

# Segmentación de Usuarios de Twitter

---

Agustín Capello  
FaMAF - UNC

# Problema

Segmentación (o división en clases), de usuarios de twitter en base al texto de sus tweets o sus retweets.

# Datos

Se obtuvieron tweets/retweets entre las 18 a 00 hs del día 22 de noviembre con hashtags relacionados a las elecciones/ballottage en Argentina

A partir de estos tweets, se obtuvo toda la información relacionada a usuarios.

El objetivo es lograr descubrir inclinaciones políticas de usuarios, como así también grupos de usuarios relacionados por la misma temática.

# Approach de la solución

Se tomó como un problema capaz de ser solucionado con machine learning, usando el algoritmo no supervisado de clustering K-Means.

Se vectorizaron los usuarios, cada uno como vector de features.

Los features utilizados fueron bag of words, bag of hashtags (ambos en base al texto de los tweets de un usuario) y relación con otros usuarios por retweets

Cada "feature" es una característica propia que posee cada usuario.

# Algunos procesos que mejoraron la segmentación

- Tokenizar/normalizar los tweets de todos los usuarios (en minúsculas, sin acentos, etc.).
- Quitar "function words" (artículos, preposiciones, pronombres, símbolos, etc.) de los tokens principales de cada usuario.
- Quitar palabras/tokens poco significativos y que se repiten mucho (ejm: argentina, balotaje2015, etc.).
- Dividir los usuarios en una mayor cantidad de clusters: con siete u ocho clusters se visualizaron mejor las diferencias que con tres o cuatro.
- Usar sólo features relevantes, y quitar algunos como la ubicación del usuario, su nombre o si está verificado.
- Utilización de diferentes métodos para evaluar la clasificación y analizar los resultados:
  - a. Palabras más repetidas de cada cluster, y cantidad de usuarios que la usaron.
  - b. Examinar los usuarios más relevantes de cada cluster (más cercanos a los centroides), y sus tweets.
  - c. Examinar los features más representativos de cada cluster.
  - d. Word-clouds de cada cluster.

# Conclusiones - Resultados

El siguiente es un ejemplo de la evaluación de un modelo entrenado, es decir, el análisis de la clasificación.

En "Reporte.txt"

dato oficial provisorio numeros dentro imbecil triunfo escucharlo conferencia ceres claudio  
telefe com daniel scioli dentro imbecil triunfo escucharlo conferencia ceres claudio  
propio escrutadas fpv mesa vera  
elegiste lt3 capital lider debe participacion em acta gana pringles cordoba ahora  
puedas presidente respetarlo ahí ganaba rosario verdadero dar pro espero minutos  
eltrece oficial seria gran ratón mauricio macri  
sergiomassa podran soy primeros laburo mauricio presidente eleger rating radiosuquia  
governistas local falam abrazo mauricio jbonifacio america tv escuela laplata penalosa  
votos transitocordoba macri lafano albertofernandez  
cpugliese 1968  
todonoticias

resentimiento  
gustan  
reencuentros  
si  
daniel  
formen  
pais  
li  
partido  
gente  
cambiemos  
mauricio  
macri

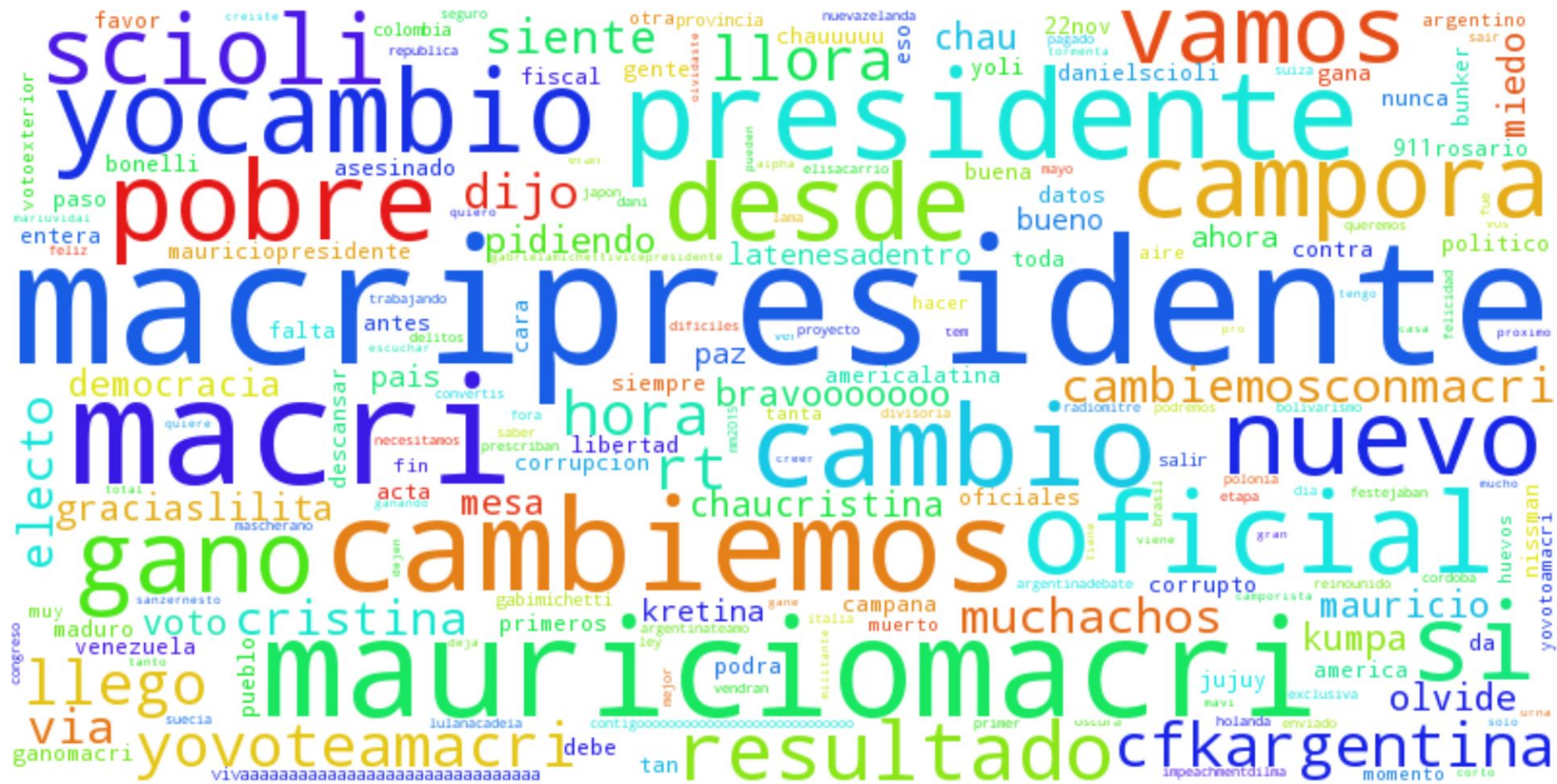
reconciliacion votar serfiscal impunidad chauuu eso afanarse  
cuba poderen central socialismo tiene nelson bocaranda vido maximo caedoguidepadua mejor cristina pibes proargentina  
lucha ir gran fraude electoral fin nopreguntes lee rt banco traste figuemin cristobal caro veces  
lazarro scioli quedan atentos nuestro historia lopez amia anos juntos timerman cuentan lucio quincioc trabajar chau  
presidencial fraude inflacion desde era exclusivo aflojemos bunker kirchner cuidemos voto laburar gusel chepereyra  
boudou eligio mundo nunca guita kicillof presidente aislados marupita mauricio



urna bien presidenciales favor presente primer pueblo nisman  
decoracion lilita dato oficial victoria  
acompanando argentino jubilo falkland yovoteamagri generales chnee expulsado  
margarita aqui devolver estado dio padron hermano  
solo venezolano sanzernesto boletin republica diciembre desde anunciando escrutadas mariuvidal participacion  
gobierno eres alak adios nuestro via vivacasa dara mucho habia ano vuelta wolffwaldo votar dime saber muy  
saenzricardo america reacciono tan gane interesa dire golpe historiconuevo  
fotodemocrata sera vota socialismo si vamo hecho urgente justicia oficial ganen historia  
minuto triunfo esperanza tienen esperando hora electoral mucha tiene mencion gusta gente cfkargentina felicidades medio  
paz julio mesa despues duron argentinos tenemos voto junto momento cristina boca primera vuelta felicitaciones  
jose nacional aplastante fue cambio fin espero bueno rosada mejor entre partido scioli  
fue cambio fin espero bueno rosada mejor entre partido scioli







[illegible]



quieren militopresidente electoral nacer sean bunker jornada acabaron seguir nacional mientras aquello bien justicia  
estoy fue gano momento san alak semiedo paz comienzo muy cambiemos presidente  
fiscales participacion diferencia quieramos viene glorioso dicen parte cara centro de computos daniel  
alcanzo fracaso solo espera tercero ojala poder caripelas bueno casi sueno militante parte  
punto llegar miliciania padron mesa retwitea tenemos atencion gente nuevo periodoista f5  
tanto c5n fpv ganen vivo importante cfkargentina dan ultimo mucho comienza odio voto  
gana actualiza atentaron quiero nunca saber salio mauricio ver desde diego branca fin pais minuto  
futura quiero siente mayo cristina sobres dame mauro ver desde diego branca fin pais minuto  
victoria mejor otra terminamos tanta votaron reaparecio final extrano continuidad ganador  
sanguches estamos termina politica responde afuera dia escuela estaran argentino acta hubo  
venezuela votaste trabajar perdidos linda tenes favor scioli ahora disponibles  
computo alegria deportista festejo chau prensa ano politico campora palabra si cayetano vamo  
escrutado ano politico brancatelli dio muchachos caos

[illegible]