

Boosting Eventus ID through Natural Language Processing

Pedro Armengol

2/6/2018

Abstract

The Boosting Eventus ID project is a Natural Language Processing Pipeline designed to apply Name Entity Recognition (NER) to news texts in Spanish. The software uses label data elaborated by the Researcher Javier Osorio to fit neuronal models (using the SpaCy package) to predict 7 categories (Target, source, location and so on). The accuracy of the model, defined as exact match of class as well as word start and end location, is around 14% for all the classes and around 45% for the location class. Further extensions of this research will include Entity Relationship Extraction from the news.

Introduction

More than 2 million articles are published every year on the web. Those articles possess a universe of information about facts and trends, many times, not available in structured datasets (that are prone to be used for statistical analysis). Also, those articles are in thousands of languages, formats and styles making a very difficult task the standardization of their content into structured information. However, through natural language processing and statistical learning techniques, I used the features of each word to predict their class (target, source, location and so on): where every class could be composed of one or more words.

One first attempt to systematize the information in news was the creation of the Conflict and Mediation Events Observations (CAMEO) catalog. CAMEO is a coding scheme that catalogs authors and events, specialized for automatized coding and for the register of sub-national authors (previous coding schemes were capable just to deal with national level authors (they were unable to catalog asymmetrical warfare for example). The cited catalog has been used in several projects including the Integrated Conflict Early Warning System a DARPA – Lockheed Martin Project joint venture project.

One step further, the Global Database of Events, Language and Tone (GDELT) is an effort to systematize the global media through the implementation of worldwide collector of news in more than 28 languages and machine learning techniques to systematize their information using the CAMEO labels. The GDELT website (now powered by Jigsaw - Google) claim to have systematized more than 500 million news since start operating.

However, for some languages, including Spanish, the GDELT database seems to perform particularly poorly (Osorio 2013). This is the reason that led the researcher Javier Osorio to develop techniques that could systematize news texts in Spanish. Particularly, Osorio (2018) did an extraction of classes based on exact matches, looking for concepts related to drug trafficking in Mexican Army reports. Despite of a high level of accuracy reached in the exercise, around 70 percent, there are concerns that this matching rules could overfit models for the particular task such that in other contexts they perform poorly.

As a contribution in the same direction, this project uses statistical learning techniques, particularly neuronal networks implemented with SpaCy, to classify the CAMEO classes that were labeled by the Osorio's research team.

Data

This research uses 724 news texts in Spanish coming from the Linguistic Data Consortium (LDC) Gigaword project. Each news text was labeled with the CAMEO classifications (particularly from the Protest sub-

Input 1

```
[In [9]: training_data[0][0]
Out[9]: 'HUELGAS EN ECUADOR\n\nQUITO, Jul 12\n\nEste martes fue un día marcado por las huelgas en Ecuador, donde los hospitales estatales ecuatorianos paralizaron sus actividades, los maestros anunciaron un paro de dos días para la próxima semana y en la población de La Concordia, en el litoral, sus habitantes dejaron de trabajar durante dos días.\nLos cambios que se generan en el Ecuador, con el inicio de las privatizaciones, la reducción de la burocracia y una cerrada ofensiva a la inflación a costa de una supuesta recesión, han generado descontento ciudadano reflejado por las huelgas hasta el punto que comentaristas locales subrayaron la "extraña" inexistencia de protestas en las últimas dos semanas.\nLos empleados de los hospitales estatales reclaman los ofrecimientos de mejoras económicas y sus condiciones de trabajo, una vez que desde principios de mes se inició un sistema de cobro a los usuarios. La atención era antes gratuita, aunque los pacientes debían adquirir medicamentos cuando faltaban en las casas de salud estatales.\nLos profesores acordaron una huelga el martes y miércoles de la próxima semana, como una acción preventiva para reclamar el pago de los atrasos salariales, lo que fue rechazado por la ministra del ramo, Rosalía Arteaga, quien sostuvo que se está pagando dentro de los diez primeros días del mes subsiguiente.\nLos maestros paralizaron sus actividades en noviembre y diciembre pasados, causando estragos en la educación estatal.\nLa ministra Arteaga, por otro lado, apoya la oposición de los maestros a una ley que se tramita en el Congreso para incluir la educación religiosa en los establecimientos estatales, si la mayoría de padres lo decide, pero el Gobierno apoya decididamente esta iniciativa eclesiástica.\nAnte los reclamos de los maestros y sectores políticos que aseguran que el restablecer la educación religiosa sería un retroceso de un siglo y se contradice con la libertad de culto, la iglesia pidió la suspensión del trámite legislativo o del proyecto hasta que se concrete un debate nacional más profundo.\nMientras tanto, en La Concordia, cuyos pobladores piden la cantonización –una nominación de la división geopolítica que le aseguraría recursos estatales–, las autoridades militarizaron la zona para evitar eventuales desmanes.\n'
```

Input 2

```
[In [10]: training_data[0][1]
Out[10]:
{'entities': [(0, 7, 'material_conflict'),
(11, 18, 'location'),
(20, 25, 'location'),
(75, 82, 'material_conflict'),
(86, 93, 'location'),
(101, 125, 'source'),
(139, 166, 'material_conflict'),
(172, 180, 'source'),
(181, 199, 'verbal_conflict'),
(256, 268, 'location'),
(276, 283, 'location'),
(289, 299, 'source'),
(300, 319, 'material_conflict'),
(700, 741, 'source'),
(742, 750, 'material_conflict'),
(1034, 1048, 'source'),
(1049, 1069, 'verbal_conflict'),
(1215, 1235, 'source'),
(1237, 1252, 'source'),
(1197, 1210, 'verbal_conflict'),
(1343, 1355, 'source'),
(1356, 1383, 'material_conflict'),
(1461, 1480, 'source'),
(1497, 1515, 'verbal_cooperation'),
(1519, 1531, 'target'),
(1770, 1778, 'source'),
(1754, 1762, 'verbal_conflict'),
(1781, 1799, 'source'),
(1804, 1812, 'verbal_cooperation'),
(1820, 1844, 'verbal_cooperation'),
(1928, 1938, 'source'),
(1939, 1958, 'verbal_cooperation'),
(1963, 1995, 'verbal_cooperation'),
(2070, 2082, 'location'),
(533, 554, 'material_conflict'),
(573, 580, 'material_conflict'),
(1146, 1188, 'verbal_conflict'),
(2210, 2221, 'source'),
(2222, 2243, 'material_conflict'),
(2090, 2100, 'target')]]
```

Figure 1: Input: tuples structure

category). The categories are “location”, “material cooperation”, “material cooperation”, “source”, “target”, “verbal cooperation” and “verbal conflict”. There are around 19 labels for each news with a standard deviation of around 16 labels (there is an important variation in the number of labels for each text).

In **Figure 1** depicts the inputs of the prediction model fitted in the next section. Basically, each text is extracted and stored as the first position of a tuple where the second position of the tuple are the labels (with class and starting/ending position of the word or words that compose the class).

Pipeline

The pipeline or process of prediction is composed of four components:

- Preprocessing
- Splitting
- Training
- Test

Preprocessing basically is divided in two functions: one takes the texts with labels provided by Osorio and creates tuples with them (one tuple for each news text) and the second takes the first element of the tuple and processes the text to create first sentences, then tokens (one for every word) and then assigns Part-of-speech (POS), entity recognition (ER) and dependency parsing (DP) (see **Figure 2**).

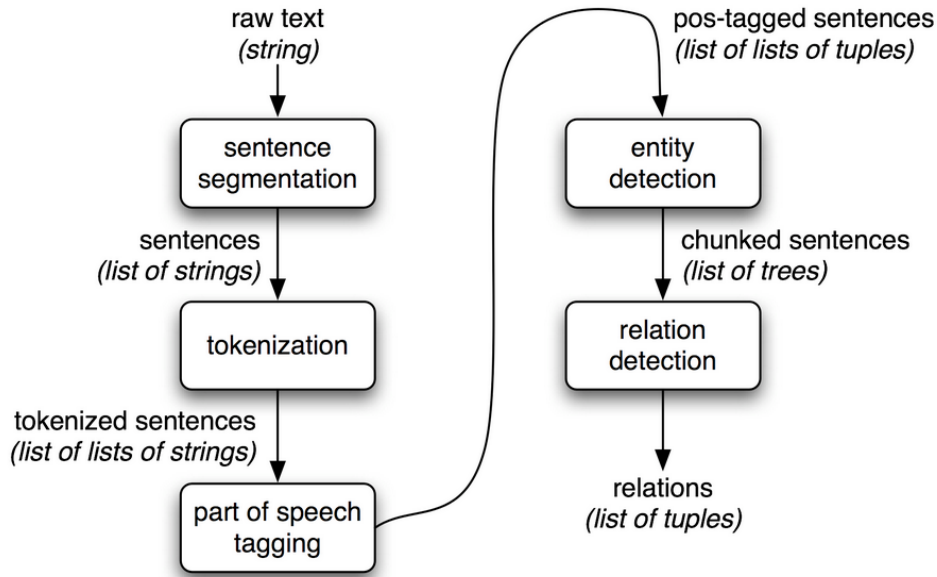


Figure 2: Preprocessing: second function funnel - source: NLTK

Spanish

MODEL		SPACY	TYPE	UAS	NER F	POS	WPS	SIZE
es_core_news_sm	2.0.0	2.x	neural	89.8	88.7	96.9	<i>n/a</i>	35MB

Figure 3: Preprocessing: accuracy, SpaCy tokenization - source: SpaCy

The implemented pipeline uses the SpaCy package for the second function of preprocessing. SpaCy, uses annotated corpuses in Spanish to predict each one of the token values (POS, ER or DP). According to the results reported in their website, the accuracy rates of those predictions are above 89 % for each category (see *Figure3**).