

# The TAGFACT annotator and editor: A versatile tool

Ana Fernández-Montraveta<sup>a</sup> – Hortènsia Curell<sup>a</sup> – Glòria Vázquez<sup>b</sup> – Irene Castellón<sup>c</sup>  
Universitat Autònoma de Barcelona<sup>a</sup> / Spain  
Universitat de Lleida<sup>b</sup> / Spain  
Universitat de Barcelona<sup>c</sup> / Spain

**Abstract** – The multifunctional tool this paper presents has been developed within the TAGFACT project, a project that aims to automate the annotation of factuality –understood as the degree of commitment with which the writer presents situations– in Spanish journalistic texts. In what follows, the tool, which allows the compilation of the texts and the manual annotation of predicates, is described. The corpus created using it has been extracted in groups of three pieces of news covering the same event from newspapers with different ideologies (left wing, right wing and centrist). It is made up of 176 different pieces of news, containing 1,359 sentences and 46,947 words. The tool has been used so far to manually annotate a section of the ‘Gold Standard’ (approximately 10,000 words). It has proved to be versatile in that it allows for both the creation and management of corpora and corpus annotation, using any tags the user wants depending on the purpose of each corpus.

**Keywords** – annotation tool; corpus creation; corpus edition; Spanish journalistic texts

## 1. INTRODUCTION

The categorization of events with respect to their factual status is an area of growing interest in the field of Corpus Linguistics and Natural Language Processing. In recent years, several projects dealing with the annotation of corpora, either manual or automatic, with this type of information have been developed. So far, the most common approach has been the annotation of the degree of certainty with which the author of a message presents an event (Saurí 2008).

The objective of our project (TAGFACT), which is two years into its development, is to create a system for the automatic annotation of the degree of certainty implicit in the situations narrated in Spanish journalistic texts, an annotation

solely grounded on linguistic knowledge (Alonso *et al.* 2018).<sup>1</sup> In Spanish, this issue has not been dealt with in much depth, and what little has been done is based primarily on statistical processes (Wonsever *et al.* 2016).

One of the first steps in our project was the creation of a corpus of Spanish journalistic texts (the TAGFACT corpus) and then a portion of this corpus, which will constitute the ‘Gold Standard’, is being annotated manually. In order to perform these two tasks, the tool presented here was created. Before presenting the tool, it is necessary to describe briefly the main aspects of the project. Thus, Section 2 presents a brief state of the art and sets the framework for our annotation scheme –described in Vázquez and Fernández-Montraveta (in press)– required to fully understand the tool. Section 3 describes the design of the corpus and the ‘Gold Standard’ and, finally, Section 4 presents the tool and how it can be used to collect corpora and carry out the manual annotation.

## 2. THE ANNOTATION OF FACTUALITY

One of the groundbreakers in the annotation of factuality in texts is *FactBank* (Saurí and Pustejovsky 2009), which constitutes an innovative proposal for the representation of this semantic category in English. *FactBank* contains 9,488 events manually annotated with factuality information, and it also takes into account the source of information.

Various authors have drawn on Saurí and Pustejovsky (2009) for the annotation of different corpora, with the factuality values established using exclusively information from the text. In this respect, some projects worth mentioning are Diab *et al.* (2009), Soni *et al.* (2014), Tonelli *et al.* (2014), van Son *et al.* (2014) and Lee *et al.* (2015) for English; Matsuyoshi *et al.* (2010) and Narita *et al.* (2013) for Japanese; Minard *et al.* (2016) for Italian; Wonsever *et al.* (2016) for Spanish; and Velupillai (2011) for Swedish. Other authors, contrary to the framework used in *FactBank*, have considered factuality as linked to the knowledge of the world (Marneffe *et al.* 2012).

In our project, we basically follow Saurí and Pustejovsky (2009) and Diab *et al.* (2009), although we propose some innovations in the annotation scheme. The first decision is whether a predicate will be annotated or not. If it is decided not to, the

---

<sup>1</sup> The authors would like to acknowledge the support from the Ministerio de Economía, Industria y Competitividad: Research Project ‘Del texto al conocimiento. Factualidad y grados de certeza en español –TAGFACT’ (Grant number FFI2017–84008– P).

predicate is disregarded altogether. If it is annotated, four categories are used: ‘Polarity’, ‘Degree of commitment’, ‘Time’ and ‘Dynamicity.’ Following Diab *et al.* (2009), we prefer the term ‘commitment’ rather than ‘certainty’ (Saurí and Pustejovsky 2009), since it reflects better that we are describing the author’s view of the event.

Regarding ‘Time’ –following van Son *et al.* (2014), Wonsever *et al.* (2016) and Matsuyoshi *et al.* (2010)– we assign one of the following values: ‘Present’, ‘Past’ or ‘Future.’ Future situations are different from present and past ones, since they can never denote facts that have happened at the point of narration. It could be argued, hence, that certainty does not apply to them. However, we claim that the writer can present a future situation with commitment or with lack of it, and this is one of the innovations of our project. Another important novelty is the inclusion of ‘Dynamicity’, in which not only do we distinguish between states and events, but we also provide a fine-grained annotation of states. Following Tonelli *et al.* (2014) and van Son *et al.* (2014), we treat absolute truths –as in *The Earth is round*– and habits –as in *In our country we usually have lunch at 2 p.m.*– differently from other types of stative situations. In the former, commitment is always stronger since they represent knowledge commonly agreed upon by a community. The latter are of interest because they always include more than one situation and some of the events are in the past, some in the present and some in the future.

Furthermore, following Saurí (2008), van Son *et al.* (2014) and Prabhakaran *et al.* (2015), among others, it was decided to signal the authorship or source of each commitment, considering that there is not such a thing as ‘reality’ and that facts are always narrated from a given perspective. In addition, our annotation includes the predicate and all the entities involved in it, since we believe that a fact necessarily contains all the participants in the situation.

Regarding the automation of the annotation process, the systems currently available follow two different approaches: those using machine-learning techniques and those based, at least partially, on linguistic information. Among the former, Mullick *et al.* (2019) present the development of a deep neural network based on the ‘Factuality Judgment Model’, while Huang *et al.* (2019) use ‘Bi-directional Long Short-Term Memory’ (BiLSTM), that is, neural networks to learn contextual information about the event in sentences. The latter consider that annotating factuality at sentence level provides an incomplete picture and their unit of analysis is the document.

On the other hand, *De Facto* (Saurí 2008) automates part of the annotation of factuality using knowledge extracted from a corpus, that is, linguistic information. It is not fully automatic, though, since some knowledge modules were created manually. Different kinds of automatic tools have also been developed for various languages, for example Minard *et al.* (2006), Narita *et al.* (2013) and Lee *et al.* (2015). In TAGFACT, only linguistic information is used to tag factuality.

As regards edition and annotation tools, there are various tools available nowadays. Some of these are designed with a general purpose and are tools for the creation, annotation and edition of corpora at different levels –*UAM Corpus Tool* (O'Donnell 2008). Some of them include the functionality of defining the categories ('Tagset'), together with the possibility to annotate at different linguistic levels. Other tools incorporate the automatic treatment of certain aspects, such as the segmentation into sentences (sentence split) or the identification of words (tokens)– *ANNIS* (Krause and Zeldes 2016). Still other annotators aim at a specific type of corpus or level of analysis, such as *Knowtator* (Ogren 2006) and *DART* (Weisser 2006). Most annotators include the production of output in XML format, the possibility of conducting complex searches and statistical tools. Our tool is versatile since, on the one hand, it allows the user to organize and manage the texts compiled for the corpora, to interact with the database created (in MySQL), to extract the final data in XML and to create tabs, while still permitting the automatic processes of tagging and parsing.

### 3. THE TAGFACT CORPUS

The TAGFACT corpus includes news articles from several Spanish newspapers. Specifically, three pieces of news, narrating the same event, were collected from newspapers with different ideologies: right wing (*La Razón*), left wing (*El Diario*) and centrist (*El Periódico*).<sup>2</sup> This will eventually allow the analysis of the author's stance and the role of ideology in journalism. The articles were mainly chosen from two genres: politics and sports. Those two genres offer the possibility of finding news and describing facts more than opinions (as opposed to, for example, op-ed columns). As for politics, since the newspapers represent clearly different ideologies, it is to be expected that the perspective from which certain events are narrated will be distinct. In sports, the

---

<sup>2</sup> When it was not possible to find one of the pieces in these media, articles were taken from another one with similar characteristics (*ABC*, *Público*, *La Vanguardia* or *20 Minutos*, among others).

well-known rivalry between the football teams Barcelona and Madrid will guarantee varied points of view. At present, the corpus includes 176 different pieces of news, containing 1,359 sentences and 46,947 words.

For each piece of news, metadata is saved in order to facilitate the access to the information about the source: name of the newspaper, section, date, author, news URL and geographical location. The data is structured in several fields, following the structure of the newspaper article. Any extra information, which is an informative part of the piece, is also included in this structure: namely the title and subtitle, the text and images and TWITTER comments, as shown in Figure 1. At present, we are manually annotating a part of the corpus, which will constitute the ‘Gold Standard.’ The total volume of words in the ‘Gold Standard’ corpus is approximately 10,000.

<b>Group</b> 1 - Máster Cifuentes		<b>Item creator</b>	
<b>Newspaper</b> 20 Minutos		<b>Section</b> Technology	
<b>Date</b> 19/03/2019	<b>Author</b> Judith Vives		
<b>Title</b> Así alimenta Youtube las teorías que afirman que la Tierra es plana			
<b>Subtitle</b> Una investigación sugiere que la plataforma de vídeos ayuda a los teóricos de la conspiración tierraplanista			
<b>Text</b> Youtube podría estar desempeñando un papel importante para convencer a algunas personas de que la Tierra es plana. Así lo sugiere un estudio de la Texas Tech University que se ha basado en realizadas a personas que han participado en conferencias sobre el tema. Las entrevistas realizadas a estas personas demuestran, según el estudio de la Texas Tech University, que la mayoría basan sus creencias en los vídeos que han visto en Youtube. Estos videos trat demuestran que la tierra no es redonda.			

Figure 1: Partial structure of the data for each item in the corpus

#### 4. THE TAGFACT TOOL

A multi-purpose tool was created to compile and annotate corpora. This tool has two main functionalities: first, to compile and manage large collections of text and, second, to facilitate the annotation of any specific corpus. The first functionality allows corpus creation, edition and management through a highly user-friendly interface. Data is collected and saved in a MySQL Database. The interface offers the possibility of querying the database and allows the downloading of all the information in either Excel or XML formats.

#### 4.1. Corpus creation, edition and management

The corpus creation tool permits the compilation of one or several collections of texts, each of which can be saved as an independent database and can be independently named, edited and modified. In addition, once a particular corpus has been collected, the tool allows for the edition, modification and management of each item in the collection, regardless of whether it has already been annotated.

The tool includes a default administrator that has the capacity to add any number of users able to interact with the database in different ways, depending on the role assigned to them (administrator or annotator). Besides assigning roles to the collaborators in the project, the administrator is in charge of assigning the sections that each annotator has to deal with, and can view all the annotations, as shown in Figure 2.

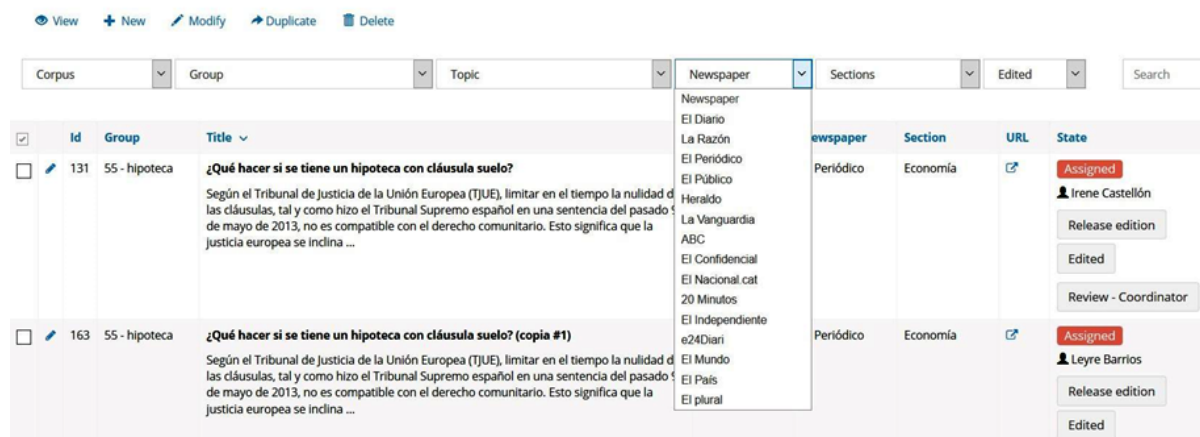


Figure 2: News items and newspapers in the corpus

The tool is connected to a parser that analyzes the texts. The administrator sends the text to the parser and has access to the pieces of news at all stages in the process. Annotators have access exclusively to items assigned to them by the administrator, as illustrated in Figure 3.

Id	Title	Text	Date	Newspaper	Section	URL	State
87	Los bomberos de Lesbos: "Estamos condenados a salvar vidas, no por salvarlas"	Manuel Blanco, Enrique Rodríguez y Julio Latorre esperan poner fin este lunes a una pesadilla "surrealista" que comenzó en enero de 2016. Estos tres bomberos sevillanos miembros de la asociación Proem-Aid, que participaban de forma voluntaria en los peores días de la crisis de refugiados en Grecia,	06/05/2018	El Periódico	Sociedad		Assigned Gloria Vázquez Release edition Edited Review - Coordinator
4	Los bomberos españoles que ayudaron a inmigrantes en Grecia llegan a juicio: "Nuestra condena sería un aviso a navegantes"	Onio Reina recuerda como si fuera ayer el momento en el que pisó por primera vez la arena de las playas de Lesbos. Era diciembre de 2015 y este bombero sevillano acababa de llegar a la isla griega junto a otros tres compañeros tras un largo trayecto en coche desde España. Uno de los vecinos les ...	06/05/2018	El Diario	Política		Processed 14/01/2019 12:24:56 Assign to...
5	Tres bomberos españoles, juzgados en Grecia por tráfico de personas	Manuel Blanco, José Enrique Rodríguez y Julio Latorre, los tres bomberos españoles detenidos en 2016 en Lesbos por la Guardia Costera Griega cuando realizaban tareas de rescate de refugiados con la ONG Proem-AID, han llegado esta mañana al tribunal que les va a juzgar por un presunto delito de ...	07/05/2018	La Razón	Política		Processed 14/01/2019 12:25:43 Assign to...

Figure 3: Corpus management

Texts undergo several stages through the process: 'Initial', 'Pending', 'Processed', 'Assigned' and 'Edited.' 'Initial' state means that an item in a collection has been introduced and documented (metadata); then, it becomes 'Pending.' It is subsequently sent to the parser, which sends it back as 'Processed', that is, segmented into sentences and with the list of predicates and their corresponding arguments, as shown in Figure 4. When the administrator assigns an item to an annotator, the stage will change to 'Assigned', and the process of annotation can start. While annotating the factual values, the syntactic structure can be corrected or 'Edited', if required.

**Predicates:**

El candidato de el PSOE a la Presidencia\_de\_el\_Gobierno , Pedro\_Sánchez , llegó el jueves a el debate sobre la moción de censura contra Mariano\_Rajoy con el acuerdo previo con el PNV de respetar los Presupuestos aprobados por el Gobierno y , con ello , garantizar se el apoyo mayoritario de la Cámara .

+ New ← Back

Main (beg. > end)	Trig	Categories	Arguments	Trig/Voice	Main	Probl
llegó (t1.12) *	✓	Applies *	El candidato de el PSOE a la Presidencia_de_el_Gobierno , Pedro_Sánchez ,	A 11	✓	✗
		Past *	a el debate sobre la moción de censura contra Mariano_Rajoy	A 10		
		Commitment	el jueves	A 2		
		Positive *	con el acuerdo previo con el PNV de respetar los Presupuestos aprobados por el Gobierno y , con ello , garantizar se el apoyo mayoritario de la Cámara	A 28		
		Event *				

Figure 4: Editor interface – predicate *llegó* 'arrived'

#### 4.2. Corpus tagging

As mentioned above, the first step in the annotation process is to send the text to an external analyzer. The most complete tools for Spanish, in terms of the different levels of analysis provided, were considered in order to choose the most adequate parser for our project. According to Soroa *et al.* (2017), these are *Freeling* (Padró and Stanilovsky 2012) and *Ixa Pipe* (Agerri *et al.* 2014). Both offer a level of document representation

and resolve co-referencing, in addition to providing a morphological and syntactic analysis. *Freeling* is the parser that has been more thoroughly evaluated and has obtained an optimal index for syntactic parsing, more specifically 84% of accuracy in the analysis of dependencies (Lloberes *et al.* 2015). No evaluation of syntactic performance of the *Ixa Pipe* parser was found, even though other levels, such as co-reference, have been assessed, achieving 55% of accuracy. In addition, some tests using our corpus were performed and the final decision was to use *Freeling* as a basic working tool.

An important problem presented by the *Freeling* output is the identification of predicates, more specifically, the recognition of eventive nouns. Regretfully, the recognition of this type of element does not seem to work well, as shown in (1), and the decision to deactivate all eventive nouns has been made.

- (1) *Esto significa que la justicia europea se inclina porque la banca tenga que devolver a sus clientes lo cobrado de más.*

‘This means that the European justice favors that the banks return the money overcharged to their clients.’

The analysis of compound verbs –both complex tenses and verb periphrases– is problematic in *Freeling* as well, since it separates the verbs in a complex as two (or more) independent predicates. This problem has been easily overcome with a pre-process that rewrites them as one single predicate. In the case of complex tenses, a simple rule identifies those structures in which *haber* ‘have’ is followed by a past participle. As for verb periphrases, another rule identifies periphrastic verbs and unites them with the corresponding main verb. The only problematic issues that cannot be solved automatically are the cases in which there is one or more lexical items placed between the auxiliary and the main verb, as shown in (2). These cases are dealt with manually.

- (2) ... *ya estaba en el agua rescatando a inmigrantes que partían de la costa.*


‘... he was already in the water rescuing immigrants leaving the coast...’

#### 4.3. Corpus annotation

Once an item of the collection has been returned from the parser one can proceed to the manual annotation. The user can validate the structure sent by the parser through the



interface and categorize each predicate regarding its factual status, according to the categories proposed in the scheme (Section 2). In our project, we propose four categories, but the number can be increased or decreased by the administrator. Figure 5 shows the first layer of annotation.



Categories of predicate			
<a href="#">View</a> <a href="#">+ New</a> <a href="#">✎ Modify</a> <a href="#">🗑 Delete</a>			
<input checked="" type="checkbox"/>	Name	Default	Order ▾
<input type="checkbox"/>	<a href="#">✎</a> NA (Does not apply)	✗	↑ 10 ↓
<input type="checkbox"/>	<a href="#">✎</a> Applies	✓	↑ 20 ↓
<input type="checkbox"/>	<a href="#">✎</a> Eventive noun	✗	↑ 30 ↓
<input type="checkbox"/>	<a href="#">✎</a> Error: no predicate	✗	↑ 40 ↓

Figure 5: Tag creation and management

This layer allows the annotator to make the first decision: whether the annotation of the factual status of an event is relevant (‘Applies’ vs. ‘(NA) Does not apply’). A predicate is only labeled as ‘Does not apply’ when the clause describes a wish or a conjecture, such as *deseara* ‘wished’, as illustrated in (3).

(3) ... *podría, si lo deseara, poner fin a la investigación o incluso ejercer su poder de perdón.*

‘... he could, if he wished, end the investigation or even exercise his power of forgiveness.’

The other options in this layer of annotation are: ‘Eventive noun’ –as explained above, these nouns are not annotated at the current stage of the project– and ‘Error: no predicate’ for words identified as eventive which are, in fact, not eventive, as shown in (4).

(4) *Hace unos meses Mongolia lanzó una campaña de apoyo para recaudar fondos.*

‘A few months ago Mongolia launched a fundraising support campaign.’

Finally, if the annotators consider that the predicate has to be tagged with respect to factuality, they use ‘Applies.’ The next step is to determine the following aspects: the time referred to by the predicate (present, past or future), the degree of the writer’s commitment towards the truth or falsehood of the predicate (‘Commitment’ or ‘Non-commitment’), polarity (‘Positive’ or ‘Negative’) and, finally, dynamicity (‘Event’, ‘Mental Predicate’ or ‘Property’).

Regarding temporal information, tense and time do not always correspond, as can be seen in (5), where a present tense indicates a past time.

(5) *Landrum alerta que el algoritmo que sugiere nuevos vídeos a las personas que buscan información sobre este tema les acaba llevando a un pozo de información incorrecta.*

‘Landrum warns that the algorithm that suggests new videos to people looking for information on this topic ends up leading them to a reservoir of incorrect information.’

Future situations are labeled differently from uncertain past and present situations because they are radically different in nature. Only in the first case is uncertainty absolute since future situations have not happened yet. The author can only express (non)-commitment towards the possibility of situations happening in the future.

Regarding polarity, one value, ‘Positive’ or ‘Negative’, is applied to the whole sentence, as in (6). Although polarity can have different scopes, at the present stage the tool only allows to assign polarity to the whole predicate. This is a limitation of the project that can be addressed in the future.

(6) *Estos vídeos tratan de mostrar evidencias que demuestren que la tierra no es redonda.*

‘These videos try to present evidence that proves that the Earth is not round.’

The ‘Dynamicity’ tag accounts for the internal structure of predicates, differentiating between stative, dynamic situations (events) and mental processes. Stative situations express properties of individuals or events, whereas events refer to actions or processes that happen in the world and have the capacity to modify it, as illustrated in (7). Mental predicates describe cognitive processes, as shown in (8). As for stative situations, if the property refers to individuals (both people and objects), the tag used is ‘Non-eventive Property’, and when it refers to events, ‘Property Event’, as can be seen in (9) and (10) respectively. Finally, ‘Property-Absolute Truth’ is used for properties considered as such by culture or scientific proof, as shown in (11).

(7) *El estudio se ha realizado a partir de las entrevistas con 30 asistentes a dos conferencias sobre teorías de la Tierra plana.*

‘The study was carried out based on interviews to 30 attendees at two conferences about flat Earth theories.’

(8) *En los últimos tiempos han proliferado las personas que no aceptan la idea de que el planeta Tierra es redondo.*

‘In recent times, people who do not accept the idea that planet Earth is round have proliferated.’

(9) *No es un ataque político, es un ataque personal.*

‘It is not a political attack; it is a personal attack.’

(10) *Landrum alerta que el algoritmo que sugiere nuevos vídeos a las personas...*

‘Landrum alerts that the algorithm that suggests new videos to people...’

(11) *.... el planeta Tierra es redondo*

‘... planet Earth is round.’

As pointed out above, the final goal of the project is to develop an automatic tool for the recognition and annotation of factuality. To this aim, the editor permits the annotation of linguistic cues (triggers), either morphological or lexical, that justify the choice of a tag so that they can be used in the automatic tool. For example, in (12), the clause containing the verb *tiene* ‘has’ is annotated, whereas the conditional clause is not since it is a condition. The trigger for its interpretation is the word *si* ‘if.’ Similarly, in (13), the verb *explica* ‘explains’ would be annotated as a trigger for the interpretation of the clause as a commitment for two reasons: first, the verb tense used (present indicative) and, second, the semantic class that the verb belongs to (verb of communication).

(12) *Pero al margen de ese posible acuerdo existen algunas opciones en función de si se tiene la cláusula suelo en la hipoteca...*

‘But apart from this possible agreement there exist some options depending on whether your mortgage has a base clause...’

(13) *... explica Óscar Serrano, abogado del “Col·lectiu Ronda”, exigiendo la devolución íntegra y retroactiva de los intereses pagados de más.*

‘... explains Óscar Serrano, a lawyer with the “Col·lectiu Ronda”, demanding the full and retroactive return of interest paid in excess.’

The possibility of annotating any relevant voices in the narration other than the writer’s is also considered, because they might modify the interpretation of the event. The author of the piece of news is always considered the main narrator, presenting events and situations from a particular perspective. When the main author provides the name of a different narrator, as in (14), and explicitly states the source of the information, it is understood that the author is somehow moving away from it.

(14) *Algunas de las personas consultadas aseguran que al principio solo miraban los videos para criticarlos...*

‘Some of the people consulted claim that at first they only watched the videos to criticize them...’

The tool provides a field where any problems encountered during the annotation process can be recorded and then discussed before the next stage in the project, as shown in Figure 6. Keeping a log of doubts allows the creation of lists of problematic configurations with the view of the systematization of the annotation.

<input type="checkbox"/>		3	Pero , ahora se abre un proceso que aún no está de el todo claro : ¿ devolverán de oficio los bancos el dinero cobrado de más o los afectados tendrán que acudir a los tribunales para reclamar lo ?		A 40	8
<input type="checkbox"/>		4	La sentencia europea no deja claro cómo tiene que aplicar se la resolución .	1	A 14	3
<input type="checkbox"/>		5	Puede que esta duda se resuelva políticamente si el Gobierno y la oposición se ponen de acuerdo sobre cómo hacer que los bancos devuelvan el dinero .	1	A 27	7
<input type="checkbox"/>		6	Pero a el margen de ese posible acuerdo existe algunas opciones en función de si se tiene la cláusula suelo en la hipoteca y no se ha reclamado ante los tribunales , si se ha planteado una demanda que ya ha sido resuelta o si se ha alcanzado un acuerdo con la entidad financiera para cambiar la hipoteca de tipo variable con suelo a tipo fijo .	1	A 67	20

There are problematic predicates

Figure 6: Mark for problematic predicates

## 5. SUMMARY AND CONCLUSIONS

In this paper we have presented the tool created in the TAGFACT project, whose main objective is to create a tool to automatically annotate factuality in Spanish. This task has become especially relevant in the last few years in the field of Natural Language Processing. The multifaceted tool presented allows for corpus creation, management and annotation and has been used to create the TAGFACT corpus and the ‘Gold Standard.’

The corpus includes texts extracted from different Spanish newspapers, belonging to different political ideologies. The extraction was carried out in groups of three pieces of news, each from a different newspaper, covering the same event, which can provide information about how facts are accounted for in each of the papers. At present, the corpus contains 46,947 words in 176 pieces of news and the ‘Gold Standard’ consists of around 10,000 words.

With respect to corpus creation, the tool greatly facilitates inputting both the text and the metadata required for text identification. In addition, it presents a user-friendly interface to edit and manage the corpora created. Regarding the annotation of the corpus, the fact that it is linked to *Freeling* permits the automatic segmentation of texts into sentences and clauses. In this way, it can be used to annotate corpora at different levels, from whole texts to just words. The tool allows users to create their own labels, so it is possible to annotate linguistic information relevant to varied projects. Another relevant feature offered by the editor is the possibility of marking the voice of a predicate, that is, the narrator of the situation, or any other word that might trigger a decision for the various levels of annotation.

Currently, we are manually annotating the pieces of news of the ‘Gold Standard’ with regard to how events are presented with respect to author’s commitment. In TAGFACT, factual information is very rich and is inferred taking into account the four layers described in this paper that cover the different aspects taken into consideration in the project.

The two functionalities of the tool, corpus creation and corpus annotation, are versatile resources that can be freely used by any researcher working in Corpus Linguistics. The tool will be made available on the Internet under a *GNU General Public License*. In the future, we aim to complete the implementation of the system for the automatic annotation of factuality for Spanish.

#### REFERENCES

- Agerri, Rodrigo, Josu Bermúdez and German Rigau. 2014. Ixa pipeline: Efficient and ready to use multilingual NLP tools. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation*. Reykjavik: European Language Resources Association, 3823–3828.
- Alonso, Laura, Irene Castellón, Hortènsia Curell, Ana Fernández-Montraveta, Sònia Oliver and Glòria Vázquez. 2018. Proyecto TAGFACT: Del texto al conocimiento. Factualidad y grados de certeza en español. *Procesamiento del Lenguaje Natural* 61: 151–154.
- Diab, Mona, Bori Levin, Teruko Mitamura, Owen Rambow, Vinodkumar Prabhakaran, Vinodkumar and Weiwe Guo. 2009. Committed belief annotation and tagging. In Manfred Stede, Chu-Ren Huang, Nancy Ide and Adam Meyers eds. *Proceedings of the Third Linguistic Annotation Workshop*. Singapur: Association for Computational Linguistics, 68–73.
- Huang, Rongtao, Zou Bowei, Wang Hongling, Li Peifeng and Zhou Guodong. 2019. Event factuality detection in discourse. In Jie Tang, Min-Yen Kan, Dongyan

- Zhao, Sujian Li and Hongying Zan eds. *Natural Language Processing and Chinese Computing*. NLPCC 2019. Lecture Notes in Computer Science. Vol. 11839. Springer, Cham, 404–414.
- Krause, Thomas and Amir Zeldes. 2016. ANNIS3: A new architecture for generic corpus query and visualization. *Digital Scholarship in the Humanities* 31/1: 118–139.
- Lee, Kenton, Yoav Artzi, Yejin Choi and Luke Zettlemoyer. 2015. Event detection and factuality assessment with non-expert supervision. In Lluís Màrquez, Chris Callison-Burch and Jian Su eds. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Lisbon: Association for Computational Linguistics, 1643–1648.
- Lloberes Marina, Irene Castellón, Lluís Padró. 2015. Suitability of ParTes test suite for parsing evaluation. *Proceedings of the 14<sup>th</sup> International Conference on Parsing Technologies*. Bilbao: Association for Computational Linguistics, 61–65.
- Marneffe, Marie-Catherine, Christopher D. Manning and Christopher Potts. 2012. Did it happen? The pragmatic complexity of veridicality assessment. *Computational Linguistics* 38/2: 301–333.
- Matsuyoshi, Suguru, Megumi Eguchi, Chitose Sao, Koji Murakami, Kentaro Inui and Yuji Matsumoto. 2010. Annotating event mentions in text with modality, focus and source information. In Nicoletta Calzolari, Khalid Choukri, Bente Maegaard, Joseph Mariani, Jan Odijk, Stelios Piperidis, Mike Rosner and Daniel Tapias eds. *Proceedings of the Seventh International Conference on Language Resources and Evaluation*. Valetta: European Language Resources Association, 1456–1463.
- Minard, Anne-Lyse, Manuela Speranza and Tommaso Caselli. 2016. Event factuality annotation task (FactA). In Pierpaolo Basile, Anna Corazza, Franco Cutugno, Simonetta Montemagni, Malvina Nissim, Viviana Patti, Giovanni Semeraro and Rachele Sprugnoli eds. *Proceedings of the Fifth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian*. Napoli: Open Edition Books, 32–39.
- Mullick, Ankan, Sourav Pal, Projjal Chanda, Arijit Panigrahy, Anurag Bharadwaj, Siddhant Singh and Tanmoy Dam. 2019. D-FJ: Deep neural network based factuality judgment. *TrueFact, Truth Discovery and Fact Checking: Theory and Practice workshop*.
- Narita, Kazuya, Junta Mizuno and Kentaro Inui. 2013. A lexicon-based investigation of research issues in Japanese factuality analysis. In Ruslan Mitkov and Jong C. Park eds. *Proceedings of the Sixth International Joint Conference on Natural Language Processing*. Nagoya: Asian Federation of Natural Language Processing, 587–595.
- O'Donnell, Mick. 2008. The UAM CorpusTool: Software for corpus annotation and exploration. In Carmen Bretones, José Francisco Fernández, José Ramón Ibáñez, M. Elena García, M. Enriqueta Cortés, Sagrario Salaberri, M. Soledad Cruz, Nobel Perdu and Blasina Cantizano eds. *Applied Linguistics Now: Understanding Language and Mind*. Almería: Universidad de Almería, 1433–1447.
- Ogren, Philip V. 2006. Knowtator: A protégé plug-in for annotated corpus construction. In Alex Rudnicky, John Dowding and Natasa Milic-Frayling eds. *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Demonstrations*. New York: Association for Computational Linguistics, 273–275.
- Padró, Lluís and Evgeny Stanilovsky. 2012. FreeLing 3.0: Towards wider multilinguality. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Mehmet

- Uğur Doğan, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk and Stelios Piperidis eds. *Proceedings of the Eight International Conference on Language Resources and Evaluation*. Istanbul: European Language Resources Association, 2473–2479.
- Prabhakaran, Vinodkumar, Tomas By, Julia Hirschberg, Owen Rambow, Samira Shaikh, Tomek Strzalkowski, Jennifer Tracey, Michael Arrigo, Rupayan Basu, Micah Clark, Adam Dalton, Mona Diab, Louise Guthrie, Anna Prokofieva, Stephanie Strassel, Gregory Werner, Yorick Wilks and Janyce Wiebe. 2015. A new dataset and evaluation for belief/factuality. *Proceedings of the Fourth Joint Conference on Lexical and Computational Semantics*. Denver: Association for Computational Linguistics, 82–91.
- Saurí, Roser. 2008. *A Factuality Profiler for Eventualities in Text*. Massachusetts: Brandeis University dissertation.
- Saurí, Roser and James Pustejovsky. 2009. FactBank: A corpus annotated with event factuality. *Language Resources and Evaluation* 43/3: 227–268.
- Soni, Sandeep, Tanushree Mitra, Eric Gilbert and Jacob Eisenstein. 2014. Modeling factuality judgments in social media text. In Kristina Toutanova and Hua Wu eds. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. Volume 2: Short Papers*. Baltimore: Association for Computational Linguistics, 415–420.
- Soraa, Aitor, German Rigau, Jordi Porta, Jordi Atserias, Xavier Gómez Guinovart and Horacio Saggion. 2017. *Plataformas y Sistemas de Procesamiento Lingüístico de Alto Rendimiento*. Plan de impulso de las tecnologías del lenguaje: Ministerio de Energía Turismo y la Agenda Digital.
- Tonelli, Sara, Rachele Sprugnoli and Manuela Speranza. 2014. NewsReader guidelines for annotation at document level. In extension of deliverable D3. *Technical Report NWR-2014-2*. Trento.
- van Son, Chantal, Marieke van Erp, Antske Fokkens and Piek Vossen. 2014. Hope and fear: Interpreting perspectives by integrating sentiment and event factuality. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asunción Moreno, Jan Odijk and Stelios Piperidis eds. *Proceedings of the Ninth International Conference on Language Resources and Evaluation*. Reykjavik: European Language Resources Association, 26–31.
- Vázquez, Gloria and Ana Fernández-Montraveta. In press. Annotating factuality in the TAGFACT corpus. Comares.
- Velupillai, Sumithra. 2011. Automatic classification of factuality levels. A case study on Swedish diagnoses and the impact of local context. In Anne Moen, Stig Kjaer Andersen, Jos Aarts and Petter Hurlen eds. *User Centred Networked Health Care Proceedings of the European Federation of Medical Informatics*. Amsterdam: IOS Press, 559–563.
- Weisser, Martin. 2016. DART – The dialogue annotation and research tool. *Corpus Linguistics and Linguistic Theory* 12/2: 355–388.
- Wonsever, Dina, Aiala Rosá and Marisa Malcuori. 2016. Factuality annotation and learning in Spanish texts. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Sara Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Helene Mazo, Asunción Moreno, Jan Odijk and Stelios Piperidis eds. *Proceedings of the Tenth Conference on Language Resources and Evaluation*. Portoroz: European Language Resources Association, 2076–2080.

*Corresponding author*

Ana Fernández-Montraveta  
Autonomous University of Barcelona  
Facultat de Filosofia i Lletres  
08193 Bellaterra  
Cerdanyola del Vallès · Barcelona  
Spain  
e-mail: Ana.Fernandez@uab.cat

received: January 2020

accepted: April 2020