

# MinaR los discursos pResidenciales

Juan Pablo Ruiz Nicolini , Camila Higa , Lucas Enrich

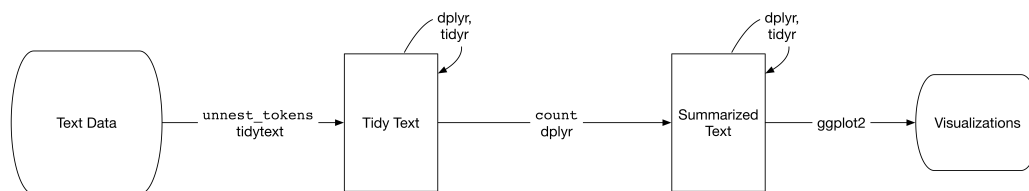
**Palabras clave:** discurso - text minning - política

## Abstract

El primero de marzo de cada año las cámaras de Diputados y Senadores de la Nación Argentina se reúnen en asamblea para dar comienzo al año legislativo y el presidente de turno encabeza el acto con un discurso<sup>1</sup>. Éstos suelen girar en torno a los ejes de gobierno o promesas y objetivos del año. Es notorio que estos mensajes tiene un estilo y contenido marcado por quien ejerce el gobierno. En este trabajo analizamos el texto contenido en los discursos presidenciales desde el primero en 1854 por Justo José De Urquiza hasta el último en 2020 por Alberto Fernández.

La *minería de texto* como estrategia de investigación es de utilidad para un rápido y eficiente análisis exploratorio del gran volumen de información contenida en los discursos presidenciales. Dentro del ecosistema de R este campo ha ido creciendo sostenidamente. Liberías como **tm** y **topicmodels** son herramientas poderosas para el procesamiento, manipulación y modelado de la información contenida en el texto. Siguiendo la filosofía de **tidyverse**, Silge y Robinson (2016) desarrollaron **tidytext** que facilita una primera introducción a esta técnica de investigación y su integración con otras como **ggplot2** para la visualización.

Un flujo de trabajo como el descripto anteriormente puede ilustrarse siguiendo el esquema propuesto por Silge y Robinson (2020):



1. Se encuentran digitalizados los 114 discursos emitidos por los 31 presidentes que dieron lugar a la apertura de las sesiones legislativas. Debe mencionarse que no hay un discurso por año debido, principalmente, a las interrupciones institucionales cuando el congreso no sesionó. Entre todos suman alrededor de 1,358,792 palabras con un promedio de 11,919 y picos mínimo de 258 (Hipolito Yrigoyen 1917) y máximo de 44,415 (Ramon Castillo en 1942).
2. Con esa información construimos una única base de datos siguiendo el principio *datos de texto ordenados* (*tidy text*) propuesto por Silge y Robinson (2016) como extensión de los *datos ordenados* (*tidy*) de Wickham (2014):
  - Cada variable debe tener su propia columna.
  - Cada observación debe tener su propia fila.
  - Cada valor debe tener su propia celda.

Silge y Robinson (2016) definen entonces a los *datos de texto ordenados* cuando están en una tabla compuesta por “un token por fila”. Un token es una unidad de texto significativa, como una palabra (o un bigrama), que estamos interesados en usar para el análisis, y la tokenización es el proceso de dividir el texto en tokens<sup>2</sup>.

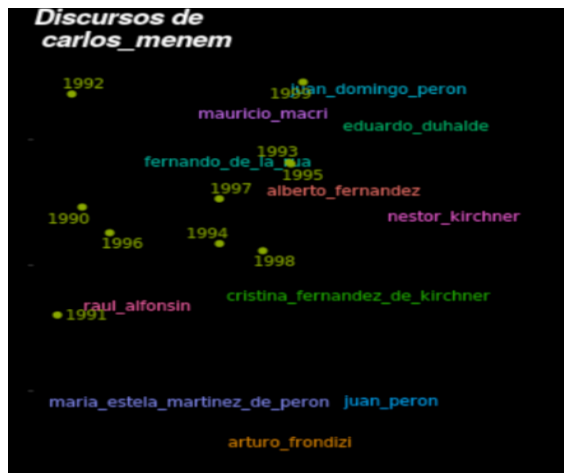
3. Trabajamos con **dplyr** para calcular frecuencias de palabras, **tidytext** para identificar las más relevantes comparadas entre discursos (*tf-idf*) y **ggplot2** para las visualizaciones.

<sup>1</sup>La fuente original de todos los discursos puede consultarse en línea en [https://www.hcdn.gob.ar/secparl/dgral\\_info\\_parlamentaria/dip/documentos/mensajes\\_presidenciales.html](https://www.hcdn.gob.ar/secparl/dgral_info_parlamentaria/dip/documentos/mensajes_presidenciales.html). Los mismos fueron posteriormente digitalizados mediante proceso de OCR y se encuentran disponibles para descargar desde R a través del paquete {polAr} (Ruiz Nicolini 2020).

<sup>2</sup>Traducción propia de *The tidy text format* (Silge and Robinson 2016).

### Ejemplo: Trayectoria en el discurso

Mediante las técnicas *TF-IDF* y *Principal Component Analysis (PCA)* es posible embeber numéricamente los discursos considerando el uso de palabras en cada uno ponderados por su longitud y así visualizar los presidentes de mayor y menor variabilidad discursiva en comparación con los promedios de los demás.



Las limitaciones de técnicas basadas en la frecuencia de palabras independientemente del orden (como son *Bag of Words* y *TF-IDF*) están vinculadas con el vocabulario por un lado, y con la semántica por el otro. Así, para trabajar mejor con textos históricos (el más antiguo en este caso tiene 166 años) es recomendable usar técnicas que hagan uso del contexto como puede ser *Doc2Vec*.

### Referencias

- Ruiz Nicolini, Juan Pablo. 2020. "PolAr: Argentina Political Analysis." <https://github.com/electorArg/polAr>.
- Silge, Julia, and David Robinson. 2016. "Tidyttext: Text Mining and Analysis Using Tidy Data Principles in R." *Journal of Open Source Software* 1 (3): 37. <https://doi.org/10.21105/joss.00037>.
- . 2020. *Text Mining with R*. O'Reilly. <https://www.tidyttextmining.com/>.
- Wickham, Hadley. 2014. "Tidy Data." *Journal of Statistical Software, Articles* 59 (10): 1–23. <https://doi.org/10.18637/jss.v059.i10>.

Juan Pablo Ruiz Nicolini  
Universidad Torcuato Di Tella  
[juan.ruiznicolini@mail.utdt.edu](mailto:juan.ruiznicolini@mail.utdt.edu)

Camila Higa  
menta Comunicación  
[chiga1226@gmail.com](mailto:chiga1226@gmail.com)

Lucas Enrich  
Universidad Nacional de La Matanza  
[lucas.a.enrich@gmail.com](mailto:lucas.a.enrich@gmail.com)