

RESEARCH

Agenda diversity and coverage bias: A quantitative approach to the agenda-setting theory

Sebastián Pinto^{1,2*†}, Claudio O Dorso^{1,2} and Pablo Balenzuela^{1,2}

*Correspondence: spinto@df.uba.ar

¹Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428 Buenos Aires, Argentina

Full list of author information is available at the end of the article

†Equal contributor

Abstract

The agenda-setting theory is a practical framework in order to understand the role of mass media on a society. The theory treats mass media as a very important actor that is able to make people think about, and in many cases how to think about certain topics. When the media succeeds in this task, we say that the media *set the agenda*. In this work we study the agenda of Argentinian newspapers in comparison with public's interests through a quantitative approach by performing topic detection over the news, identifying the main topics covered and their evolution over time. We look for characterizing the differences and similarities over time between what we call the Media Agenda and the Public Agenda. On the other hand, we aim to detect coverage bias among the newspapers involved in the analysis in the emerging topics.

Keywords: agenda-setting; mass media influence

Background

Agenda Setting Theory

In the famous study performed in Chapel Hill during the US presidential elections in 1968 [1], Maxwell McCombs and Donald Shaw found that those aspects of public affairs that are prominent in the news become prominent among the public. This study is considered the founding of the agenda-setting theory, which focus in the influence of mass media in public opinion. From [2], *“The media agenda is the pattern of news coverage over a period of days, weeks (...) for a set of issues or other topic. In other words, the media agenda is a systematic compilation of the issues or topics presented to the public that identifies the degree of emphasis on these topics.”*

Since the Chapel Hill research, several directions of agenda-setting were established [3]. In its basic stage, the theory is focus on the comparison between the topics coverage by the media and the public agenda, i.e. the topics that the public consider as priority. It looks for answering the question if the media is able to set the public agenda, which would transform the media as an important actor in the formation of public opinion. The interaction between media and public is rather complex, for instance, [4] shows that not necessarily the journalists and public preferences coincide.

With the emergency of the Internet, the end of agenda-setting were predicted due to the audience fragmentation onto multiple sources, which would virtually lead to

a highly individualized agenda. However, it is based on two assumptions that not necessarily are true: that the public spreads its attention in an homogeneous way across the multiple sources, and that the agendas of that sources are different [3].

The basic agenda-setting sometimes is called the “first level agenda-setting”. The very often quoted phrase of Bernard C. Cohen “*The press may not be successful much of the time in telling people what to think, but it is stunningly successful in telling its readers what to think about.*” illustrates its object of study. On the other hand, the “second level agenda-setting”, sometimes called *attribute agenda-setting*, studies the *objects* (in a social psychology way, where an *attitude object* designate a thing that an individual has an attitude or opinion about) present in the media agenda. When the media talks about an object some attributes are emphasized, and others not.

The “second level agenda-setting” is linked with *framing* [5] [6]. To frame is to *select some aspects of a perceived reality and make them more salient in a communicating text, in such a way as to promote a particular problem definition, causal interpretation, moral evaluation and/or treatment recommendation* (Robert Entman) [3]. This stage of agenda-setting theory can be summary in the phrase “*the media not only can be successful in telling us what to think about, they also can be successful in telling us how to think about it.*”

Other interesting stage of agenda-setting concerns with the sources of media agenda, i.e. if the media set the public agenda, *who sets the agenda media?* Within this framework, *intermedia agenda-setting* observes the competition between different media and how they influence each other. The competition between mass media for the same audience can lead to a homogenization of the agendas [7], which is in the opposite direction of one of the assumptions that predict the end of agenda-setting, as was mentioned before.

Social Media

The interaction between media and public is rather complex, and the advance of social media gives new direction in the exploration of this aspect of agenda-setting theory.

In [8] the question of *causality* is faced up. Their study shows that sometimes the traditional media set the agenda and sometimes, the social media does. They show that social media is always more interested in social issues than the traditional one, and despite the existence of correlation, the social media agenda can not be seen as a *slave* of the traditional media. [9] shows that the newspapers and Twitter have an opposite reaction to the changes of the unemployment rates, and for instance, in Argentina, [4] shows that not necessarily the journalists and public’s preferences coincide, by studying the most viewed articles and the home page articles in online news sites.

Other works study various aspects of social media, such as the danger of selective exposure [10][11][12] and the role of media organization in Twitter discussions [13][14].

0.1 Topic detection

To cite some works on topic detection

Argentinian works

In Argentina, the role of Mass Media has been particularly discussed during the administration of Argentina's president *Cristina Fernández de Kirchner* (from 2007 to 2015), where the government started to confront with news organizations that were critical to its management as they were an opposition party [15]. Several works by E. Zunino, like [16] and [17], studied the coverage given by the main newspapers of particular events which have happened during that period. On the other hand, Sagarzazu studied the bias present in the argentinian media respect to *Hugo Chávez Frias* as a reflection of the ideology of the respective newspapers [18].

Our contribution

In this work we propose a simple method to the study of Mass Media and audience response through topic detection algorithms. Our work intends to contribute another quantitative approach which complements the agenda-setting theory describe above.

On the other hand, we aim to take an insight about Media dynamics and Public response in order to create useful tools at the time of constructing mathematical models about the interaction between Media and Public, investigation that we started in [19].

Results

Description of the work

We analyzed a corpus of news' articles that were published between July 31th and November 5th in the section *Politics* of the electronic editions of the Argentinian newspapers *Clarín*, *La Nación*, *Página12*, and the news portal *Infobae*. The first two lead the sale of printed editions in *Buenos Aires* city, but *Clarín* reaches roughly two times the readers of *La Nación*, and ten times the readers of *Página 12* ^[1]. On the other hand, *Infobae* has the website with more visitors, above the websites of *Clarín* and *La Nación* ^[2].

The corpus analyzed is constituted by 2908 politics articles of *Clarín*, 3565 of *La Nación*, 3324 of *Página 12*, and 2018 of *Infobae*. Except *Página 12*, all articles were taken from the section *Política* of the respective news portal, while the articles which belong to *Página 12* were taken from the section *El país*.

The analysis made basically consist of topic detection over the corpus' articles in order to describe it a set of topics which evolve over time. Topic detection is a powerful computational technique that allows us to analyze a big amount of texts that can be impossible otherwise [20]. For a careful description of the methodology implemented please see section 0.1. This methodology not only gives us the evolution over time of the topics, but also a set of keywords that allow us to interpret and understand what the topics are talking about. We take advantage of this keywords in order to make queries to the *Google Trends* tool and have an insight of what the public interests are.

After all this proceedings, which we are going to give more details during the description and discussion of the work, we obtain two objects of study which we

^[1]www.ivc.org.ar

^[2]<https://www.alexa.com/topsites/countries/AR>

call the **Media Agenda (MA)** and the **Public Agenda (PA)**. The first one is the set of topics detected in the corpus, and the second one belongs to do the respective Google Trends queries. Both are normalized so the Agendas are described as a time dependent distribution over the topic's space, where the time scale is day by day. Therefore the Agendas give us the relative importance of a given topic respect the other (i.e. it does not give us absolute values such as the number of titles associated).

As was mentioned above, the starting point of our analysis is the topic decomposition of the corpus. We chose to decompose the corpus in 10 topics, which four of them we interpret to be talking about the same macro-topic, so we joined them in a topic called *Elections* as can be seen below. So the final description is made with only 7 topics. We opted for choosing this arbitrary number of topics in order to describe the corpus with as least information as possible, but it is not more than an arbitrary decision, validated in some manner by our knowledge of the corpus. After the proceedings describe in section 0.1, we construct the **Media Agenda (MA)** and the **Public Agenda (PA)** as time-dependent distributions.

Media and Public agenda: a qualitative approach

In figure 1 we show a bump chart of the **MA** and **PA**. A bump chart is a very useful visualization tool for displaying the relative weight of the topics and at the same time their ranking, putting on the top the most important topic at a particular date. The topics' names are introduced in the figure, where we also point out some important events related to them.

In figure 2, we show the wordclouds of the keywords that define each topic, where the size of the word belongs to the importance in the topic's definition given by the topic detection algorithm. In green color, we point out the words involved in the Google Trends queries in order to construct the Public Agenda. The queries employed in Google Trends are specifically shown in table ?? together with the linear correlation between the topics' temporal profiles that form the Public Agenda and their counterparts in the Media Agenda.

In figure 2, we also show the radar plots of the distributions made up by averaging over time the Agendas of figure 1, as a way to reduce that information which would help in farther discussions. We use radar plots as an alternative of histograms because its facility in distributions comparison.

The figures introduced above show in a qualitative way the differences between the agendas, and the dynamics of the topics, i.e. the range of dates in which a given topic was an important one, for which topic it was replaced, and so on. We can see, for instance in the radar plot of figure 2, a more interest of the audience in the topic *Missing person* than the Media, or inversely in the topic *Prosecutor's death*, when we see all the period analyzed as a whole.

On the other hand, the linear correlations of table ?? are in all cases positive and statistically significant, which we interpret as a form of validation of the topics found in the corpus and the keywords that describe it. We expected that the Media's and public's interest should generally follow a similar a pattern due to the external events, although the periods where those differ are of particular interest for us. A non positive (or a non significant) correlation may imply that we are not properly detecting the keywords or features that describe a particular topic, so the Google

Trends' pattern would not be able to reflect a similar behavior that its counterpart in the Media.

A quantitative approach

Agenda diversity

In order to quantify the similarities and differences between the Media Agenda and the Public Agenda, we start by asking how is the distribution of each agenda among the topic's space. In particular, we measure how diverse is each agenda. Following [21], we calculate the normalized Shannon's entropy (H , see eq.8) in order to measure the diversity of the **MA** and **PA**.

In figure 3 we can see the value of H as a function of time. We can see that there are periods where the diversity is lower than the usual. In particular, we pay attention to four dates in the Public Agenda which we detect as outliers of the typical behavior. The small value in the diversity is due to the fact that the most important topic attracts practically all the attention of the public, as can be seen in the radar plots included also in figure 3.

Two of the points belong to the topic *Elections* and coincide with the primary and general legislative elections that took place in August 13th and October 22th. In the other hand, the other two point belong to the topic *Missing person*: The first one a month after the disappearance of Santiago Maldonado, and the second one, when the Santiago Maldonado's body was found, a few days before the general legislative election (see section 0.1).

From the measure of H we can also see that the median of the Public Agenda diversity ($H_{PA} = 0.70$) is lower than the each the Media Agenda's one ($H_{MA} = 0.85$) being this difference statistical significant ($p < 10^{-9}$). We conclude that its an important fact about audience behavior: given a finite set of topics, **the Public Agenda is less diverse than the Media Agenda**, because the public seems to focus more in the most important topic than the Media can do.

Public Agenda's distance

The measurement of the Shannon's entropy made above is an independent property of each distribution. Here we directly compare the Agendas by computing the Jensen-Shannon distance. We again identify outliers and aim to interpret them. In figure 4 we show the Jensen-Shannon distance as a function of time. We inspect three points that seems to be of particular interest. In all cases, the radar plots shows that a greater distance is associated with a more interest of public in the topic *Missing person*.

Points (b) and (c) shows that both the public and the Media are interest in that topic, but the Media have to cover other topics, so the distance value can be seen as a derivation of the diversity effect discussed in the last section. However point (a) seems to show an interest of the public in the topic *Missing person* that it is not reflected in the Media. In figure 1 we can see that this topic reached the first place in public's interest before that in the Media. We associate this fact with a campaign made in social media like Facebook and Twitter in August 26th, that paid for the appearance of Santiago Maldonado and had a great repercussion, maybe at first underestimated by the Media (see section 0.1).

It is important to recall that it is our interpretation based on the knowledge of the context, and that we are not studying causality, i.e. we can't say, for instance, that in this case the Public Agenda set the Media Agenda. However, the Jensen-Shannon distance, in conjunction with the measurement of the agenda diversity given by the Shannon entropy, give an insight of independent behavior of the Public and the Media, and its identification can be a starting point to study the Media reaction to a change in audience's interests.

Media Agenda: differences between newspapers

In this section we leave aside the Public Agenda and we study how the Media agenda varies when one consider the newspapers separately. In figure 5 we show the bump charts which belongs to each newspaper analogously to figure 1. The topics are the same introduced in the wordclouds of figure 2, but at computing the topics' weights the articles are separated by newspaper. We also show the radar plots showing the average distribution, as made in figure 2.

In figure 5 we can see in a qualitative way the slightly differences between the newspapers' agendas. For instance, we can see how *Página 12* gave more importance to the topics *Missing person* and *Social leader*, while it did not pay too much attention to the *Former Planning minister* as the others did.

Independent behavior

In other to detect an independent behavior of a newspaper respect the others we again calculate the Jensen-Shannon distance between the newspapers agenda and the Media Agenda. Note that this is the distance between the distributions of figure 5 and the top panel of figure 1.

In figure 6 we show the Jensen-Shannon distance is a function of time. We detect three points as outliers, although we discard the point **(b)** due to the low information of *Infobae* in that period. The other two points belongs to a difference between *Página 12* and the other newspapers. We interpret that the coverage of the topic *Missing person* is the principal cause of the outliers. *Página 12* paid more attention to it than the Media Agenda at point **(a)** when the first notices of the Santiago Maldonado's disappearance before the primary elections, and in point **(c)** when a march two months after the disappearance took place, and had not the same coverage as the other march a month before (see section 0.1). Also, in point **(c)**, it can be seen a greater coverage of *Página 12* in the topic *Social leader* while the others seemed to be more interested in the topic *Former Vice-President*.

Coverage bias

The greater coverage in the topic *Missing person* by *Página 12* is even more clear if we inspect the temporal profile of the topic and compare the coverage given by each newspaper. A difference in the coverage is what it is called *coverage bias*. In figure 7 we show the temporal profile of the topic *Missing person* (panel (a)) and the topic *Former Planning minister* in panel (b), as an example where the behavior is the opposite, as can be seen below.

From panel (a) of figure 7, we can see a more coverage of *Página 12* respect the other newspapers at the first of the period. We can for instance quantify this

difference by the median of the signals. If we focus in the period between July 31th and August 27th, the median of the topic relative weight for *Página 12* is roughly 0.14 and this is statistically significant larger ($p < 10^{-7}$) than other medians, which are lower than 0.05. Analyzing the same period, but in panel (b), we again can show that the median in *Página 12*, which is roughly 0.01, is lower than the others, which oscillate around 0.05 ($p < 10^{-3}$). This quantification is proposed as method of studying coverage bias in the context of the methodology implemented in our work.

Finally, in figure 7 we also show wordclouds of topic's keywords but separate those that are most frequent mentioned by each newspaper, filtering the words that are common to all and basically define the factual details of the topic. Although most words are less informative, some are interesting to inspect, for instance, the word *represión* (repression) when *Página 12* talks about the topic *Missing Person* and the word *Kirchner* which is employed by all newspapers except *Página 12* when talk about the topic *Former Planning minister* (see section 0.1). We don't go farther but we think that a more deeply study of topics' keywords is a first approximation in the study of framing, which will constitute the core of futures works.

Conclusions

The study of Mass Media, and in particular the agenda setting theory, can be empowered by the used of data mining or machine learning algorithms. In this work, through the implementation of a topic detection algorithm we could describe the Agenda of the Media as a distribution which evolves over time and which is defined in a topic's space which emerges from the analysis of the corpus. This gave us an insight of how we can construct and follow the Public's interests, the Public Agenda, in order to compare with the Media Agenda, i.e. Media interests.

Given the Agendas, we found that the Public one is usually less diverse than the Media, showing that when there is a very attractive topic, the audience focus on this one, when the Media has to cover the other too. On the other hand, the measurement of distances between Agendas can be employed to rapidly detect periods when the Public may have an independent behavior respect to the Media.

The methodology implemented allow us to detect coverage bias in newspapers and gave us a first approximation in the theory of framing.

Competing interests

The authors declare that they have no competing interests.

Author's contributions

Text for this section ...

Acknowledgements

Text for this section ...

Author details

¹Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428 Buenos Aires, Argentina. ²Instituto de Física de Buenos Aires (IFIBA), CONICET, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428 Buenos Aires, Argentina.

References

1. McCombs, M.E., Shaw, D.L.: The agenda-setting function of mass media. *Public opinion quarterly* **36**(2), 176–187 (1972)
2. McCombs, M., Valenzuela, S.: *Agenda-setting theory: The frontier research questions*. Estados Unidos: Oxford handbooks online (2014)

3. McCombs, M.: A look at agenda-setting: Past, present and future. *Journalism studies* **6**(4), 543–557 (2005)
4. Mitchelstein, E., Boczkowski, P.J., Wagner, C., Leiva, S.: La brecha de las noticias en argentina: factores contextuales y preferencias de periodistas y público. *Palabra Clave* **19**(4) (2016)
5. Guggenheim, L., Jang, S.M., Bae, S.Y., Neuman, W.R.: The dynamics of issue frame competition in traditional and social media. *The ANNALS of the American Academy of Political and Social Science* **659**(1), 207–224 (2015)
6. Tsur, O., Calacci, D., Lazer, D.: A frame of mind: Using statistical models for detection of framing and agenda setting campaigns. In: *ACL* (1), pp. 1629–1638 (2015)
7. Vargo, C.J., Guo, L.: Networks, big data, and intermedia agenda setting: An analysis of traditional, partisan, and emerging online us news. *Journalism & Mass Communication Quarterly*, 1077699016679976 (2017)
8. Russell Neuman, W., Guggenheim, L., Mo Jang, S., Bae, S.Y.: The dynamics of public attention: Agenda-setting theory meets big data. *Journal of Communication* **64**(2), 193–214 (2014)
9. Soroka, S., Daku, M., Hiaeshutter-Rice, D., Guggenheim, L., Pasek, J.: Negativity and positivity biases in economic news coverage: Traditional versus social media. *Communication Research*, 0093650217725870 (2017)
10. Feezell, J.T.: Agenda setting through social media: The importance of incidental news exposure and social filtering in the digital era. *Political Research Quarterly*, 1065912917744895 (2017)
11. Messing, S., Westwood, S.J.: Selective exposure in the age of social media: Endorsements trump partisan source affiliation when selecting news online. *Communication Research* **41**(8), 1042–1063 (2014)
12. Bakshy, E., Messing, S., Adamic, L.A.: Exposure to ideologically diverse news and opinion on facebook. *Science* **348**(6239), 1130–1132 (2015)
13. Calvo, E., Aruguete, N.: Time to# protest: Polarization and time-to-retweet in argentina (2016)
14. Malik, M.M., Pfeffer, J.: A macroscopic analysis of news content in twitter. *Digital Journalism* **4**(8), 955–979 (2016)
15. Mitchelstein, E., Boczkowski, P.J.: Information, interest, and ideology: Explaining the divergent effects of government-media relationships in argentina. *International Journal of Communication* **11**, 20 (2017)
16. Zunino, E., Aruguete, N.: La cobertura mediática del conflicto campo-gobierno. un estudio de caso. *Global Media Journal* **7**(14) (2010)
17. Koziner, N., Zunino, E.: La cobertura mediática de la estatización de ypf en la prensa argentina: un análisis comparativo entre los principales diarios del país. *Global Media Journal* **10**(19) (2013)
18. Sagarzazu, I., Mouron, F.: Hugo chavez's polarizing legacy: Chavismo, media, and public opinion in argentina's domestic politics. *Revista de Ciencia Política* **37**(1) (2017)
19. Pinto, S., Balenzuela, P., Dorso, C.O.: Setting the agenda: Different strategies of a mass media in a model of cultural dissemination. *Physica A: Statistical Mechanics and its Applications* **458**, 378–390 (2016)
20. Griffiths, T.L., Steyvers, M.: Finding scientific topics. *Proceedings of the National academy of Sciences* **101**(suppl 1), 5228–5235 (2004)
21. Boydston, A.E., Bevan, S., Thomas, H.F.: The importance of attention diversity and how to measure it. *Policy Studies Journal* **42**(2), 173–196 (2014)
22. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *Journal of machine Learning research* **3**(Jan), 993–1022 (2003)
23. O'Callaghan, D., Greene, D., Carthy, J., Cunningham, P.: An analysis of the coherence of descriptors in topic modeling. *Expert Systems with Applications* **42**(13), 5645–5657 (2015)
24. Řehůřek, R., Sojka, P.: Software Framework for Topic Modelling with Large Corpora. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pp. 45–50. ELRA, Valletta, Malta (2010)
25. 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., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
26. Fuglede, B., Topsøe, F.: Jensen-shannon divergence and hilbert space embedding. In: *Information Theory, 2004. ISIT 2004. Proceedings. International Symposium On*, p. 31 (2004). IEEE

Context

Elections

Two legislative elections were celebrated during the period in great part of the Argentina: Primary elections on August 13th and the general elections on October 22th. A special focus was put on the elections in the Buenos Aires province, where the former President *Cristina Fernández de Kirchner* participated as a senator candidate representing the alliance *Unidad Ciudadana*, confronting to the alliance *Cambiamos*, which is the alliance of the current President *Mauricio Macri* and the current governor of Buenos Aires province, *Maria Eugenia Vidal*.

Current President

Mauricio Macri is the current Argentinian President, since December 2015. Most articles in political sections are logically devoted to him under different contexts. However, it is important to point out that during the period analyzed, and specially after the general elections of October 22th, a controversial labour reform promoted by the government was been discussing.

Missing Person

Santiago Maldonado vanished on August 1st after a minor clash between the Gendarmerie (Border Guards) and a group of Mapuches, which recognize as themselves the original population of an area in the Patagonia. Since that event, the *Mauricio Macri*'s administration was accused by several people as the responsible of a **forced disappearance**.

A very massive campaign in social media took place on August 26th, under the motto "Where is Santiago Maldonado?", followed by two massive protest marches to the *Plaza de Mayo* took place on September 1st and October 1st, which the first one had a great repercussion due to several incidents that took place during the march.

The body of *Santiago Maldonado* was found on 17th October in the *Chubut* river, near the place where he was seen the last time, and the autopsy report told that *Santiago Maldonado* had died from 'asphyxia after being submerged,' with no injuries on his body. However the responsibility of the current administration is still being discussed.

Former Planning minister and Former Vice-President

Julio de Vido was the Planning minister during the administration of *Nestor Kirchner* and *Cristina Fernandez de Kirchner* (2003-2015). In 2015, he was elected to integrate the Chamber of Deputies, which finally voted to strip *De Vido* of his congressional immunity over corruption allegations and was immediately jailed on October 27th. *Amado Boudou* was the Vice-President of the *Cristina Kirchner's* administration. *Boudou* was arrested on November 3th on charges including money-laundering and hiding undeclared assets.

Social leader and Prosecutor's death

Milagro Sala is an indigenous leader's. She has been incarcerated under pre-trial detention ever since she was first detained in January 2016. She faces allegations of embezzlement related to government funding for housing projects managed by *Tupac Amaru*, her social organization. *Sala* accused the government of "violating her human rights", and several people think that she is a political prisoner of the *Mauricio Macri* administration.

Alberto Nisman was a special prosecutor who were investigating the 1994 terror attack on the Argentine Israeli Mutual Association (AMIA), until his suspicious death on January 2015. During the period analyzed in this work, a team of experts led by the Gendarmerie (Border Guard) concluded that late prosecutor's death may have been a case of murder, not suicide.

Comparison with LDA

In this section we apply other topic model, Latent Dirichlet Allocation [22] (LDA), to our corpus and compare its results to the shown in this paper. Due to the increasingly use of LDA, we think that a few words about the performance of LDA in our work is necessary.

Naturally the topics found with LDA may not coincide with the NMF ones. However, one expects that the corpus under studying should be in some manner robust to the election of the topic model. On the other hand, as was discussed in [23], NMF can be a more suitable topic modeling method in certain domains, in the way that it produces more coherent topics, while LDA tends to return higher levels of generality and redundancy. Topic coherence is defined as the semantic interpretability of the terms used to describe a particular topic, although the coherence of a topic may depend on the end user's expectations.

We define a simple coherence measure defined in equation 1, where d_{ij} is the number of documents where the term i and term j appear simultaneously, and d_x is the number of documents where appears the term x . The summation is over the N top terms of the topic. It's important to note that if two terms have no co-occurrences, the contribution to the summation is zero, and if these ones appear only together the contribution is one. A topic with higher coherence is a topic where the terms that define it co-occurrence frequently.

$$TC = \sum_{i < j}^N \frac{2d_{ij}}{d_i + d_j} \quad (1)$$

We perform a decomposition into 10 topics using LDA with the python module *gensim* [24], which allow us to modify the number of times the corpus is read, improving the coherence of the topics. Unlike to what we see with NMF, the LDA's performance depends strongly on the initial condition of the algorithm. After 10 iterations, we chose the one with highest mean topic coherence, and compared this with the NMF results.

In figure 8 we show the temporal profiles of topics *Elections* and *Missing Person* for both NMF and LDA. The association between topic models was simply made by looking at the topics which share common keywords. As can be seen from the figure and the table 2, those LDA topics which can be linked to NMF ones or to a combination of these, show a temporal profile highly correlated.

Nevertheless, LDA returns other topics which can not be directly associated, some of them composed of very general words. By keeping only those topics which can be associated with NMF and re-defining the Media Agenda over this topic space with reduced dimension, we observed similar results by both methods. The same procedure is proposed in absence of an alternative topic model to which make the comparison: Keep only those topics easily interpretable and define the Agendas over this reduced space.

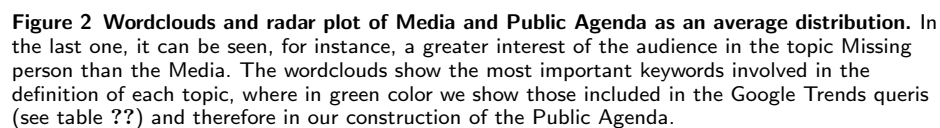
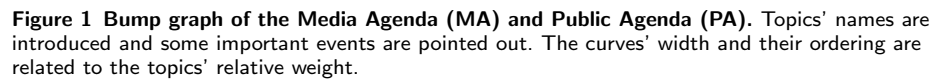
Figures

Tables

Methodology

Numerical representation of the corpus

In order to perform the analysis of the articles in the corpus, we describe them as numerical vectors through the *term frequency - inverse document frequency (tf-idf)* representation. *Tf-idf* gives greater values to terms that appear in less documents of the corpus (i.e more specific terms) and/or to those which appear more frequent in a document. Given the set of terms made up by all the corpus' words after removing the non-informative ones, such as prepositions and conjunctions, the *tf-idf* algorithm represents the i -document as a vector $v_i = [x_{i1}, x_{i2}, \dots, x_{it}]$, where the component x_{ij} is computed by the equation 2, where tf_{ij} is the number of times the j -term appears in the i -document, d is the number of documents in the corpus, and n_j is the number of documents where the j -term appears. Each document's vector is normalized to unit Euclidean length. Once the documents' vector are constructed, we join them in a document-term matrix (M), which has dimensions of number of documents in the corpus (d) per number of terms selected (t).



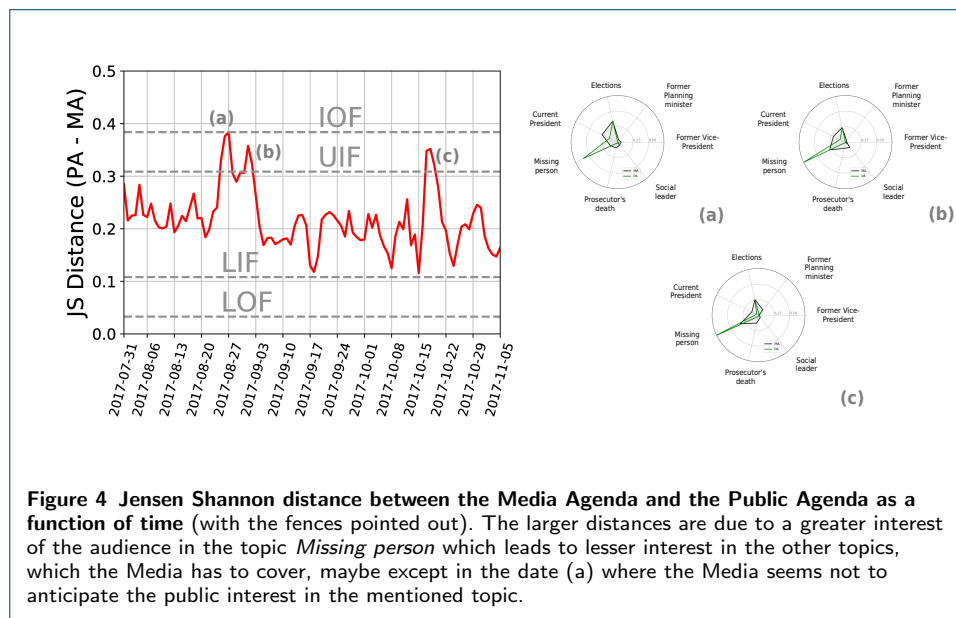
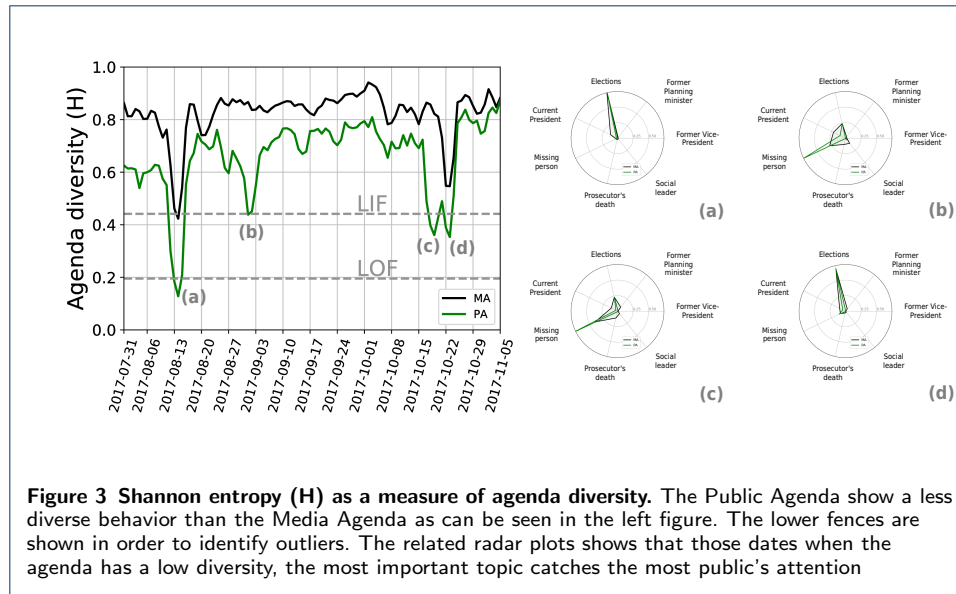


Table 1 Queries performed in Google Trends in order to made up the Public Agenda. We also shown the correlation between the topics' temporal profiles of the Public Agenda and their counterpart in Media Agenda. All correlation values are statistical significant ($p < 10^{-4}$).

Topic	Google Trends query	Correlation MA and PA
Elections	elecciones + cambiemos + cristina kirchner + maría eugenia vidal + unidad ciudadana	0.83
Missing person	santiago maldonado + juez otranto + mapuche + gendarmería + desaparición forzada	0.76
Former Planning minister	julio de vido + río turbio + tragedia de once + desafío + minnicelli	0.91
Current President	mauricio macri + cgt + reforma laboral + peña + triaca	0.80
Social leader	milagro sala + cidh + tupac amaru + pullen llermanos + morales	0.40
Prosecutor's death	alberto nisman + amia + memorándum con irán + timerman + juez bonadio	0.54
Former Vice-President	amado boudou + ciccone + juez lijo + vandenbroele + núñez carmona	0.91

$$\text{idf}_j = 1 + \log\left(\frac{1 + d}{1 + n_j}\right)$$

$$x_{ij} = \text{tf}_{ij} \cdot \text{idf}_j$$

(2)

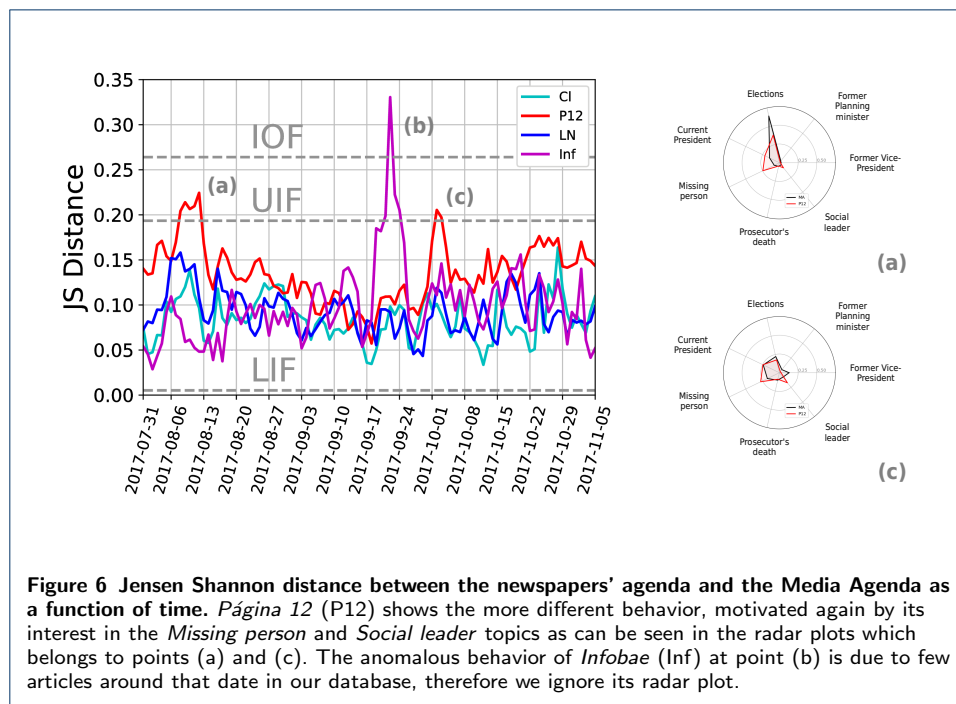
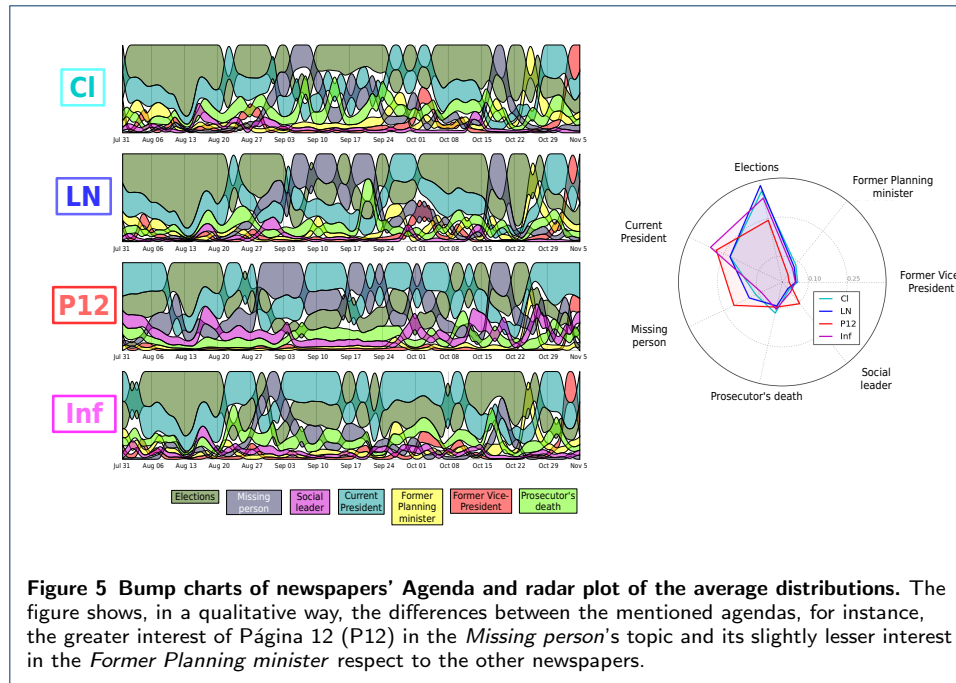
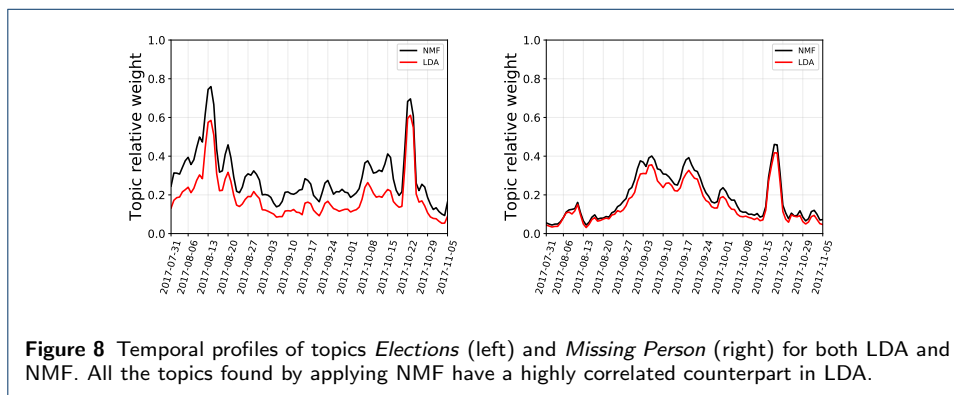
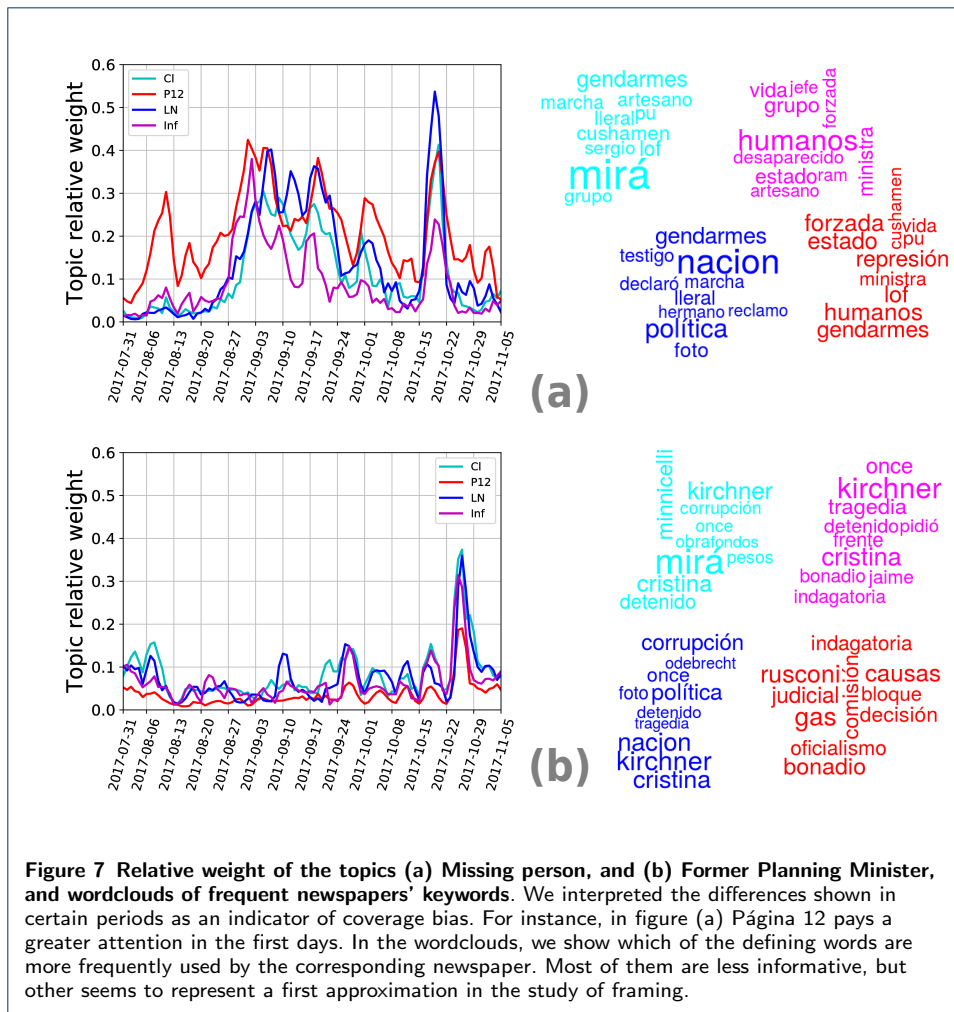


Table 2 Correlation between the temporal profiles of the topics found in NMF and associated topics in LDA.

	Correlation between NMF and LDA
Elections	0.98
Missing person	0.99
Former Planning minister + Former Vice-President	0.89
Current President	0.94
Social leader	0.94
Prosecutor's death	0.83



Topic detection

We perform *non-negative matrix factorization (NMF)* on the document-term matrix (M) in order to detect the main topics in the corpus. *NMF* is an algorithm which factorize a matrix M into two matrices W and H (eq.3), with the property that all three matrices have no negative elements. This non-negativity makes the resulting matrices easier to inspect, and very suitable for topic detection.

$$M^{(d \times t)} \sim H^{(d \times k)} \cdot W^{(k \times t)} \quad (3)$$

As can be seen in eq.3, the resulting matrix H has dimensions of number of documents per k , and matrix W has dimensions of k per number of terms. The number k is interpreted as the number of topics in the documents and it is a parameter that must be set before the factorization. The matrix H is interpreted as the representation of the documents in the topic-space, and the matrix W as the topics represented in the original term-space. This factorization usually can not be made exactly so it is approximated by minimizing the reconstruction error, i.e. the distance between matrix M and its approximated form $\tilde{M} = H \cdot W$. We performed *NMF* through the python module *scikit-learn* [25].

Topic interpretation and temporal profiles

The matrix H obtained by *NMF* gives the representation of the documents in the topic-space. In order to improve its interpretation we normalize each document's vector described in that space to unit l_1 -norm. Therefore the components of these vectors can be viewed as a degree of membership of a given document in the set of topics. The index of the largest component tells us which is the most representative topic of the document.

On the other hand, each row of matrix W represent the topic over the term-space. Therefore, the terms associated with the largest components of the i -row are the most representative ones and give an insight of what the topic is talking about.

We define the temporal profile of the topic i , $W_i(day)$ by the eq. 4 where $l(j)$ is the number of words of the document j , and h_{ji} is the degree of membership of document j on topic i . This definition allows all documents to contribute to any topic weight, providing by the fact that each document's vector can have non-zero components in more than one topic.

As last steps of the data's preprocessing, we filter the topics' weight in order to reduce the noise but keeping the most details as possible, by redefining $W_i(day)$ as the mean value of a 3 days width window, centered on the day, as described in equation 5.

Finally we normalize again all the temporal profiles in order to describe each newspaper as dsitribution over the topics' space which evolves over time. This last normalization prevent us against the differences in the number of articles that each newspaper publishes.

$$W_i(day) = \sum_j^d l(j) \cdot h_{ji} \cdot \delta_{d,day} \quad (4)$$

$$\tilde{W}_i(day) = \frac{1}{3}(W_i(day) + W_i(day - 1) + W_i(day + 1)) \quad (5)$$

Outliers identification

We identify outliers values in a data set of N observations by following the box plot construction. The quantities (called fences) in 6 are used, where $Q1$ is the lower quartile (range of the distribution where lies the 25th percent of the data), $Q3$ is the upper quartile (where lies the 75th percent of the data) and $IQ = Q3 - Q1$.

$$\begin{aligned} \text{lower inner fence (LIF)} &= Q1 - 1.5IQ \\ \text{upper inner fence (UIF)} &= Q3 + 1.5IQ \\ \text{lower outer fence (LOF)} &= Q1 - 3IQ \\ \text{upper outer fence (UOF)} &= Q3 + 3IQ \end{aligned} \quad (6)$$

A point above the upper inner fence considered a mild outlier and a point above an outer fence is considered an extreme outlier. The same holds for the lower fences. [3]

Jensen shannon divergence

In probability theory and statistics, the Jensen–Shannon divergence is a method of measuring the similarity between two probability distributions. It is based on the Kullback–Leibler divergence (D_{KL}), but have useful properties such as it is symmetric and it is always a finite value. The square root of the Jensen–Shannon divergence is a metric often referred to as Jensen–Shannon distance. [26]

$$\begin{aligned} D_{KL}(P||Q) &= - \sum P(i) \log\left(\frac{Q(i)}{P(i)}\right) \\ \text{JS Divergence}(P||Q) &= \frac{1}{2}[D_{KL}(P||M) + D_{KL}(Q||M)] \\ \text{JS Distance (JSD)} &= \sqrt{\text{JS Divergence}} \end{aligned} \quad (7)$$

[3] <http://www.itl.nist.gov/div898/handbook/prc/section1/prc16.htm>

Normalized Shannon's H Information Entropy

The normalized Shannon's entropy is a way to measure how spread is a distribution, taking the maximum value where all outcomes are equally probable in the case of a discrete distribution.

$$H[p] = \frac{-\sum_{i=1}^N p(x_i) * \ln(p(x_i))}{\ln(N)} \quad (8)$$

Additional Files

Additional file 1 — Sample additional file title

Additional file descriptions text (including details of how to view the file, if it is in a non-standard format or the file extension). This might refer to a multi-page table or a figure.

Additional file 2 — Sample additional file title

Additional file descriptions text.