

# **PRAC1: Web Scraping**

Máster Universitario en Ciencia de Datos (UOC)

Tipología y ciclo de vida de los datos

Daniel Núñez Sanz

## **Contexto**

El año 2016 fue un año muy convulso en la política a nivel internacional. Contra todos los pronósticos y las encuestas: el Brexit gana las elecciones en Gran Bretaña, gana el no en el plebiscito sobre los acuerdos de paz en Colombia y gana Trump las elecciones. En ese momento se desconocía que había pasado o como se pudieron suceder estos tres sucesos a priori tan poco probables.

En el año 2018 aparece el escándalo de Cambridge Analytica que arrojará luz sobre estos sucesos y otros que se producirán justo después como la victoria de Bolsonaro en Brasil o el auge del populismo de derechas en Europa.

El escándalo de Cambridge Analytica vino por el uso de información de millones de usuarios de Facebook para crear anuncios políticos personalizados en las elecciones presidenciales de EE. UU. de 2016. Tras esto, uno de sus trabajadores revela que la empresa había creado una maquinaria para manipular las decisiones de los votantes. Pero ya no solo es que usaran las redes sociales para crear anuncios, sino que van más allá. En un video con cámara oculta nombran algunas tácticas que han usado; como el uso de prostitutas, sobornos o usar antiguos espías del M16 o el Mossad israelí para conseguir información de los rivales<sup>1</sup>.

Además, crean mentiras y las difunden por redes para beneficiar a sus clientes. Pero detrás de esto hay mucho más. No solo un encargo de un partido hacia una empresa o una consultora de marketing para ganar las elecciones como hacen todos. Detrás de esta empresa hay un grupo de personas dedicadas a manipular las decisiones de los votantes en América y Europa a favor de sus intereses, en favor de una ideología concreta.

Este entramado lo podemos simplificar en una persona clave. Steve Bannon, Vicepresidente de Cambridge Analytica, asesor de la campaña de Trump en 2016 y director de Breitbart (medio de extremaderecha). Bannon ha declarado en numerosas entrevistas su deseo de crear un movimiento político populista de derechas<sup>2</sup>. Creó la agrupación The Movement con la idea de unir a los partidos populistas de derechas en Europa y Latinoamérica. Además de su conexión directa con Trump, Bannon ha tenido relación o ha asesorado a partidos de extremaderecha como el Frente Nacional francés, Alternativo Por Alemania, Liga Norte en Italia o a Vox en España entre muchos otros.

A pesar de que Cambridge Analytica cerró, luego crearon una nueva empresa. A parte de esta, hay más empresas que han cogido el relevo, como CLS Estrategies. Entre Cambridge Analytica,

---

<sup>1</sup> <https://www.elperiodico.com/es/internacional/20180320/ejecutivos-cambridge-analyticafilmados-alardeando-prostitutas-sobornos-influencia-elecciones-channel-4-6702488>

<sup>2</sup> <https://www.eltiempo.com/mundo/eeuu-y-canada/entrevista-con-steve-bannon-sobre-el-populismo-y-el-gobierno-de-donald-trump-294662>

CLS, Bannon y los partidos populistas de derechas se puede ver una estrategia común: creación de miles de bots en redes sociales, difusión de fakes news o anuncios de fraude electoral en caso de que vayan a perder las elecciones como pasó en Bolivia y el posterior Golpe de estado. Es más, lo estamos viendo en estos momentos con Trump en las actuales elecciones a la presidencia de EE. UU., pidiendo que no se tenga en cuenta el voto por correo, que se pare el recuento del voto o su recurso a la justicia y sus acusaciones de fraude ante la OEA.

El 10 de abril de 2018 Vox anuncia que se reunirá con Bannon<sup>3</sup> para impulsar su agenda política y establecer estrategias de propaganda de cara a las siguientes elecciones<sup>4</sup>.

Y si nos fijamos en los resultados electoral de Vox en todos los niveles, vemos que el partido a pesar de presentarse a elecciones desde 2014, no es hasta diciembre de 2018 que logra sus primeros escaños. Tras esto, logra escaños en casi todos los parlamentos autonómicos y en el Congreso.

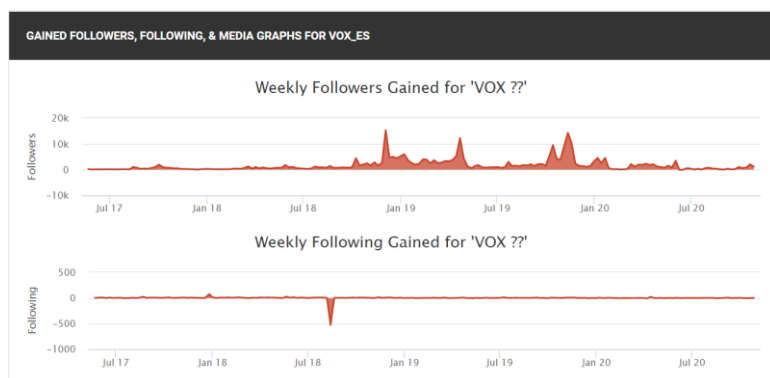
En este contexto, resulta interesante analizar las redes sociales de Vox. ¿Cómo han cambiado antes y después del encuentro de Bannon? ¿Ha tenido algún efecto el asesoramiento de Bannon en las redes sociales de Vox?

La principal red social donde se puede enviar su mensaje sin filtros y hacer más trampas sin pagar, es Twitter. Esto se debe a que Facebook, Instagram o YouTube son vistos principalmente entre quien sigue las cuentas, pero en Twitter puede lograr a poner un Hashtag en Trending Topic y lo verá todo el mundo.

Para ellos vamos a extraer el histórico de los seguidores y seguidos de la cuenta oficial de Vox en esta red social y analizarla. También cogeremos los de los otros tres principales partidos (PSOE, PP y Podemos) para facilitar un futuro análisis.

Dado que no podemos acceder al historial de los seguidores de una cuenta de Twitter, tenemos que recurrir a otra plataforma que los haya recopilado anteriormente. Esta información la podemos encontrar en Social Blade, web que proporciona información y métricas de usuarios de numerosas redes sociales, entre ellas Twitter. En el final de la web se puede ver un gráfico que tiene la información que necesitamos:

[https://socialblade.com/twitter/user/vox\\_es/monthly](https://socialblade.com/twitter/user/vox_es/monthly)



<sup>3</sup> Esta reunión se da gracias a Fajardí, ex del PP que se pasó a Vox, pasó por FAES y ahora está en Atlantic Council, vinculada con CLS Estrategies.

<sup>4</sup> <https://www.voxespana.es/noticias/bannon-el-artifice-de-la-victoria-de-trump-apuesta-por-vox-para-espana-20180410>

También, indagando encontramos a otras personas que han realizado web scraping en la página de Social Blade. Por ejemplo, Anya Vastava<sup>5</sup> que lo hace para recopilar datos de YouTube, y en su caso recurre a Web Archive para obtener datos de hace más tiempo. En nuestro caso, lo intentamos, pero no había ninguna captura en Web Archive de los datos que necesitábamos, por lo que lo hicimos directamente desde Social Blade.

## Proceso de obtención de datos

Para su obtención hemos raspado con Python la web Social Blade. Primero intentamos coger la información de la web directamente, pero nos da error. Sabe que estamos accediendo desde Python:

```
In [6]: import requests
        from bs4 import BeautifulSoup

        page = requests.get("https://socialblade.com/twitter/user/vox_es/monthly")
        soup = BeautifulSoup(page.content)
```

```
print(soup.prettify())
```

```
<html>
<head>
<style>
/*!
 * Font Awesome 4.2.0 by @davegandy - http://fontawesome.io - @fontawesome
 * License - http://fontawesome.io/license (Font: SIL OFL 1.1, CSS: MIT License)
 */@font-face{font-family:FontAwesome;src:url(data:application/vnd.ms-fontobject;base64,xt0AAODZAAACAAIABAAAAAAAAAAAAAAAAABAB
JABAAAAEAEQAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAFMQAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABYARgBvAG4AdABBAHCAZQBzAG8AbQBIAAAADgBSAGUAZwB1AGwAYQ
ByAAAAJABWAGUAcgBzAGkAbwBuACAANAAUADIALgAwACAAAGAwADEAMwAAACyARgBvAG4AdABBAHCAZQBzAG8AbQBIAAAUgB1AGcAdQBzAGEAcgAAAAAAAAQ1NHUAA
AAAAAAAAAAAAAAAAAADAX7AANnYANncALFgFM3pJM/SEVmjRApN5UYeAHIZoKQ+sJvend4XEADGjy1bHBQxybVI0e2miS1BKUbg1dmcMFLNLtMvC1M3u1UZJ+sm
YQqc02ETIzcxE14Thr13UpW5IKhQ87qRcpd9W8S29TPYVaHut2PNpK1xBSFbduw19Vc0EhvjQJBMr5KcViHZYKIMFuANcsN3r0q1UT70vU1xrw9L/ZXAp1/5XR
hv5KA6RyVZII88jIPmrFvFqpjLUzlwUfOQ765UJ1qfxx9qjbu520FfnvRb/yJM1g6NkNMfYBzJIGFBEHLuRvvuF0yq/pp2ZQAOFNwDR2GNMUMXIGGbyu6gMKj7r9
Y4e1NBisEtI1eRyMjqrA9Nw6kRXj5mEny7KCAAtgZMDaCLBrwUn4+oBRjStP/UyqOGmA1ChNuhShUZWZYOMHmEDPyb58EgkQTIKU1D2nQ6E9yuh1pzzuJ2bysL1
zSAMnGwUJaS8ieN2JINvnj1habA5XjCF1GZk4UoU4piq7WHDZUZVa0yswYs1xUm7BxoTJx1a1Z+jfSMBwNBA7tJj0jNIXLq+ZXAVR0jA4MzwCQKaQF1pQ8HaYCR
U6F3hTCo0I7vYqznQFhvebeae5e4Ztaq9q9o4hVw4SzaLKKJTJ0aDCoXBauhkcIGtpkgTqUxewMLowkpZQxUqYiuvjBSBRxu0BeDCeQvU9NjH1BwK+B6NiVmF08x
E0erIQmr4Ld8FeF+icaLJ5Bu8nsiwTz9pVAUSW4+66wRshg2FRou10cr9IzD6HbYGINCLAyFZaPubPbWtcJ10iTMd4EBv0LA23kCe3zfuUmNKc305i8FLHqm0F6x
WyuGrfD4VbthK4FzBBHLuN3USN9yQ8m71Y6W0L4X0WVGziJIIdrgYeCSaVQzu+nxQ1TZqq4XI7G0zDhFRk0yKpuEsOHLKySp5de3KT6TN+rumJY99agZK67gnWK
po3RfKTrOtD5SzkjzVbc2Ysz58cjNW62kf1k1LBczahXFzZ8+ZFeZRwLQ6h2WnBgXi8IVBS1T0pD36AWxDiXH3zIb6isSk5S0867xHgc02TrqkaIAW1ljZA6erc
vBqBV85g+SYMoujmeZh7thicgaOPDTkwHckZGhJk4sLbS0FRMwaj5scrjiOJYucdxvaIiAKIQIP5iZPg9CXLiHPVe6JDHijj0kux/c2aDNVsN2WqeEGG83qc+9h65
```

```
print(soup.prettify())
```

```
</head>
<body>
<div class="cf-error-details cf-error-1020">
  <h1>
    Access denied
  </h1>
  <p>
    This website is using a security service to protect itself from online attacks.
  </p>
  <ul class="cferror_details">
    <li>
      Ray ID: 5ec5f359af5814f9
    </li>
    <li>
```

<sup>5</sup> <https://medium.com/swlh/how-to-scrape-socialblade-for-youtube-subscription-data-ec7c4bde6933>

```
print(soup.prettify())

<br/>
<ol>
  <li>
    1) You're not actually a human (i.e. a crawler/scraper)
  </li>
  <li>
    2) You're accessing the website from a computer network not typically used by humans (i.e. a data center)
  </li>
  <li>
    3) You're accessing the website from a network that has a high level of abuse (i.e. a vpn provider)
  </li>
</ol>
<br/>
If you believe you're receiving this in error please take a screen shot of this entire page and open a support ticket at
<a href="https://support.socialblade.com">
  support.socialblade.com
</a>
<br/>
If you are trying to access Social Blade from a data center via automated means you probably want to look at using our
<a href="https://socialblade.com/business-api">
```

Para solucionar este problema modificamos el user agent y algunas cabeceras HTTP:

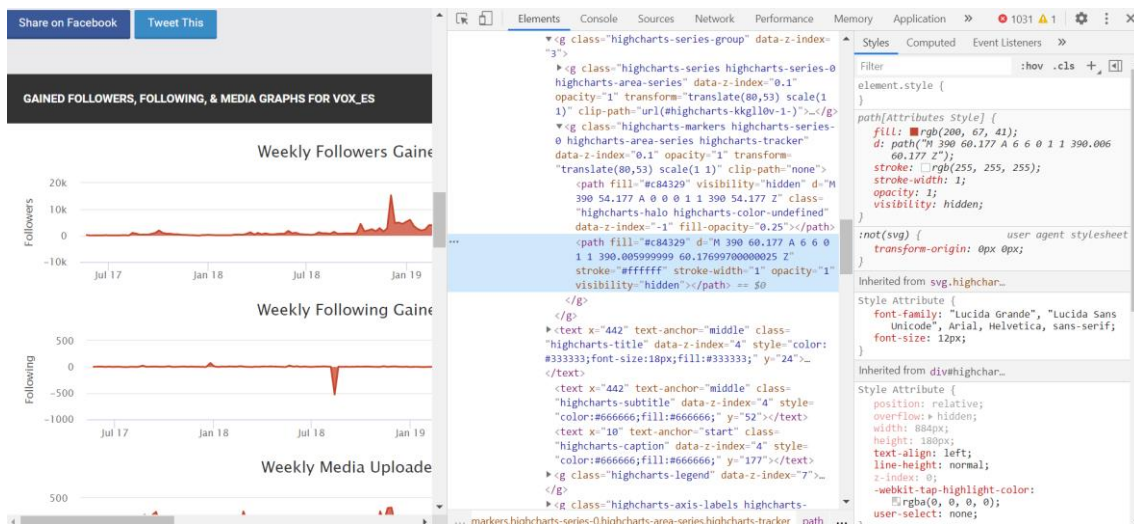
**jupyter** Prac1-Webscrap Last Checkpoint: hace 5 horas (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help Trusted
+ % Run C Code
import time
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

%matplotlib inline

In [3]: #Modificamos las cabeceras HTTP y el user agent al detectar la web que accedemos desde python
headers = {
    "Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,\
    /*;q=0.8",
    "Accept-Encoding": "gzip, deflate, sdch, br",
    "Accept-Language": "en-US,en;q=0.8",
    "Cache-Control": "no-cache",
    "dnt": "1",
    "Pragma": "no-cache",
    "Upgrade-Insecure-Requests": "1",
    "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_12_3) AppleWebKit/5\
    37.36 (KHTML, like Gecko) Chrome/56.0.2924.87 Safari/537.36"
}
```

Analizamos la web:



Vemos donde está la información que queremos y que su código es en JavaScript

DATE		FOLLOWERS	FOLLOWING
2020-10-06	Tue	-	422,764
2020-10-07	Wed	+25	422,789
2020-10-08	Thu	+16	422,805
2020-10-09	Fri	+86	422,891
2020-10-10	Sat	+187	423,078
2020-10-11	Sun	+278	423,356
2020-10-12	Mon	+236	423,592
2020-10-13	Tue	+457	424,049
2020-10-14	Wed	+78	424,127
2020-10-15	Thu	+93	424,220
2020-10-16	Fri	+93	424,313
2020-10-17	Sat	+66	424,379
2020-10-18	Sun	+30	424,409
2020-10-19	Mon	+44	424,453

TWITTER STATS SUMMARY / USER STATISTICS FOR VOX\_ES (2020-10-06 - 2020-10-19)

```

<div style="width: 860px; padding: 0px; float: left; margin: 10px 0px;"></div>
<div style="clear: both;"></div>
<script type="text/javascript"></script> == $0
<div style="width: 860px; background: #333; padding: 20px; text-transform: uppercase; color:#fff; margin-top: 30px;"></div>
<div style="width: 885px; padding: 15px 10px 15px 5px; background: #fff; float: left; margin: 0px 0px 10px 0px;">
  <div id="graph-twitter-weekly-change-followers-container" style="min-width: 800px; height: 180px; margin: 0px auto; overflow: hidden;" data-highcharts-chart="0">
    <div id="highcharts-kgllyw-0" dir="ltr" class="highcharts-container" style="position: relative; overflow: hidden; width: 884px; height: 180px; text-align: left; line-height: normal; z-index: 0; webkit-tap-highlight-color: rgba(0, 0, 0, 0); user-select: none;">
      <svg version="1.1" class="highcharts-root" style="font-family:'Lucida Grande','Lucida Sans Unicode', Arial, Helvetica, sans-serif;font-size:12px;" xmlns="http://www.w3.org/2000/svg" width="884" height="180" viewBox="0 0 884 180">
        <desc>Created with Highcharts 8.2.2</desc>
        <defs>...</defs>
        <rect fill="#ffffff" class="highcharts-background" x="0" y="0" width="884" height="180" rx="0" ry="0"></rect>
        <rect fill="none" class="highcharts-plot-background" x="80" y="53" width="794" height="90"></rect>
      </svg>
    </div>
  </div>
  
```

Styles Computed Event Listeners  
 Filter :hov .cls +  
 element.style { }  
 script { user agent stylesheet }  
 display: none;  
 Inherited from body  
 body, html { main.css?v=45679023:19 }  
 margin: 0 pxpx; padding: 0 pxpx; color: #232323; font-size: 11pt; font-family: 'Roboto', sans-serif; display: block;  
 Inherited from html,fa-even...  
 html { main.css?v=45679023:28 }  
 color: #222; font-size: 9pt;

Tras esto, cogemos con BeautifulSoup la información de la web indicando es de tipo JavaScript

```
soup = bs(soup.text, 'lxml')
script_divs = soup.find_all('script', {'type': 'text/javascript'})
```

[illegible]

Primero cogemos solo la fila 5 que es donde está la información que queremos.





```

lista = [test.replace(']', ',').replace('\"', '\"') for test in lista]
lista

['149542560000,209',
'149603040000,69',
'149663520000,107',
'149724000000,105',
'149784480000,97',
'149844960000,165',
'149905440000,111',
'149965920000,109',
'150026400000,106',
'150086880000,138',
'150147360000,177',
'150216480000,89',
'150268320000,1027',
'150328800000,791',
'150389280000,337',
'150449760000,406',
'150518880000,357',
'150570720000,665',
'150631200000,1017',
'150691680000,1055']

```

Con esta función obtendremos una lista con la fecha y separado por una coma, el dato según las filas que se pongan.

Función final:

```

# Definimos una función para coger solo la parte que nos interesa
def lista_scrap(x, y):
    # Nos quedamos solo con la fila 5 de la lista ya que es donde esta la información que queremos
    cadena_ori = script_divs[5]
    # Creamos una lista separando por [ que es por donde empieza nuestra información
    cadena_sep = str(cadena_ori).split('[')
    # Seleccionamos las filas que queremos según si es seguidores, menos seguidores...
    lista = cadena_sep[x:y]
    # eliminamos ,]
    lista = [test.replace(']', ',').replace('\"', '\"') for test in lista]
    return lista

```

Ahora hacemos otra función para depurar la lista y hacer un data frame. En ella usamos la función anterior para obtener la lista. Explicamos paso a paso la función:

Creamos un data frame separando la lista por comas. Indicando que la posición 0 teniendo en cuenta la coma como separación, obtenemos la Semana...

```

df = pd.DataFrame()
df['Semana'] = [x.split(',')[0] for x in lista]
df

```

	Semana
0	149542560000
1	149603040000
2	149663520000
3	149724000000
4	149784480000
5	149844960000
6	149905440000
7	149965920000
8	150026400000

y la posición 1 son los seguidores.

```
df['Seguidores'] = [x.split(',')[1] for x in lista]
df
```

	Semana	Seguidores
0	1495425600000	209
1	1496030400000	69
2	1496635200000	107
3	1497240000000	105
4	1497844800000	97
5	1498449600000	165
6	1499054400000	111

Además, tenemos que depurar el formato de la fecha que está en epoch con tres ceros. Primero le quitamos los tres ceros, lo pasamos a numérico y hacemos un bucle para pasarlo a un formato de fecha normal.

```
df['Semana'] = df['Semana'].map(lambda x: str(x)[:3])
df['Semana'] = pd.to_numeric(df['Semana'])
for n in range(len(df)):
    df['Semana'][n] = datetime.datetime.fromtimestamp(df['Semana'][n]).strftime('%d-%m-%Y')
df
```

C:\Users\Dany\Anaconda2\lib\site-packages\ipykernel\_launcher.py:4: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html>  
after removing the cwd from sys.path.

C:\Users\Dany\Anaconda2\lib\site-packages\pandas\core\indexing.py:190: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html>  
self.\_setitem\_with\_indexer(indexer, value)

	Semana	Seguidores
0	22-05-2017	209
1	29-05-2017	69
2	05-06-2017	107
3	12-06-2017	105
4	19-06-2017	97

Por último, obtenemos el data frame depurado. Función final es la siguiente:

```
# Hacemos otra funcion para crear un dataframe con los datos de la anterior funcion
def lista_dataf(x, y):
    df = pd.DataFrame()
    #Usamos la funcion anterior para coger los datos
    lista = lista_scrap(x, y)
    # Señalamos la primera columna como la fecha y sacamos los datos de la lista
    df['Semana'] = [x.split(',')[0] for x in lista]
    # Señalamos la segunda como los seguidores
    df['Seguidores'] = [x.split(',')[1] for x in lista]
    #Ahora hay que pasar la fecha a un formato mas legible
    #Convertir epoch a fecha normal
    # tenemos que eliminar los tres ultimos ceros para pasarlo
    df['Semana'] = df['Semana'].map(lambda x: str(x)[:3])
    #pasamos la columna a numeros
    df['Semana'] = pd.to_numeric(df['Semana'])
    #pasamos a formato fecha
    for n in range(len(df)):
        df['Semana'][n] = datetime.datetime.fromtimestamp(df['Semana'][n]).strftime('%d-%m-%Y')
    return df
```

Ahora inspeccionamos la descarga de datos que hemos hecho en el scraping y vemos la fila 5 separándola por “[“ obteniendo una lista. Si analizamos los datos obtenidos vemos que de la fila 3 a la 182 son los datos de los nuevos seguidores, de la 185 a la 364 son los seguidos y de la 367 a la 546 es la media. Con esto ya podemos aplicar la función a cada una de las filas para obtener un data frame para cada uno de los datos.



```
# Si inspeccionamos Los datos vemos de donde a donde van cada uno de Los graficos
cadena_ori = script_divs[5]
cadena_sep = str(cadena_ori).split('[')

seguidores_ganados = lista_dataf(3, 182)
seguidos_nuevos = lista_dataf(185, 364)
media = lista_dataf(367, 546)
```

Ahora creamos el data frame definitivo para un partido juntando los datos de las tres bases.

```
# Juantamos Las tres variables
Base_conjunta = pd.DataFrame()
Base_conjunta['Semana'] = seguidores_ganados['Semana']
Base_conjunta['Seguidores'] = seguidores_ganados['Seguidores']
Base_conjunta['Seguidos'] = seguidos_nuevos['Seguidores']
Base_conjunta['Media'] = media['Seguidores']

#Vemos La cabecera de La base
Base_conjunta.head(5)
```

	Semana	Seguidores	Seguidos	Media
0	22-05-2017	209	5	113
1	29-05-2017	69	7	84
2	05-06-2017	107	9	121
3	12-06-2017	105	2	71
4	19-06-2017	97	8	127

Por último, guardamos la base:

```
# Guardamos La base en csv
Vox_Twitter = Base_conjunta
Vox_Twitter.to_csv('Vox_Twitter.csv', index=False)
```

Ahora tenemos que hacer el mismo proceso con el PP, PSOE y Podemos para tener la base completa con los cuatro partidos más relevantes y así poder comparar.

Se decide eliminar la media ya que no parece aportar mucha información.

```
# PSOE
soup = requests.get("https://socialblade.com/twitter/user/PSOE/monthly", headers=headers)
soup = bs(soup.text, 'lxml')
script_divs = soup.find_all('script', {'type': 'text/javascript'})

seguidores_ganados = lista_dataf(3, 230)
seguidos_nuevos = lista_dataf(233, 460)
media = lista_dataf(463, 690)

Base_conjunta = pd.DataFrame()
Base_conjunta['Semana'] = seguidores_ganados['Semana']
Base_conjunta['Seguidores'] = seguidores_ganados['Seguidores']
Base_conjunta['Seguidos'] = seguidos_nuevos['Seguidores']
Base_conjunta['Media'] = media['Seguidores']

PSOE_Twitter = Base_conjunta

PSOE_Twitter.to_csv('PSOE_Twitter.csv', index=False)
```

Una vez aplicado a los cuatro partidos, solo tenemos que juntar las cuatro bases y guardarla

```
# Juntamos las cuatro bases y dejamos solo seguidores y seguidos
Partidos_Twitter = pd.DataFrame()

# Eliminamos la media ya que no nos aporta
PSOE_Twitter.drop(['Media'], axis = 'columns', inplace=True)
Podemos_Twitter.drop(['Media'], axis = 'columns', inplace=True)
Vox_Twitter.drop(['Media'], axis = 'columns', inplace=True)
PP_Twitter.drop(['Media'], axis = 'columns', inplace=True)

# Al juntarlos vemos que hay algunas semanas que no coincide exactamente y no se juntan bien.
PSOE_Twitter.Semana[77] = '11-12-2017'
PP_Twitter.Semana[58] = '11-12-2017'
Vox_Twitter.Semana[11] = '07-08-2017'
Vox_Twitter.Semana[16] = '11-09-2017'
Vox_Twitter.Semana[29] = '18-12-2017'
Vox_Twitter.Semana[44] = '26-03-2018'

# Juntamos las bases en una conjunta
Partidos_Twitter = pd.merge(PSOE_Twitter, Podemos_Twitter, on='Semana', how='outer', suffixes=( '_PSOE', '_Podemos'))
Partidos_Twitter = pd.merge(Partidos_Twitter, PP_Twitter, on='Semana', how='outer', suffixes=( '_prueba', '_PP'))
Partidos_Twitter = pd.merge(Partidos_Twitter, Vox_Twitter, on='Semana', how='outer', suffixes=( '_PP', '_Vox'))

Partidos_Twitter.to_csv('Partidos_Politicos_Twitter.csv', index=False)
```

## Dataset

Al dataset se le da el nombre de “Partidos\_Politicos\_Twitter” y queda de la forma siguiente. Se muestra el final ya que en el inicio solo hay datos del PSOE que es el único partido que tiene datos desde 2016, mientras que el resto empiezan a tener datos en 2017.

Partidos_Twitter.tail(5)									
	Semana	Seguidores_PSOE	Seguidos_PSOE	Seguidores_Podemos	Seguidos_Podemos	Seguidores_PP	Seguidos_PP	Seguidores_Vox	Seguidos_Vox
223	21-09-2020	315	-3	224	-1	3277	2671	255	2
224	28-09-2020	355	-5	84	2	409	-11	1021	4
225	05-10-2020	435	-6	454	2	613	-2	623	0
226	12-10-2020	472	-1	493	3	652	-33	817	-3
227	19-10-2020	988	-6	992	6	-387	-234	2048	-1

El data set está compuesto por nueve variables:

- Semana que hace referencia a la fecha semanal de los datos de las otras variables.
- Seguidores del partido (hay una por cada partido), que son el número de nuevos seguidores que ha conseguido el partido esa semana.
- Seguidos del partido (hay una por cada partido), que indica el número de cuentas nuevas que sigue el partido en Twitter, este puede ser negativo en caso de dejar de seguir cuentas.

El data set contiene los datos del PSOE, PP, Podemos y Vox. En el momento del Scraping, la base tiene 228 filas. Temporalmente, comprenden desde junio de 2016 hasta octubre de 2020 para el PSOE, y desde mayo de 2017 para el resto de los partidos.

Como último detalle, publicamos el data set en Zanodo con una licencia de **Creative Commons Attribution 4.0 International** ya que los datos en si son de dominio público y no somos quien para poner otra licencia que no sea libre. También es importante que los datos sean públicos para una mayor difusión de conocimiento.

<https://zenodo.org/record/4247917#.X6RguVCCHcd>

Publication date:

November 5, 2020

DOI:

DOI 10.5281/zenodo.4247917

Keyword(s):

política, partidos políticos, Twitter

License (for files):

[Creative Commons Attribution 4.0 International](#)

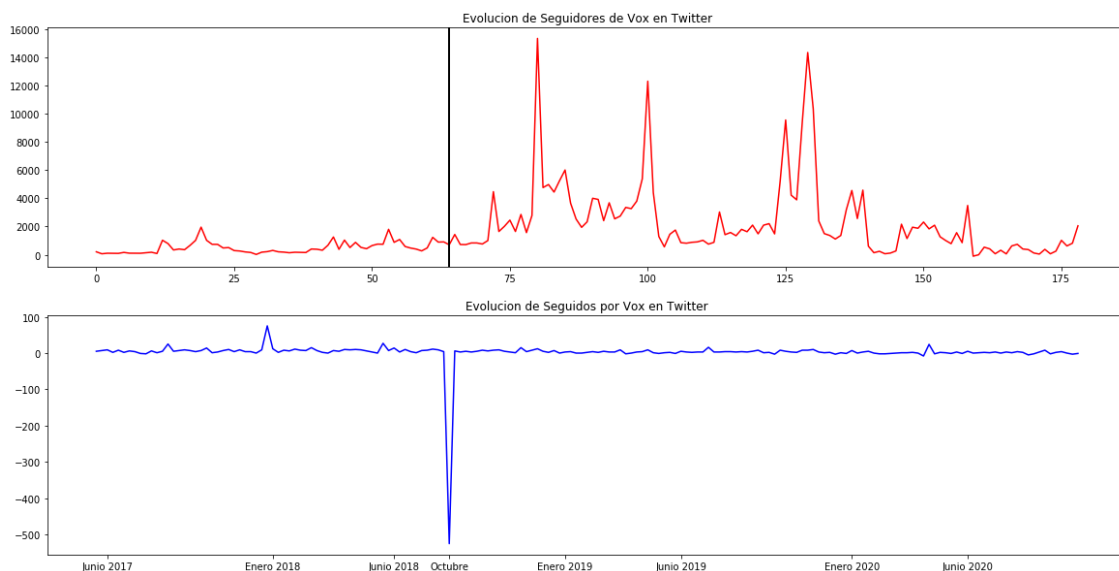
También cuando subimos el código y la base a Github le hemos puesto una licencia de **GNU General Public License v3.0**

[https://github.com/dnunezs/Partidos\\_Politicos\\_Twitter](https://github.com/dnunezs/Partidos_Politicos_Twitter)

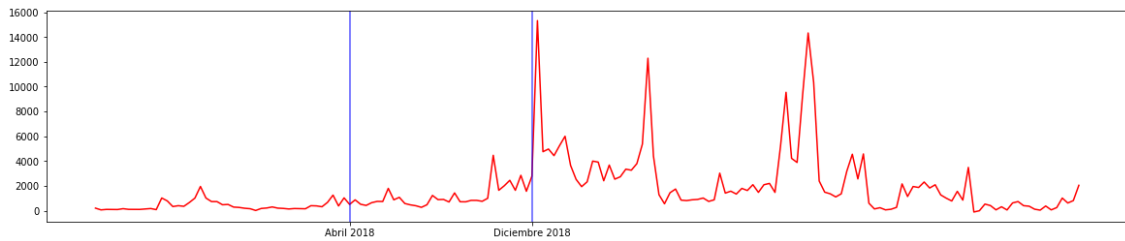
## Interés del conjunto de datos y análisis preliminar.

De un primer análisis podemos sacar las siguientes conclusiones que responden a las preguntas planteadas al principio.

Para ello, lo vamos a hacer con gráficos (el código de los gráficos está en otro archivo en GitHub). En el primero se puede ver la evolución histórica de los seguidores y seguidos de la cuenta oficial de Vox en Twitter desde 2017 hasta ahora. Está marcado con una línea negra un pico descendiente que tienen de cuentas que sigue Vox en octubre de 2018. Recordemos que la primera reunión de Vox y Bannon se produjo en abril de 2018. Antes de esto no hay movimientos reseñables ni en los seguidores ni en los seguidos. Y el único movimiento que vemos grande es que la cuenta deja de seguir a 500 cuentas en una misma semana de octubre de 2018.



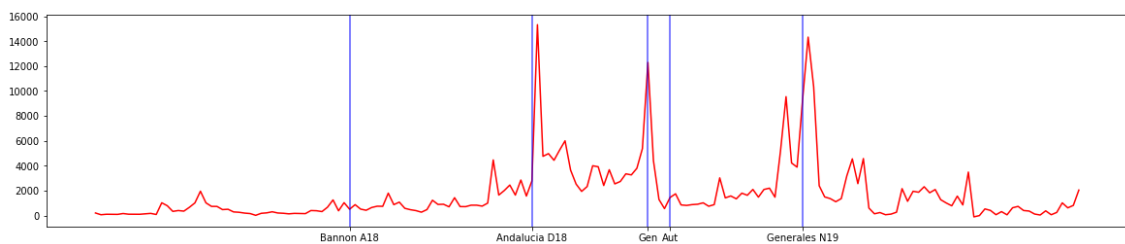
Ahora nos centraremos en los nuevos seguidores. Marcamos en el gráfico una línea en la semana del 10 abril de 2018 que es cuando se hizo pública la primera reunión con Bannon y otra en la semana del 2 de diciembre de 2018 que son las elecciones al Parlamento andaluz, momento en el cual Vox logra sus primeros escaños.



Desde 2017 hasta la llegada oficial de Bannon (abril 2018) no hay muchas variaciones en los seguidores que tiene Vox. Es desde marzo de 2018, cuando Vox empieza a tener picos de seguidores más continuos.

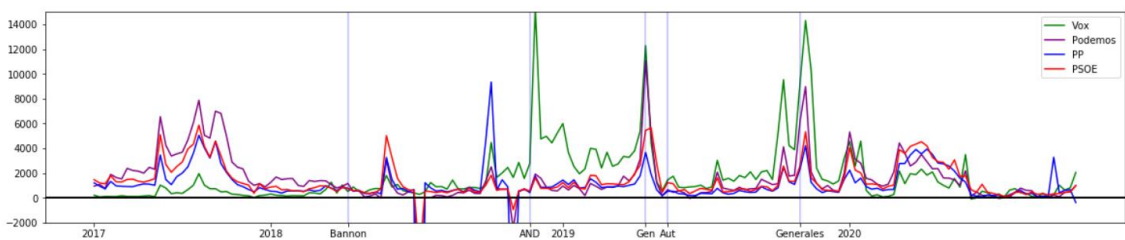
Se aprecia que, desde la llegada de Bannon hasta las elecciones en Andalucía, la cuenta de Vox empieza a tener muchos más nuevos seguidores semanales. Esto implica que, o por un lado Vox consigue en ese periodo llegar a la gente, ser conocidos, que se interesen por ellos y les empiecen a seguir, o por otro, que parte de la estrategia de Bannon sea crear miles de cuentas para difundir cuentas como las de Vox y su mensaje. Así, ampliar un mensaje y aparentar que tiene muchos seguidores. Lo más probable es que haya sido una combinación de las dos técnicas.

Nos podemos fijar más, poniendo el punto de mira en las elecciones, ya que es lógico que haya picos altos tras unas elecciones donde ganen escaños y salgan en las televisiones y les sigan.



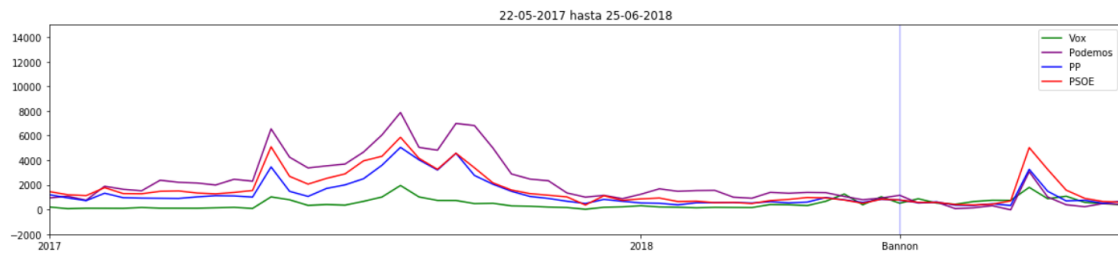
Si nos fijamos, los picos más altos se producen la misma o la siguiente semana de unas elecciones, algo que puede ser normal por la difusión en los medios.

Para analizar estos picos vamos a compararlo con el resto de los partidos. Si todos los partidos tienen picos en común, quiere decir que tienen una causa común: campaña electoral o elecciones.

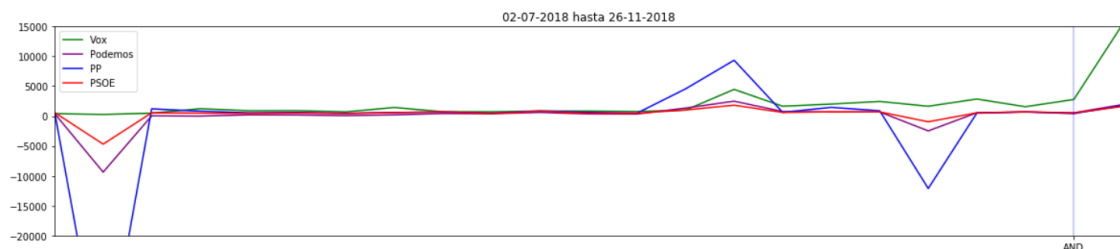


Con rasgos generales se puede ver una evolución muy similar en los cuatro partidos. Las subidas y bajadas son muy parecidas. Pero hay dos cosas que llaman la atención. La primera, desde el principio de los datos hasta un poco antes de las elecciones en Andalucía, la evolución de Vox es similar o inferior a la del resto de partidos, pero después y hasta 2020 tiene una evolución muy superior al resto. La segunda, es que hay periodos concretos donde Vox rompe la tendencia del resto de partidos.

Hagamos zoom para verlo mejor. Primero vemos desde 22-05-2017 hasta 25-06-2018:

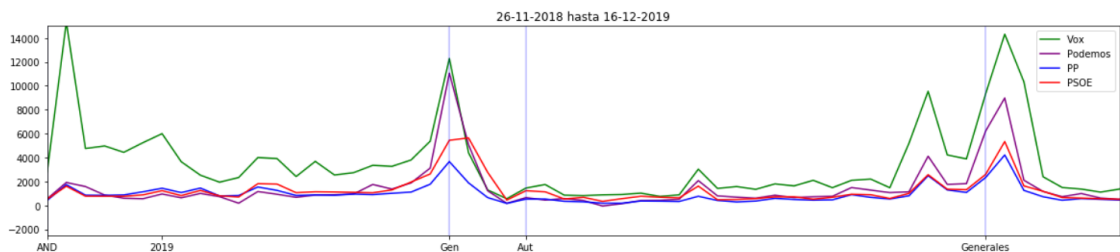


En este periodo la evolución es la misma en todos los partidos. Las subidas y bajadas siguen el mismo patrón. Pero desde entonces empiezan a haber diferencia con Vox. Vemos desde 02-07-2018 hasta 26-11-2018 justo el periodo anterior a las elecciones en Andalucía.



Aquí se ve como hay dos puntos donde todos los partidos pierden seguidores salvo Vox que no gana muchos, pero nunca llega a perder seguidores.

Por último, veamos el periodo desde 26-11-2018 hasta 16-12-2019 que comprende desde las elecciones andaluzas hasta las Generales de noviembre de 2019. Un periodo con muchas elecciones.



Podemos ver como la evolución de Vox es la misma que el resto de los partidos. Tienen las mismas subidas y bajadas de nuevos seguidores. Solo hay una diferencia significativa, la evolución de Vox es muy superior a la del resto de partidos justo antes de las elecciones y las semanas posteriores. En el periodo entre las andaluzas y justo antes de las primeras generales de 2019 hablamos de una media de 3.000 seguidores más a la semana que el partido que más gane cada semana, un total de 61.000 seguidores más que el resto.

Lo mismo ocurre los dos meses antes de las elecciones generales de noviembre de 2019. Vox gana una media de 2.500 nuevos seguidores cada semana más que el partido que más seguidores consigue. En total suma 26.000 nuevos seguidores más que el resto de los partidos.

En conclusión, aunque en un análisis preliminar si parece haber ciertos datos inflados artificialmente en comparación con el resto de los partidos, es imposible afirmarlo categóricamente. La única forma de comprobarlo sería analizar todos sus seguidores y ver hasta que punto hay bots o cuentas artificiales que no correspondan con una persona real de a pie.

# Autoría

El trabajo se ha realizado individualmente por Daniel Núñez Sanz.

Contribuciones	Firma
Investigación previa	DNS
Redacción de las respuestas	DNS
Desarrollo código	DNS