

Table 4: Some topics of pro-government and opposition journalists.

Another usual analysis of the texts in a corpus is determining the *polarity* of each document by averaging the polarity of its component words. A popular tool for this task is **TextBlob** that calculates sentiment polarity in the range of $[-1, 1]$ where 1 means *positive* sentiment and -1 means *negative* sentiment. In that way, it is usual to show some articles with the highest/lowest or even close to neutral (zero) sentiment polarity score or give some distribution of the articles according to their polarity scores. Due to space restrictions, we will only analyze the last alternative presenting in Figure 5 the sentiment polarity distributions of both pro-government (Figure 5a) and opposition (Figure 5b) journalists. There, it can be observed a greater frequency of pro-government articles on higher scores (around 0.1) than the one shown by the opposition journalists (around 0.05). Besides, the highest positive score achieved by opposition journalists (0.2) is surpassed by several pro-government articles. The other way around, the global lowest (negative) score is also obtained by pro-government journalists (less than -0.1) indicating that pro-government articles show the greatest variation range in polarity scores.



(a) Histogram of pro-government articles.       (b) Histogram of opposition articles.

Fig. 5: Sentiment Polarity Distribution in pro-government and opposition journalists.

## 4   LIWC-based Analysis

LIWC is a tool developed by the American psychologist J. Pennebaker and colleagues [9] and have been used in several studies related to psychological aspects of individuals. LIWC calculates the proportions of certain grammatical, lexical, and semantic markers, as well as markers belonging to other categories (up to 90 text features depending on the version). In our study, we used the most recent version of LIWC, LIWC2015 [10]. For each text file, LIWC2015 generates approximately 90 output variables as one line of data to an output file. This data record includes the file name and word count, 4 summary language variables (analytical thinking, clout, authenticity, and emotional tone), 3 general descriptor categories (words per sentence, percent of target words captured by the dictionary, and percent of words in the text that are longer than six letters), 21 standard linguistic dimensions (e.g., percentage of words in the text that are pronouns, articles, auxiliary verbs, etc.), 41 word categories tapping psychological constructs (e.g., affect, cognition, biological processes, drives), 6 personal concern categories (e.g., work, home, leisure activities), 5 informal language markers (assents, fillers, swear words, netspeak), and 12 punctuation categories (periods, commas, etc.).

Properties of documents generated by LIWC2015 have been used as document representations in several studies and also to analyze how these measures differ between articles of different classes. This last approach will be the one used in the present article. We will first identify what are the features/characteristics in which there are statistical differences between both classes and then we will show more detailed information on some of them. Since the distribution of the feature values is not known and we cannot make any assumption about it, we used, the same as similar works with LIWC features [11], the (non-parametric) Wilcoxon signed-rank test for comparing paired data samples with a p-value $< 0.05$ for statistical significance. The null hypothesis ($H0$) that we are trying to refute is that there is no statistically significant relationship between the mean value of a feature belonging to the pro-government class and the mean value of the same feature belonging to the opposition class. In that context, we determined significant statistical differences in 34 LIWC categories. For instance, pro-government journalists show a greater use of *verbs*, *adverbs*, *first person singular* ("yo", "mi", "mio"), *social* processes ("compañero", "hablar", "ellos") and, *perceptual* processes ("mirar", "escuchar", "sentir"). Opposition journalists, on the other hand, make a higher use of words with a *length* $> 6$ letters and make more references to expressions related to *money* ("dinero", "efectivo", "adeudar"). As an example, Figure 6 shows comparative boxplots of pro-government and opposition journalists for 2 LIWC categories with statistically significant differences: Perceptual processes and Money.
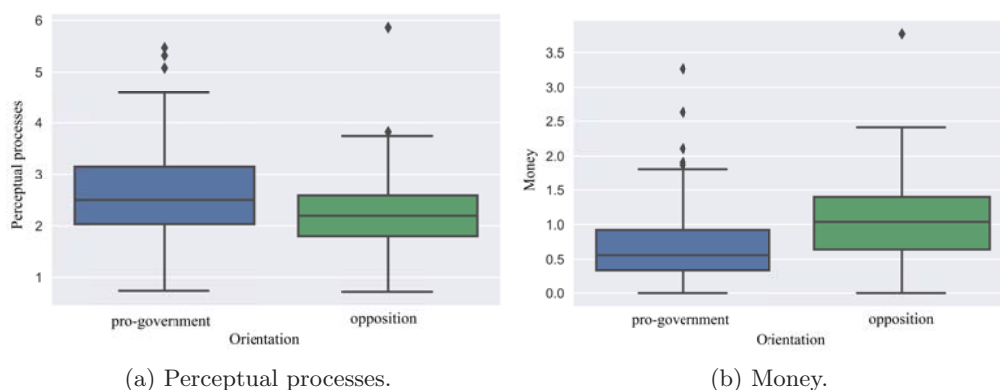


(a) Perceptual processes.            (b) Money.

Fig. 6: Comparative box-plots for Perceptual processes and Money categories.

## 5   Conclusions and Future Work

This article introduces a new corpus for political alignment identification of Argentinian Journalists. In that context, a comprehensive analysis of that corpus has been carried out which included the study of the corpus statistics, topic modelling and sentiment analysis, and a comparison of texts based on LIWC categories. As a result of this analysis, some interesting patterns were identified that reveal evident differences between the writings of pro-government and opposition journalists. For instance, they differ in how their sentiment polarities are distributed, what are the topics they talk about and, in their values for many categories of the LIWC system.

As future work, we plan to use the different types of information obtained in the present work in the representation of documents for supervised (classification) and, non-supervised (clustering) tasks. Thus, the idea is using LIWC-based and LDA/topics-based features in text classification tasks, and comparing them against classical (bag of words) and more recent approaches like deep neural-networks with word embeddings.

Finally, we will reorganize the documents of the corpus analyzed in the present work according to the *authors* of these documents. In that way, we will have ten different classes (one for each journalist) and the task will be addressed as an *authorship attribution* task. An interesting point, in this case, will be determine how the hardness of this task is incremented when the authorship attribution is constrained to journalists of the same political orientation.

# References

1. Raviv Cohen and Derek Ruths. Classifying political orientation on twitter: It's not easy! In *Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media*, 2013.
2. Michael D. Conover, Bruno Gonçalves, Jacob Ratkiewicz, Alessandro Flammini, and Filippo Menczer. Predicting the political alignment of twitter users. In *Proceedings of 3rd IEEE Conference on Social Computing (SocialCom)*, pages 192–199, 10 2011.
3. Konstantina Lazaridou and Ralf Krestel. Identifying political bias in news articles. *Bulletin of the IEEE TCDL*, 12, 2016.
4. Andranik Tumasjan, Timm O. Spenger, Philipp G. Sandner, and Isabelle M. Welpe. Predicting elections with twitter: What 140 characters reveal about political sentiment. *Proceedings of the Fourth International AAAI conference on Weblogs and Social Media (ICWSM), Washington, S. 178-185.*, (2010).
5. Richard. B. Slatcher, Cindy K. Chung, James W. Pennebaker, and Lori D. Stone. Winning words: individual differences in linguistic style among u. s. presidential and vice presidential candidates. j. *ournal of Research in Personality, vol. 41, pp.63-75.*, (2007).
6. James W. Pennebaker and Thomas C. Lay. Language use and personality during crises: analyses of mayor rudolph giuliani's press conferences. *Journal of Research in Personality, vol. 36, pp.271-282.*, (2002).
7. María Jesús Carrera-Fernández, Joan; Guárdia-Olmos, and Maribel Peró-Cebollero. Linguistic style in the mexican electoral process: Language style matching analysis. *Revista Mexicana de Psicología, 31(2), 138-152.*, (2014).
8. Mercedes Fernández-Cabana, José Rúas-Araújo, and Maria Teresa Alves-Pérez. Psicología, lenguaje y comunicación: análisis con la herramienta liwc de los discursos y tweets de los candidatos a las elecciones gallegas de 2012. *Anuario de Psicología, 44(2), pp.169-184.*, (2014).
9. James W.a Pennebaker, Roger J. Booth, and Marta E. Francis. Linguistic inquiry and word count (liwc). *[Software].*, (2001).
10. James W. Pennebaker, Ryan L. Boyd, Kayla Jordan, and Kate Blackburn. *The development and psychometric properties of LIWC2015*. University of Texas at Austin, Austin, TX, 2015.
11. José Rúas, Mercedes Fernández, and Iván Puentes. Aplicación de la herramienta liwc al análisis del discurso político. los mítines de los candidatos en las elecciones al parlamento de galicia de 2012. In *Actas del 2do Congreso Nacional sobre Metodología de la Investigación en Comunicación*, pages 47–64, 2013.