

**How to cite this article in bibliographies / References**

R Sánchez Sabaté, C del Valle, M Mensa (2019): “Method for the construction of large thematic corpora of online news articles. Towards a corpus of food-related news”. *Revista Latina de Comunicación Social*, 74, pp. 594 to 617.

<http://www.revistalatinacs.org/074paper/1347/30en.html>

DOI: [10.4185/RLCS-2019-1347en](https://doi.org/10.4185/RLCS-2019-1347en)

# Method for the construction of large thematic corpora of online news articles. Towards a corpus of food-related news

**Rubén Sánchez Sabaté** [CV] [ORCID] [GS] Universidad de La Frontera (UFRO), Scientific and Technological Research Group in Social Sciences and Humanities, Centre of Excellence in Economic and Consumer Psychology (CEPEC), Temuco, Chile. [ruben.sanchez@ufrontera.cl](mailto:ruben.sanchez@ufrontera.cl)

**Carlos del Valle** [CV] [ORCID] [GS] Universidad de La Frontera (UFRO), Scientific and Technological Research Group in Social Sciences and Humanities, PhD in Communication, Temuco, Chile. [carlos.delvalle@ufrontera.cl](mailto:carlos.delvalle@ufrontera.cl)

**Marta Mensa** [CV] [ORCID] [GS] Institute of Social Communication, Universidad Austral de Chile, Inequality and Human Rights Interdisciplinary Research Group, Valdivia, Chile. [marta.mensa@uach.cl](mailto:marta.mensa@uach.cl)

## Abstract

**Introduction.** Food-related news disseminated by the mass media have increased over the last two decades. These contents are powerful structuring elements of Western contemporary society. This article presents a method for the selection of a corpus of thousands of news published in the digital press –without having to resort to paid news databases– that can be subsequently processed with quantitative analysis software. **Methods.** Food-related content is operationalised through a series of keywords that allow for the retrieval of news on the subject in question using the Google search engine and the Import.io service. **Results.** More than 2500 news items published in 2016 by three Chilean newspapers were retrieved and stored for their subsequent analysis with big data techniques. **Discussion and conclusions.** The differences and similarities of the results obtained for each newspaper are discussed from different theoretical and methodological points of view. The different research possibilities offered by the data obtained are also mentioned.

**Keywords:**

Methodology; Big Data; digital press; food and journalism.

**Contents**

1. Introduction. 1.2. Communication and food. 1.3. Theoretical framework. 1.4. Objectives. 2. Methods. 2.1. Operationalisation of food. 2.2. Procedure to retrieve news. 3. Results. 4. Discussion. 5. Conclusion. 6. Notes. 7. References.

Translation by **CA Martínez-Arcos**  
(PhD in Communication, University of London)

**1. Introduction**

Due to the negative value that Western philosophy in general, from Plato to Karl Marx, has bestowed on food and physiological pleasures (Heldke, 2006; Korthals & Kooymans, 2004; Telf, 2002), it is easy not to realise the true social, cultural and political relevance of food and the –private and public– discourses that are built around them. Fortunately, thanks to anthropology and sociology, since the beginning of the 20<sup>th</sup> century, academic work has been developed on the psychosocial and political dimensions of food. It has been found that food is such a complex and primordial element for human life that since the dawn of the 21<sup>st</sup> century an interdisciplinary field of study has been configured in English-speaking countries under the name of “Food Studies” (Poulain, 2017).

The presence of food-related content in the mass media has increased over the last two decades (Contreras & Gracia Arnaiz, 2005, p. 16; Koldobsky, Daniela, 2011, p. 15; Narváez, 2013). These discourses are powerful structuring elements of modern societies not only because they deal with a human/social phenomenon, but also because they are disseminated by the mass media, the primary agents in the configuration of the “public sphere” (Habermas, Domènech & Grasa, 1981). Despite its socio-political importance, the study of food-related content in the mass media is still scarce. The editors of *Food as communication/Communication as food*, a book published in 2011 that collects the contributions of more than 20 scholars on food and communication, explicitly recognise that although food has been widely analysed from an anthropological, sociological and historical-cultural perspective, the communication sciences have barely looked into this subject (Cramer, Greene, & Walters, 2011).

With the purpose of contributing to the field of food studies from the communication sciences, this article offers theoretical and methodological tools to perform content analysis on the online press in a way that allows us, at least, to explore, describe and analyse the news coverage of food, understood in a broad and integral way as we will see later. We agree with other scholars who see content analysis as a prerequisite and as an essential step to study the relationship between this coverage and 1) its production by the media; 2) the effects that such coverage may have on audiences (2005); and finally, the relationship between the food-related content disseminated by the mass media and the food cultures, ideologies and practices of a given society.

Thus, the reader will find a brief literature review on food and communication, a theoretical framework on food from the social sciences, and a detailed description of the procedure designed to form large corpora of news on food –or any other subject deemed appropriate– that can be analysed with software. The article also presents the results of the application of this method in three Chilean online newspapers, and a discussion on the results and the research possibilities that emerge once the desired corpus has been formed.

## 1.2. Communication and food

Academic research in non–advertising food-related content disseminated by the media is still incipient. Generally, it takes food consumption as a subsidiary of two subjects that do receive regular attention by the academic community: health, for example food, obesity and eating disorders (Evans, Rich, Davies, & Allwood, 2008; Lawrence, 2004; Menezes Ferreira, de Castro Oliveira, & Terrón Blanco, 2015; Square, 2012; Sandberg, 2007; Westall, 2011), and climate change, such as food and global warming (Roig, 2013). The few research works on the media that take food as the only object of study have to do mainly with gastronomy (Blanco Hernández, 2015; Martínez & Poyatos, 2015; Sánchez Gómez, 2010), although there is a study on food information from the point of view of nutrition (Bernabeu-Peiró, 2015), and another one that quantitatively analyses food news in the Madrid press (Fúster, Ribes, Bardón, & Marino, 2009).

French semiologist Roland Barthes was one of the first to study food discourses in the media and especially in advertising (Contreras & Gracia Arnaiz, 2005, p. 90). With regards to the latter, Barthes (1997) detected three motives present in food advertising: 1) the remembrance of a past and a rural space that takes the consumer back in the history of the country; 2) health for nutrition (functional nutrition); and 3) the association of food and food practices with specific situations in modern life, for example, a cup of coffee to relax or a snack to refer to people with lot of work and little time to eat.

A broader and more current look at the food discourses that appear in the mass media allows us to distinguish at least four types of discourses (Koldobsky, Daniela, 2011): 1) didactic, used in cooking programmes; 2) advertising; 3) scientific–medical, especially focused on the nutritional impact on people’s health; and 4) gastronomic critique. We add a fifth discourse that obviously could go unnoticed: the journalistic discourse in the strictest sense. Journalistic food discourses are those that report and analyse food production, distribution, consumption and disposal, especially considering the political-economic and environmental viewpoints.

A historical review of food journalism (Albala, 2013) shows that, at least in the English-speaking world, until 1957, the press published many news on food production and distribution but very few news, or none at all, on the pleasures of eating. The gastronomic critique and culinary education emerged strongly in the last third of the twentieth century. Finally, with the advent of the Internet there were two major changes: 1) food became a topic of high interest. Blogs and websites specialised in food emerged to respond to the great demand of the audience; and 2) the paper crisis in the written press caused by the Internet led to the development of online journalism. Therefore, food journalism also “moved” to the Internet.

The internet basically offers three types of journalistic publications –at least in English– that cover food issues (Albala, 2013): 1) online editions of the general-interest press. These websites in general reproduce the same content as their paper versions, although they increasingly offer extra content exclusively for online readers. The sections dedicated to food can be expanded in the online version. 2) web-only publications specialised in food; and 3) personal food blogs, which, on occasions, have earned high ratings that have allowed their publishers to do this job exclusively and professionally, thanks to online advertising.

### 1.3. Theoretical framework

The scarce literature on food discourses in the mass media is not justified if we take into consideration the fact that currently these discourses are not only broadening the knowledge of ordinary citizens about the complex food phenomenon but are also reconfiguring our relationship with the planet and with our bodies (Frye & Bruner, 2013). From the field of semiotics in particular, authors like Roland Barthes (1997), Jean Soler (1997), Gastón Gaínza (2003) and Fabio Parsecoli (2011) have shown that food practices and eating constitute a communication system that, like verbal language, reflects a certain sociocultural world. Food and eating, therefore, are signs; texts that are located in the particular semiotic sphere of each community and play a fundamental role in the biological and cultural reproduction of society. For socio-semiotics, food texts communicate ideologies that configure power relations.

The discourses that are built around food are as important as food's ability to symbolise. The reason is obvious: if food "says a lot" about the society that produces, cooks and eats it, when society talks about food, after all, it is talking about itself. Thus, food discourses are another mechanism of socialisation by which we understand and build our culture (Cramer et al., 2011). Therefore, we can also refer to the ways of speaking and representing food consumption as powerful ideological signifiers that build power relations (LeBesco & Naccarato, 2008). Thompson (2012) equates the emerging food discourses to the human rights discourses that emerged in the second half of the twentieth century and considers that food discourse is becoming a language of political and cultural struggle that transcends food matters. It is in the discourses about food where power is negotiated and where identities and agency are sought in today's globalised world.

Food content disseminated by the mass media is especially relevant because the latter are paradigmatic and cultural institutions of modern societies (J. B. Thompson, 1998) that determine, through the contents they produce and disseminate, the limits and horizons of the possible worlds (Alsina, 1989; Farré, 2004). These contents structure and promote socio-imaginary meanings that guide, order, classify and organise the social events in which individuals are involved (Dijk, 2009). Thus, the Anthropology of Communication says that the mass media are the modern "host structure" –education and culture agent– that today complements, and even to supplants, the three traditional host structures: the family, the city and the church (Duch & Chillón, 2012). Therefore, the power of the media as institutional producers of symbolic forms to shape and transform society is evident (Velázquez, 1992). So much that, for Thompson (1998) and Bourdieu (2001), the mass media are the agent, the institution, that owns and manages the symbolic power that circulates in mediated societies like ours.

To study the food content disseminated by the mass media, we propose to conceptualise food consumption in a holistic way by following the theoretical contributions made mainly by anthropology and sociology. Today, food is considered a totally social phenomenon (Contreras & Gracia Arnaiz, 2005), in the Maussian sense of the term, and a total human phenomenon, in the words of French philosopher and sociologist E. Morin (2014). “Totally social” because all areas of culture and all types of institutions (economic, legal, political, religious, etc.) (Elias, 1989; Goody & Willson, 1995) find in food a simultaneous expression; and “Totally human” because of its biological and ecological components, and its power to organise society, “being at the same, or even greater, level of importance than sexuality and kinship” (Contreras & Gracia Arnaiz, 2005). On the latter, the works of Claude Lévi-Strauss (1981, 2013; 1970), Mary Douglas (1972, 1980) and Pierre Bourdieu (1984) are particularly noteworthy. The latter, in his well-known work on social distinction, postulated that eating habits are a way to naturalise not only social differentiation but also ideological differentiation.

This broad conception of food allows us to address the food content disseminated by the mass media beyond nutrition and gastronomy. Understanding food as a “totally social practice” also allows us to combine two key concepts on human nutrition: “food system” and “food culture”. This way, we can work with a comprehensive definition of food that includes both its symbolic and material dimension. Thus, “food system” must be necessarily understood at least [1] as the set of relationships that are established between the production, distribution, preparation, consumption and disposal of food (Goody & Willson, 1995). And, “food culture” must be understood as “the set of representations, beliefs, knowledge and practices, inherited and/or learned, that are associated with food and are shared by individuals of a given culture or a determined social group within a culture” (Contreras & Gracia Arnaiz, 2005). Only in this way, considering the multiple dimensions of food, one can characterise quantitatively and qualitatively, and in an integral way, the media coverage of food-related contents.

#### **1.4. Objectives**

The objective of this article is to explain how to define a sufficiently broad and representative corpus of food-related news from the online press, in a way that allows us to carry out, at least, exploratory and descriptive studies on the coverage of food fact based on content analysis and supported by big data techniques. Bernard Berelson (1952) defines content analysis as a research technique for the objective, systematic and quantitative description of the manifest content of communication. Manifest content refers to the meaning of the symbol in question, on which both sender and receiver would spontaneously agree. A more current definition of content analysis is the one provided by Riffe et al. (2005): it is “the systematic allocation of categories to communicational content according to pre-established rules, to be able to analyse the relationships between categories by means of statistical methods”. This with the aim of making “reproducible and valid inferences that can be applied to their context” (Krippendorf, 1990).

#### **2. Methods**

The increase in food content that the Internet has allowed, together with the computing possibilities developed to analyse big data, invites us to consider corpora of hundreds or thousands of news (units

of analysis) to obtain a general characterisation of the food coverage carried out by the media in general, beyond specific events, vicissitudes or crises. This type of analysis allows us to observe *what* the media is talking much or little about (for example, the agenda setting), before focusing on a discursive analysis that is more interested in *how* the media is talking about the food system and food culture (for example, framing theory).

Generating a corpus of hundreds or thousands of news (or blog entries) published in digital format can be achieved with relative ease by using the Google search engine, although the process, as we will see, is somewhat arduous. It is true that, for some media, it is possible to obtain news pieces that meet previously established criteria through news aggregators such as LexisNexis and Factiva. However, these aggregators have two drawbacks: 1) they are not free to use and therefore not all universities or research centres provide this service to their respective researchers; and 2) not all news media, especially those from Latin America, are present in such aggregators.

For us, the best alternative to these aggregators is Google because it is a free service and is available to the entire academic community; and, most importantly, because, as it has been documented, it is reliable when it comes to searching for news –the scientific literature has established that the search for news in English in Google and LexisNexis have a degree of coincidence of over 80% (Weaver & Bimber, 2008).

## 2.1. Operationalisation of food

Whether the researcher uses Google or a database like LexisNexis to form the research corpus, it is necessary to operationalise the category “food” mentioned in the introduction through a series of keywords. To this end, we have departed from the theoretical discussion on the food system performed by food anthropologists Jesús Contreras and Mabel Gracia in their book “Food and culture” (2005), and two quantitative works conducted in Spain on food content in the media (Fúster et al., 2009; Marín-Murillo, Flora, Armentia-Vizueté, José-Ignacio, & Olabarri-Fernández, Elena, 2016). This selection of keywords to retrieve relevant news pieces to form the desired corpus is consistent with the best practices recommended for content analysis (Lacy, Watson, Riffe, & Lovejoy, 2015).

Taking the five fundamental processes of the food system described by Goody (1995), production, distribution, preparation, consumption and waste disposal, we generated a matrix with five columns, each of which corresponds to one process. Then, based on the theoretical discussion of Contreras et al. (2005), we added subcategories to some processes in order to recognise their complexity. Finally, we identified the keywords that best denote the aforementioned food processes with the help of: a) previous studies on food content in the media; b) scientific literature on the anthropology and sociology of food consumption; c) the anecdotal experience provided by the reading of food content in the three selected Chilean digital newspapers that will be identified later; and d) the triangulation between the researcher and co-researchers of this project.

After the identification of the keywords that operationalise “food”, we removed those words that were repeated in more than one process. This way, we simplified the search without risking the exclusion

of relevant contents. For example, instead of asking Google to show news containing the terms “food” and “healthy eating”, we simply searched for “food”, which would also show results of news containing the term “healthy eating”.

The list of keywords that are proposed to get relevant news to the desired corpus are as follows:

Food/s; food consumption; biofuel; agribusiness; transgenic; seed; agricultural exploitation; fishing; fish farming; aquaculture; fishing ground; fruticulture; wine; livestock; agriculture; supermarket; kitchen; cooking; culinary; diet; nutrition; nutrients; malnutrition; nutritional; meal; eating; hunger; beverage/s; drinking; gastronomy; gastronomic; tasting; restaurant; cookery; snacks; chef; patisserie; pastry.

The terms and expressions related to food garbage and waste are not included in the list of search terms because “garbage” is too generic and would lead Google to show irrelevant results and, second, because from our experience, whenever we talk about food waste, words like “meals” and “food”, which are already in our search list, show up.

## 2.2. Procedure to retrieve news

In order to illustrate the procedure used to retrieve news using Google, below we describe the steps followed to get the news published in 2016 by the most visited Chilean online newspaper [emol.com](http://emol.com) [2]. This newspaper would be the context unit (Colle, Raymond, 2011) from where we obtained the units of analysis that will end up configuring the news universe to be analysed.

While some online newspapers have their own search engine, Google is the best choice for two reasons: 1) it is not possible to know the reliability of an online newspaper’s search engine; and 2) if we want to compare the contents of different online newspapers, their search engines will likely have differences in terms of reliability and performance, which would compromise the reliability of the data to. Google, on the other hand, can act as referee between online newspapers when it comes to establishing what contents have been published by each written medium.

In order for Google to provide results that are not personalised –i.e., not conditioned by the digital print on the researcher– and are close as possible to the reality of the online newspaper in question, it is necessary to follow the following steps. First, it is necessary to create a new user account in the computer operating system that will be used for this purpose. Once created, it is necessary to restart the computer and sign in with this new account. Then, the researcher opens the preferred browser (we recommend Firefox), and opens a private browsing window (in Firefox, “new private window” and in Chrome “new incognito window”).

In the search bar, type “Google search” and the country where the online media under study are based. In our case, we write “Google search Chile”. Once the researcher is in the Google search engine of the corresponding country, it is necessary to configure it correctly for the search work. To do this, the researcher clicks in the “preferences” menu located at the right of the bottom of the page to make sure “SafeSearch” is not activated. Then in “Results per page”, choose “100”; then click on “Search

History” and make sure it is disabled; then select the region in which the newspaper is based. This is imperative to force Google to yield results stored on the servers of the country where the newspaper is located, because otherwise Google delivers results based on the location of the computer that generates the search. Following these steps will ensure we will obtain results that are very similar to those obtained by the citizens of the country of the online newspaper under study.

Once this configuration is done, it is necessary to access “advanced search”. Go to the “Preferences” menu and click on “Advanced Search”. This service enables Boolean search and restricts the search to a single web domain, among other options. In the field “all these words”, write, to begin with, the first keyword of our listing: “food”. Introduced the domain of the medium you want to study (in our case, [emol.com](http://emol.com)) in the “site or domain” field. Then is click on “Advanced search” and Google will provide a page with the first 100 results. At the top of the results page, immediately below the search bar, there is a menu with several options. Click on “Tools” to filter the results and show only the news published between 1 January 2016 and 31 December 2016.

If the researcher has followed the steps described above, the Google results page would display the following message in the search bar: “Food site: [www.emol.com](http://www.emol.com)”. The next step is the most important to ensure Google yields depersonalised results. According to marketing specialists (Nedelko, 2013), this step is so effective that some of the previous step would be unnecessary. It consists of the following: once the first page with the 100 first results is displayed –there will be keywords that yield less than 100 results– in the browser bar (do not mistake it with the Google search bar) write the following at the end of the URL without quotation marks: “&pws=0” [3] And press “Enter”. This expression forces Google to deliver depersonalised results. As way of example, below are the two URLs of the same search query, the first without the depersonalisation parameter, and the second with that parameter in bold:

URLs for food site: <a href="http://www.emol.com">www.emol.com</a> between 1-1-2016 and 31-12-2016
<a href="https://www.google.es/search?q=alimento+site%3Awww.emol.com&amp;num=100&amp;lr=&amp;hl=es-419&amp;as_qdr=all&amp;source=Int&amp;tbs=cdr%3A1%2Ccd_min%3A1%2F1%2F2016%2Ccd_max%3A12%2F31%2F2016&amp;tbm=">https://www.google.es/search?q=alimento+site%3Awww.emol.com&amp;num=100&amp;lr=&amp;hl=es-419&amp;as_qdr=all&amp;source=Int&amp;tbs=cdr%3A1%2Ccd_min%3A1%2F1%2F2016%2Ccd_max%3A12%2F31%2F2016&amp;tbm=</a>
<a href="https://www.google.es/search?q=alimento+site%3Awww.emol.com&amp;num=100&amp;lr=&amp;hl=es-419&amp;as_qdr=all&amp;source=Int&amp;tbs=cdr%3A1%2Ccd_min%3A1%2F1%2F2016%2Ccd_max%3A12%2F31%2F2016&amp;tbm=&lt;b&gt;&amp;pws=0&lt;/b&gt;">https://www.google.es/search?q=alimento+site%3Awww.emol.com&amp;num=100&amp;lr=&amp;hl=es-419&amp;as_qdr=all&amp;source=Int&amp;tbs=cdr%3A1%2Ccd_min%3A1%2F1%2F2016%2Ccd_max%3A12%2F31%2F2016&amp;tbm=<b>&amp;pws=0</b></a>

It is important to note that the search with the parameter that asks Google for depersonalised results, is only valid for the first page of results. It is necessary to re-paste the parameter “&pws=0” in the browser bar in the second page of results for the same keyword.

From now on, just by changing the search term, Google will show the news published in the selected period and domain [emol.com](http://emol.com) in groups of 100. Thus, it would not be necessary to repeat all the

configuration steps explained for each of the keywords listed above, but it would be necessary to add the parameter “&pws=0” in every case to ensure we get depersonalised results.

Ideally, a search should be done for each keyword. However, in order to lighten the search process, it is possible to use more than one keyword in each search. If this is done, the following considerations must be taken into account. First, make sure that we use the Boolean operator OR and not AND –the latter is used by default in Google when we write several words in the search bar. To do this, on the Advanced search page, instead of filling the field “All these words”, we have to write the keywords that we want to search separated by comma in the “any of these words” field. Another way of doing this is directly in the search bar of the results page, separating each word with a comma and the Boolean OR. For example, “food, OR feeding”.

The second consideration is that Google offers at most 500 Results. That is, if for example, the search for the keyword “food” among the news published by [emol.com](http://emol.com) in 2016 displayed more than 500 news containing that word, Google would only display 500. That is why, that if we choose to search for more than one word at a time, we need to make sure that Google displays less than 500 results, otherwise we might be losing news. If the search for a keyword yielded more than 500 results, it would be necessary to divide the period considered into several searches. For example, instead of asking Google to display results from 2016, we can request results for the first half of 2016. Finally, it should be noted that, whether we carry out a new search for each keyword, or searches with more than one keyword, when gathering the results of all searches we will get a significant number of duplicated news. As it will be shown later, obtaining duplicated results is a minor problem that is easily resolved in a more advanced phase of the process.

Once Google displays the results page with the 100 news stories published in [emol.com](http://emol.com) in 2016, there are two options to visualise each of the news listed by Google, although we can also click on each of them and copy the contents into a spreadsheet or applications that can automate this process. The second option is much more agile and reliable, and that is why we should opt for it.

There are different options with regards to computer apps to get the news. Our recommendation requires a low-cost online service that allows us to get the URLs and content of the news results displayed by Google.

This service we used is [import.io](http://import.io), which offers a free trial period with limited options [4]. Its main disadvantage is that it is only available in English. The service works as follows: the user inputs X number of URLs into the system and the latter extracts the desired elements from each of these URLs. In the case that concerns us, we delivered the URLs of each of the results pages provided by Google for each search. That is, if a given search on Google with X keywords yielded 200 to 300 news results, were collected the three URLs displayed in the browser bar –not the Google search bar– of the three results pages. Then, we asked [import.io](http://import.io) to pick up the headline, URL, publication date and description of each result and organise the collected data in rows and columns, where each row is a result or case, and each column is one aspect of the result (headline, URL, date, description, etc.).

Once this process has been completed for all Google searches, all the spreadsheets produced by [import.io](http://import.io) were gathered in a single file using appropriate software such as Microsoft Excel or LibreOffice. Then, it is necessary to debug the results using a double procedure. First, select the URL column and ask the software to remove duplicates. Second, order results by date to confirm all news are actually from 2016. Third, it requires human intervention because it consists of detecting false positives, i.e., URLs that contain relevant keywords but refer to news that have nothing to do with food, or simply do not deal with the issue of food. Below are three examples of [emol.com](http://emol.com) to illustrate false positives:

**Example 1:**

Headline of result according to Google: “San Fermín 2016: Third day leaves 13 wounded...”

URL of result:

[Http://www.emol.com/noticias/Internacional/2016/07/10/811755/San-Fermin-2016-Tercera-jornada-deja-a-13-heridos-dos-de-ellos-por-asta-de-toro.html](http://www.emol.com/noticias/Internacional/2016/07/10/811755/San-Fermin-2016-Tercera-jornada-deja-a-13-heridos-dos-de-ellos-por-asta-de-toro.html)

Publishing date according to Google: 10 July 2016.

Description of result according to Google: “10 July 2016- This Sunday, at 02 AM local time, we have the participation of Salamanca’s Pedraza de Yeltes stockbreeding, who makes its debut in this version of San...”

This is an example of false positive because although the news contains keywords, the theme is the bullfighting festivities held in Spain, and not human’s food consumption.

**Example 2:**

Headline of result according to Google: “Contest: Do you like “simple” cooking? We have a winner...”

URL of result:

[Http://www.emol.com/noticias/Tendencias/2016/05/09/801927/CONCURSO-Te-gusta-la-cocina-simple-Comparte-tu-receta-mas-facil-y-gana-un-robot-Moulinex.html](http://www.emol.com/noticias/Tendencias/2016/05/09/801927/CONCURSO-Te-gusta-la-cocina-simple-Comparte-tu-receta-mas-facil-y-gana-un-robot-Moulinex.html)

Publishing date according to Google: 9 May 2016

Description of the result according to Google: 9 May 2016-For those who seek to make their life easier, Moulinex’s Multicook is the ideal alternative to save time and cook easily...”

This is another example of false positive because it is a “cooking news story” that is actually covert advertising for a kitchen appliance.

**Example 3:**

Headline of result according to Google: “Star of “Scandal” TV series is honoured at Harvard | Emol...”

URL of result:

[Http://www.emol.com/noticias/Espectaculos/2016/01/29/785909/Protagonista-de-la-serie-Scandal-es-homenajeada-en-Harvard.html](http://www.emol.com/noticias/Espectaculos/2016/01/29/785909/Protagonista-de-la-serie-Scandal-es-homenajeada-en-Harvard.html)

Publishing date according to Google: 29 January 2016

Description of the result according to Google: “29 January 2016-Besides, Washington took the Hasty Pudding prize, a golden pot-shaped trophy to cook pudding. The theatre group, the oldest in the country...”

This is also an example of false positive because while the news contains the keyword “cooking”, food is not the main or secondary topic in this news: it is simply anecdotal.

To carry out the identification of false positives, we trained three coders who could evaluate each news according to the concepts of “food”, “food system” and “food culture” as presented in the introduction of this article. Afterwards, they were asked to judge the relevance or lack thereof to the corpus based on the following exclusion criteria: 1) “false news”: covert advertising, contests, etc., in which there is no intention to inform the public; 2) news in which food is not a primary or secondary topic, i.e., food is anecdotal; 3) news in which alcoholic beverages are the main topic, with the exception of wine and beer; 4) news in non-written format such as audio, photography and video. And finally, with the intention of covering as many news as possible, coders were told to consider as relevant all those news items whose relevance was unclear for them.

Two encoders were assigned to independently assess each newspaper. Evaluation should be exclusively based on the spreadsheet produced by [import.io](http://import.io), which contained the elements of each case. If it was not possible to make a decision with this information, coders had to copy the URL of the news in the browser to read the whole story and make a final decision. If considered relevant, the row for that news was left blank. If considered irrelevant, the row for that news was marked in red. In cases of doubt, the row was marked in yellow. Once this process was completed, a third encoder and one of the authors of this article reviewed the news highlighted in red and yellow and made a final decision on their relevance.

Once all the URLs were debugged for all newspapers, [import.io](http://import.io) was used to retrieve the elements of the selected news. In our case, were collected the following items: newspaper, date, author, place from which news is written, headline, subtitle, news text, and section. We input the URLs of every online newspaper to [import.io](http://import.io) and configured it to retrieve the required data from each news item to get an Excel file (data can also be download in CSV and NDJSON formats) in which each row is a unit of analysis (a news item) and each column is a different variable from each unit of analysis. It should be noted that in the extraction process, some URLs can be “lost” for different reasons, for example, because some URLs have expired, or because some URLs leads to a PDF or DOC document, which obviously could no longer be considered a unit of analysis. To avoid the first case, we recommend that to focus the search on news published as recently as possible.

At this point, we have formed the corpus or universe of study, and can proceed to the sampling stage or directly to the analysis if we have specialised software for the analysis of large amounts of textual data. More information about big data textual analysis software can be found on the second edition of *The Content Analysis*, written by K. A. Neuendorf (2016, 304 and subsequent pages). It is important to note that the procedure developed here only requires basic knowledge of office automation. It is therefore within the reach of every researcher familiar with the use of today’s daily technology.

### 3. Results

In order to obtain a representative corpus of the food content that potentially has more audience, we applied the aforementioned method to the following three Chilean online newspapers: [emol.com](http://emol.com), [latercera.com](http://latercera.com) and [elmostrador.cl](http://elmostrador.cl). The first is the online version of *El Mercurio* newspaper, founded in 1827 in the city of Valparaiso. It currently has several editions in the country. The edition of the capital, Santiago de Chile, was first published in 1900. *El Mercurio* is part of *El Mercurio Group*, one of the two main media conglomerates in Chile. According to audience data provided by Alexa and SimilarWeb in July 2016, its online version is the most-read news website in Chile. Meanwhile, [latercera.com](http://latercera.com) is the online version of *La Tercera* newspaper, founded in 1950 by the Copesa group, the other large Chilean media conglomerate. According to Alexa and SimilarWeb, [latercera.com](http://latercera.com) was the second or third most visited online medium in July 2016. Finally, [elmostrador.cl](http://elmostrador.cl) is an independent newspaper published exclusively online by La Plaza SA. According to Alexa and SimilarWeb, in July 2016, it was the seventh most-visited news website.

<b>Results of the proposed method applied in three Chilean newspapers</b>	
<b>Emol.com</b>	<b>Total</b>
Total number of news items retrieved with Google	4318
Total number of news items without duplicated URLs	3032
Total number of news items without duplicated headline	3024
Total number of news items without false positives	961
Total number of news items retrieved with <a href="http://import.io">import.io</a>	912
<b>Latercera.com</b>	
Total number of news items retrieved with Google	4462
Total number of news items without duplicated URLs	2739
Total number of news items without duplicated headline	2715
Total number of news items without false positives	973
Total number of news items retrieved with <a href="http://import.io">import.io</a>	957
<b>Elmostrador.cl</b>	

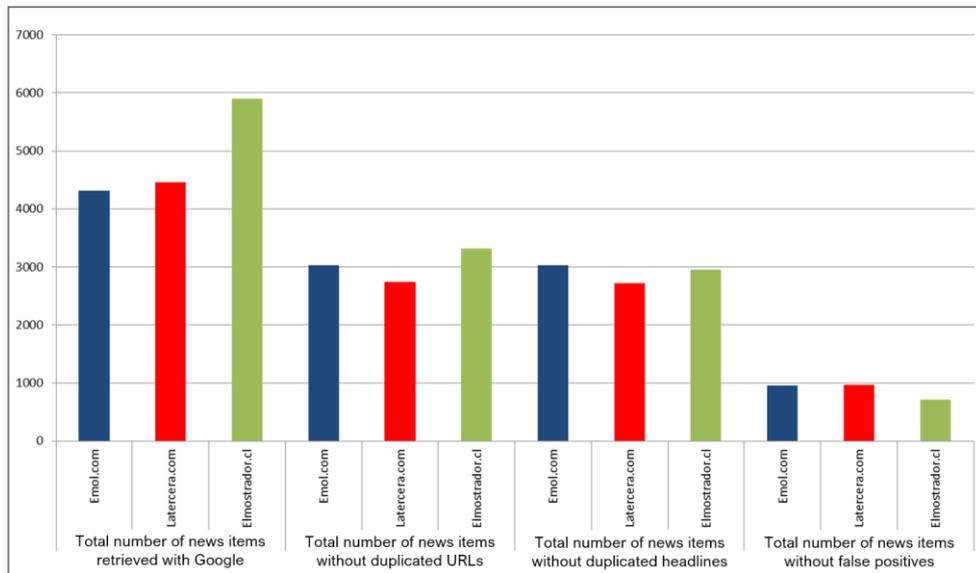
Total number of news items retrieved with Google	5897
Total number of news items without duplicated URLs	3313
Total number of news items without duplicated headline	2949
Total number of news items without false positives	695
Total number of news items retrieved with <a href="http://import.io">import.io</a>	682
Source: Authors' own creation	

The results of the application of the method proposed to form the corpus of online news on food are the following. The Google search engine yielded 4,318 cases for [emol.com](http://emol.com), 4,462 for [latercera.com](http://latercera.com), and 5,897 for [emol.com](http://emol.com). The elimination of URLs duplicates left [emol.com](http://emol.com) with 3,032 cases, [latercera.com](http://latercera.com) with 2,739, and [elmostrador.cl](http://elmostrador.cl) with 3,313. The elimination of cases based on duplicated headline (same news with two URLs) left [emol.com](http://emol.com) with 3,024 cases, [latercera.com](http://latercera.com) with 2,715, and [emol.com](http://emol.com) with 2,949. After the evaluation process carried out by the reviewers, the total numbers of news that deal with food and were published in 2016 was: [emol.com](http://emol.com), 961; [latercera.com](http://latercera.com), 1,029; and [elmostrador.com](http://elmostrador.com), 695. Finally, the total numbers of news items that were retrieved with [import.io](http://import.io) is the following: [emol.com](http://emol.com), 912; [latercera.com](http://latercera.com), 957; [elmostrador.cl](http://elmostrador.cl), 682. Thus, the corpus of this research was made up of a total of 2,551 news items published by the three news media outlets in question during 2016 (see Table 1).

“Total of news items retrieved by Google” refers to the total number of entries displayed by Google on the results page for the list of aforementioned keywords. “Total number of news items without URLs” refers to the total number of cases resulting from the detection and elimination of duplicated URLs carried out in Excel. “Total number of news items without duplicated headlines” is the total number of cases resulting from the process of detection and elimination of duplicated headlines carried out in Excel. “Total number of news items without false positives” refers to the total number of news considered relevant and valid by encoders. Finally, “total number of news items retrieved with [import.io](http://import.io)” refers to the number of retrieved news that are part of research corpus.

Figure 1 shows and compares quantitatively the evolution of this process for the three online newspapers under study. It is observed that [emol.com](http://emol.com) and [latercera.com](http://latercera.com) have very similar numbers and an almost identical evolution. This is not the case of [elmostrador.cl](http://elmostrador.cl). Even though the Google search for this newspaper provided approximately 33% more results than the rest of the newspapers, the total number of news that finally became part of the corpus was about 30% less. The discussion section offers some hypothesis that explain this significant difference.

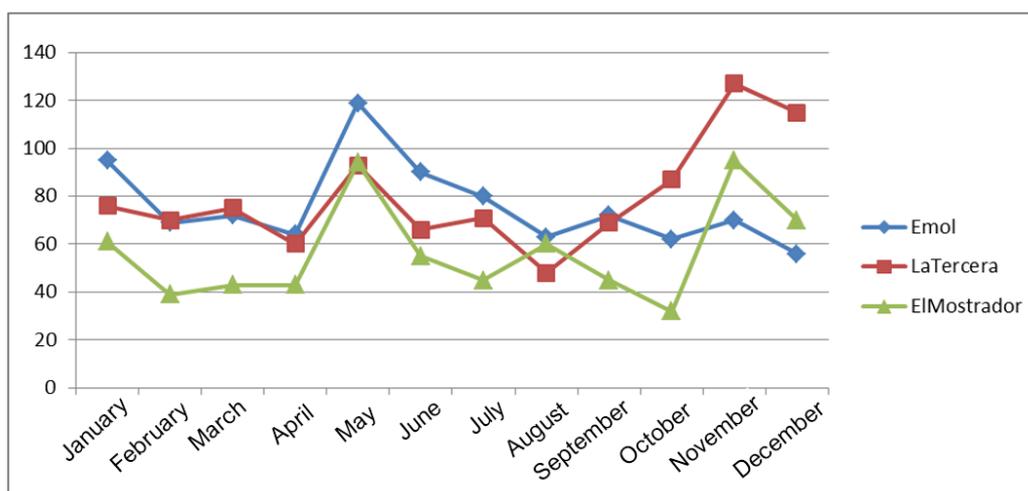
**Figure 1. Evolution of debugging of Google search results**



Source: Authors' own creation

The results also allow us to count and compare how many food news have been published in these three newspapers during each month in 2016. Figure 2 shows the similar evolution of the three newspapers, although May and October stand out for different reasons. While in May both newspapers significantly increase their news on food, offering a similar amount of news, October shows the largest quantitative disparity, according to our search and procedure, with [elmostrador.cl](http://elmostrador.cl) publishing the lowest number of news items (38), and [latercera.com](http://latercera.com) the highest (116), leaving [emol.com](http://emol.com) in the middle (67).

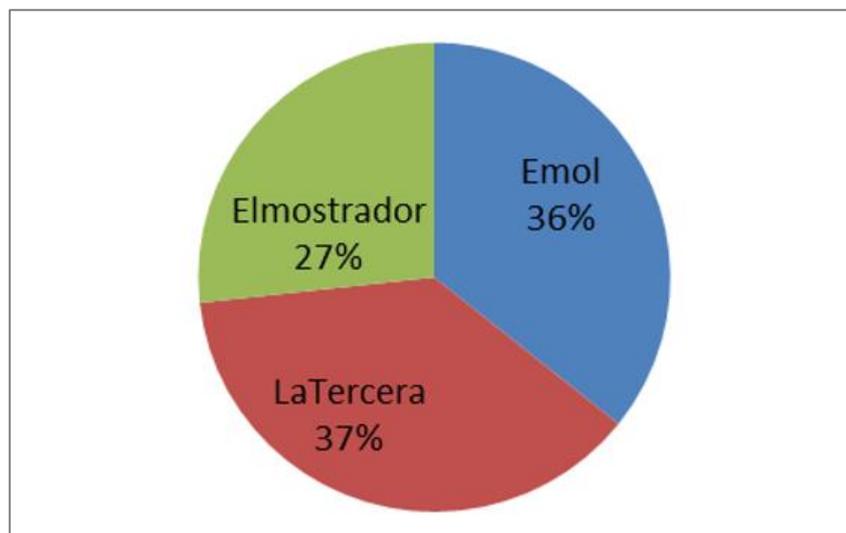
**Figure 2. Monthly evolution of the number of news items that are part of the research corpus**



Source: Authors' own creation

Figure 3 shows the percentage contribution of each online newspaper to the research corpus. The distribution is not homogeneous given that [elmostrador.cl](http://elmostrador.cl) only contributed with 26%, while [latercera.cl](http://latercera.cl) provided 37% of the corpus.

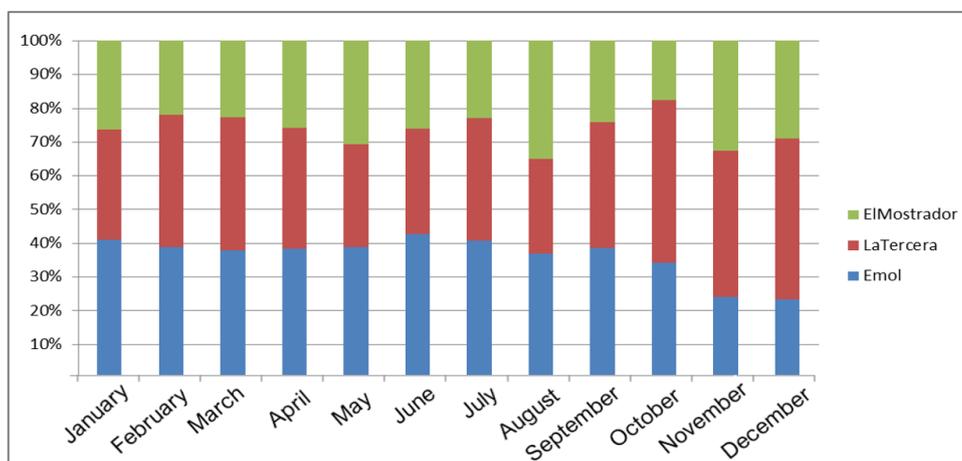
**Figure 3. Percentage contribution of each newspaper to the research corpus**



Source: Authors' own creation

As shown in Figure 4, this percentage distribution tends to be maintained throughout 2016, although in May and August [latercera.com](http://latercera.com) is relegated to the third place in terms of the percentage of news it contributed to the corpus, while in November and December [emol.com](http://emol.com) contributed the least.

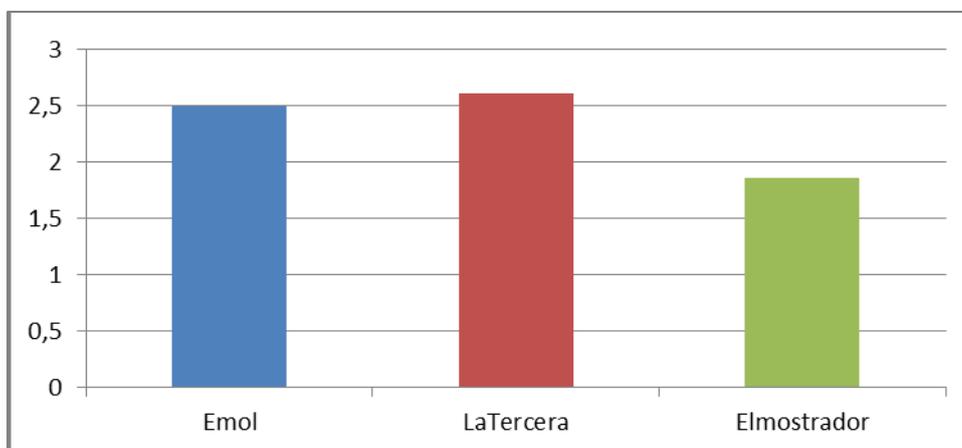
**Figure 4. Monthly percentage contribution of each newspaper to the research corpus**



Source: Authors' own creation

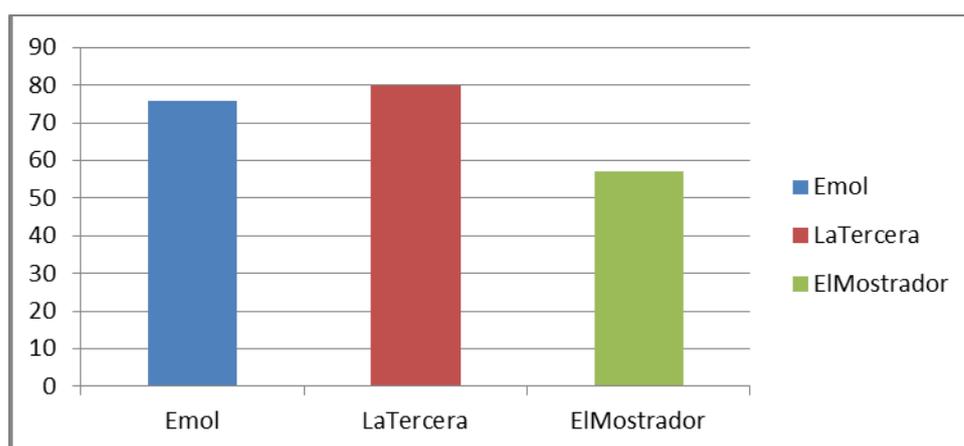
Figure 5 shows the monthly average of total number of news items that the three newspapers would have published in 2016, while Figure 6 shows the same average but by day. Quantitatively speaking, the differences between the two most-visited online newspapers in Chile are negligible. On the other hand, [elmostrador.cl](http://elmostrador.cl) has a significantly lower monthly and daily average (-26% approximately).

**Figure 5. Monthly average number of food news published by each newspaper**



Source: Authors' own creation

**Figure 6. Daily average number of food news published by each newspaper**



Source: Authors' own creation

#### 4. Discussion

Below is a discussion of the aspects of the results that can generate more doubts and, second, of the possibilities of exploiting the data that we have generated.

It is striking that [elmostrador.cl](http://elmostrador.cl) provided the largest number of URLs at the beginning of the process but ended up contributing the least to the corpus. It is unknown why Google displayed more results for this newspaper than for the rest. In any case, this initial difference disappears when the duplicated URLs and headlines are eliminated, leaving its number of URLs at the same level of the that of the other two newspapers.

The explanation for the large number of duplicated URLs in the case of the three newspapers is simple: the same news can contain more than one keyword. Therefore, Google displays the same news item in different searches because it contains some of the keywords requested. It is estimated that if a search is carried out separately for each of the keywords, without using the Boolean OR operator, the number of duplicates would be greater.

It might also be surprising that false positives add up to two thirds of the URLs retrieved by Google (not counting duplicates). However, it should be taken into account that the search that was carried out in Google with the keywords could equated to bottom trawling. In other words, Google is asked to retrieve all the URLs that contain any of the keywords. Thus, using the keyword “kitchen”, Google will offer URLs in which the news is actually about cooking or is related to it, but also URLs in which the word “kitchen” is used metaphorically or anecdotally, or news that contain links to other content that contain the word kitchen. Hence manual coding is very important to be able to separate the true positives from false positives.

Another aspect that needs to be discussed is the small difference between the number of URLs coded as pertinent and the number of news eventually retrieved with import.io. We have observed that some URLs were rejected by encoders because they were news in non-textual format, such as audio, video, or photography, or because they were links to Doc and PDF documents.

Having discussed the partial results of each phase of the process to retrieve the news that make up the corpus, we will discuss the final results, i.e., the number of news that each newspaper contributes to the corpus. To begin with, there is a remarkable quantitative similarity between [emol.com](http://emol.com) and [latercera.com](http://latercera.com), and the difference between these two newspapers and [elmostrador.cl](http://elmostrador.cl). It should not be surprising that the latter contributed 25% less news than the other two if one analyses these three online media from the perspective of political economy. As it has been documented, in Chile there is a strong media concentration (Diaz, 2017) that can be consider as duopoly or oligopoly (Gronemeyer et al. 2013). The two main media groups are Mercurio SAP and Copesa, the publishers of [emol.com](http://emol.com) and [latercera.com](http://latercera.com), respectively. In 2008, these two media groups concentrated almost 95% of daily print newspaper sales and received 65% of the advertising investment in the press (Mellado, 2012). The current situation, at least in terms readership, is very similar (ACHAP AG, 2017). So, while [emol.com](http://emol.com) and [latercera.com](http://latercera.com) are the two main newspapers of Chile’s media duopoly, with offline and online versions, [elmostrador.cl](http://elmostrador.cl) is an independent newspaper published by La Plaza S.A., which only owns an online version and whose readership is lower than the other two newspapers. Therefore, from a quantitative point of view, it would be logical for the two mainstream newspapers with similar structures and economic power to produce, in absolute terms, a similar amount of news about food and a higher amount than that produced by the independent newspaper that has less economic resources.

Although the amount of news shows differences, it is observed that the quantitative trend throughout 2016 of the three newspapers is quite homogeneous. In eight of the 12 months, all three newspapers follow the same trend, and in the remaining four months, at least two of them follow a similar trend. Perhaps the most striking case occurs in May, when the three newspapers experience a significant increase in food news. In the absence of a content analysis, a look at the headlines of the news published by the three newspapers indicates that they covered widely a critical event, namely, the ‘red tide’ crisis in the island of Chiloe, Chile, in May 2016. The quantitative similarities may be another argument in favour of the reliability of the method proposed here to form a large corpus of food news. The very similar tendencies, similar news numbers, and the fact that the three newspapers increased their publication of news in May, shows that Google has been able to retrieve the news for every month with similar efficacy. Therefore, Google would prove to be a good referee when it comes to locating the food news published by different newspapers. On the other hand, and again in the absence of content analysis, it is very likely that this similarity is also a quantitative proof of the homogeneity of the news agenda in Chile (Díaz & Mellado, 2017).

Figures 3 and 4 show the total and monthly percentage contribution of each newspaper to the research corpus. It is striking that [elmostrador.cl](http://elmostrador.cl) contributes the least –with a difference of 9.5% with respect to each of the other two newspapers– to the annual total, but in four months this same newspaper makes the second largest contribution. The explanation to this difference should be sought in the journalistic production of each newspaper rather than in the method used to locate and identify food news. It could also be another argument in favour of the consideration of long periods in exploratory studies interested in characterising the coverage of food.

Finally, figures 5 and 6 show the monthly and daily average numbers of food news published by each newspaper. In the absence of previous studies on the amount of food content published, there is only the possibility of comparing those averages with the anecdotal experience of the authors. An daily average of two to three news is no less of what one would expect from mainstream press, and probably a little more we would expect a priori. Therefore, the procedure presented here allows us to assume that the sample obtained is exhaustive and highly representative, at least in quantitative terms.

Beyond these results, it is necessary to close this section talking about the different ways in which you can exploit the 2551 news items identified and stored on a spreadsheet. First of all, we can analyse them using different software programs designed specifically for quantitative text analysis. In *The Content Analysis Guidebook* (Neuendorf, 2016), there is a section dedicated to the description of different applications that can extract frequencies, detect coincidences, identify statistically significant relationships between variables, terms and expressions; as well perform content analysis based on dictionaries built by the researcher. This possibility of analysis opens at least other three more.

The first and most obvious possibility is to carry out mix methods research. Combining quantitative and qualitative analysis, supported by software that allows for the quantitatively analysis of qualitative categories built by the researcher, allows doing research with large samples and with an important degree of depth.

The second possibility would take a first exploratory quantitative analysis with a smaller sample but a much deeper content analysis. In other words, rather than base research in communication and food on previous theory or critical or emergent events, these works could be justified in the results of a first quantitative analysis. For example, if a first quantitative analysis detects a very high percentage of food and health news and a very low percentage of news on the food distribution sector, it would be necessary to formulate research questions and objectives that could explain the difference.

Finally, a third possibility would be to use quantitative analysis to define and delimit samples to study a specific food dimension or topic. In our case, we are analysing the coverage that these three newspapers made in 2016 on the relationship between food and mental health. Assisted by software we can identify the news that dealt with this topic and create a new smaller corpus that will be analysed quantitatively and, if deemed relevant, qualitatively. Obviously, this analysis could also be complemented with critical discourse analysis.

## 5. Conclusion

This article has shown how a corpus of thousands of news published in online newspapers can be formed cheaply, effectively and technically within the reach of anyone with basic computer skills and minimal English skills, and that it can be subsequently analysed with big data techniques. The application of this method has been exemplified through the formation of a corpus of food news published by three Chilean newspapers. To do this, we departed from social theory to conceptualise food, and with the help of the same theory along with recent research studies, we operationalised food through a set of keywords. Finally, Google, [import.io](http://import.io) and three encoders were used to obtain a corpus of more than 2,500 food news items.

The authors of this article consider that the Spanish-speaking academic community should increase research on food and communication because the food content circulating in the mass media contributes significantly to people's knowledge of food. In addition, these same contents are rearranging the relationship of humans with the planet and their own bodies. Analysing the food content published by the online press is a first step to analyse their possible effects on the relationships that humans establish with food. Only if food is conceptualised as a totally social and human phenomenon, and only if it is operationalised as a food system and food culture, we will be able to cover all the contents related to the human food consumption that are disseminated by the mass media.

Finally, the results show that from a quantitative point of view, the coverage of the three Chilean newspapers under study is quite homogeneous, especially similar when we compare the two newspapers that have the largest online audience and represent the two media conglomerates that make up the Chilean duopoly, which has been described by the scientific literature. These two newspapers contributed to the research corpus 10% more news each, compared to the third newspaper, which belongs to an independent Chilean media group. This difference would be explained, on the one hand, by the differences in economic power that exist between this latter newspaper and the two newspapers belonging to the Chilean duopoly, and on the other hand, the homogeneity of the media agenda in Chile.

- **Funded research:** This article is a product of the research project titled “Food and the online press: Quantitative analysis of the food content published by Chile’s four most-visited digital newspapers” (DI17-0089), financed by the University of the Frontier/ Universidad de La Frontera, a Chilean public university that grants research funds through internal competition. The evaluation of the research projects participating in the competition is carried out by external scholars.

Dates:

-Start of research: June 2017

-End of research: November 2018

## 6. Notes

1. Although it is not incomplete, the food system model presented by Goody is the most basic. For a discussion on the different definitions of the food system that have been proposed up until 2005, see Contreras, 2005. For our purposes, we consider Goody’s scheme to be more than enough.
2. <http://www.alexa.com/topsites/countries/CL> (visited on July 2016 and January 2019)
3. See: <https://moz.com/blog/google-personalized-search> and <https://honeypotmarketing.com/remove-personalized-google-results/>
4. As of December 2018, this service offers a 7-day trial during which users can make a maximum of 500 queries. This would be enough to make tests and confirm that the service is useful to retrieve the news from the selected newspaper(s). Users can also request an estimate for academic projects to guarantee the necessary number of queries for the lowest price.

## 7. References

- ACHAP AG. (2017). Valida - Boletín de Circulación y Lectura 1º semestre 2017. Retrieved from <https://www.dropbox.com/s/xh7lk39edaaojro/Boletin%20de%20Circulaci%C3%B3n%20y%20Lectura%201%C2%B0%20semestre%202017.pdf?dl=0>
- Albala, K. (2013). *Routledge International Handbook of Food Studies*. Routledge.
- Alsina, M. R. (1989). *La construcción de la noticia*. Paidós España.
- Barthes, R. (1997). Toward a psychosociology of contemporary food consumption. *Food and Culture: A Reader*, 2, 28–35.
- Berelson, B. (1952). *Content analysis in communication research*. New York, NY: Free Press.

- Bernabeu-Peiró, À. (2015). La divulgación radiofónica de la alimentación y la nutrición. El ejemplo de Radio 5 Todo Noticias. *Revista de Comunicación y Salud*, 5, 36–53.
- Blanco Hernández, N. (2015). La nota gastronómica y el artículo de costumbres. *Estudios Sobre El Mensaje Periodístico*, 21(2), 953–967.
- Bourdieu, P. (1984). *Distinction: A social critique of the judgement of taste*. Harvard University Press.
- Bourdieu, P., Inda, A. G., Beneitez, M. J. B., Ordovás, M. J. G., & Lalana, D. O. (2001). *Poder, derecho y clases sociales* (Vol. 2). Desclée de Brouwer Bilbao.
- Colle, Raymond. (2011). *El Análisis de Contenido de las Comunicaciones*. Sociedad Latina de Comunicación Social.
- Contreras, J., & Gracia Arnaiz, M. (2005). *Alimentación y cultura: perspectivas antropológicas*. Barcelona: Editorial Ariel.
- Cramer, J. M., Greene, C. P., & Walters, L. M. (Eds.). (2011). *Food as communication: Communication as food*. Peter Lang New York.
- Díaz, M., & Mellado, C. (2017). Agenda y uso de fuentes en los titulares y noticias centrales de los medios informativos chilenos. Un estudio de la prensa impresa, online, radio y televisión. *Cuadernos. Info*, (40), 107–121.
- Dijk, T. A. van. (2009). *Discurso y poder. Contribuciones a los estudios críticos del discurso*.
- Douglas, M. (1972). Deciphering a meal. *Daedalus*, 61–81.
- Douglas, M. (1980). Las abominaciones del Levítico. *Pureza y Peligro. Un Estudio de Contaminación y Tabú*, 63–81.
- Duch, L., & Chillón, A. (2012). *Un ser de mediaciones*. Herder.
- Elias, N. (1989). *El proceso de la civilización. Investigaciones sociogenéticas y psicogenéticas*. fce.
- Evans, J., Rich, E., Davies, B., & Allwood, R. (2008). *Education, disordered eating and obesity discourse: Fat fabrications*. Routledge.
- Farré, M. (2004). *El noticiero como mundo posible: estrategias ficcionales en la información audiovisual*. La Crujía.
- Frye, J., & Bruner, M. S. (2013). *The rhetoric of food discourse, materiality, and power*. New York: Routledge. Retrieved from [http://www.novanet.ebib.com/EBLWeb/patron/?target=patron&extendedid=P\\_1039373\\_0](http://www.novanet.ebib.com/EBLWeb/patron/?target=patron&extendedid=P_1039373_0)

- Fúster, F., Ribes, M. Á., Bardón, R., & Marino, E. (2009). Análisis cuantitativo de las noticias de alimentación en la prensa madrileña en 2006. *Revista Española de Documentación Científica*, 32(1), 99–115.
- Gaínza, G. (2003). La práctica alimentaria y la historia. *LOTMAN DESDE AMÉRICA*.
- Goody, J., & Willson, P. (1995). *Cocina, cuisine y clase. Estudio de sociología comparada*. Editorial Gedisa.
- Habermas, J., Domènech, A., & Grasa, R. (1981). *Historia y crítica de la opinión pública: la transformación estructural de la vida pública*. Gustavo Gili Barcelona.
- Heldke, L. (2006). The Unexamined Meal is Not Worth Eating: Or, Why and How Philosophers (Might/Could/Do) Study Food. *Food, Culture and Society: An International Journal of Multidisciplinary Research*, 9(2), 201–219. <https://doi.org/10.2752/155280106778606035>
- Koldobsky, Daniela. (2011). La gastronomía en el discurso crítico actual. *DeSignis*, (18).
- Korthals, M., & Kooymans, F. (2004). *Before dinner: philosophy and ethics of food*. Dordrecht: Springer.
- Krippendorff, K. (1990). *Content Analysis Method: Theory and practice*. Buenos Aires: Paidós.
- Lacy, S., Watson, B. R., Riffe, D., & Lovejoy, J. (2015). Issues and best practices in content analysis. *Journalism & Mass Communication Quarterly*, 92(4), 791–811.
- Lawrence, R. G. (2004). Framing obesity, the evolution of news discourse on a public health issue. *The Harvard International Journal of Press/Politics*, 9(3), 56–75.
- LeBesco, K., & Naccarato, P. (Eds.). (2008). *Edible ideologies: representing food and meaning*. Albany: State University of New York Press.
- Lévi-Strauss, C. (1981). *El origen de las maneras de mesa* (Vol. 3). Siglo XXI.
- Lévi-Strauss, C. (2013). The Culinary Triangle. *Food and Culture: A Reader*, 40.
- Marín-Murillo, Flora, Armentia-Vizuite, José-Ignacio, & Olabarri-Fernández, Elena. (2016). Alimentación y Salud: Enfoques predominantes en la prensa española. *Revista Latina de Comunicación Social*.
- Martínez, A. A., & Poyatos, M. D. F. (2015). Gastronomy in the Spanish Press during the 19th century. *Estudios Sobre El Mensaje Periodístico*, 21(1), 17.
- Menezes Ferreira, C., de Castro Oliveira, V., & Terrón Blanco, L. (2015). Una temática de peso: el tratamiento de la obesidad en los periódicos brasileños. In *La pantalla insomne*. Revista Latina de

Comunicación Social. Retrieved from <http://www.revistalatinacs.org/15SLCS/libro-colectivo-2015.html>

Morin, E. (2014). *Le paradigme perdu. La nature humaine*. Seuil.

Narváez, T. A. (2013). Consejos dietéticos y nutricionales en la prensa española. *Revista Española de Comunicación En Salud*, 4(1), 17–26.

Nedelko, D. (2013, July 9). See How Your Website is Really Ranking - Unbiased Google Results. Retrieved January 2, 2019, from <https://honeypotmarketing.com/remove-personalised-google-results/>

Neuendorf, K. A. (2016). *The content analysis guidebook*. Sage.

Parasecoli, F. (2011). Savoring semiotics: food in intercultural communication. *Social Semiotics*, 21(5), 645–663. <https://doi.org/10.1080/10350330.2011.578803>

Plaza, J. F. (2012). Medios de comunicación, anorexia y bulimia. La difusión mediática del ‘anhelo de delgadez’: un análisis con perspectiva de género. *Revista Icono14. Revista Científica de Comunicación y Tecnologías Emergentes*, 8(3), 62–83.

Poulain, J.-P. (2017). *The sociology of food: eating and the place of food in society*. London: New York: Bloomsbury Academic.

Riffe, D., Lacy, S., & Fico, F. (2005). Analyzing Media Messages-2/E.: Using Quantitative Content Analysis in Research.

Roig, N. A. (2013). Alimentación y calentamiento global: “La larga sombra del ganado” en la prensa española. *Estudios Sobre El Mensaje Periodístico*, 19(1), 17–33.

Sánchez Gómez, F. (2010). La función didáctica del periodismo gastronómico en Internet. In *Alfabetización mediática y culturas digitales* (p. 172). Universidad de Sevilla.

Sandberg, H. (2007). *A matter of looks: the framing of obesity in four Swedish daily newspapers*.

Soler, J. (1997). The semiotics of food in the Bible. *Food and Culture: A Reader*, 55–66.

Strauss, C. L., Verón, E., & Menéndez, E. L. (1970). *Antropología estructural*. Eudeba.

Telfer, E. (2002). *Food for thought philosophy and food*. London; New York: Routledge. Retrieved from <http://search.ebscohost.com/login.aspx?direct=true&scope=site&db=nlebk&db=nlabk&AN=68639>

Thompson, J. B. (1998). *Los media y la modernidad: una teoría de los medios de comunicación*.

Thompson, J. R. (2012). Food talk: Bridging power in a globalizing world. *The Rhetoric of Food: Discourse, Materiality, and Power*, 58–70.

Velázquez, T. (1992). Los políticos y la televisión. Aportaciones de la teoría del discurso al diálogo televisivo. *Ariel, Barcelona*.

Weaver, D. A., & Bimber, B. (2008). Finding news stories: a comparison of searches using LexisNexis and Google News. *Journalism & Mass Communication Quarterly*, 85(3), 515–530.

Westall, D. (2011). La obesidad infantil en la prensa española. *Estudios Sobre El Mensaje Periodístico*, 17(1), 215–224.

related paper:

A C Yemsi-Paillissé, Y Acosta Meneses, M Martínez, E Calvo Gutiérrez (2018): “Aplicación de la crítica de los dispositivos a la cena performativa “El Somni” de El Cellar de Can Roca y Fran Aleu”. *Revista Latina de Comunicación Social*, 73, pp. 1267 a 1283.

---

### How to cite this article in bibliographies / References

R Sánchez Sabaté, C del Valle, M Mensa (2019): “Method for the construction of large thematic corpora of online news articles. Towards a corpus of food-related news”. *Revista Latina de Comunicación Social*, 74, pp. 594 to 617.

<http://www.revistalatinacs.org/074paper/1347/30en.html>

DOI: [10.4185/RLCS-2019-1347en](https://doi.org/10.4185/RLCS-2019-1347en)

Paper received on 7 December. Accepted on 23 February.  
Published on 28 February