

---

---

# Predicting persuadable individuals using text mining and graph machine learning on Twitter

---

---

— Federico Albanese —

(falbanese@dc.uba.ar)

 @f\_albanese

---

---



Esteban Feuerstein<sup>1</sup>



Leandro Lombardi<sup>2</sup>



Pablo Balenzuela<sup>3</sup>

*"En el panóptico digital no existe ese Big Brother que nos extrae informaciones contra nuestra voluntad. Por el contrario, nos revelamos, incluso nos ponemos al desnudo por iniciativa propia."*

***Byung-Chul Han***

"En el pa  
ex  
contra

**DATA**



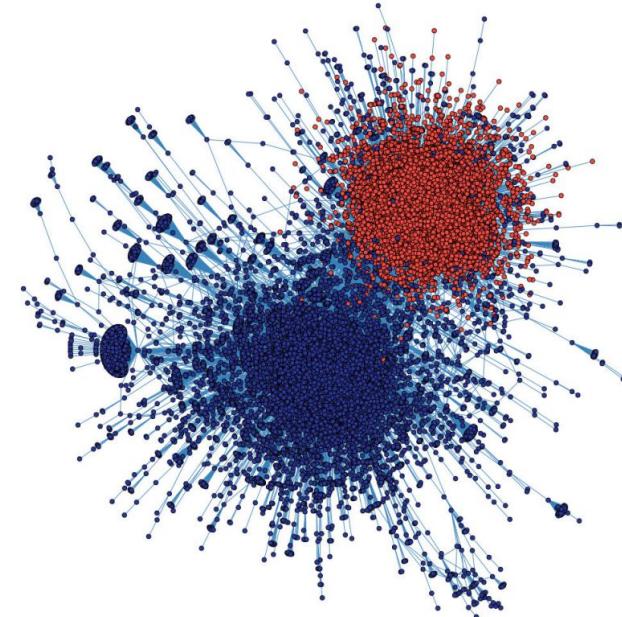
que nos  
. Por el  
desnudo  
propia."

*Chul Han*

# Trabajos previos

# Trabajos previos

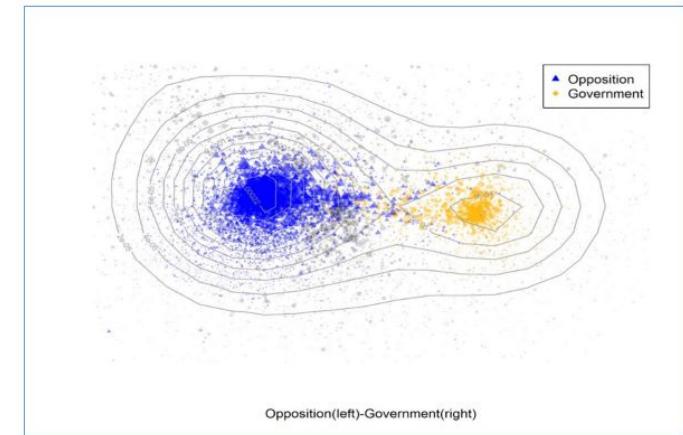
- Conover et al. (2011): Usan detección de comunidades y muestran la estructura polarizada durante las elecciones del 2010 en Estados Unidos.



Red de Retweets (RT):  
En azul Demócratas, en rojo Republicanos.

# Trabajos previos

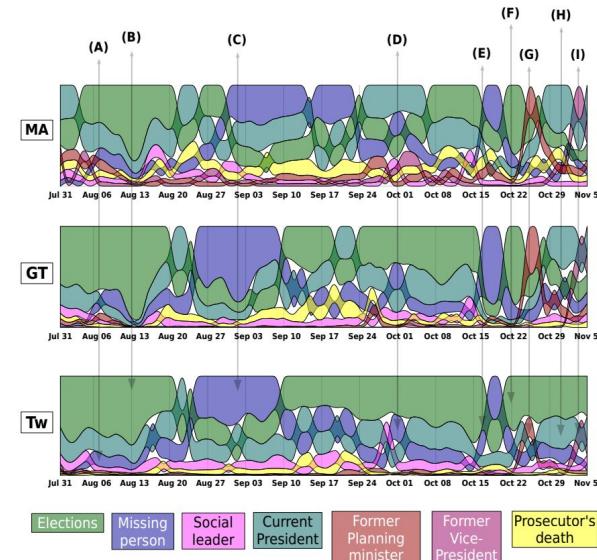
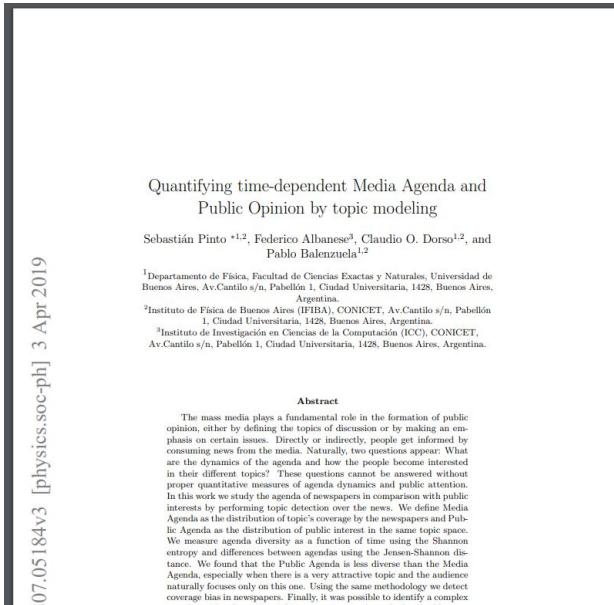
- Aruguete et al. (2018): muestran cómo los usuarios “like-minded” se agrupan e interactúan entre sí en la red de RT formando comunidades cerradas.



Red de RT:  
**Azul:** oposición; **Amarillo:** oficialismo; **Gris:** otros.

# Trabajos previos

- Pinto et al. (2019): Detectan y analizan los tópicos de discusión para entender el rol de los medios de comunicación en el proceso de formación de opinión pública.



Bump Graph de los temas:  
MA: Medios, GT: Google Trends, Tw: twitter

# ¿Qué hicimos en este trabajo?

- Armamos un Machine Learning Framework

# ¿Qué hicimos en este trabajo?

- Armamos un Machine Learning Framework
  - Objetivo  
Detectar *cambios de agrupamiento político* de las personas

# ¿Qué hicimos en este trabajo?

- Armamos un Machine Learning Framework
  - Objetivo  
Detectar ***cambios de agrupamiento político*** de las personas
  - Herramientas
    - Natural Language Processing techniques.
    - Graph Machine Learning.
    - Gradient Boosting Models.
  - Datos
    - Twitter

# ¿Qué hicimos en este trabajo?

- Armamos un Machine Learning Framework
  - Objetivo  
Detectar cambios de agrupamiento político de las personas
  - Herramientas
    - Natural Language Processing techniques.
    - Graph Machine Learning.
    - Gradient Boosting Models.
  - Datos
    - Twitter
- Evaluamos el Framework

# ¿Qué hicimos en este trabajo?

- Armamos un Machine Learning Framework
  - Objetivo  
Detectar cambios de agrupamiento político de las personas
  - Herramientas
    - Natural Language Processing techniques.
    - Graph Machine Learning.
    - Gradient Boosting Models.
  - Datos
    - Twitter
- Evaluamos el Framework
- Caracterizamos a los usuarios que cambian el agrupamiento político.
  - Análisis de feature importance.

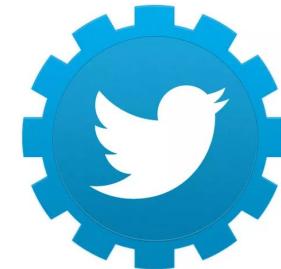
# ¿Qué hicimos en este trabajo?

- Armamos un Machine Learning Framework
  - Objetivo  
Detectar cambios de agrupamiento político de las personas
  - Herramientas
    - Natural Language Processing techniques.
    - Graph Machine Learning.
    - Gradient Boosting Models.
  - Datos
    - Twitter
- Evaluamos el Framework
- Caracterizamos a los usuarios que cambian el agrupamiento político.
  - Análisis de feature importance.
- Identificación de los temas de discusión más importantes y persuasivos.

# Data Collection

# Data Collection

- ¿Dónde?  
Argentina.
- ¿En qué idioma?  
Español.
- ¿Cuándo?  
Durante la semana de las elecciones.
- ¿Queries?  
Nombres de candidatos, partidos y políticos.

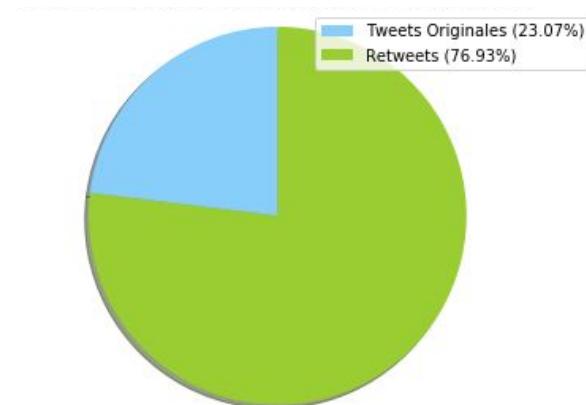


Streaming API

# Datos

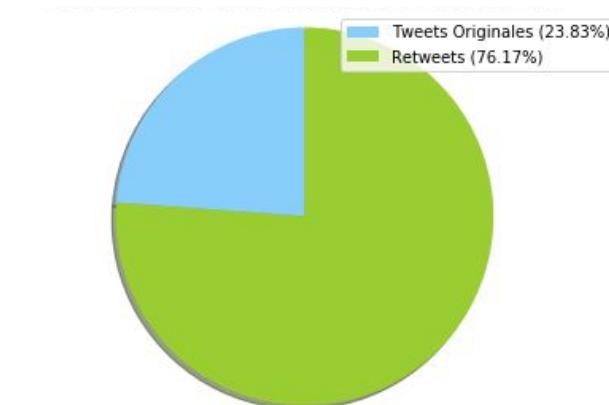
## Elecciones PASO 2017:

- 2.117.708 tweets
- 1.629.200 RT
- 86.361 usuarios



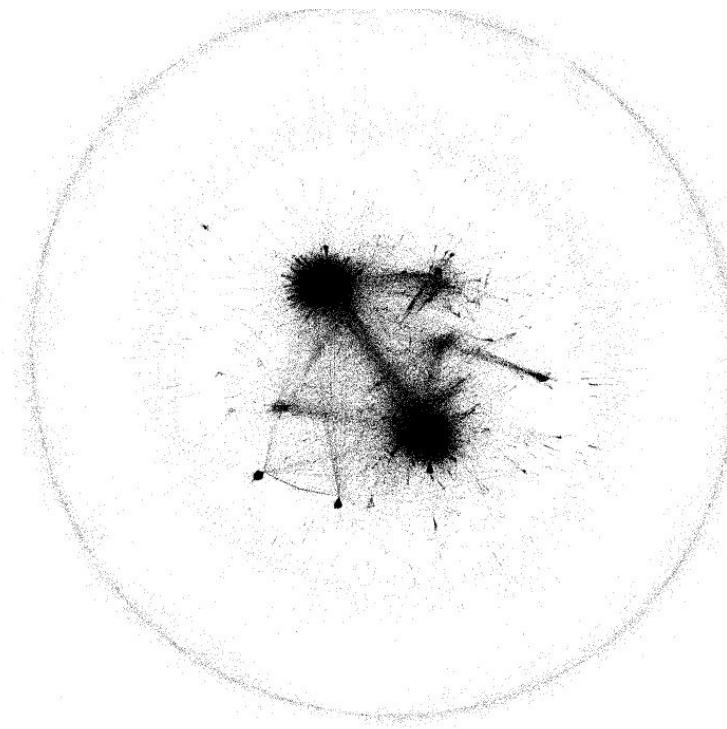
## Elecciones Generales 2017:

- 1.751.813 tweets
- 1.334.294 RT
- 298.866 usuarios



# Análisis de los datos

# Red de RT

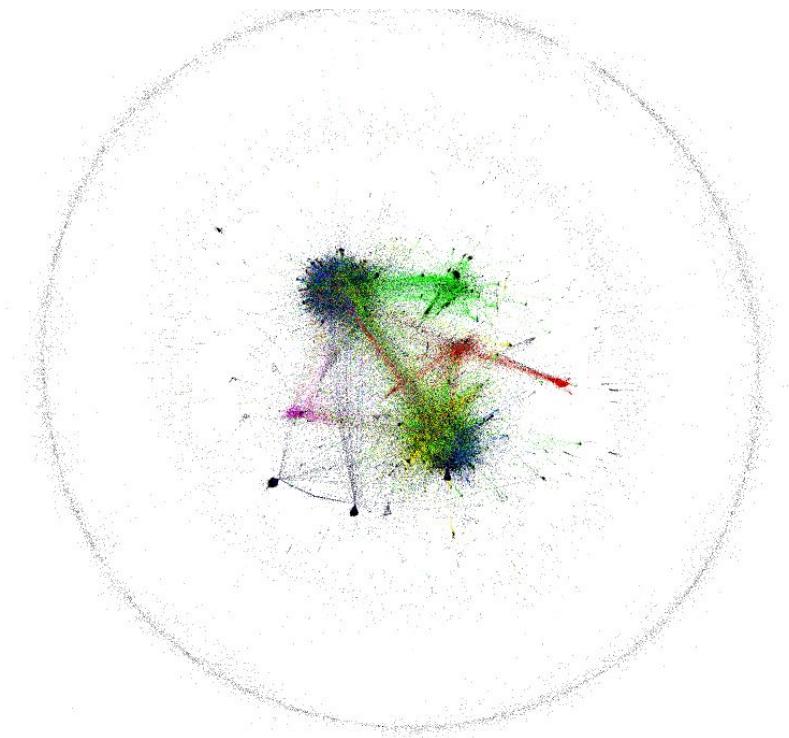


Visualización hecha con *Force Atlas 2* con los datos de las elecciones PASO 2017.

# Red de RT

¿De quién hablan?

- *Unidad Ciudadana*
- *Cambemos*
- *Frente Justicialista*
- *1 País*
- *Macri*

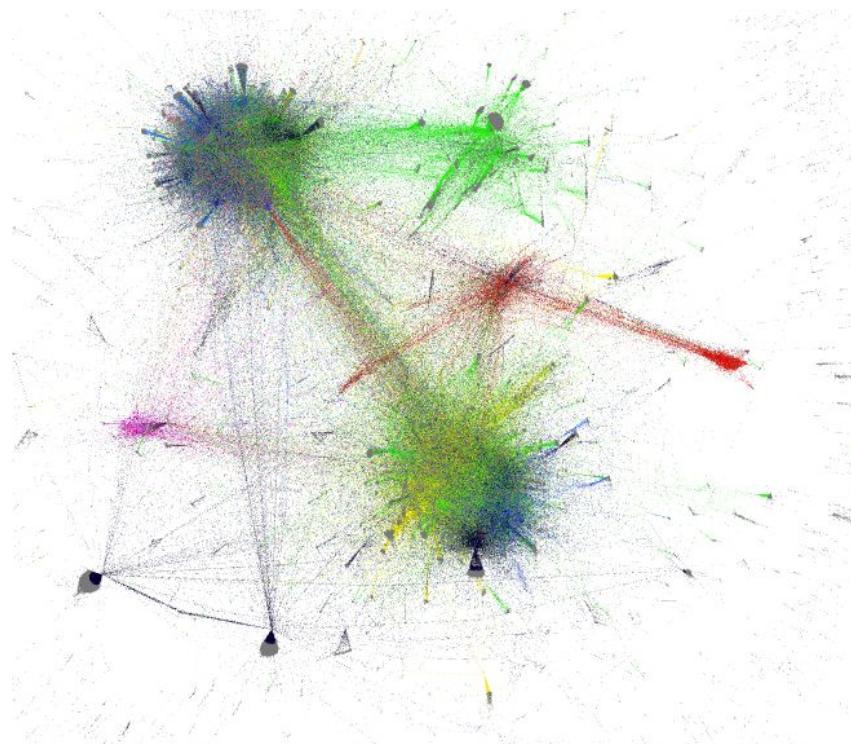


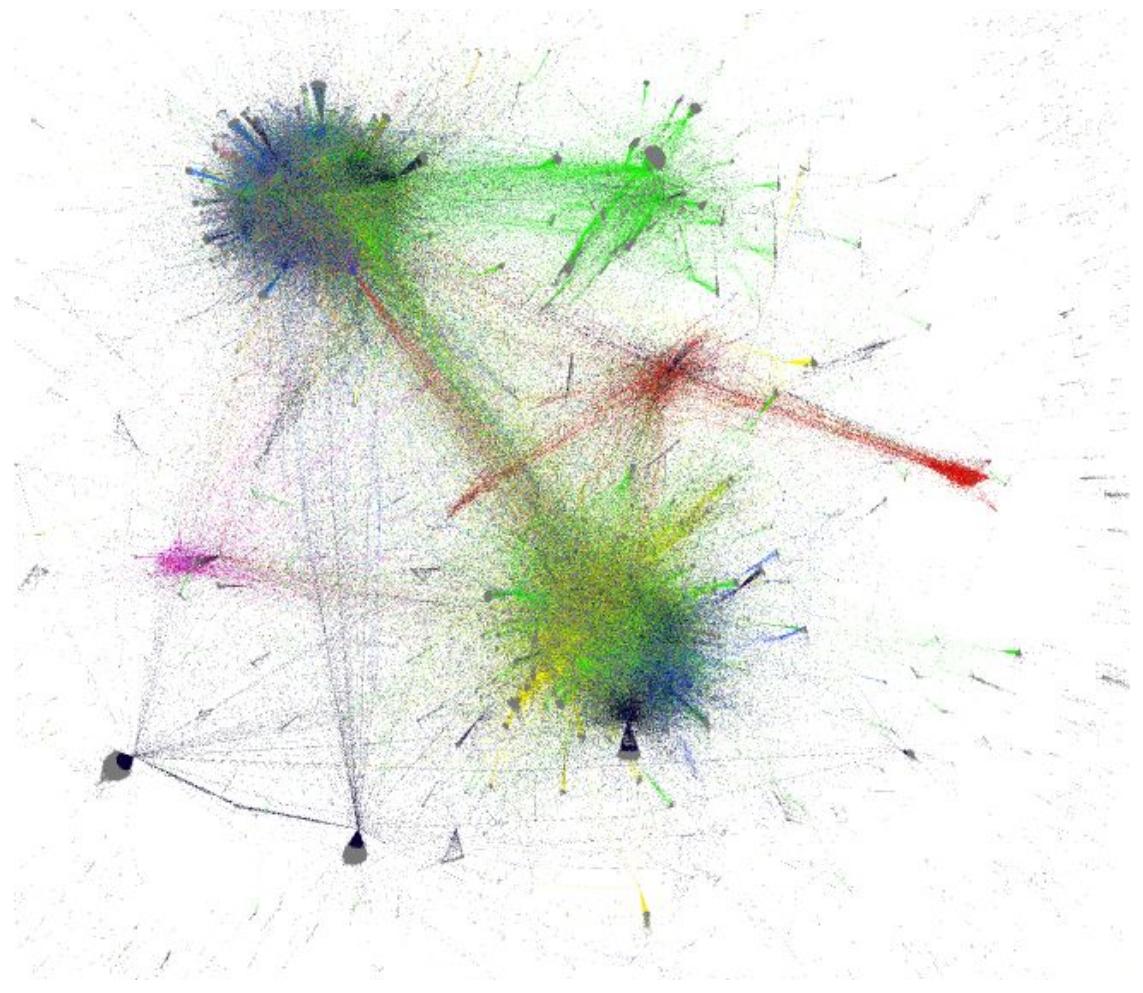
Visualización hecha con *Force Atlas 2* con los datos de las elecciones PASO 2017.

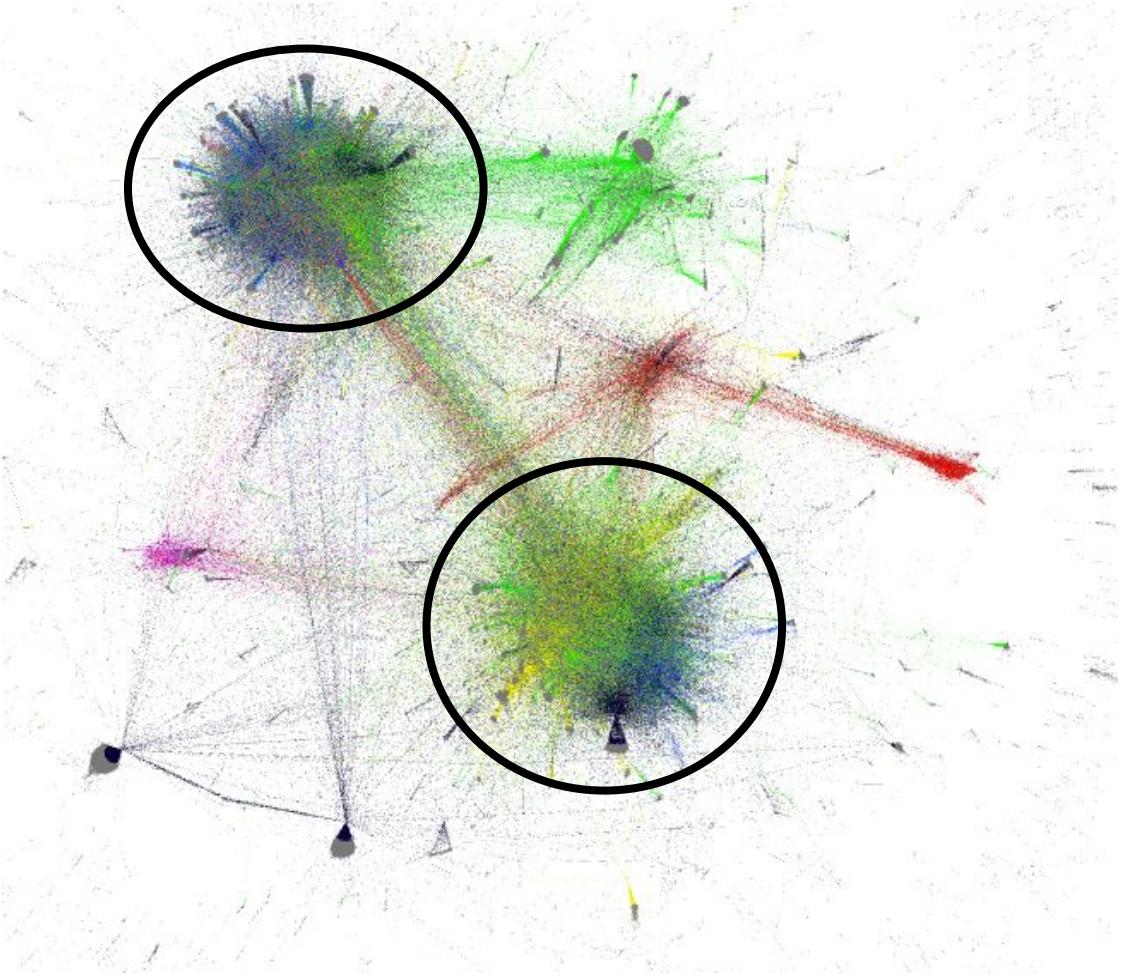
# Red de RT

¿De quién hablan?

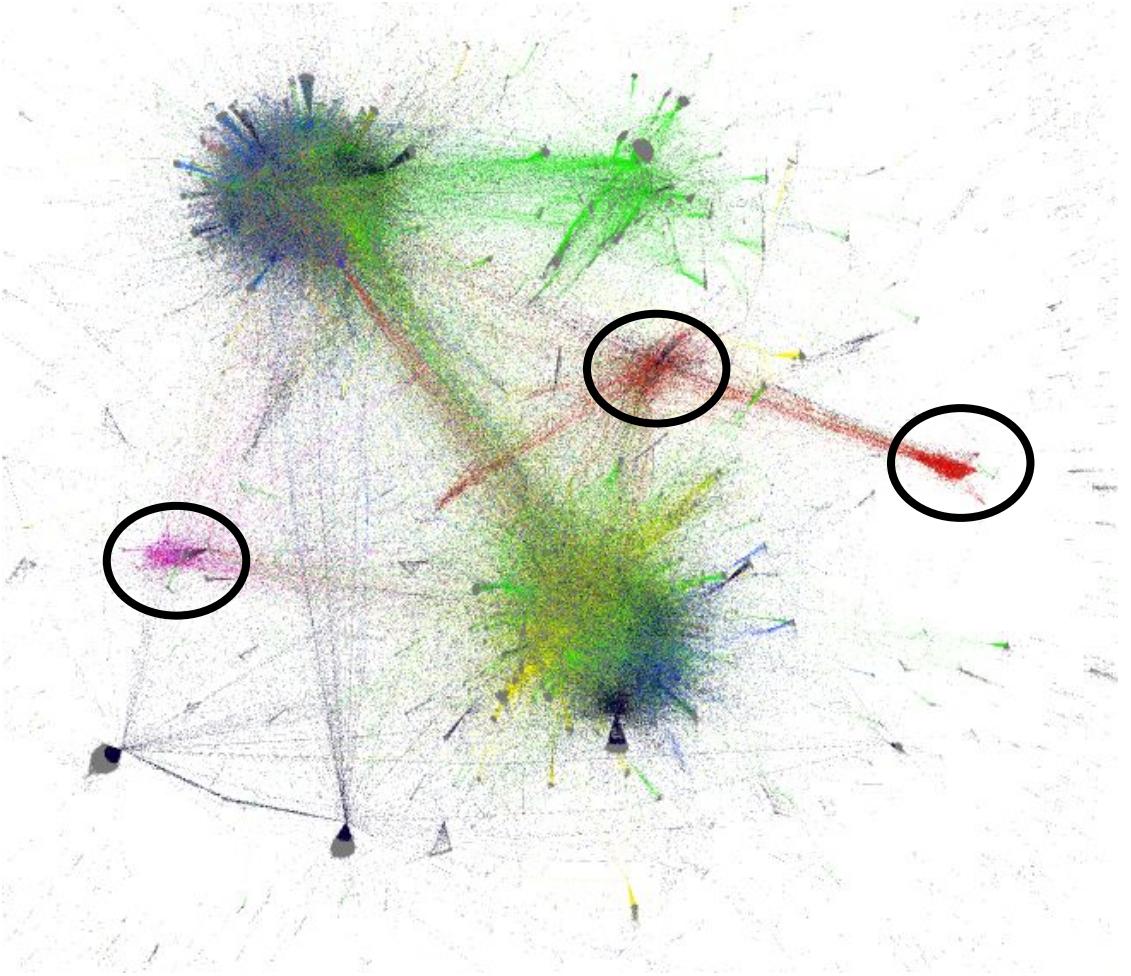
- *Unidad Ciudadana*
- *Cambios*
- *Frente Justicialista*
- *1 País*
- *Macri*



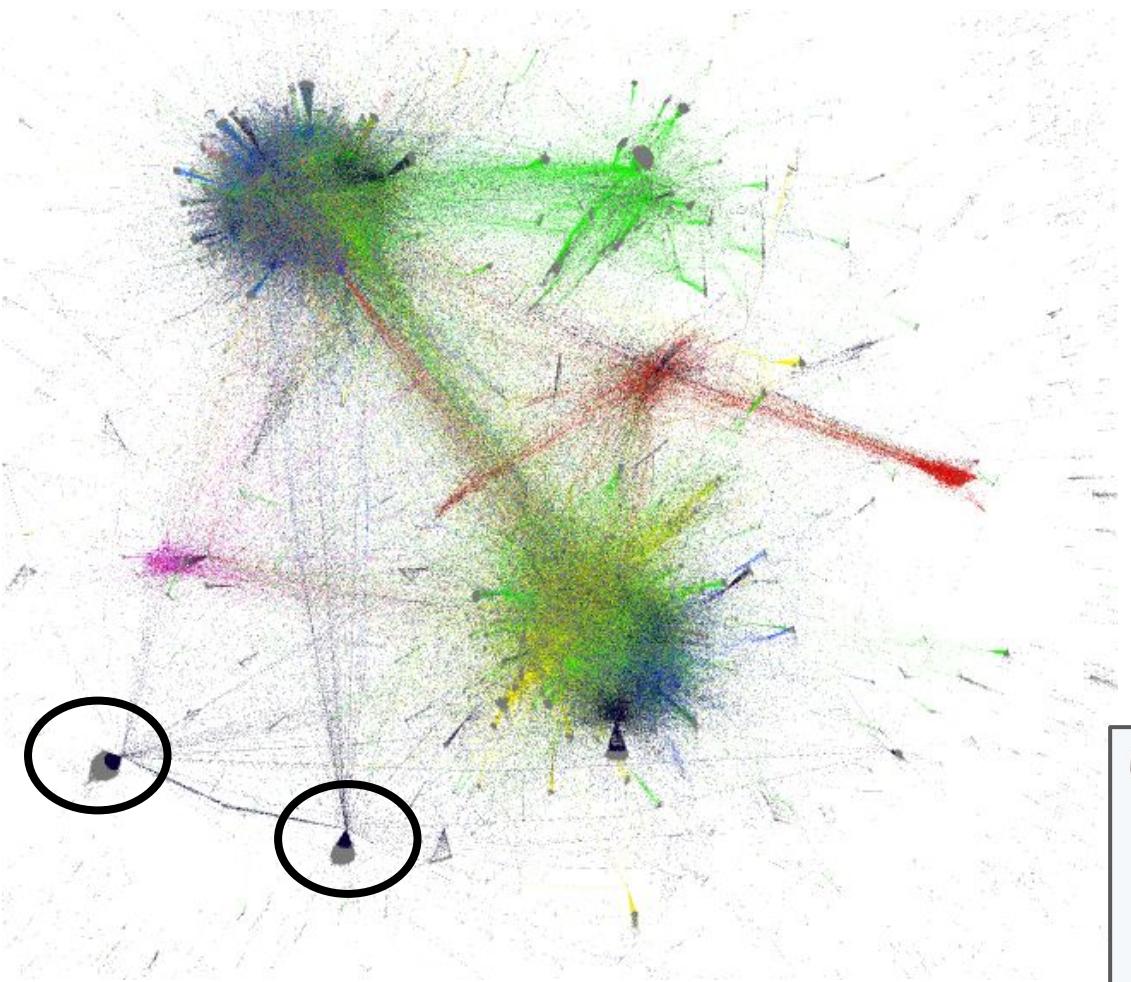




- Dos comunidades principales:  
***"la grieta"***



- Dos comunidades principales:  
***"la grieta"***.
- Comunidades pequeñas de los otros partidos.

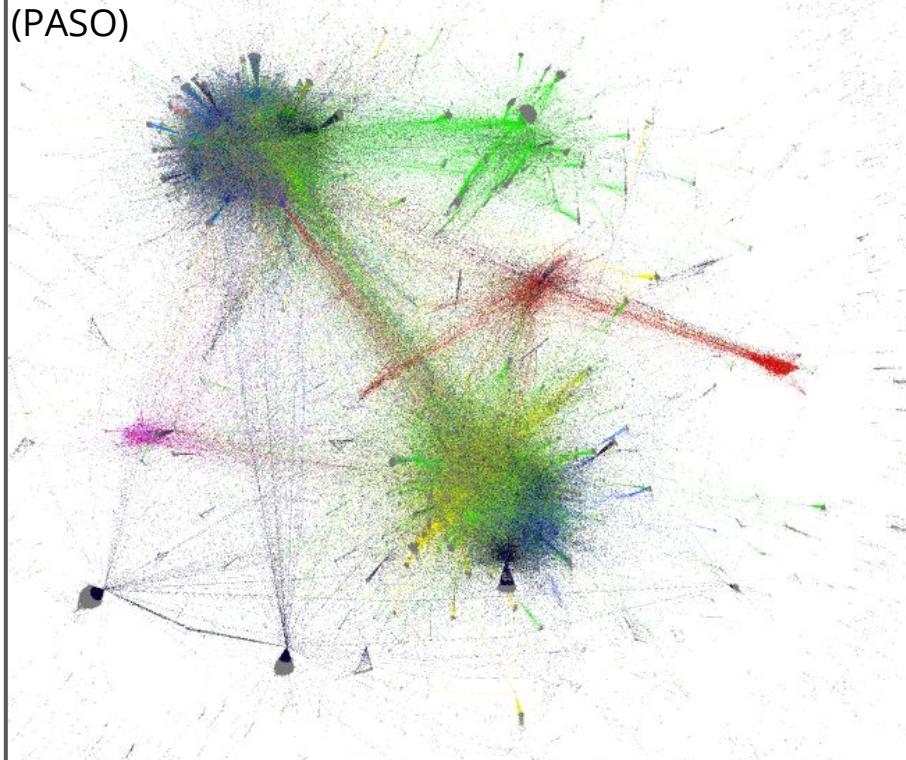


- Dos comunidades principales:  
***"la grieta"***.
- Comunidades pequeñas de los otros partidos.
- Famosos por un día.

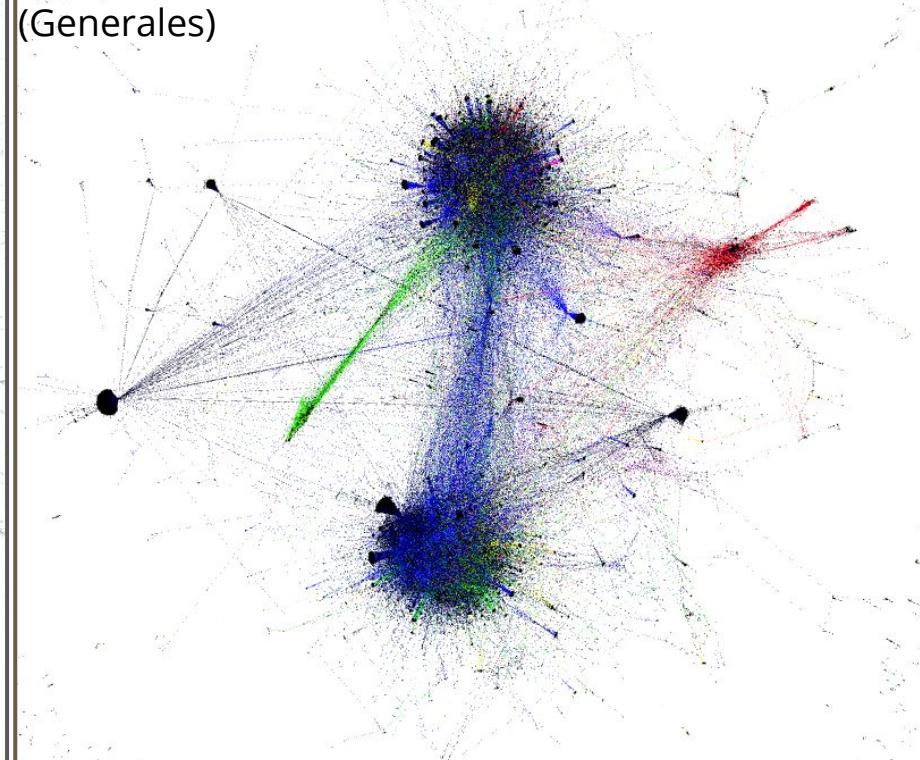


# Redes de RT

(PASO)

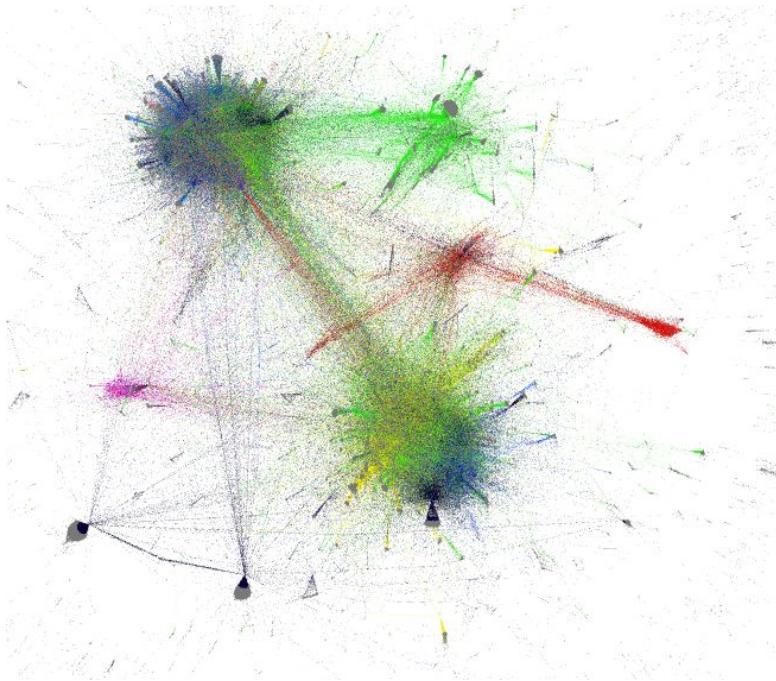


(Generales)

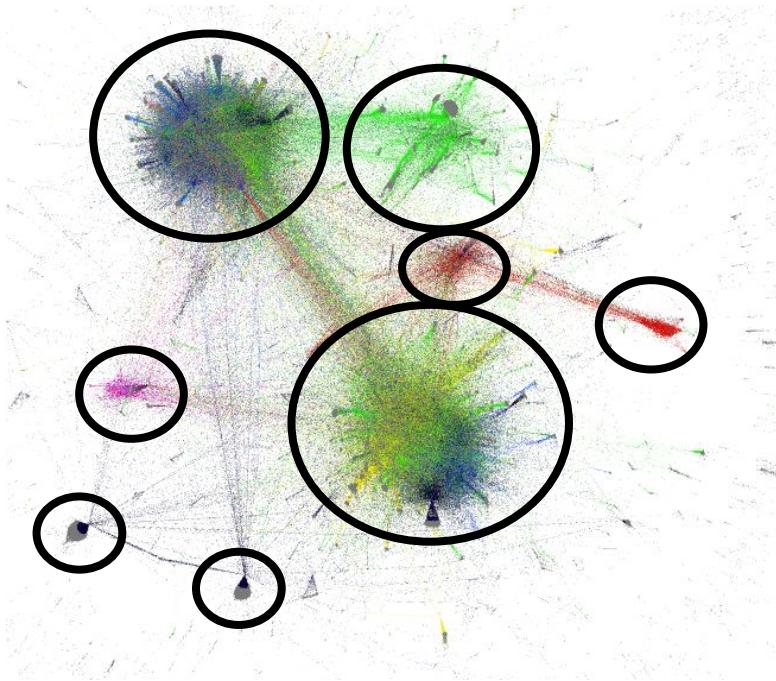


# Unsupervised Machine Learning: Clustering

# Unsupervised Machine Learning: Clustering

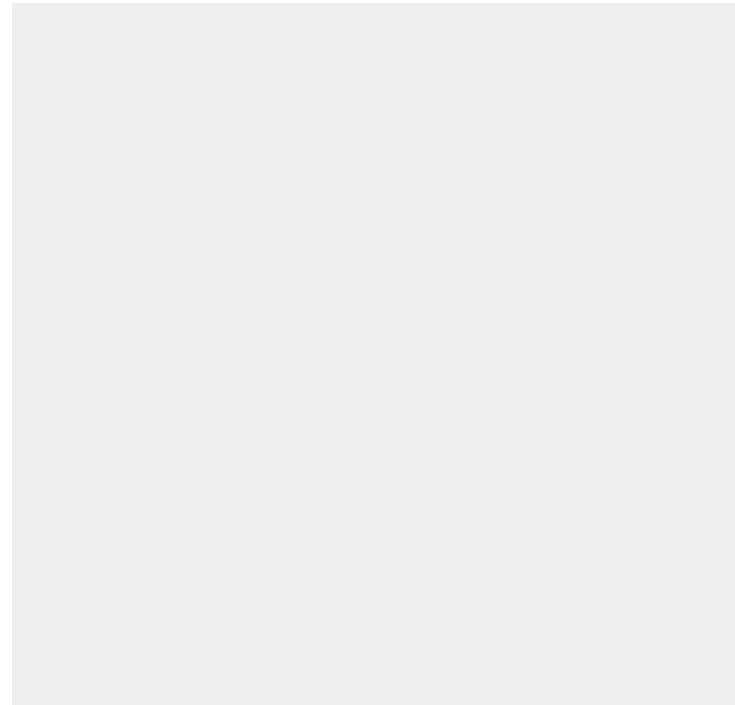


# Unsupervised Machine Learning: Clustering



# Método de Louvain

1. A todo nodo en la red se le asigna a una comunidad.



# Método de Louvain

1. A todo nodo en la red se le asigna a una comunidad.
2. Para cada nodo, se le cambia iterativamente la comunidad por la de un vecino. Uno se queda con el cambio que más aumenta la modularidad.

Modularidad:

$$Q = \frac{1}{2m} \sum_{i,j} [A_{i,j} - \frac{k_i k_j}{2m}] \delta(c_i, c_j)$$

Donde:

- $A$  es la matriz de adyacencia.
- $k_i$  es la suma de los pesos de las aristas del nodo  $i$
- $2m$  es la suma de los pesos de todas las aristas de la red.
- $c_i$  es la comunidad del nodo  $i$ .

# Método de Louvain

1. A todo nodo en la red se le asigna a una comunidad.
2. Para cada nodo, se le cambia iterativamente la comunidad por la de un vecino. Uno se queda con el cambio que más aumenta la modularidad.
3. Se repite el paso iterativo hasta llegar al máximo de modularidad.

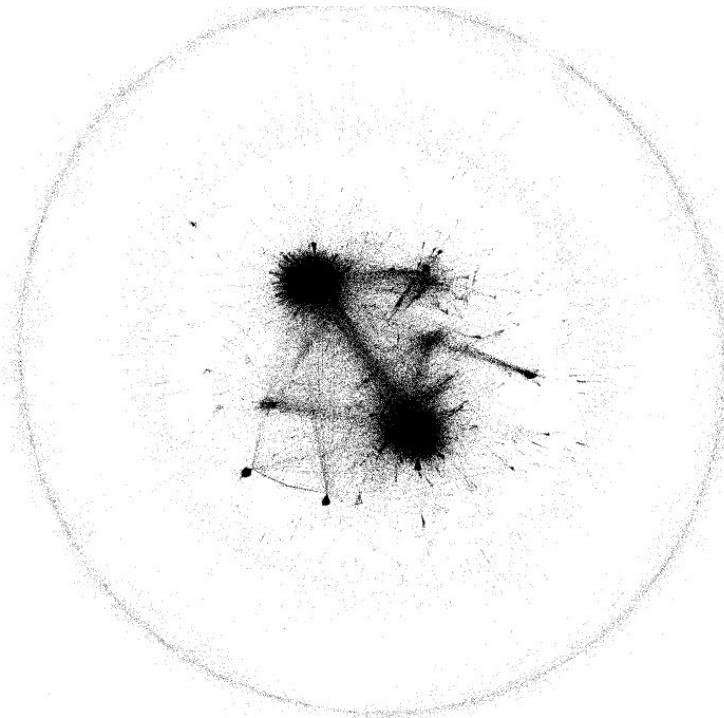
Modularidad:

$$Q = \frac{1}{2m} \sum_{i,j} [A_{i,j} - \frac{k_i k_j}{2m}] \delta(c_i, c_j)$$

Donde:

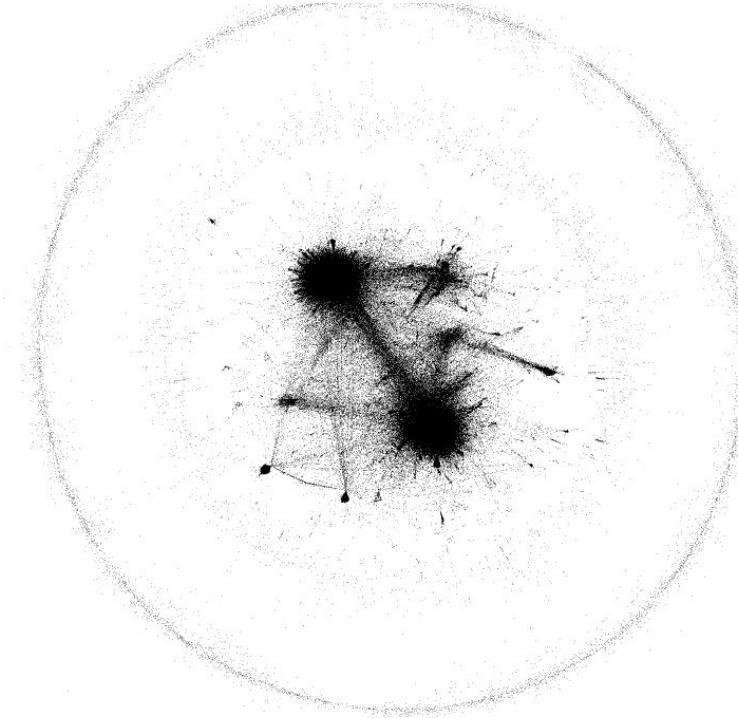
- $A$  es la matriz de adyacencia.
- $k_i$  es la suma de los pesos de las aristas del nodo  $i$
- $2m$  es la suma de los pesos de todas las aristas de la red.
- $c_i$  es la comunidad del nodo  $i$ .

# Graph Features



# Graph Features

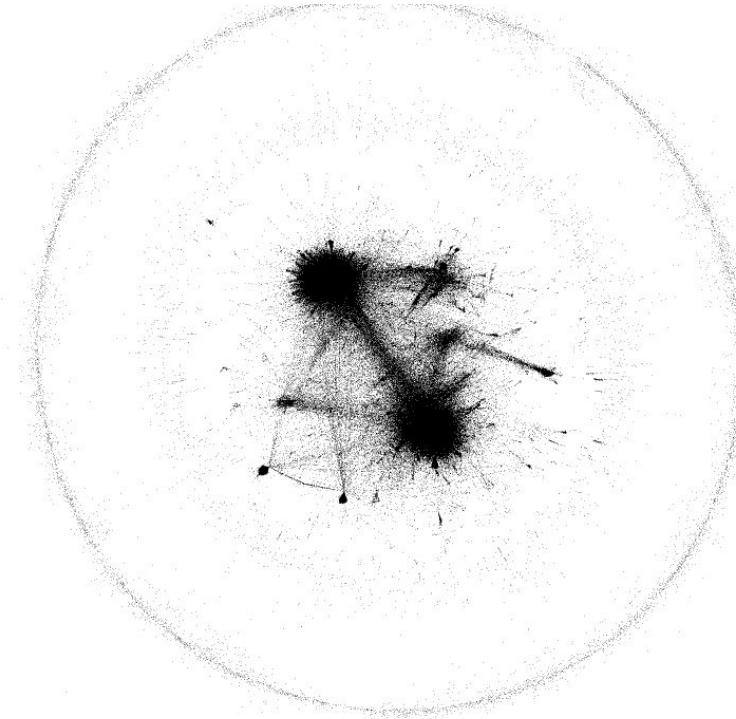
- Degree



- [1] Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab
- [2] Brandes, U. (2001). A faster algorithm for betweenness centrality. Journal of mathematical sociology, 25(2), 163-177.
- [3] Saramki, J., Kivel, M., Onnela, J. P., Kaski, K., & Kertesz, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. Physical Review E, 75(2), 027105.
- [4] Fast unfolding of communities in large networks, Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008 (12pp)
- [5] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. nature, 433(7028), 895.

# Graph Features

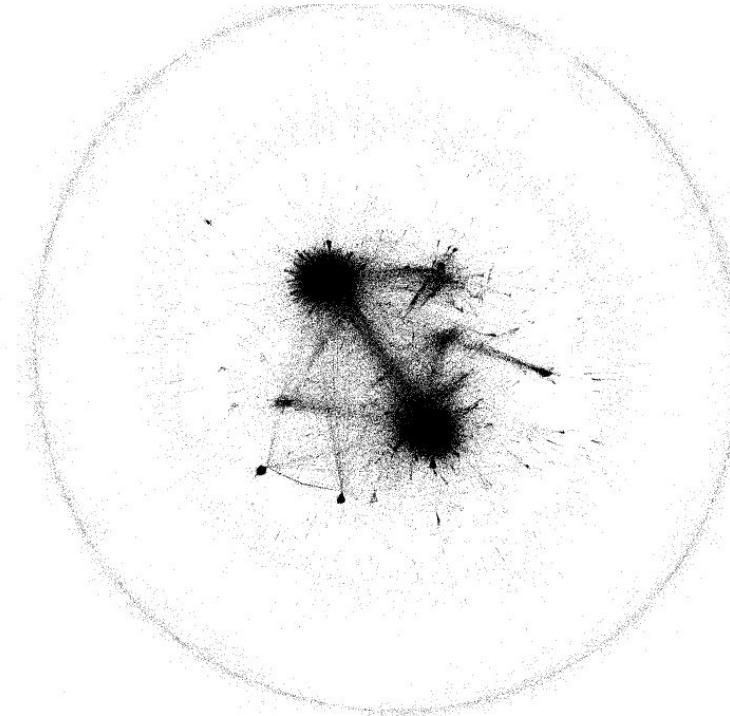
- Degree
- PageRank [1]



- [1] Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab
- [2] Brandes, U. (2001). A faster algorithm for betweenness centrality. Journal of mathematical sociology, 25(2), 163-177.
- [3] Saramki, J., Kivel, M., Onnela, J. P., Kaski, K., & Kertesz, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. Physical Review E, 75(2), 027105.
- [4] Fast unfolding of communities in large networks, Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008 (12pp)
- [5] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. nature, 433(7028), 895.

# Graph Features

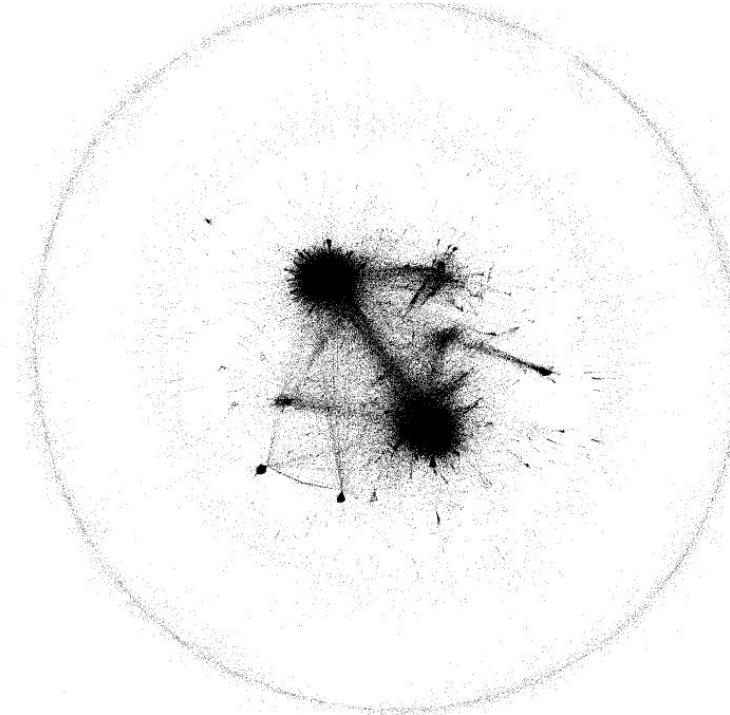
- Degree
- PageRank [1]
- Betweenness centrality [2]



- [1] Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab
- [2] Brandes, U. (2001). A faster algorithm for betweenness centrality. Journal of mathematical sociology, 25(2), 163-177.
- [3] Saramki, J., Kivel, M., Onnela, J. P., Kaski, K., & Kertesz, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. Physical Review E, 75(2), 027105.
- [4] Fast unfolding of communities in large networks, Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008 (12pp)
- [5] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. nature, 433(7028), 895.

# Graph Features

- Degree
- PageRank [1]
- Betweenness centrality [2]
- Clustering coefficient [3]



[1] Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab

[2] Brandes, U. (2001). A faster algorithm for betweenness centrality. Journal of mathematical sociology, 25(2), 163-177.

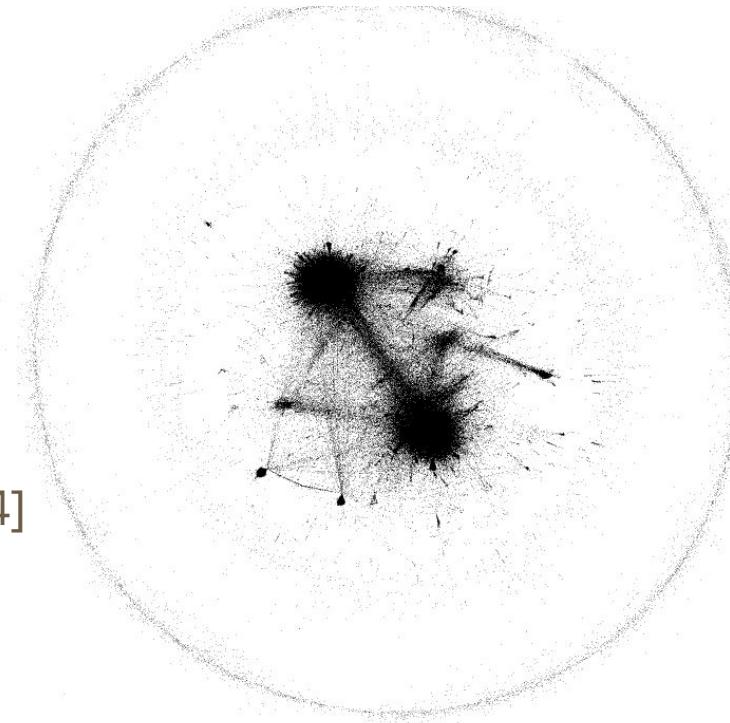
[3] Saramki, J., Kivel, M., Onnela, J. P., Kaski, K., & Kertesz, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. Physical Review E, 75(2), 027105.

[4] Fast unfolding of communities in large networks, Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008 (12pp)

[5] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. nature, 433(7028), 895.

# Graph Features

- Degree
- PageRank [1]
- Betweenness centrality [2]
- Clustering coefficient [3]
- Clustering affiliation (Louvain) [4]



[1] Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab

[2] Brandes, U. (2001). A faster algorithm for betweenness centrality. Journal of mathematical sociology, 25(2), 163-177.

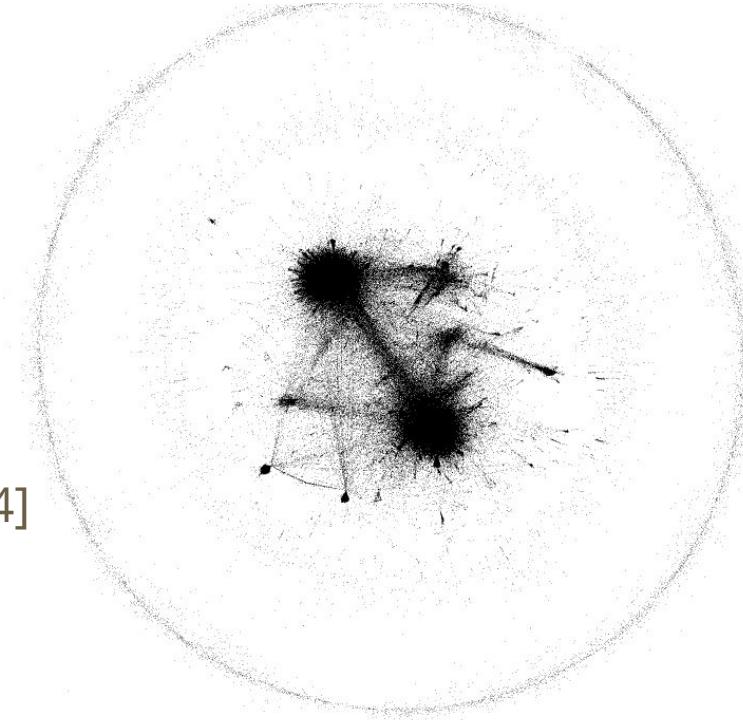
[3] Saramki, J., Kivel, M., Onnela, J. P., Kaski, K., & Kertesz, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. Physical Review E, 75(2), 027105.

[4] Fast unfolding of communities in large networks, Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008 (12pp)

[5] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. nature, 433(7028), 895.

# Graph Features

- Degree
- PageRank [1]
- Betweenness centrality [2]
- Clustering coefficient [3]
- Clustering affiliation (Louvain) [4]
- z-score del grado interno [5]
- Coeficiente de participación [5]



[1] Page, L., Brin, S., Motwani, R., & Winograd, T. (1999). The PageRank citation ranking: Bringing order to the web. Stanford InfoLab

[2] Brandes, U. (2001). A faster algorithm for betweenness centrality. Journal of mathematical sociology, 25(2), 163-177.

[3] Saramki, J., Kivel, M., Onnela, J. P., Kaski, K., & Kertesz, J. (2007). Generalizations of the clustering coefficient to weighted complex networks. Physical Review E, 75(2), 027105.

[4] Fast unfolding of communities in large networks, Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P10008 (12pp)

[5] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. nature, 433(7028), 895.

# Graph Features: cartografía de la red

NIH Public Access  
Author Manuscript

Nature. Author manuscript; available in PMC 2008 January 7.  
Published in final edited form as:  
*Nature*. 2005 February 24; 433(7028): 895–900.

**Functional cartography of complex metabolic networks**

Roger Guimera and Luis A. Nunes Amaral  
NICo and Department of Chemical and Biological Engineering, Northwestern University, Evanston, Illinois 60208, USA

**Abstract**  
High-throughput techniques are leading to an explosive growth in the size of biological databases and creating the opportunity to revolutionize our understanding of life and disease. Interpretation of these data remains, however, a major scientific challenge. Here, we propose a methodology that enables us to extract and display information contained in complex networks<sup>1–3</sup>. Specifically, we demonstrate that we can find functional modules<sup>4,5</sup> in complex networks, and classify nodes into universal roles, similar to the notion of ‘universal proteins’<sup>6</sup>. The method that yields a ‘cartographic representation’ of complex networks. Metabolic networks<sup>6–8</sup> are among the most challenging biological networks and, arguably, the ones with most potential for immediate applicability<sup>9</sup>. We use our method to analyse the metabolic networks of twelve organisms from three different superkingdoms. We find that, typically, 80% of the nodes are only connected to other nodes within their respective modules, and that nodes with different roles are affected by different evolutionary constraints and pressures. Remarkably, we find that metabolites that participate in only a few reactions but that connect different modules are more conserved than hubs whose links are mostly within a single module.

If we are to extract the significant information from the topology of a large, complex network, knowledge of the role of each node is of crucial importance. A cartographic analogy is helpful to illustrate this point. Consider the network formed by all cities and towns in a country (the nodes) and all the roads that connect them (the links). It is clear that a map in which each city and town is represented by a circle of fixed size and each road is represented by a line of fixed width is hardly informative. In contrast, a map in which each city and town are represented by dots of different sizes and the lines connecting them have different widths is one that can obtain scale-specific information about the network. Similarly, it is difficult, if not impossible, to obtain information from a network with hundreds or thousands of nodes and links, unless the information about nodes and links is conveniently summarized. This is particularly true for biological networks.

Here, we propose a methodology, which is based on the connectivity of the nodes, that yields a cartographic representation of a complex network. The first step in our method is to identify the functional modules<sup>4,5</sup> in the network. In the cartographic picture, modules are analogous to countries or regions, and enable a coarse-grained, and thus simplified, description of the network. Then we classify the nodes in the network into a small number of system-independent ‘universal roles’.

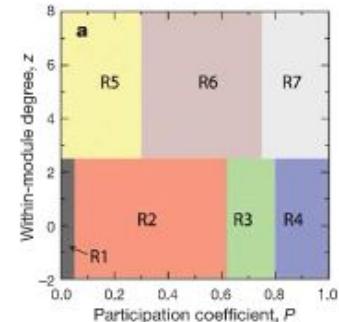
It is common that social networks have communities of highly interconnected nodes that are less connected to nodes in other communities. Such modular structures have been reported not only in social networks<sup>9,10–12</sup>, but also in food webs<sup>13</sup> and biochemical networks<sup>4,14–16</sup>. It is widely believed that the modular structure of complex networks plays a critical role in

- z-score del grado interno:

$$Z_i = \frac{(K_i - K_s)}{S_s}$$

- Coeficiente de participación:

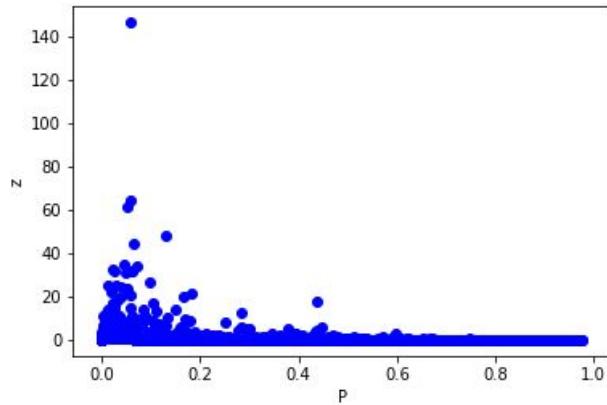
$$P_i = 1 - \sum_s (K_{i,s} - K_i)^2$$



- [1] Raghavan, U. N., Albert, R., & Kumara, S. (2007). Near linear time algorithm to detect community structures in large-scale networks. *Physical review E*, 76(3), 036106.
- [2] Guimera, R., & Amaral, L. A. N. (2005). Functional cartography of complex metabolic networks. *nature*, 433(7028), 895.

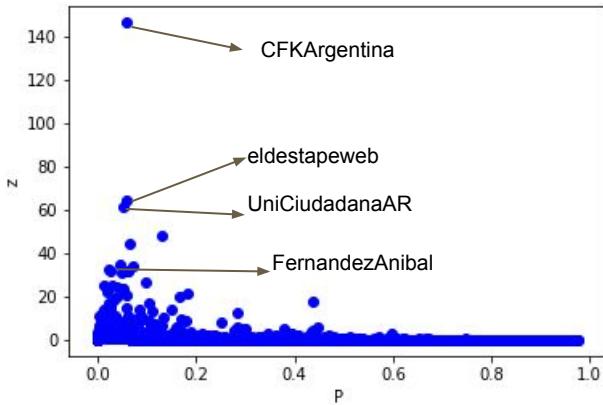
# Graph Features: cartografía de la red

# Comunidad: Unidad Ciudadana



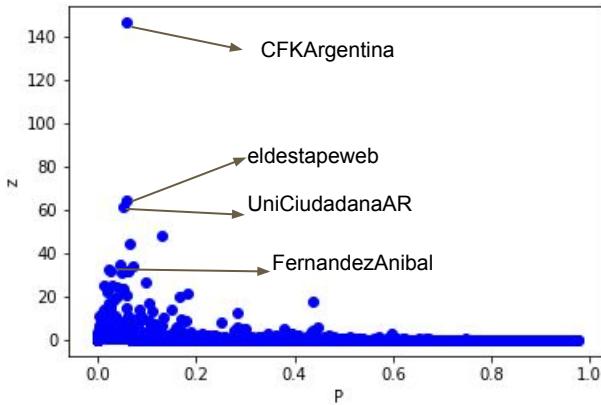
# Graph Features: cartografia de la red

Comunidad:  
Unidad Ciudadana

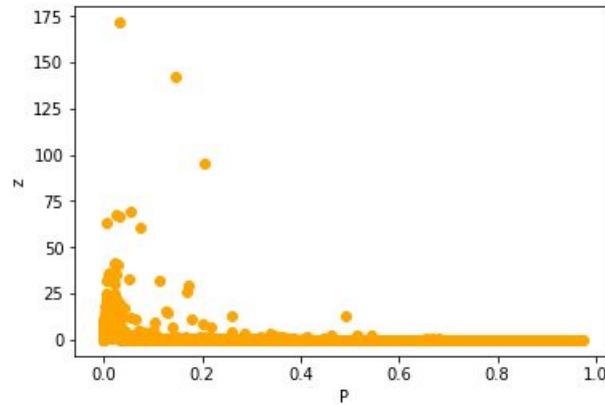


# Graph Features: cartografia de la red

Comunidad:  
Unidad Ciudadana

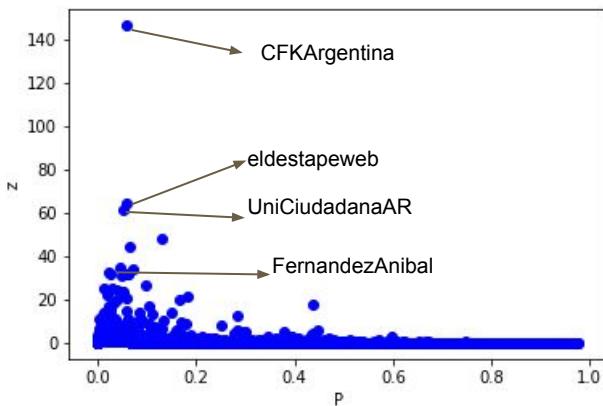


Comunidad:  
Cambiemos

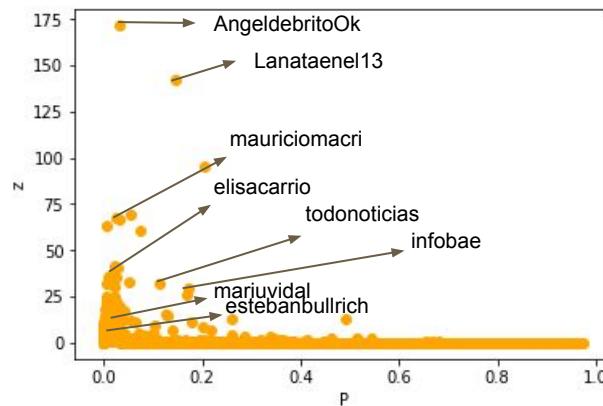


# Graph Features: cartografía de la red

Comunidad:  
Unidad Ciudadana

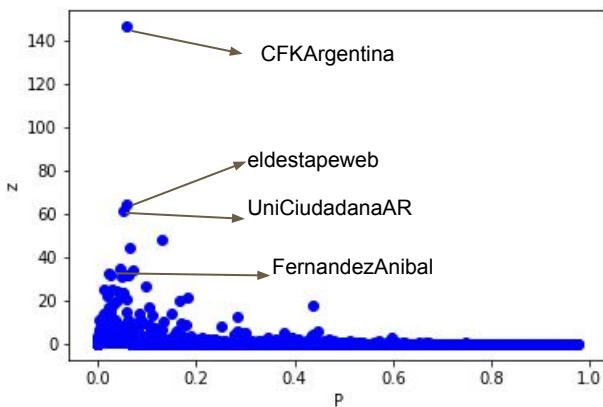


Comunidad:  
Cambiemos

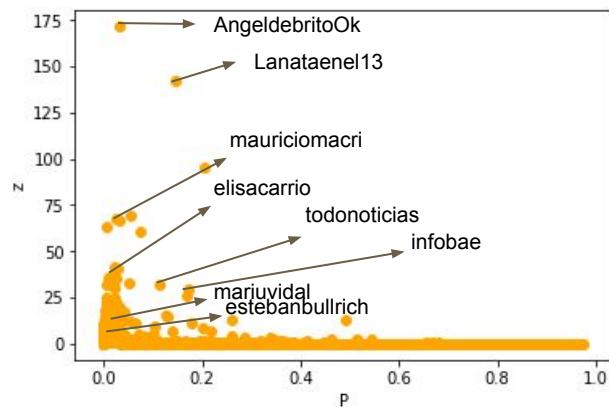


# Graph Features: cartografia de la red

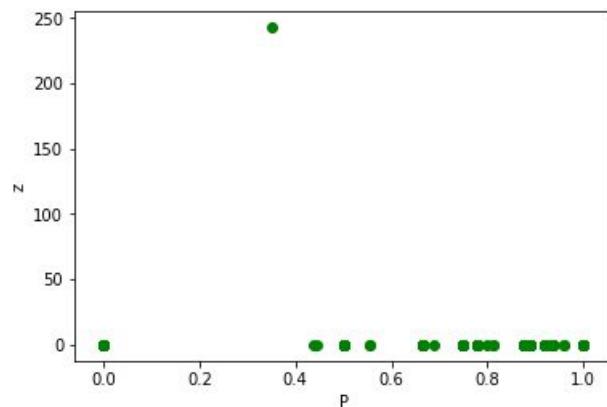
Comunidad:  
Unidad Ciudadana



Comunidad:  
Cambiemos

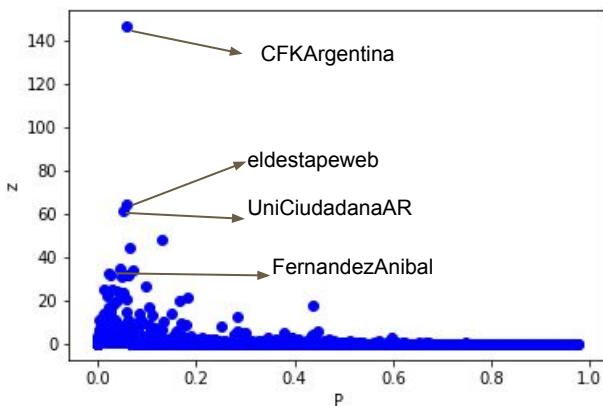


Comunidad:  
Famoso por un dia

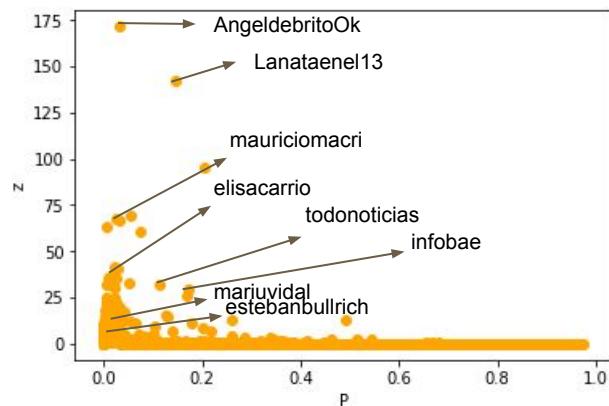


# Graph Features: cartografía de la red

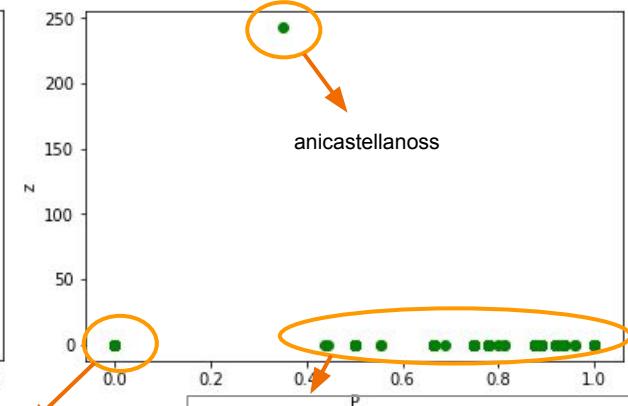
Comunidad:  
Unidad Ciudadana



Comunidad:  
Cambiemos



Comunidad:  
Famoso por un dia



Usuarios con un único RT a anicastellano ( $z=0$  y  $p=0$ )

A diferencia de los otros casos, este cluster no está definido ideológicamente y la mayoría de los usuarios están muy conectados a otros clusters ( $p>0.4$ ) pero poco dentro del mismo cluster ( $z=0$ )

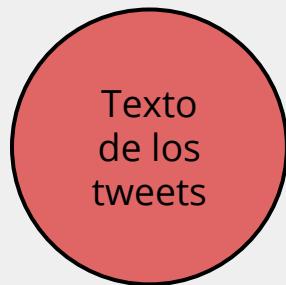
# Text mining

# Text mining



Texto  
de los  
tweets

# Text mining



 **Mauricio Macri**  @mauriciomacri · 21 oct. 2019

Queda una semana. Todos a trabajar para darla vuelta. Vamos que sí se puede!!!

3,7 mil  13,6 mil  64,3 mil  

 **Roberto Lavagna**  @RLavagna · Oct 13, 2019

No hay oportunidad de progresar si en los primeros años no se recibe la alimentación necesaria para crecer sano. Tampoco un joven puede cumplir sus aspiraciones si no cuenta con los medios para estudiar correctamente

#DebatePresidencial #DebatanPropuestas #DebateAr2019

4  37  80  

[Show this thread](#)

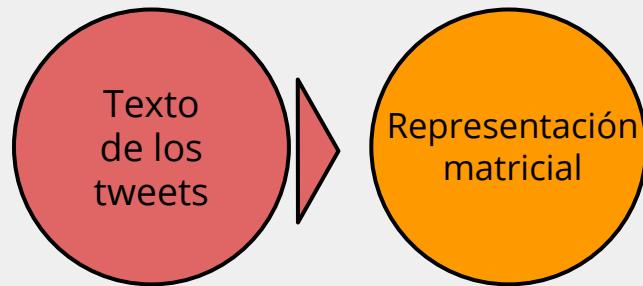
 **Alberto Fernández**  @alferdez · 21 oct. 2019

Gracias nuevamente a todas las personas que siguieron el debate en sus casas. El próximo domingo vamos a demostrar en las urnas que los argentinos decidimos dejar atrás nuestros desencuentros y trabajar entre todos para poner a la Argentina de pie. #DebateAr2019

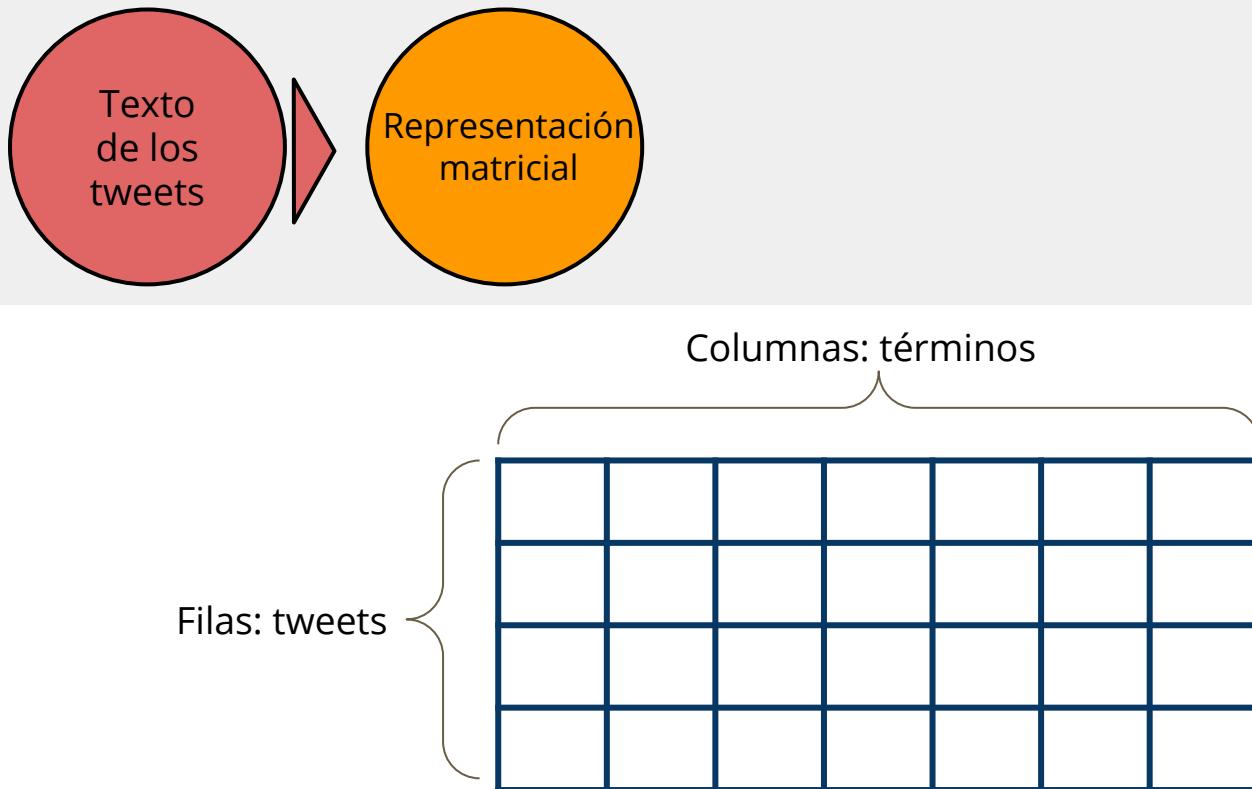


774  2,1 mil  11,4 mil  

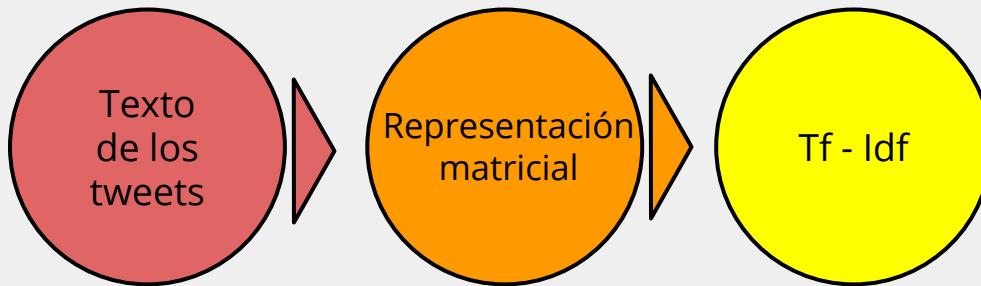
# Text mining



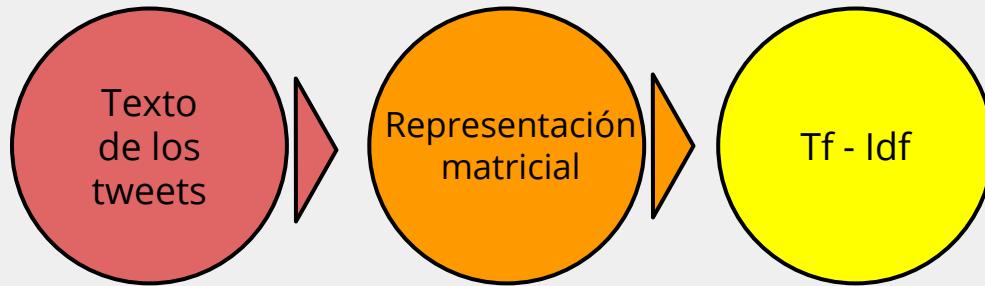
# Text mining



# Text mining



# Text mining



Using TF-IDF to Determine Word Relevance in Document Queries

Joan Ramos  
Department of Computer Science, Rutgers University, 2351 BRPO Way, Piscataway, NJ, 08854

**Abstract**  
In this paper, we examine the results of applying TF-IDF to determine what words in a corpus of documents are relevant to a query. As the term implies, TF-IDF calculates word relevance based on the frequency of the word in each document and the inverse proportion of the frequency of the word in all documents. We show how the words in the documents the word appears in. Words with higher frequencies in the document have a stronger relationship with the document they appear in, making them more relevant to the query. If a query, the document could be of interest to the user. This approach is very efficient and can also be used to quickly categorize relevant words from a large query dataset.

**1. Introduction**  
Before proceeding in depth into our experiments, it is useful to describe the nature of the query retrieval problem. In general, the problem is to find relevant information to answer a query. There are many different approaches to solving this problem, including TF-IDF.

**1.1 Query Related Problems**  
The field of text mining has a well-defined area that has become as common and natural in recent years that some might say it is a discipline in its own right. The growing use of query retrieval variants, content-based filtering, and other related topics is a testament to the problems.

Information retrieval can be described as the task of searching a collection of data, by that data documents, documents, or other items. First, we will focus ourselves to searching a collection of documents, which is often referred to as text mining. This becomes the task of searching this corpus of

JOANRAMOS@RUTGERS.EDU  
PDF DOI

**1.2 Algorithm for Ad-Hoc Retrieval**  
Let us briefly examine other approaches used for information retrieval. One of the most common and intuitive we present for the problem, the use of statistical methods, is the use of TF-IDF. This approach is based on responding to the problem (Bergé & Lafferty, 1999) for document retrieval. They propose two main ways to calculate the relevance of a document to a query. Both are aimed to enhance their approach. They suggest that the relevance of a document to a query is calculated as a sequence of words  $t$  in the actual document. They then calculate the relevance of the document to a query using Bayes' Law in equation (1), they then calculate the remaining approach.

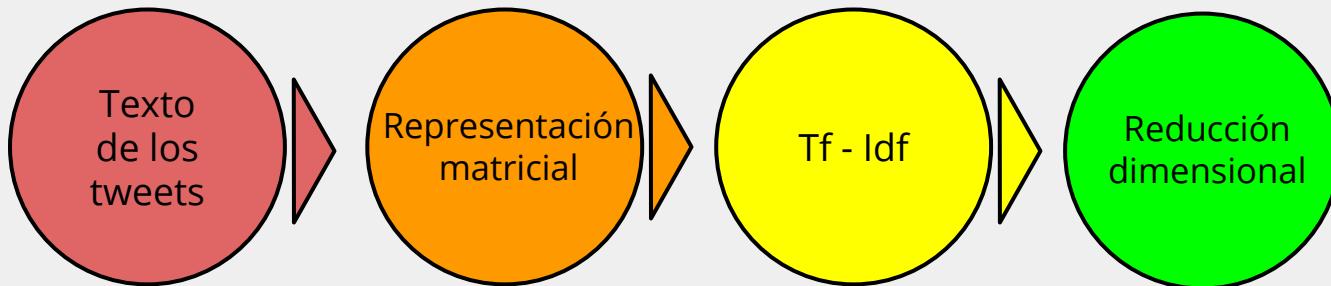
**1.3 Likelihood Methods**  
Likelihood methods for performing ad-hoc retrieval also known as probabilistic models (Bergé & Lafferty, 1994) have been used for document retrieval. These algorithms called Latent Semantic Indexing (LSI). In this approach, the documents are represented as vectors in a vector space that captures an n-dimensional representation of the documents. The distance between the numerical representations is compared to the cosine similarity of the documents. The closer the distance, the more similar the documents. This approach is highly effective in query retrieval, even when the documents are written in different languages (Lemire &

$$TF = N_t$$

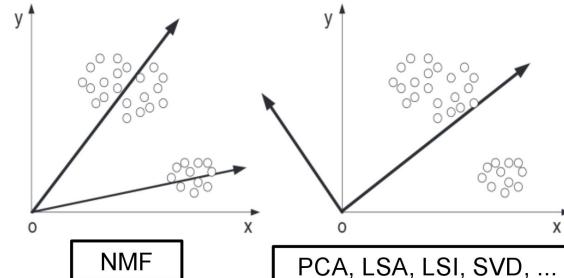
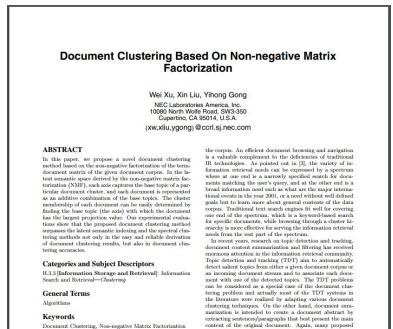
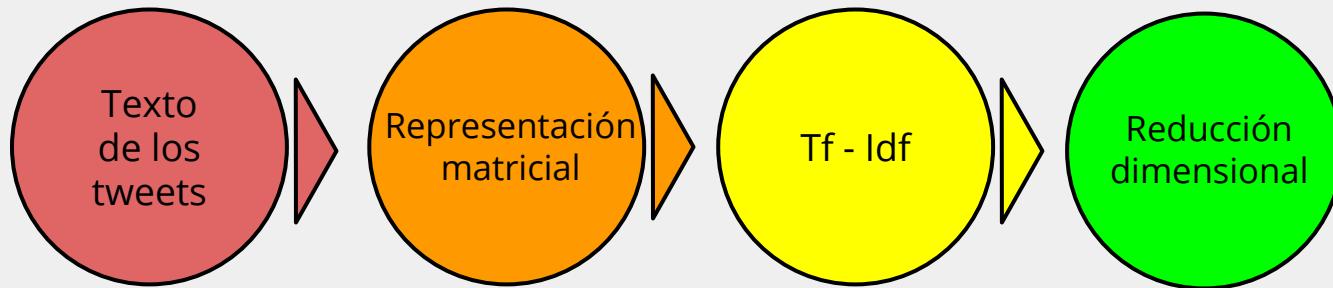
$$IDF = 1 + \log\left(\frac{1 + N}{1 + n_t}\right)$$

Donde  $N$  es el # total de documentos y  $N_t$  el # de documentos donde aparece el término  $t$ .

# Text mining



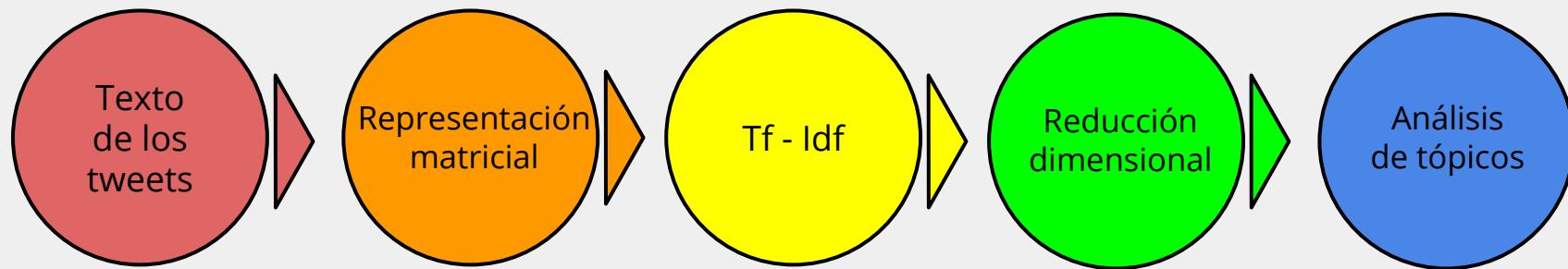
# Text mining



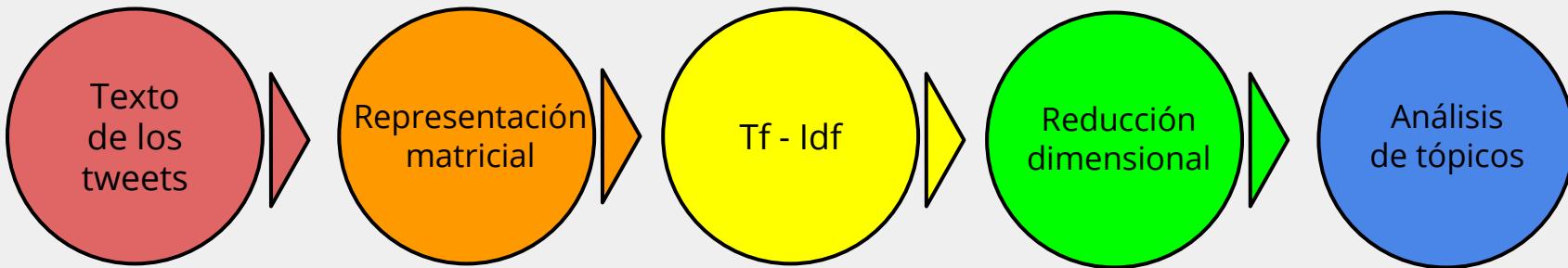
$$M \simeq H * W$$

Donde  $H$  es la matriz de documentos por tópico y  $W$  una matriz de tópicos por término.

# Text mining



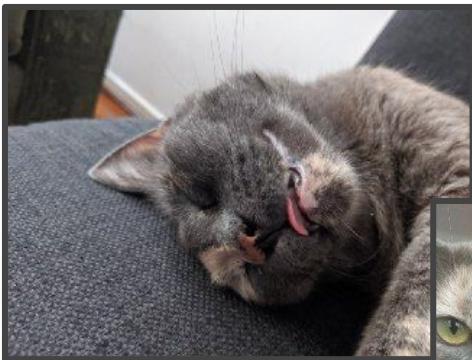
# Text mining



- **Santiago Maldonado:** Desaparecido en agosto del 2017.
- **Faltante de Boletas.**
- **Autoridades de Mesa:** Ausencia en las escuelas.
- **Economía:** centrado en pobreza, desempleo y programas de ajuste.
- **Tragedia de "Once":** En 2012 hubo un accidente de trenes donde murieron 51 personas.
- **Encuestas.**
- **Provincia de Santa Cruz:** Una de las principales provincias opositoras gobernada por Alicia Kirchner.
- **Provincia de Buenos Aires:** La provincia de mayor superficie y población de la Argentina
- **Venezuela.**

# Modelo predictivo

# Clasificación Supervisada



¿Gato o perro?

Clasific



perro?

# XGBoost

## XGBoost: A Scalable Tree Boosting System

Tianqi Chen  
University of Washington  
tqchen@cs.washington.edu

Carlos Guestrin  
University of Washington  
guestrin@cs.washington.edu

### ABSTRACT

Tree boosting is a highly effective and widely used machine learning method. In this paper, we describe a scalable end-to-end tree boosting system called XGBoost, which is used widely by data scientists to achieve state-of-the-art results on many machine learning challenges. We propose a novel sparsity-aware algorithm for sparse data and weighted quantile sketch for approximate tree learning. More importantly, we provide insights on cache access patterns, data compression and sharding to build a complex tree boosting system. By combining these insights, XGBoost scales beyond billions of examples using far fewer resources than existing systems.

### Keywords

Large-scale Machine Learning

### 1. INTRODUCTION

Machine learning and data-driven approaches are becoming very important in many areas. Smart spam classifiers protect our email by learning from massive amounts of spam data and user feedback; advertising systems learn to match the right ads in the right context; fraud detection systems protect banks from malicious attackers; anomaly event detection systems help experimental physicists to find events that lead to new physics. There are two important factors that drive these successful applications: usage of effective (statistical) models that capture the complex data dependencies and scalable learning systems that learn the model of interest from large datasets.

Among the machine learning methods used in practice, gradient tree boosting [10] is one technique that shines in many applications. Tree boosting has been shown to give state-of-the-art results on many standard classification benchmarks [16]. LambdaMART [5], a variant of tree boosting for ranking, achieves state-of-the-art result for ranking

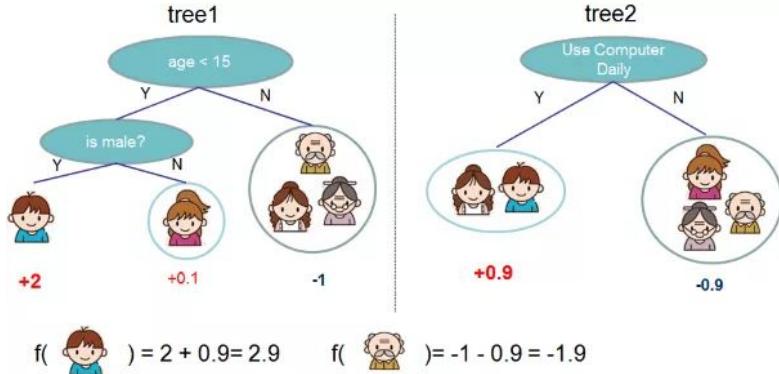
problems. Besides being used as a stand-alone predictor, it is also incorporated into real-world production pipelines for ad click through rate prediction [15]. Finally, it is the de-facto choice of ensemble method and is used in challenges such as the Netflix prize [3].

In this paper, we describe XGBoost, a scalable machine learning system for tree boosting. The system is available as an open source package<sup>1</sup>. The impact of the system has been widely recognized in a number of machine learning and data mining challenges. Take the challenges hosted by the machine learning community Kaggle for example. Among the 29 challenge-winning solutions<sup>2</sup> published at Kaggle's blog during 2015, 17 solutions used XGBoost. Among these solutions, eight solely used XGBoost to train the model, while most others combined XGBoost with neural nets in ensembles. For comparison, the second most popular method, deep neural nets, was used in 11 solutions. The success of the system was also witnessed in KDDCup 2015, where XGBoost was used by every winning team in the top-10. Moreover, the winning teams reported that ensemble methods outperform a well-configured XGBoost by only a small amount [1].

These results demonstrate that our system gives state-of-the-art results on a wide range of problems. Examples of the problems that these winning solutions include: store sales prediction; high energy physics event classification; sentiment classification; customer behavior prediction; motion detection; ad click through rate prediction; malware classification; product categorization; hazard risk prediction; massive online course dropout rate prediction. While domain dependent data analysis and feature engineering play an important role in these solutions, the fact that XGBoost is the consensus choice of learner shows the impact and importance of our system and tree boosting.

The most important factor behind the success of XGBoost is its scalability in all scenarios. The system runs more than ten times faster than existing popular solutions on a single machine and scales to billions of examples in distributed or

## Árboles de decisión



<sup>1</sup>Gradient tree boosting is also known as gradient boosting

[1] Chen, T., & Guestrin, C. X. A scalable tree boosting system. CoRR. 2016. arXiv preprint arXiv:1603.02754.

[2] Nielsen, D. (2016). *Tree Boosting With XGBoost-Why Does XGBoost Win "Every" Machine Learning Competition?* (Master's thesis, NTNU).

# XGBoost

## XGBoost: A Scalable Tree Boosting System

Tianqi Chen  
University of Washington  
tqchen@cs.washington.edu

Carlos Guestrin  
University of Washington  
guestrin@cs.washington.edu

### ABSTRACT

Tree boosting is a highly effective and widely used machine learning method. In this paper, we describe a scalable end-to-end tree boosting system called XGBoost, which is used widely by data scientists to achieve state-of-the-art results on many machine learning challenges. We propose a novel sparsity-aware algorithm for sparse data and weighted quantile sketch for approximate tree learning. More importantly, we provide insights on cache access patterns, data compression and sharding to build a scalable tree boosting system. By combining these insights, XGBoost scales beyond billions of examples using far fewer resources than existing systems.

### Keywords

Large-scale Machine Learning

### 1. INTRODUCTION

Machine learning and data-driven approaches are becoming very important in many areas. Smart spam classifiers protect our email by learning from massive amounts of spam data and user feedback; advertising systems learn to match the right ads with the right context; fraud detection systems protect banks from malicious attackers; anomaly event detection systems help experimental physicists to find events that lead to new physics. There are two important factors that drive these successful applications: usage of effective (statistical) models that capture the complex data dependencies and scalable learning systems that learn the model of interest from large datasets.

Among the machine learning methods used in practice, gradient tree boosting [10] is one technique that shines in many applications. Tree boosting has been shown to give state-of-the-art results on many standard classification benchmarks [16]. LambdaMART [5], a variant of tree boosting for ranking, achieves state-of-the-art result for ranking

problems. Besides being used as a stand-alone predictor, it is also incorporated into real-world production pipelines for ad click through rate prediction [15]. Finally, it is the de-facto choice of ensemble method and is used in challenges such as the Netflix prize [3].

In this paper, we describe XGBoost, a scalable machine learning system for tree boosting. The system is available as an open source package<sup>1</sup>. The impact of the system has been widely recognized in a number of machine learning and data mining challenges. Take the challenges hosted by the machine learning community Kaggle as an example. Among the 29 challenge-winning solutions<sup>2</sup> published at Kaggle's blog during 2015, 17 solutions used XGBoost. Among these solutions, eight solely used XGBoost to train the model, while most others combined XGBoost with neural nets in ensembles. For comparison, the second most popular method, deep neural nets, was used in 11 solutions. The success of the system was also witnessed in KDDCup 2015, where XGBoost was used by every winning team in the top-10. Moreover, the winning teams reported that ensemble methods outperform a well-configured XGBoost by only a small amount [1].

These results demonstrate that our system gives state-of-the-art results on a wide range of problems. Examples of the problems that these winning solutions include: store sales prediction; high energy physics; sentiment analysis; text classification; customer behavior prediction; motion detection; ad click through rate prediction; malware classification; product categorization; hazard risk prediction; massive online course dropout rate prediction. While domain dependent data analysis and feature engineering play an important role in these solutions, the fact that XGBoost is the consensus choice of learner shows the impact and importance of our system and tree boosting.

The most important factor behind the success of XGBoost is its scalability in all scenarios. The system runs more than ten times faster than existing popular solutions on a single machine and scales to billions of examples in distributed or

<sup>1</sup>Gradient tree boosting is also known as gradient boosting

### Pasos (forma simplificada):

1. Ajustar un modelo a los datos  
$$F_1(X) = Y$$
2. Ajustar un modelo al gradiente de la loss del modelo anterior  
$$h_1(X) = Y - F_1(X)$$
3. Crear un nuevo modelo  
$$F_2(X) = F_1(X) + h_1(X)$$

donde  $F_2$  es la versión "boosted" de  $F_1$ .

Y así sucesivamente

$$F_m(X) = F_{m-1}(X) + \gamma h_{m-1}(X)$$

donde  $\gamma$  es un factor de escala.

[1] Chen, T., & Guestrin, C. X. A scalable tree boosting system. CoRR. 2016. arXiv preprint arXiv:1603.02754.

[2] Nielsen, D. (2016). *Tree Boosting With XGBoost-Why Does XGBoost Win" Every" Machine Learning Competition?* (Master's thesis, NTNU).

# Resultados

# Modelo predictivo

## Input:

Matriz donde cada fila es un usuario y donde cada columna es una feature del usuario.

Lista de features:

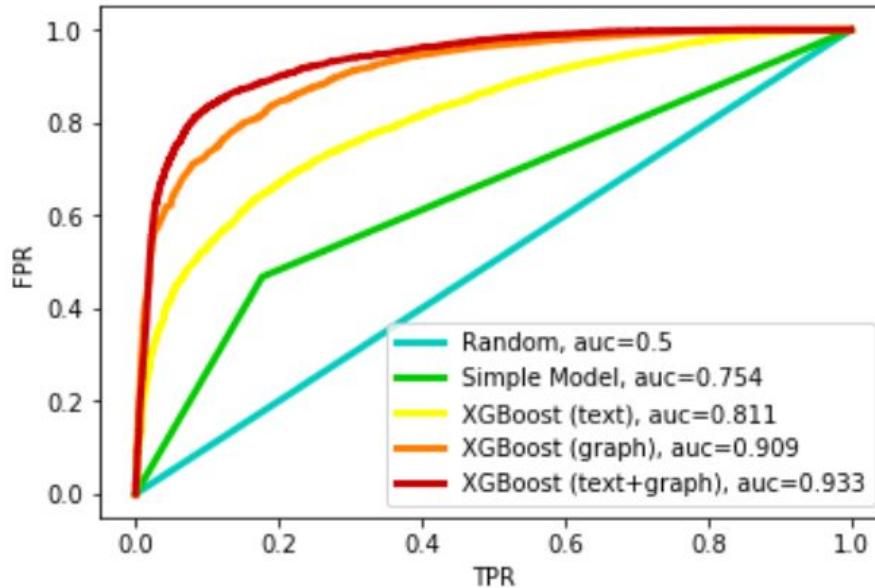
- ***Graph features:*** Degree, PageRank, Betweenness centrality, Clustering coefficient, Clustering affiliation (Louvain), z-score del grado interno, Coeficiente de participación.
- ***Text features:*** de qué tema habla.

## Output:

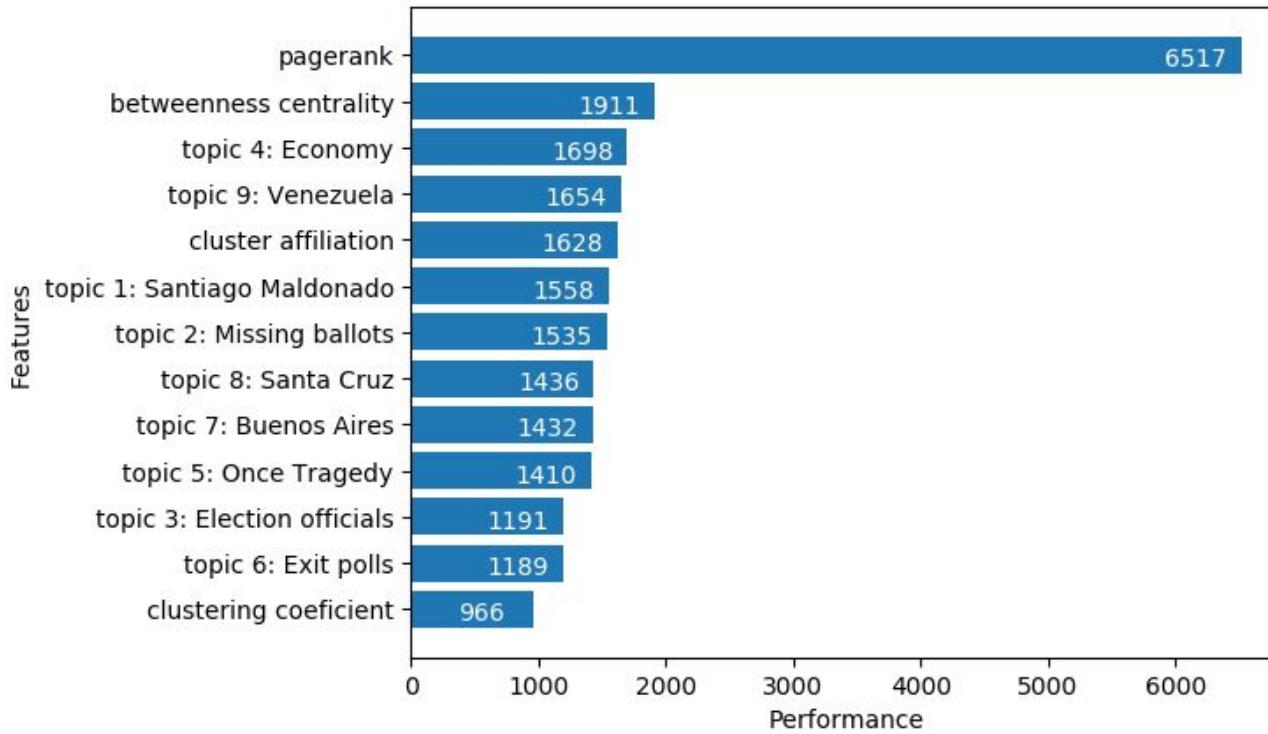
Vector de 1 y 0, cambió de comunidad entre elecciones o no respectivamente.

# Evaluación del modelo

Curvas ROC

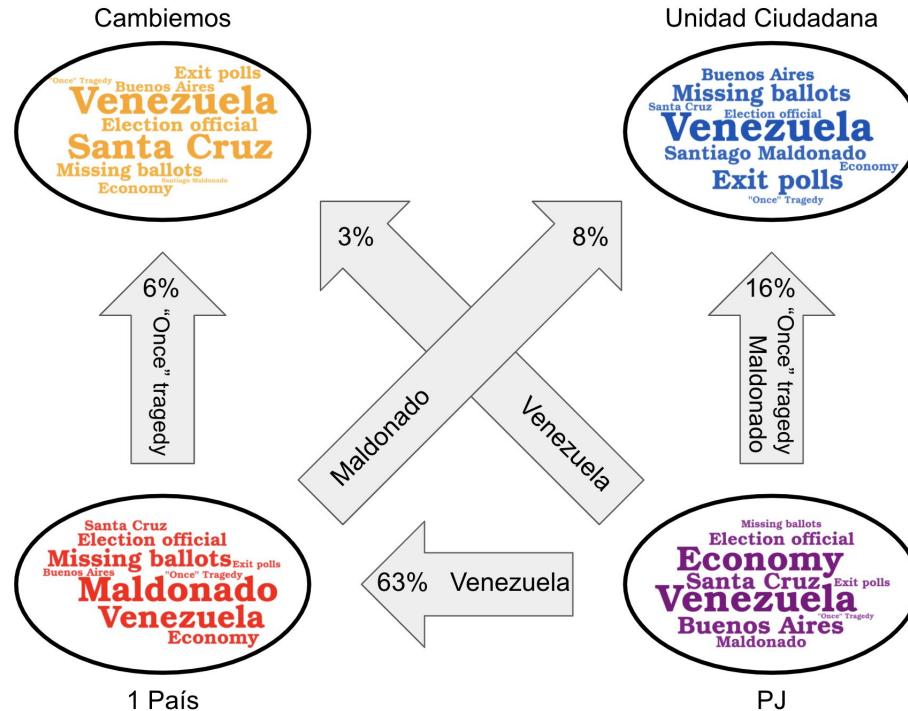


# Resultados: Feature Importance



# Tópicos más persuasivos

- ¿Quiénes cambian de grupo?
- ¿De qué habla la gente que cambia de grupo?
- ¿Y los que se quedan?



---

---

# Resultados 2019

---

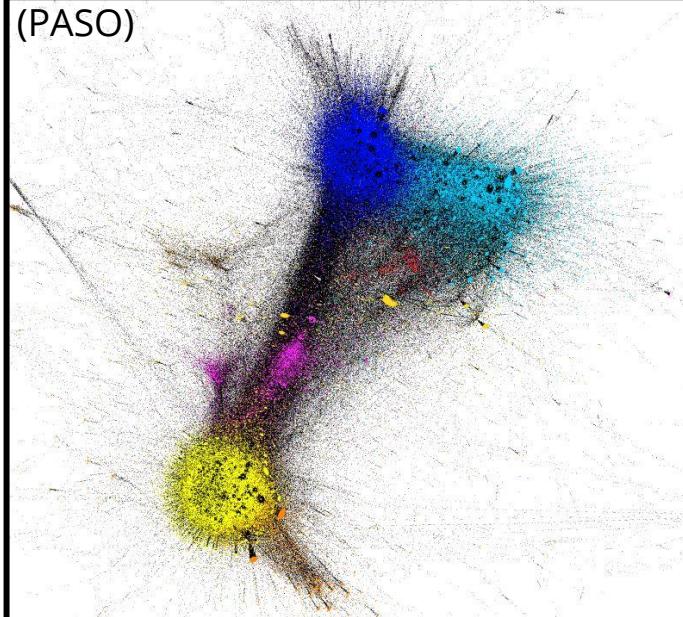
---

# Datos 2019

## Elecciones PASO 2019:

- 1.638.585 tweets
- 1.236.172 RT
- 256.860 usuarios

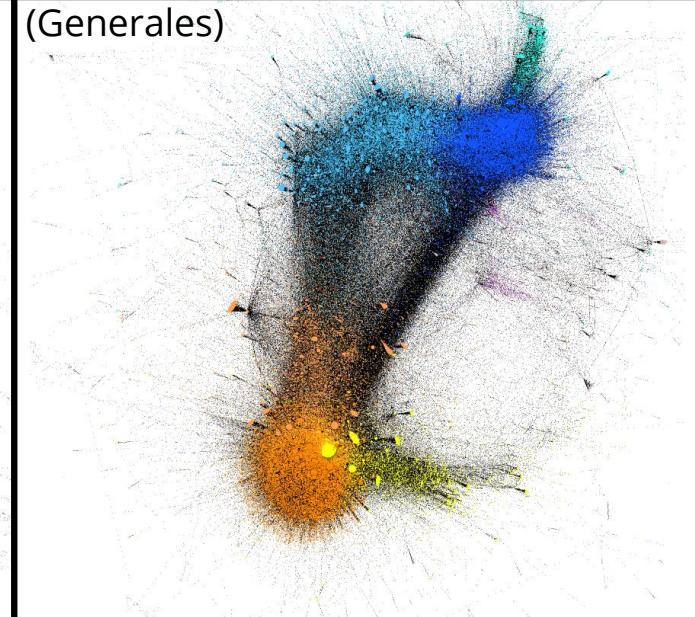
(PASO)



## Elecciones Generales 2017:

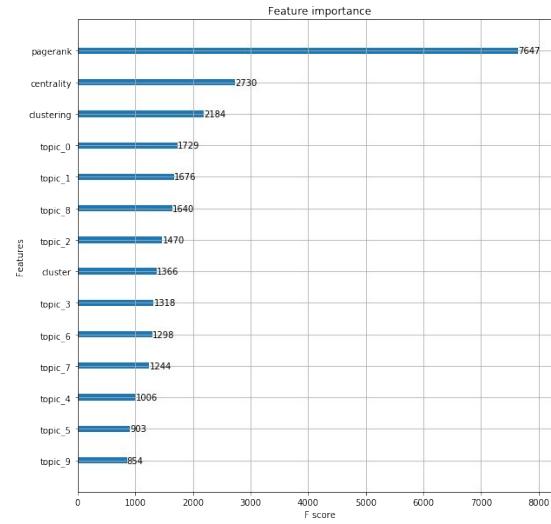
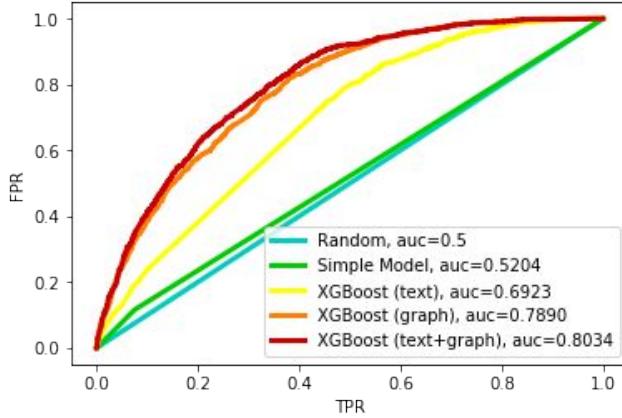
- 1.587.563 tweets
- 1.282.184 RT
- 271.910 usuarios

(Generales)



# Resultados 2019

Curvas ROC



## Tópicos más persuasivos:

Los de **Consenso Federal (Lavagna)** que se van a **Juntos por el Cambio (Macri)** hablan principalmente del "gobierno de la ciudad".

Mientras que los que se van al **Frente de Todos (Fernández)** hablan de Trabajo y Empleo (temas económicos).

# Conclusiones

# Conclusiones

# Conclusiones

- El modelo predictivo, usando datos de twitter, efectivamente **detecta cambios en la afiliación política** de las personas ( $auc= 0.933$ ) .

# Conclusiones

- El modelo predictivo, usando datos de twitter, efectivamente **detecta cambios en la afiliación política** de las personas ( $auc= 0.933$ ).
- Mostramos la importancia de las **características topológicas** de los usuarios en la red de RT (comúnmente ignoradas en la literatura).

# Conclusiones

- El modelo predictivo, usando datos de twitter, efectivamente **detecta cambios en la afiliación política** de las personas ( $auc= 0.933$ ).
- Mostramos la importancia de las **características topológicas** de los usuarios en la red de RT (comúnmente ignoradas en la literatura).
- Detectamos los **tópicos más persuasivos** e importantes. En una sociedad polarizada, esta herramienta tiene **muchas aplicaciones** para las ciencias sociales y para los partidos políticos.

# Conclusiones

- El modelo predictivo, usando datos de twitter, efectivamente **detecta cambios en la afiliación política** de las personas ( $auc= 0.933$ ).
- Mostramos la importancia de las **características topológicas** de los usuarios en la red de RT (comúnmente ignoradas en la literatura).
- Detectamos los **tópicos más persuasivos** e importantes. En una sociedad polarizada, esta herramienta tiene **muchas aplicaciones** para las ciencias sociales y para los partidos políticos.
- Los resultados son **consistentes** para distintas elecciones.

---

---

# ¡Muchas Gracias!

---

---

— Federico Albanese —

(falbanese@dc.uba.ar)

 @f\_albanese

---

---

# Preguntas

