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

Sebastián Pinto[1,2]*, Federico Albanese[3], Claudio O Dorso[1,2] and Pablo Balenzuela[1,2]

*Correspondence: spinto@df.uba.ar
[1]Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428, Buenos Aires, AR
[2]Instituto de Física de Buenos Aires (IFIBA), CONICET, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428, Buenos Aires, AR
Full list of author information is available at the end of the article
[†]Equal contributor

**Abstract**

The mass media play a fundamental role in the formation of public opinion, either by defining the topics of discussion or by making a degree of emphasis on certain issues. Directly or indirectly, people get informed by consuming news from the media. But which is the dynamics of the agenda and how the people become interested in their different topics? The Agenda Setting theory provides a conceptual framework in order to understand the role played by the mass media in public opinion formation, but the previous questions can not be answered without proper quantitative measures of agenda's dynamics and public attention. In this work we study the agenda of Argentinian newspapers in comparison with public interests through a quantitative approach by performing topic detection over the news, identifying the main topics covered and their evolution over time. We measure Agenda's diversity as a function of time using Shannon's entropy and differences between Agendas using Jensen-Shannon's distance. We found that the Public Agenda is less diverse than the Media Agenda, and we were also capable to detect periods of time where coverage of certain issues are biased (coverage bias).

**Keywords:** agenda-setting; opinion formation; topic detection; Mass Media influence

## Introduction

The understanding of the ecosystem of information flow and opinion formation is one of the challenges in complex social system research. Within this framework, many ingredients can be part of different mechanism of social influence. A major role in this ecosystem is played by the Mass Media outlets, which use to be the source of information of many people. Then, informed people tends to interact with each other, either via personal interactions or through social networks, giving rise to complex dynamics where opinions are shaping and changing in time. In this scenario, becomes essential to understand the role of Mass Media Influence in a given social group.

Naturally, how the Media affects public opinion was first explored in the area of social sciences. 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. In its basic stage, known as first-level agenda-setting [2], the theory focuses on the

comparison between the topics coverage by the media and the public agenda, i.e. the topics that the public consider as priority [3, 4, 5]. For instance, in Argentina, [6] examines the difference between public and journalists preferences and several works by E. Zunino like [7, 8], studied the coverage given by the main newspapers of particular events which have happened during the former administration of Argentinian president *Cristina Fernández de Kirchner* (from 2007 to 2015), where the government started to confront with news organizations that were critical to its management [9].

The Chapel Hill investigation also induced other several research directions [2]. One of them focused on detecting bias in the media, either by taking into account the number of mentions related to a preferred political party [10, 11] or by identifying the ideology through the position of the media regards to certain issues or actors [12, 13]. This research line can be linked with the theory of *framing* [14, 15], which is concerned with how the media emphasizes some attributes of an object, while understating others. Other investigations pay attention to the sources of the Media Agenda, theory known as *intermedia agenda-setting* [16, 17, 18], in which the competition between different media, and the influence exerted between each other, is observed.

Nonetheless, most of the works cited above follow typical methods of social science research, and fall in what we call a "static approach": They used to summarize certain issues present in news media and looking for effects in the audience. With the irruption of internet in our lives and the access to big data, a richer analysis is allowed, which may take into account the temporal dependence of the Media and Public attention. For instance, in [19] it is shown that the newspapers and Twitter have an opposite reaction to the changes of the unemployment rates; in [14], the competition of frames about gun control is explored; in [20], the authors show how fluctuations of Twitter activity in different regions depend on the location of terrorist attacks; and in [21], the complex interplay between the social media and the traditional one is followed over time on a set of predefined, but general, issues. However, one drawback of several works which perform a dynamical analysis of agenda-setting is that the research is usually made on a single issue or on a set of predefined issues. Frequently, the last are arbitrarily selected by the own researcher, and normally reflect general themes, such as "health" or "gun control".

The question about what issues should be selected can be solved by an useful tool employed in the analysis of large corpus of documents, but understated in the agenda-setting framework and its derivatives: Unsupervised topic modeling. It is an alternative to the dictionary-based analysis, which is the most popular automated analysis approach [22], and allows to work with a corpus without a prior knowledge, letting the topics emerge from the data. Although many works employ unsupervised topic modeling on news corpus, much of them emphasize the performance of the topic model over a labeled corpus, focusing on the proper detection of the topics [23, 24, 25]. In general, issues about the temporal profile of topics are embedded in the context of topic tracking [26, 27], or in the recognition of emerging topics in real-time [28], mostly applied to social media.

In this work we propose a novel method in order to study the dynamics of Mass Media Agenda, which consist in performing an unsupervised topic model on newspapers articles, and studying how the emerging topics evolve in time. We look also

at their correspondence with the Public and Social Network Agendas.Rather than focus on a single issue or a set of independent topics, we work with the agendas (the Media and the Public) as an object in their own. Our work intends to contribute another quantitative approach which complements the agenda-setting theory describe above. It mainly stands within the framework of first-level agenda setting, but issues about the comparison between media agendas and framing are also discussed.

On the other hand, in recent years new approaches to study social dynamics through several tools taken from statistical physics have proposed mathematical models to explore the interplay between Mass Media and society [29, 30, 31, 32, 33]. Much of them lack in being contrasted with real data. With our work we aim to gain a closer insight on the complex interaction between Media and Public, and provide a quantitative research that would be useful at the time of constructing better and more data-driven models.

## Materials and Methods

### The Media Agenda

In this work we analyze a three month period of the Argentinian Media's Agenda composed by a corpus of news' articles that were published between July 31th and November 5th of 2017. The articles come from the *Politics* section of the online 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* [34]. On the other hand, *Infobae* has the most visited website, much more than *Clarín* and *La Nación* [35]. The corpus analyzed is made up 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* (Politics) of the respective news portal, while the articles which belong to *Página 12* were taken from the section *El país* (The country).

The articles are described as numerical vectors through the *term frequency - inverse document frequency (tf-idf)* representation [36]. Given the set of terms contained in the corpus' words after removing non-informative ones, such as prepositions and conjunctions, the *tf-idf* algorithm represents the $i$-document as a vector $v_i = [x_{i1}, x_{i2}, ..., x_{it}]$, where the component $x_{ij}$ is computed by the eq.(refec:tfidf), where $\text{tf}_{ij}$ 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 vectors is normalized to unit Euclidean length. Once the documents' vector are constructed, we put them together in a document-term matrix ($M$), which has dimensions of number of documents in the corpus ($d$) per number of terms selected ($t$).

$$x_{ij} = \text{tf}_{ij} \cdot \text{idf}_j = \text{tf}_{ij} \cdot [1 + \log(\frac{1 + d}{1 + n_j})] \tag{1}$$

We perform *non-negative matrix factorization (NMF)* [36, 37] on the document-term matrix ($M$) in order to detect the main topics in the corpus. A topic is a group of similar articles which roughly talks about the same subject. *NMF* is an algorithm which factorize the matrix $M$ into two matrices $W$ and $H$ (eq.(2)), 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 [1].

$$M^{(d \times t)} \sim H^{(d \times k)} \cdot W^{(k \times t)} \tag{2}$$

Such as 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 therefore interpreted as the number of topics in the documents and it is a parameter that must be set before the factorization. In this work, we arbitrarily set $k = 10$. We found this value to be a suitable one based on our knowledge of the corpus.

The matrix $H$ is the representation of the documents in the topic-space. Each row is normalized to unit $l_1$-norm so that their components can be viewed as a degree of membership of a given document in the set of topics. In particular, the index of the largest component tells us which is the most representative topic of the document. On the other hand, $W$ gives the topics representation in the original term-space. The largest components of a row give the most representative words of each topic (called keywords) and therefore an insight of what the topic is talking about. Since the factorization of eq.(2) usually can not be made exactly, it is approximated by minimizing the reconstruction error, i.e. the distance between matrix $M$ and its approximated form $\tilde{M} = H \cdot W$. The *NMF* factorization through the python module *scikit-learn* [38].

After performing NMF, we represent a time-dependent Media Agenda as a time evolving distribution of topics. The daily weight of the topic $i$, $W_i(day)$, is calculated following eq.(3), where $l(j)$ is the number of words of the document $j$, $h_{ji}$ is the degree of membership of document $j$ on topic $i$, $d_j$ is the date of document $j$, and $\delta$ is the Kronecker delta. Providing by the fact that each document's vector can have all non-zero components, it is allowed that a document contributes to more than one topics' weight. In order to reduce noise, we apply a linear filter with a three day wide sliding window, and finally we normalize the temporal profiles in order to describe each newspaper agenda as a distribution over the topics' space which evolves over time.

$$W_i(day) = \sum_j l(j) \cdot h_{ji} \cdot \delta_{d_j, day} \tag{3}$$

### The Google and Twitter Agendas

The other side of time-dependent Media Agenda is to have some measure of the public interest about these topics. With this goal, we use *Google Trends* and *Twitter* as a proxy of public interest by looking for the same topics in the same period of time. We take advantage of topics' keywords to make queries in the *Google Trends*

---

[1]Even though there are other techniques in topic detection, as for instance LDA (Latent Dirichlet Allocation), the NMF decomposition suites perfect for the kind of corpus we have analyzed in this work (A detailed comparison of both, NMF and LDA methods could be found in Supplementary Material).

tool and into the advance search tool of Twitter to get the relative weight searches and tweets in both platforms.

In section **Results** we give a more detailed description of the keywords involved in the construction of the Google and Twitter Agendas.

### Normalized Shannon's H Information Entropy

In fact, one of the key of our work is the representation of the Agendas as time evolving topic distributions, and the association between the different measures that we can derive from them with the social events observed during the period. In order to do that, the first measure we emphasize is the concept of *diversity of the attention*. The *diversity* is a very important variable that must be taken into account when dealing with multiple issues [39], due to the fact that it tells us how the attention is distributed across the different topics of discussion. As was proposed in [39], we take the normalized Shannon's entropy as a very suitable way to quantify the diversity within our framework.

The normalized Shannon's entropy $H[p]$ referred in eq.(4) give us a measure of how spread is a -discrete- distribution, taking the maximum value of 1 when all outcomes are equally probable, and 0 when there is just one possible outcome. This last case -or approximated ones- will be of particular interest for us because it tells about a topic which absorbs all the attention either from the Media or the Public, while the former shows an absolutely unbiased interest towards any of the topics.

$$H[p] = \frac{-\sum_{i=1}^{N} p(x_i) * ln(p(x_i))}{ln(N)} \tag{4}$$

### Jensen-Shannon distance

While the diversity is a property of each distribution, a natural question that arise from comparing different distributions is how similar they are. For instance, we will be particularly interested in measuring the similarity between the Media and Public Agendas, because those dates when the similarity is low tell us about distant interests of the Media and the Public. We proposed here to measure the similarity between distributions by the use of the Jensen-Shannon distance ($JSD$). When the similarity of the distributions is low, the distance between them is high.

The Jensen-Shannon distance ($JSD$) is a metric between distributions based on the Jensen-Shannon divergence ($JS_{Div}$) [40], which is in turn a symmetric version of the well-known Kullback-Leibler divergence ($D_{KL}$, eq.(5)). We recall that the $JSD$ has the advantage of being symmetric and also a well-defined distance, which make it conceptually easy to deal. As can be seen in eq.(6) the $JSD$ between the distributions $P$ and $Q$ is simply the square root of the Jensen-Shannon divergence.

$$D_{KL}(P||Q) = -\sum P(i)log(\frac{Q(i)}{P(i)}) \tag{5}$$

$$JSD(P,Q) = \sqrt{JS_{Div}(P,Q)} = \sqrt{\frac{1}{2}[D_{KL}(P||M) + D_{KL}(Q||M)]} \tag{6}$$

Outliers identification

As was mentioned above, once the concepts of diversity and distance between distributions were defined, we are particularly interested in those dates when these concepts take extreme values: When the diversity is low, we know that a topic is absorbing much of the attention of any of the Agendas; on the other hand, when the distance between two distributions is high, we can conclude that they have distant interests. However, we still lack the definition of what a low diversity or a high distance are. We face this problem by treating the measures of each observable as random samples from a population with unknown distribution and identifying those extreme values as outliers of the distribution. In order to detect these outliers we follow the popular box-plot construction proposed by Tukey [41], which is a simple data-driven method and make no assumption about the distribution of the data.

The box-plot construction is as follows: First a quartile division of the $N$ observations is made, by naming $Q1$ as the lower quartile, $Q2$ the median of the distribution, and $Q3$ the upper quartile. Recall that $Q1$ ($Q3$) is defined to be the division where the 25th (75th) percent of the observations lies below (by definition, the median $Q2$ separates the distribution in two equal parts). On the other hand, the inter-quartile range $IQ$ is defined as $IQ = Q3 - Q1$. This is the range where the bulk of the data lies inside.

We are not interesting in the visualization of the box-plot in its own but instead in its procedure to identify outliers. Therefore from the identification of the quartiles, new quantities called *fences* are defined in eq.(7): The *lower inner fence* ($LIF$), the *upper inner fence* ($UIF$), the *lower outer fence* ($LOF$), and the *upper outer fence* ($UOF$). The fences can be interpreted as the limits of the distribution. We then have all the ingredients to label a point as an outlier: A point which lies above the upper inner fence is considered a *mild outlier*, while a point that lies above the upper outer fence is considered an *extreme outlier*. The same holds for the lower fences, i.e. if a point lies below the lower inner (outer) fences is considered as a mild (extreme) outlier [42].

$$
\begin{aligned}
LIF &= Q1 - 1.5IQ \\
UIF &= Q3 + 1.5IQ \\
LOF &= Q1 - 3IQ \\
UOF &= Q3 + 3IQ
\end{aligned}
\tag{7}
$$

We will indicate the proper fences in each figure either when the diversity or the distance is being analyzed. We will paid attention not only to those values labeled as outliers, but also to those that are next to any of the fences despite not being strictly defined as that.

## Results

We initially focus in the ten most important issues in the three months period corpus of news reported above. These ten topics, represented in the word clouds of figure 1, were reduced to seven (shown in the radar plot of the same figure) given the similarities found in, for instance, three of them that could be interpreted as part of

the same macro-topic called *Elections* and two others which were classified as part of a macro-topic called *Missing person*. The meaning of the topics or macro-topics is contextualized in the **Supplementary Material**.

Following the procedure described in the previous section, we construct the **Media Agenda (MA)** and the **Public Agenda (PA)**, in both its Google and Twitter derivations, as time evolving topic distributions.

### The Public and Media agendas

In figure 1 we show a ten topic decomposition of the whole corpus using radar plots for the Media **MA** and Public agendas **PA** discriminated by **GT** (Google Trends) and **Tw** (Twitter). In this figure we also show the wordclouds of the keywords that define each topic, where the size of the word reflects its importance in the topic's definition. In green color, we point out the words involved in the Google Trends and Twitter [2] queries in order to construct the Public Agenda. The queries employed are also specified in table 1. On the other hand, in table 2 we show the linear correlation between the topics' temporal profiles from the Public Agenda and their counterparts in the Media Agenda.

We can observe that both Public Agendas (*GT* and *Tw*) look similar in this representation, but they show specific differences with the Media Agenda. For instance, a greater interest of the audience in the topic *Missing person* than the Media is observed, or inversely, a lower interest in the topic *Prosecutor's death* takes place. However, this static representation is not able to show the complex dynamics of agenda's evolution and the importance of punctual and specific facts which can erase or amplify their differences.

This can be observed in figure 2 where the time evolution of the topics is shown in a bump chart of the Agendas. The bump chart provides a clear visualization of the relative weight of the topics at the same time with their ranking. In this figure we also highlight some important events related to the dynamics of the topics. It is possible to appreciate how the main topic change in time and have a glance of the qualitative differences between agendas. In particular, it can be observed some differences between Public Agendas (*GT* and *Tw*) that were not observed in the previous figure, as for instance, the persistence of main topics is longer in Twitter than in Google Trends. This is more evident at the end of the analyzed period, where topics discussed in Google Trends shows more responsive to change in Media Agenda than in Twitter, suggesting the existence that the interaction between people is stronger in the analyzed social network.

---

[2]Due to different characteristics in the search tool of Twitter, we adapted the queries employed here but preserving at least at we can the most important keywords. The queries employed by topic were:

Elections: elecciones + cambiemos + kirchner + massa + randazzo

Missing person: maldonado + otranto + gendarmería + desaparición

Former Planning minister: vido + desafuero + minnicelli + baratta

Current President: macri + cgt + laboral + triaca

Social leader: sala + cidh + tupac + amaru + pullen + llermanos + morales

Prosecutor's death: nisman + amia + memorandum + timerman + bonadio

Former Vice-President: boudou + ciccone + lijo + vandenbroele + carmona

The linear correlations between same topics of $PA$ and $MA$ were also calculated. In all cases, we found that they are positive and statistically significant, as can be shown in table 2. This can be interpreted as a validation of the topics found in the corpus and the keywords that describe it. Even though it is expected that the Media's and public's interest should generally follow a similar a pattern, mainly driven by external events, we are interested in those periods where they significantly differ. However, a non positive (or a non significant) correlation may imply, besides the obvious conclusion of agenda's disengagement, that we could eventually fail in properly detect the keywords or features that describe a particular topic, so the Google Trends' or Twitter's pattern would not be able to reflect a similar behavior that its counterpart in the Media.

## A quantitative description of the Agendas

*Agenda diversity*

How dominant is a main topic? Is the degree of dominance of a given topic in the Media Agenda reflected in the Public Agenda? In order to answer such kind of questions we quantify the diversity of the Agendas through the normalized Shannon's entropy $H$, which was introduced in section **Material and Methods**.

In figure 3 we can see the value of $H$ as a function of time for the three agendas. It is important to pay attention to those periods of time where the diversity is lower than usual. This effect is notoriously more pronounced in the Public Agenda giving by **GT** and, in particular, in four specific days where four local minima of Shannon's entropies can be detected. Three of them are outliers as defined in section **Material and Methods** , two of them from **GT** and one from **Tw**. The other one is not an outlier but a pronounced minima and therefore a point of interest in this description.

A lower value in the agenda's diversity is due to the fact that the most important topic attracts practically all the attention of the Public and the Media, collapsing the agenda to one of the issues involved. In the radar plots included in figure 3 we can see how two of these outliers (**a** and **d**) belong to the topic *Elections*. They are related to the primary and general legislative elections that took place in August 13th and October 22th respectively. In all the agendas these points were detected as outliers except point (d) in Twitter Agenda. Why is that? The radar plot of the Twitter agenda for this day displays an association between the topic *Elections* and the *Current President*, decreasing the important of this topic. Discussions in Twitter about elections appear also in point (c), when the other agendas seems to be more diverse. On the other hand, and despite not being classified as outlier, we also focus in point (b) because Shannon's Entropy of $GT$ Agenda displays a minimum (collapsing agenda) which is not corresponded neither in the Media nor in the Twitter Agendas. Crawling in the context, we see that it belongs to the topic *Missing person* and this date correspond to the rally that took place one month after the disappearance of *Santiago Maldonado* (see **Supplementary Material**). We would like to emphasize the discussion about this topic (*Missing person*) because its dynamics show interesting features, as we will show bellow.

From the measure of $H$ we have also observed that the median of the Public Agenda diversity is statistical significant lower than the Media Agenda's one. Specifically $H_{GT} = 0.73$ and $H_{Tw} = 0.74$ are statistical significant lower than $H_{MA} = 0.85$

with $p < 10^{-18}$, while there is no significant difference between the first two. However from figure 3 we can see that **GT** shows more abrupt dropouts in the diversity in response to specific events. From all this analysis we can conclude that given a finite set of topics, **the Public Agenda is less diverse than the Media Agenda**, because the public seems to focus in the most important topics than Media can do, maybe due to editorial decisions.

*Distance between Media and Public Agenda's*

Given our descriptions of the Agendas as time-evolving distributions, we can compare them by computing the Jensen-Shannon distance. In this context, outliers in selected dates will correspond to divergences between Public and Media Agenda, i.e., specific events where the public interest do not match with media offer. In figure 4 we show the Jensen-Shannon distance between Media and Publics Agendas as a function of time. We focus in three points that seems to be relevant enough. In all cases, the topic distributions at that particular dates displayed by the radar plots show that the increment in the distance between agendas is due to a greater interest of public opinion in the topic *Missing person*.

Points **(c)** and **(d)** show that both the Public and the Media highlight this topic, but the Media do not disregard other topics, so the corresponding distance between them can be interpreted as lack of the diversity in Public Agenda as discussed in the last section.

On the other hand, points **(a)** (we take this point due to be a local maximum despite not being an outlier) and **(b)** show a major interest of the public in the topic *Missing person* which it is not reflected in the Media. In figure 2 we can see that this topic becomes the most important in public's interest (both **GT** and **Tw**) days before that it happen in the Media Agenda. It can be associates with a social networks (like Facebook and Twitter) campaign in favor of the appearance of *Santiago Maldonado* ("The missing person") that took place on August 26th. This campaign was massive and initially underestimated by the main Media outlets in Argentina (see **Supplementary Material**).

Finally, it is important to say that the Jensen-Shannon distance, in conjunction with the measurement of agenda diversity given by the Shannon entropy, give an insight of independent behavior, in certain particular dates, of the Public and the Media. Its identification can be a starting point to study the Media reaction to a change in audience's interests.

Agenda bias in different Media outlets

In this section we leave aside the Public Agenda as an unified corpus and we study the composition and evolution of the Media Agenda of each media outlet. In figure 5 we show the bump charts corresponding to each of the analyzed newspaper analogously to figure 2. The topics are the same introduced in the wordclouds of figure 1, but when computing the topics' weights, the articles are discriminated by newspaper. We also show the radar plots with the average distribution, as made in figure 1.

In figure 5 we can qualitative have a glance of the differences between the newspapers' agendas. For instance, we can see how the newspaper called *Página 12* gave

more importance to the topics *Missing person* and *Social leader*, while it reduces to minimum the coverage of the topic *Former Planning minister* as the others did.

In other to detect significative bias coverage of a given newspaper, we again calculate the Jensen-Shannon distance, but between the individual newspapers agenda and the Media Agenda. Note that this is the distance between the distributions of figure 5 and the top panel of figure 2. In figure 6 we show the Jensen-Shannon distance as a function of time. We detect three points as outliers, although we finally disregarded point **(b)** due to a lack of information of newspaper *Infobae* in that period. The other two points corresponds to differences between *Página 12* and the other newspapers and correspond to differences in the coverage of the topic *Missing person*.

The point **(a)** corresponds to the first news of the Santiago Maldonado's disappearance reported by *Página 12* before the primary elections and the point **(c)** corresponds to the two months' rally after the disappearance (see **Supplementary Material**). Another singularity of point **(c)** corresponds to a greater coverage of *Página 12* in the topic *Social leader* where the others media outlets seem to be more interested in the topic *Former Vice-President*.

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* [43]. 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 the larger coverage of *Página 12* in comparison to other newspapers at the beginning of the period. We can, for instance, quantify this difference calculating the median of the signals. If we focus in the period between July 31th and August 27th, the median of the topic's 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 highlighting the more frequently mentioned in each newspaper and filtering the common words to all newspapers. Although most of the words are not relevant enough, some of them are quite interesting, as for instance the word *represión* (repression) when *Página 12* talks about the topic *Missing Person* and the word *Cristina* (Fernández de Kirchner, former president) which is employed by all newspapers except *Página 12* when they talk about the topic *Former Planning minister* (see **Supplementary Material**). We think that a more deeply study of topics' keywords could be a first approximation in the study of framing, which will constitute the core of futures works.

### A brief discussion about Agenda-Setting

In a world where social media exists and the feedback between a Media and its audience is common currency, nowadays the idea that the Media sets the agenda and

the audience blindly follows it (as it's seemed to be suggested in the original work of McCombs) is too naive. Based in the data analyzed above, the behavior of the Media and Public Agendas, either by looking at Google Trends or Twitter, shows periods of strong similarity (specially in the presence of an unexpected event) and periods of disengagement. Therefore, it is not trivial to establish a causal relationship between agendas, specially when they are represented as evolving in time topic distribution as we did in this work. However, it is possible to discuss agenda setting if we focus in a single topic. We think that the *Missing person* topic is the most adequate topic to be discussed because:

- It caused a great impact in either the Media and the audience
- its coverage fully deploys along the time lapse analyzed in this manuscript (see **Supplementary Material**).

In figure 8 we show the topic's relative weight from the Public and Media Agendas. After the initial coverage, the agendas seem to differentiate around August 15th, when topic started to became more important in the Public Agenda than in the Media one. Around August 24th, the topic abruptly increase in public's interest while the reaction in the Media is slower, showing a significant peak in the plot of the discrete difference (panel (b)). This date is very close to August 26th, when a campaign in social media took place. After that, the Media increase its coverage about the topic.

Is this a case of reverse agenda-setting, i.e, when the audience set the Media Agenda? After all, the audience get involved about this topic by the Media, so how was the coverage before those events? We can answer these questions by calculating the cumulative coverage of this topic since the first events took place. This measure can be seen as the numerical integrate of the temporal profiles of figure 7 between the initial date and the current date. It is interesting to note that the newspaper responsible to accumulate coverage during the initial stage was the minority one: *Página 12*. Our interpretation about the setting agenda dynamics of this particular topic is the following: A small newspaper (*Página 12*) gives great coverage to this topic; it is amplified by public in the social networks and also expressed by reiterative Google searches. Then the rest of the Media pay attention to this subject and it becomes an important topic in the Media Agenda. This interpretation try to catch in a qualitative way how the information flow was in this specific topic. On the other hand, behind this interpretation there are two important facts that must be mentioned: First, the disappearance of a person is a very sensible theme for the Argentinian society, and second, there are political reasons in why *Página 12* was particularly interested in cover this topic while the other Media did not follow this interest (see **Supplementary Material**).

Even tough we face the question about causality only in a qualitative way and just for a specific topic, it was possible to highlight the complex feedback dynamics that take place between public and media agendas.

## Conclusions

The Mass media play a fundamental role in opinion formation and therefore it's of vital importance to have an accurate quantitative description of the Media and Public Agenda and their relationship in the framework of Agenda Setting theory. In

this work, through the implementation of a topic detection algorithm we describe the Agenda of the Media as a distribution which evolves in 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, meanwhile the Media keeps certain degree of diversity. 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 here also allow us to detect coverage bias in newspapers and gave us a first approximation in the theory of framing.

We hope that some of the elements studied here will give us insights at the time of proposing a mathematical model about Mass Media and Public interaction. Future works may include a more systematic study and its extension to international Media, a deeper study of framing through topic detection and sentiment analysis, and a more quantitative analysis about causality.

## Supplementary Material

Context

We provide here a more detailed explanation of the topics discussed among this work. The ideology of the media in Argentina expresses the highly polarized political climate observe in Argentinian society. During the administration of *Cristina Fernández de Kirchner* (2007-2015), the government maintained a conflict with several news organizations. It led media such as *Clarín, La Nación* and the portal news *Infobae* to be very critics of the *Fernández*'s administration, emphasizing the allegations of corruption related to it as can be seen in the importance given to the topics *Former Planning minister* and *Former Vice-President.* On the other hand, *Página 12* has a opposite ideological inclination, supporting the former administration of *Cristina Kirchner* and therefore being very critical with the current *Mauricio Macri*'s administration, doing special emphasis on issues related to human rights, as can be again observed in the coverage given to the topics *Social leader* and *Missing person.*

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 *Cambiemos*, 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 analyzed period, and specially after the general elections of October $22^{th}$, 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 (Patagonian native population), 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 $26^{th}$, under the motto "Where is Santiago Maldonado?", followed by two massive protest marches to the *Plaza de Mayo* that 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 $17^{th}$ 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 *Néstor Kirchner* and *Cristina Fernández 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 immediatly jailed on October $27^{th}$.

*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. 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 Túpac 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 between NMF and LDA
In this section we apply other topic model, Latent Dirichlet Allocation [44] (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 study displays some degree of robustness when considering different topic model. On the other hand, as was discussed in [45], 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 8, 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.

$$TC = \sum_{i<j}^{N} \frac{2d_{ij}}{d_i + d_j} \tag{8}$$

We perform a decomposition into 10 topics using LDA with the python module *gensim* [46], 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 9 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 3, 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.
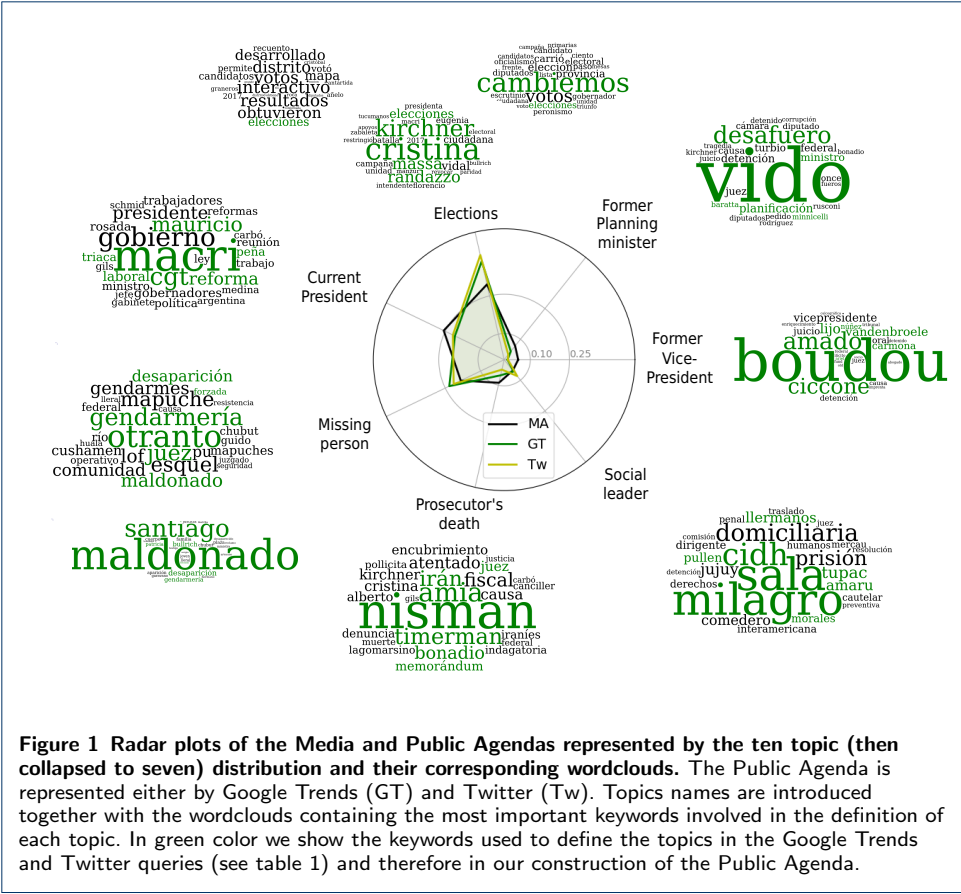
**Author details**
[1]Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428, Buenos Aires, AR. [2]Instituto de Física de Buenos Aires (IFIBA), CONICET, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428, Buenos Aires, AR. [3]Instituto de Investigación en Ciencias de la Computación (ICC), CONICET, Av. Cantilo s/n, Pabellón 1, Ciudad Universitaria, 1428, Buenos Aires, AR.

**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.: A look at agenda-setting: Past, present and future. Journalism studies **6**(4), 543–557 (2005)
3. Brians, C.L., Wattenberg, M.P.: Campaign issue knowledge and salience: Comparing reception from tv commercials, tv news and newspapers. American Journal of Political Science, 172–193 (1996)
4. Gerber, A.S., Karlan, D., Bergan, D.: Does the media matter? a field experiment measuring the effect of newspapers on voting behavior and political opinions. American Economic Journal: Applied Economics **1**(2), 35–52 (2009)
5. Coleman, R., McCombs, M.: The young and agenda-less? exploring age-related differences in agenda setting on the youngest generation, baby boomers, and the civic generation. Journalism & Mass Communication Quarterly **84**(3), 495–508 (2007)
6. 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)
7. Zunino, E., Aruguete, N.: La cobertura mediática del conflicto campo-gobierno. un estudio de caso. Global Media Journal **7**(14) (2010)
8. 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)
9. 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)
10. Lazaridou, K., Krestel, R.: Identifying political bias in news articles. Bulletin of the IEEE TCDL **12** (2016)
11. Baumgartner, F.R., Chaqués Bonafont, L.: All news is bad news: Newspaper coverage of political parties in spain. Political Communication **32**(2), 268–291 (2015)
12. Elejalde, E., Ferres, L., Herder, E.: On the nature of real and perceived bias in the mainstream media. PloS one **13**(3), 0193765 (2018)
13. 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)

14. 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)
15. 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)
16. 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)
17. Harder, R.A., Sevenans, J., Van Aelst, P.: Intermedia agenda setting in the social media age: How traditional players dominate the news agenda in election times. The International Journal of Press/Politics, 1940161217704969 (2017)
18. Guo, L., Vargo, C.J.: Global intermedia agenda setting: A big data analysis of international news flow. Journal of Communication **67**(4), 499–520 (2017)
19. 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)
20. Ali, A.E., Stratmann, T.C., Park, S., Schöning, J., Heuten, W., Boll, S.C.: Measuring, understanding, and classifying news media sympathy on twitter after crisis events. arXiv preprint arXiv:1801.05802 (2018)
21. 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)
22. Guo, L., Vargo, C.J., Pan, Z., Ding, W., Ishwar, P.: Big social data analytics in journalism and mass communication: Comparing dictionary-based text analysis and unsupervised topic modeling. Journalism & Mass Communication Quarterly **93**(2), 332–359 (2016)
23. Dai, X.-Y., Chen, Q.-C., Wang, X.-L., Xu, J.: Online topic detection and tracking of financial news based on hierarchical clustering. In: Machine Learning and Cybernetics (ICMLC), 2010 International Conference On, vol. 6, pp. 3341–3346 (2010). IEEE
24. Po, L., Rollo, F., Lado, R.T.: Topic detection in multichannel italian newspapers. In: Semanitic Keyword-based Search on Structured Data Sources, pp. 62–75 (2016). Springer
25. Brun, A., Smaïli, K., Haton, J.-P.: Experiment analysis in newspaper topic detection. In: String Processing and Information Retrieval, 2000. SPIRE 2000. Proceedings. Seventh International Symposium On, pp. 55–64 (2000). IEEE
26. Hu, X.: News hotspots detection and tracking based on lda topic model. In: Progress in Informatics and Computing (PIC), 2016 International Conference On, pp. 248–252 (2016). IEEE
27. Li, W., Joo, J., Qi, H., Zhu, S.-C.: Joint image-text news topic detection and tracking by multimodal topic and-or graph. IEEE Transactions on Multimedia **19**(2), 367–381 (2017)
28. Cataldi, M., Di Caro, L., Schifanella, C.: Emerging topic detection on twitter based on temporal and social terms evaluation. In: Proceedings of the Tenth International Workshop on Multimedia Data Mining, p. 4 (2010). ACM
29. Crokidakis, N.: Effects of mass media on opinion spreading in the sznajd sociophysics model. Physica A: Statistical Mechanics and its Applications **391**(4), 1729–1734 (2012)
30. González-Avella, J.C., Cosenza, M.G., San Miguel, M.: A model for cross-cultural reciprocal interactions through mass media. PloS one **7**(12), 51035 (2012)
31. Moussaïd, M.: Opinion formation and the collective dynamics of risk perception. PLoS One **8**(12), 84592 (2013)
32. Rodríguez, A.H., Moreno, Y.: Effects of mass media action on the axelrod model with social influence. Physical Review E **82**(1), 016111 (2010)
33. 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)
34. Instituto Verificador de Circulaciones. http://www.ivc.org.ar
35. Alexa. https://www.alexa.com/topsites/countries/AR
36. Xu, W., Liu, X., Gong, Y.: Document clustering based on non-negative matrix factorization. In: Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Informaion Retrieval, pp. 267–273 (2003). ACM
37. Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. Nature **401**(6755), 788–791 (1999)
38. 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)
39. Boydstun, 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)
40. Fuglede, B., Topsoe, F.: Jensen-shannon divergence and hilbert space embedding. In: Information Theory, 2004. ISIT 2004. Proceedings. International Symposium On, p. 31 (2004). IEEE
41. Tukey, J.W.: Exploratory Data Analysis vol. 2. Reading, Mass., ??? (1977)
42. Natrella, M.: Nist/sematech e-handbook of statistical methods (2010)
43. Dallmann, A., Lemmerich, F., Zoller, D., Hotho, A.: Media bias in german online newspapers. In: Proceedings of the 26th ACM Conference on Hypertext & Social Media, pp. 133–137 (2015). ACM
44. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. Journal of machine Learning research **3**(Jan), 993–1022 (2003)
45. 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)
46. 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)

**Figures**

**Figure 1 Radar plots of the Media and Public Agendas represented by the ten topic (then collapsed to seven) distribution and their corresponding wordclouds.** The Public Agenda is represented either by Google Trends (GT) and Twitter (Tw). Topics names are introduced together with the wordclouds containing the most important keywords involved in the definition of each topic. In green color we show the keywords used to define the topics in the Google Trends and Twitter queries (see table 1) and therefore in our construction of the Public Agenda.

**Tables**

**Table 1** Queries used in Google Trends in order to build the Public Agenda.

| Topic's name | Google Trends query |
|---|---|
| Elections | elecciones + cambiemos + cristina kirchner + massa + randazzo |
| Missing person | santiago maldonado + juez otranto + patricia bullrich + gendarmería + desaparición forzada |
| Former Planning minister | de vido + desafuero + ministro de planificación + minnicelli + baratta |
| Current President | mauricio macri + cgt + reforma laboral + peña + triaca |
| Social leader | milagro sala + cidh + tupac amaru + pullen llermanos + morales |
| Prosecutor's death | nisman + amia + memorándum con irán + timerman + juez bonadio |
| Former Vice-President | amado boudou + ciccone + ariel lijo + vandenbroele + núñez carmona |

**Table 2** Correlation between the topics' temporal profiles of the Public Agenda and their counterpart in Media Agenda. All correlation values are statistical significant ($p < 10^{-9}$), except (*) which is significant with $p < 0.05$.

| Topic's name | Correlation MA and GT | MA and Twitter | GT and Twitter |
|---|---|---|---|
| Elections | 0.81 | 0.59 | 0.75 |
| Missing person | 0.68 | 0.76 | 0.89 |
| Former Planning minister | 0.92 | 0.82 | 0.87 |
| Current President | 0.77 | 0.75 | 0.63 |
| Social leader | 0.49 | 0.25(*) | 0.57 |
| Prosecutor's death | 0.56 | 0.59 | 0.75 |
| Former Vice-President | 0.90 | 0.92 | 0.97 |

**Figure 2 Bump graph of the time-dependent Media (MA) and Public Agendas extracted from Google Trends (GT) and Twitter (Tw)**. Widths and rankings of the curves encode topic's relative weight. Also, some important events related to the topics are pointed out: **A**: First news about Santiago Maldonado's disappearance (Missing person); **B**: Primary elections; **C**: March one month after Santiago Maldonado's disappearance; **D**: March two months after Santiago Maldonado's disappearance; **E**: Appearance of Santiago Maldonado's body; **F**: General elections; **G**: Julio De Vido's detention; **H**: Debates on labor reform; **I**: Amado Boudou's detention (Former vice-president). A more detailed explanation is given in section .

**Table 3** Correlation between the temporal profiles of the topics found in NMF and associated topics in LDA.

| Topic's name | 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 |

**Figure 3 Shannon entropy (H) as a measure of agenda diversity.** The Public Agenda show a less diverse behavior than the Media Agenda as can be seen in the top figure. The horizontal lines correspond to the lower inner fences of each signal in order to identify outliers. The related radar plots show the agenda at the selected days where the time series exhibit dropouts (points a-d), indicating that the most important topic catches the most public's attention. **E**: Elections; **FPm**: Former Planning minister; **FVP**: Former Vice-President; **Sl**: Social leader; **Pd**: Prosecutor's death; **Mp**: Missing person; **CP**: Current President.
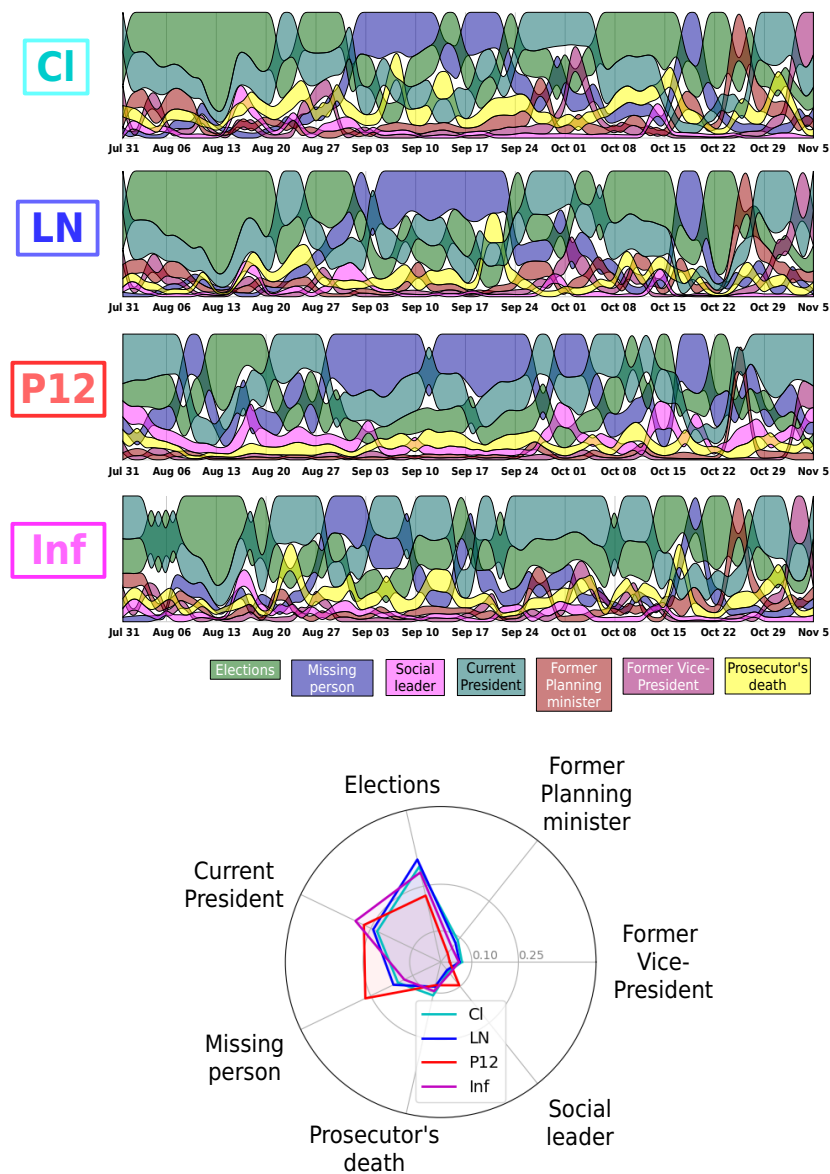
**Figure 4 Jensen Shannon distance between the Media and Public Agendas as a function of time** (with upper inner fences pointed out). Larger distances are due to a greater interest of the audience in the topic *Missing person* which decrease the interest in other topics. On the other side the Media Agenda still keeps certain degree of diversity. **E**: Elections; **FPm**: Former Planning minister; **FVP**: Former Vice-President; **Sl**: Social leader; **Pd**: Prosecutor's death; **Mp**: Missing person; **CP**: Current President.

**Figure 5 Bump charts of newspapers' Agenda and radar plot of the average distributions.** The figure shows, in a qualitative way, the bias in the different newspaper's agendas. For instance, the greater interest of *Página 12 (P12)* in the *Missing person* topic and its slightly lower coverage in the *Former Planning minister* respect to the other newspapers.
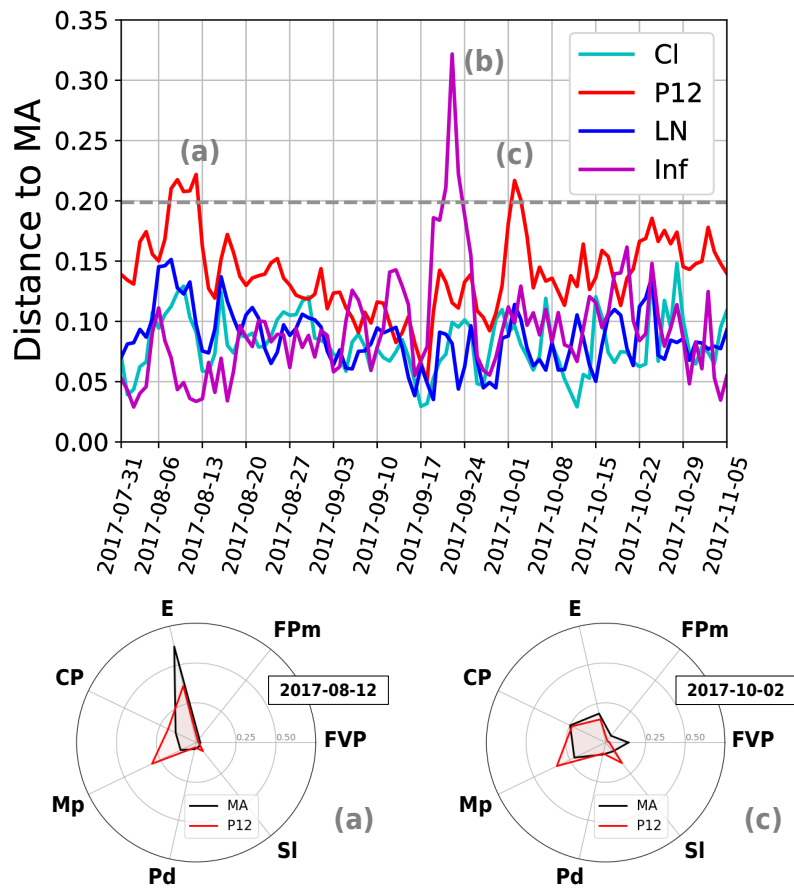
**Figure 6 Jensen-Shannon distance between the newspapers agenda and the Media Agenda as a function of time**. *Página 12* shows the more different behavior, motivated again by its interest in the *Missing person* and *Social leader* topics as can be seen in the radar plots which belongs to points (a) and (c). The anomalous behavior of *Infobae* at pint (b) is due to few articles around that date in our database, therefore we ignore its radar plot. **E**: Elections; **FPm**: Former Planning minister; **FVP**: Former Vice-President; **Sl**: Social leader; **Pd**: Prosecutor's death; **Mp**: Missing person; **CP**: Current President.
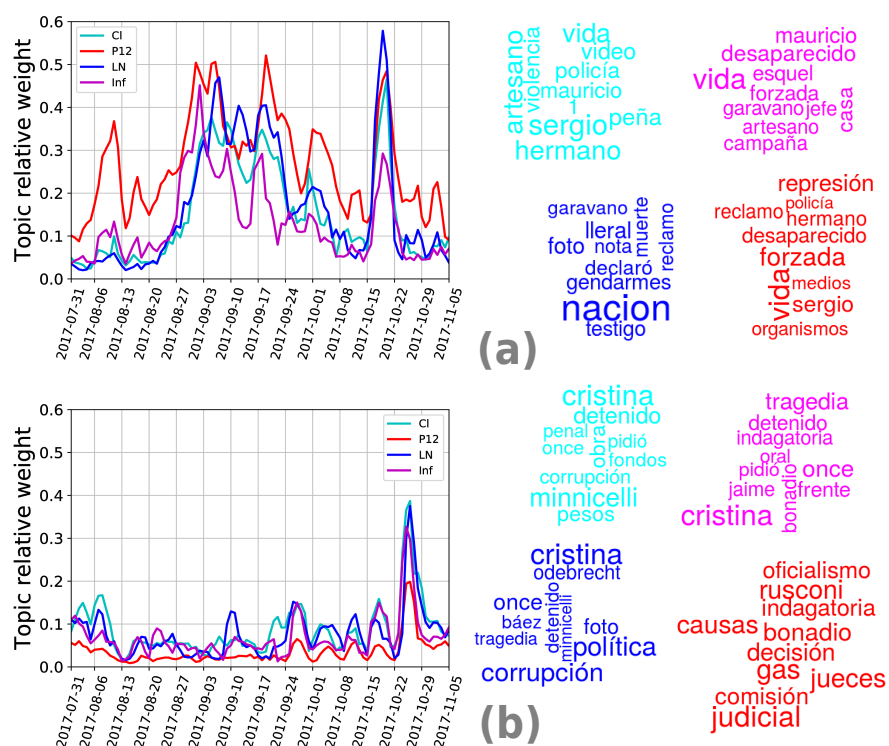
**Figure 7 Relative weight of the topics (a) Missing person, and (b) Former Planning Minister, and their corresponding word-clouds of frequent newspapers' keywords**. We interpreted the differences shown in given periods as an indicator of coverage bias. For instance, in figure (a) *Página 12* pays a greater attention in the first days. In the word-clouds, we show which of the defining words are more frequently used by the corresponding newspaper. Most of them are less informative, but other seems to represent a first approximation in the study of framing.
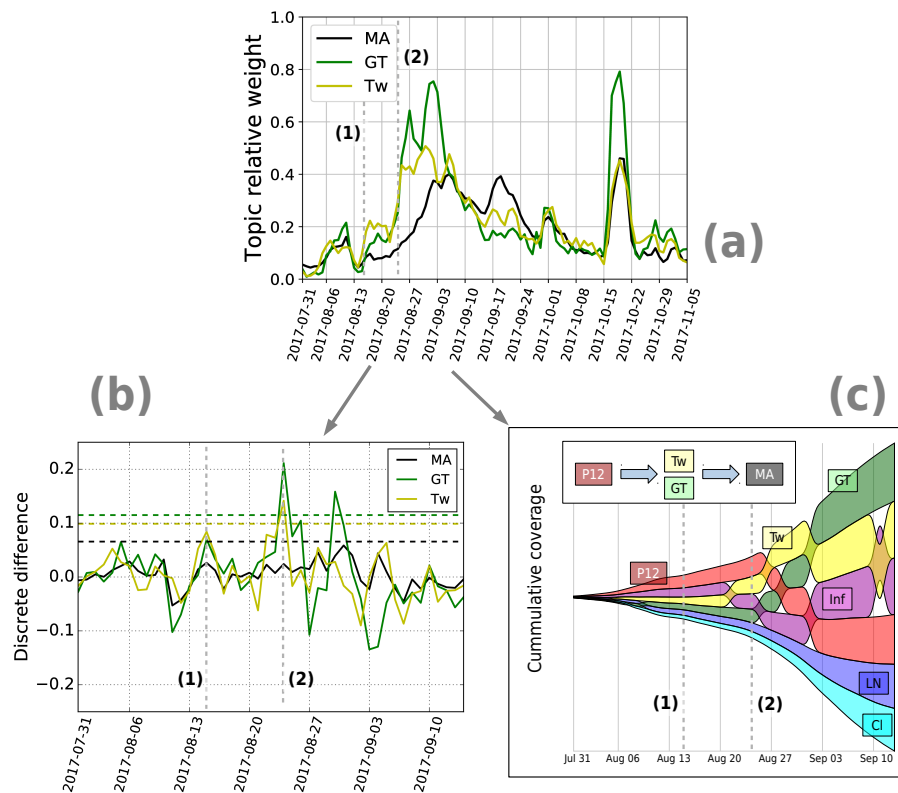
**Figure 8 Agenda setting direction in Missing person's topic?** The temporal profiles of figure (a) show that the Public and Media agenda seem to differentiate around August 15th (vertical grey line (1)) and the Public increase abruptly its interest in the topic around August 24th (grey line (2)). It can be seen also in figure (b), where the discrete differences were computed. With the computing of the cumulative sum of figure 7 and figure (a), represented as a bump chart in figure (c), we suggest that the topic was first set by *Página 12* and then the Public's interest cause the coverage of the other Media.
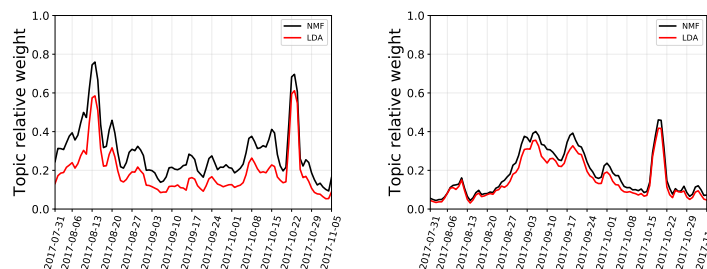


**Figure 9** Temporal profiles of topics *Elections* (left) and *Missing Person* (right) for both LDA and NMF. All the topics found by applying NMF have a highly correlated counterpart in LDA.