

Categorización de noticias skl NJ

December 19, 2019

1 Categorización de noticias

1.1 Preguntas a responder con el análisis

- ¿Se pueden catalogar las noticias con la descripción y los titulares? Compare su clasificación con las categorías incluidas en el set de datos.
- ¿Existen estilos de escritura asociados a cada categoría?
- ¿Qué se puede decir de los autores?
- ¿Qué información útil se puede extraer de los datos?

1.1.1 Análisis exploratorio y carga de los datos

```
[1]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
from tensorflow.contrib import learn
import tensorflow as tf
import string
from random import randint

[2]: df = pd.read_json("../Data/News_Category_Dataset_v2.json", lines = True,
                     convert_dates = True)

[3]: df.head()
```

```
authors      category      date \
0  Melissa Jeltsen      CRIME  2018-05-26
1    Andy McDonald  ENTERTAINMENT  2018-05-26
2      Ron Dicker  ENTERTAINMENT  2018-05-26
3      Ron Dicker  ENTERTAINMENT  2018-05-26
4      Ron Dicker  ENTERTAINMENT  2018-05-26

                                               headline \
0  There Were 2 Mass Shootings In Texas Last Week...
1  Will Smith Joins Diplo And Nicky Jam For The 2...
2  Hugh Grant Marries For The First Time At Age 57
3  Jim Carrey Blasts 'Castrato' Adam Schiff And D...
4  Julianna Margulies Uses Donald Trump Poop Bags...
```

```

link \
0 https://www.huffingtonpost.com/entry/texas-ama...
1 https://www.huffingtonpost.com/entry/will-smit...
2 https://www.huffingtonpost.com/entry/hugh-gran...
3 https://www.huffingtonpost.com/entry/jim-carre...
4 https://www.huffingtonpost.com/entry/julianna-...

short_description
0 She left her husband. He killed their children...
1 Of course it has a song.
2 The actor and his longtime girlfriend Anna Ebe...
3 The actor gives Dems an ass-kicking for not fi...
4 The "Dietland" actress said using the bags is ...

```

[4]: df.describe()

	authors	category	date	headline	\
count	200853	200853	200853	200853	
unique	27993	41	2309	199344	
top		POLITICS	2013-01-17 00:00:00	Sunday Roundup	
freq	36620	32739	100	90	
first	NAN	NAN	2012-01-28 00:00:00	NAN	
last	NAN	NAN	2018-05-26 00:00:00	NAN	

	link	short_description
count	200853	200853
unique	200812	178353
top	https://www.huffingtonpost.comhttps://www.wash...	
freq	2	19712
first	NAN	NAN
last	NAN	NAN

1.1.2 Variable category

[5]: df_temp = pd.DataFrame(pd.value_counts(df['category']))
df_temp['Categoria'] = df_temp.index # organizar top categorias
df_temp = df_temp.rename(columns={'category': 'Frecuencia', 'Categoria': 'Categoria'})
df_temp = df_temp.reset_index(drop = True)
df_temp.head(10)

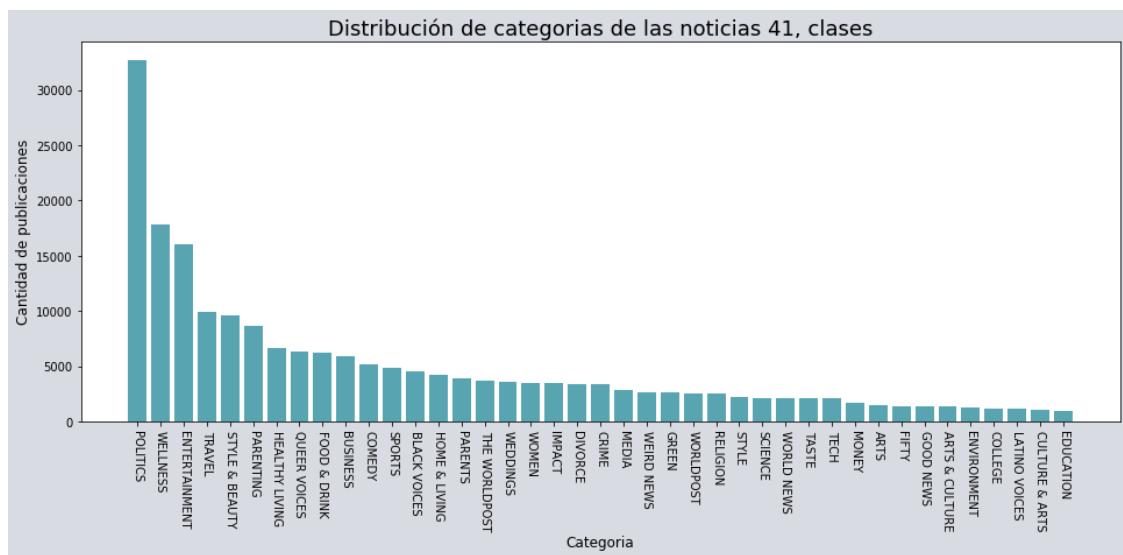
	Frecuencia	Categoria
0	32739	POLITICS
1	17827	WELLNESS
2	16058	ENTERTAINMENT
3	9887	TRAVEL
4	9649	STYLE & BEAUTY

```

5      8677      PARENTING
6      6694  HEALTHY LIVING
7      6314    QUEER VOICES
8      6226  FOOD & DRINK
9      5937      BUSINESS

```

```
[6]: plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.bar(df_temp.Categoría, df_temp.Frecuencia, color = '#58a4b0')
plt.xlabel("Categoría", fontsize = 12)
plt.xticks(rotation=270)
plt.ylabel("Cantidad de publicaciones", fontsize = 12)
plt.title("Distribución de categorías de las noticias {}, clases".format(len(df_temp)), fontsize = 18)
plt.show()
```



```
[7]: uni_c = df["category"].nunique()
len_c = len(df["category"])
len(df)
print("Se tienen {} categorías, y no tienen datos faltantes {} observaciones".
      format(uni_c, len_c))
```

Se tienen 41 categorías, y no tienen datos faltantes 200853 observaciones

Se observan algunos nombres similares por tanto se limpian para tener menos categorías y que el dataset este un poco mas balanceado

```
[8]: def category_group(x):
    if x == 'THE WORLDPOST':
        return 'WORLDPOST'
    elif x == 'PARENTING':
```

```

        return 'PARENTS'
    elif x == 'ARTS' or x == 'CULTURE & ARTS':
        return 'ARTS & CULTURE'
    elif x == 'STYLE':
        return 'STYLE & BEAUTY'
    elif x == 'COLLEGE':
        return 'EDUCATION'
    elif x == 'TASTE':
        return 'FOOD & DRINK'
    elif x == 'WORLD NEWS':
        return 'WORLDPOST'
    else:
        return x

df['category'] = df.category.apply(category_group)

```

[9]:

```

uni_c = df["category"].nunique()
len(df)
print("Al limpiar las categorias se tiene {} categorias".format(uni_c))

```

Al limpiar las categorias se tiene 33 categorias

[10]:

```

df_temp_2 = pd.DataFrame(pd.value_counts(df['category']))
df_temp_2['Categoria'] = df_temp_2.index # organizar top categorias
df_temp_2 = df_temp_2.rename(columns={'category': 'Frecuencia', 'Categoria': 'Categoria'})
df_temp_2 = df_temp_2.reset_index(drop = True)
df_temp_2.head(10)

```

	Frecuencia	Categoría
0	32739	POLITICS
1	17827	WELLNESS
2	16058	ENTERTAINMENT
3	12632	PARENTS
4	11903	STYLE & BEAUTY
5	9887	TRAVEL
6	8420	WORLDPOST
7	8322	FOOD & DRINK
8	6694	HEALTHY LIVING
9	6314	QUEER VOICES

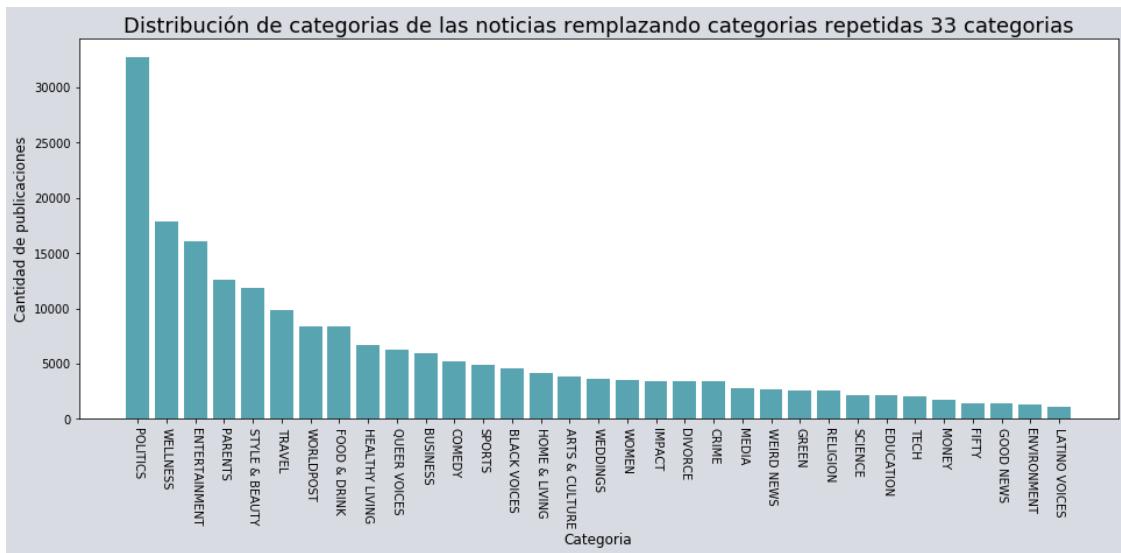
[11]:

```

plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.bar(df_temp_2.Categoría, df_temp_2.Frecuencia, color = '#58a4b0')
plt.xlabel("Categoria", fontsize = 12)
plt.xticks(rotation=270)
plt.ylabel("Cantidad de publicaciones", fontsize = 12)
plt.title("Distribución de categorias de las noticias remplazando categorias\u2192repetidas {} categorias".format(len(df_temp_2)),
          fontsize = 18)

```

```
plt.show()
```



```
[12]: c1 = len(df[df['category'] == 'ENTERTAINMENT'])
c2 = len(df[df['category'] == 'POLITICS'])
c3 = len(df[df['category'] == 'WORLDPOST'])
print("ENTERTAINMENT = {} noticias, POLITICS = {} noticias, WORLDPOST = {}".
      format(c1, c2, c3))
```

ENTERTAINMENT = 16058 noticias, POLITICS = 32739 noticias, WORLDPOST = 8420

```
[13]: p = (c1 + c2 + c3)/len(df)*100
p = np.round(p, 2)
print("Las categorias ENTERTAINMENT, POLITICS y WORLDPOST son el {} % de las
      noticias,".format(p))
print("un total de {} noticias".format(c1 + c2 + c3))
```

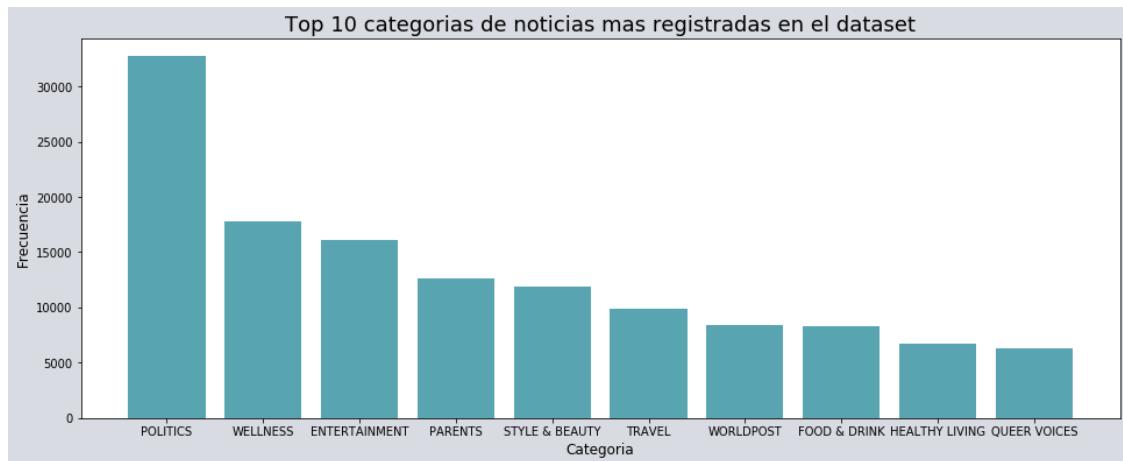
Las categorias ENTERTAINMENT, POLITICS y WORLDPOST son el 28.49 % de las noticias,
un total de 57217 noticias

La variable **category** tiene 34 categorías pero nos están distribuidas de manera uniforme ENTERTAINMENT, POLITICS y WORLDPOST tienen el 25.58% de todas las noticias

```
[14]: top_category = pd.DataFrame(pd.value_counts(df['category']))
top_category['Cat'] = top_category.index # organizar top categorias
top_category = top_category.reset_index(drop = True)
top_category.columns = ['Frecuencia', 'Categoria']
top_category = top_category.iloc[0:10]
top_category
```

```
[14]:   Frecuencia      Categoria
       0      32739      POLITICS
       1      17827      WELLNESS
       2      16058  ENTERTAINMENT
       3      12632      PARENTS
       4      11903  STYLE & BEAUTY
       5      9887       TRAVEL
       6      8420      WORLDPOST
       7      8322  FOOD & DRINK
       8      6694  HEALTHY LIVING
       9      6314  QUEER VOICES
```

```
[15]: plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.bar(top_category.Categoría, top_category.Frecuencia, color = '#58a4b0')
plt.xlabel('Categoría', fontsize = 12)
plt.ylabel('Frecuencia', fontsize = 12)
plt.title('Top 10 categorías de noticias mas registradas en el dataset', fontweight = 'bold', fontsize = 18)
plt.show()
```



1.1.3 Variable authors

```
[16]: df["authors"].nunique()
```

```
[16]: 27993
```

```
[17]: pd.value_counts(df['authors'])
```

```
[17]: 36620
```

Lee Moran
2423
Ron Dicker
1913

Reuters, Reuters
1562
Ed Mazza
1322
Cole Delbyck
1140
Andy McDonald
1068
Julia Brucculieri
1059
Carly Ledbetter
1054
Curtis M. Wong
1020
Mary Papenfuss
974
Bill Bradley
965
Dana Oliver
936
Sam Levine
893
David Moye
893
Michelle Manetti
876
Michelle Persad
875
Nina Golgowski
868
Igor Bobic
866
Ellie Krupnick
861
Dominique Mosbergen
784
Jamie Feldman
772
James Michael Nichols
764
Caroline Bologna
762
Rebecca Adams
753
Jenna Amatulli
711
Matthew Jacobs

702
Ryan Grenoble
698
Daniel Marans
669
Suzy Strutner
650

...

Cara E Bigony, Contributor\nDoctoral Candidate in Counseling Psychology at Fordham Univers... 1
Skip Shuda, ContributorDigital market strategizing, new world exploring, big picture dad 1
Siddiquah Greene, ContributorPretty girl with the fro and a mind like woah 1
Ali Reff, ContributorOwner of Alice & Wonder, a Chicago-based apparel and gifts shop 1
Rebecca Falconer and Lydia O'Connor 1
Laura Wattenberg, ContributorBaby Name Wizard 1
Idrees Ali and Megha Rajagopalan, Reuters 1
Krystal Hu, Yahoo Finance 1
Jim McKay, ContributorState Director, Prevent Child Abuse West Virginia 1
Kristin Wyatt, Associated Press 1
Liz Spikol, Contributor\nEditor-in-Chief, Curbed Philly 1
James Marshall Crotty, Contributor\nForbes Education Columnist; Author, 'How to Talk American'; Di... 1
Hiba Dlewati & Nawar Oliver, Syria Deeply 1
Brian Secemsky, M.D. and Joshua Liao, Contributors 1
Mica Rosenberg and Dan Levine, Reuters 1
Theresa Pickett, ContributorWriter 1
Donna Rockwell, ContributorClinical Psychologist, Mindfulness Teacher, Celebrity Mental H... 1
Mesfin Fekadu, AP 1
CafeMom, ContributorCafeMom is the meting place for moms 1
Melissa Sandgren, ContributorFeminist, Policy Nerd

```

1
Sarah Klein, Amanda L. Chan, Kate Bratskeir, and Gina Ryder, Contributors
1
Nicole Johnson, ContributorNicole is a fiction writer and blogger.
1
Dr. Rock Positano, Contributor\nHealth Columnist for the New York Daily News
1
Stephanie Sprenger, Contributorwriter, editor, blogger at Mommy, for Real
1
Paul Gale and John Trowbridge
1
Emily Stephenson and Luciana Lopez, Reuters
1
Jack Bouboushian, Courthouse News
1
Mariam Navaid Ottimofiore, ContributorExpat writer and blogger at 'And Then We
Moved To' 1
Maureen Greenwood-Basken, Contributor\nExecutive Director, Universal Access
Project of the United Nat... 1
Penny Phillips, ContributorUnofficial ambassador for young women striving to
create their... 1
Name: authors, Length: 27993, dtype: int64

```

[18]:

```

top_authors = pd.DataFrame(pd.value_counts(df['authors']))
top_authors['Aut'] = top_authors.index # organizar top categorias
top_authors = top_authors.reset_index(drop = True)
top_authors.columns = ['Frecuencia', 'authors']
top_authors = top_authors.iloc[0:11]
top_authors

```

[18]:

	Frecuencia	authors
0	36620	
1	2423	Lee Moran
2	1913	Ron Dicker
3	1562	Reuters, Reuters
4	1322	Ed Mazza
5	1140	Cole Delbyck
6	1068	Andy McDonald
7	1059	Julia Brucculieri
8	1054	Carly Ledbetter
9	1020	Curtis M. Wong
10	974	Mary Papenfuss

[19]:

```

plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.bar(top_authors.authors, top_authors.Frecuencia, color = '#58a4b0')
plt.xlabel('Autor', fontsize = 12)
plt.ylabel('Frecuencia', fontsize = 12)
plt.title('Top 10 autores de noticias registradas en el dataset', fontsize = 18)

```

```
plt.show()
```



```
[20]: pd.value_counts(df['authors'] == '')
```

```
[20]: False    164233  
True      36620  
Name: authors, dtype: int64
```

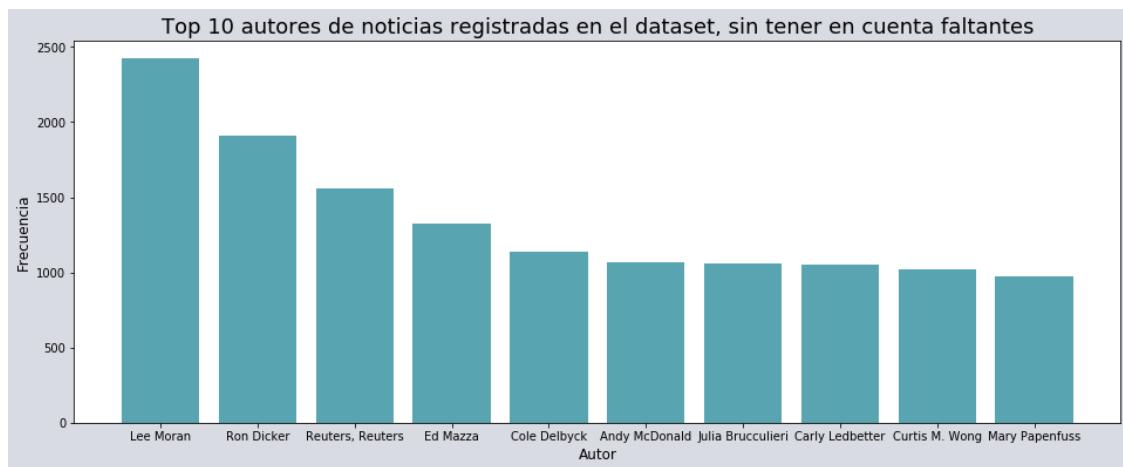
```
[ ]:
```

```
[21]: 36620/ len(df)
```

```
[21]: 0.18232239498538733
```

En el análisis se puede observar que se tienen **36620** registros sin el autor esto equivale al **18.23%** del total de los artículos por esta razón se procede a visualizar los datos sin tener en cuenta los valores que no tienen el autor

```
[22]: top_authors_nw = top_authors.iloc[1:12] # no se seleccionan los datos sin autor  
plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot  
plt.bar(top_authors_nw.authors, top_authors_nw.Frecuencia, color = '#58a4b0')  
plt.xlabel('Autor', fontsize = 12)  
plt.ylabel('Frecuencia', fontsize = 12)  
plt.title('Top 10 autores de noticias registradas en el dataset, sin tener en cuenta faltantes', fontsize = 18)  
plt.show()
```



1.1.4 Variable date

```
[23]: date = pd.DataFrame(df['date'])
date.head()
```

```
[23]: date
0 2018-05-26
1 2018-05-26
2 2018-05-26
3 2018-05-26
4 2018-05-26
```

```
[24]: type(date)
```

```
[24]: pandas.core.frame.DataFrame
```

```
[25]: date['Freq'] = 1
date.index = date['date']
date = date.drop(columns=['date'])
date.head()
```

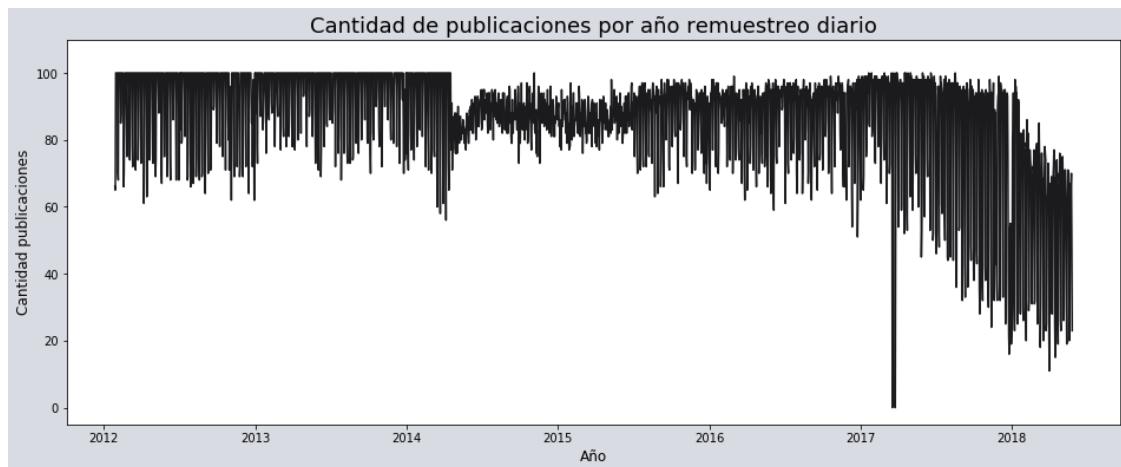
```
[25]: Freq
date
2018-05-26    1
2018-05-26    1
2018-05-26    1
2018-05-26    1
2018-05-26    1
```

```
[26]: date_res_d = date.resample('D').count() # resample por día
date_res_d.head()
```

```
[26]: Freq
date
```

```
2012-01-28    66
2012-01-29    65
2012-01-30   100
2012-01-31   100
2012-02-01   100
```

```
[27]: plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.plot(date_res_d, color = '#1b1b1e')
plt.xlabel('Año', fontsize = 12)
plt.ylabel('Cantidad publicaciones', fontsize = 12)
plt.title('Cantidad de publicaciones por año remuestreo diario', fontsize = 18)
plt.ylim([-5, 110])
plt.show()
```



```
[28]: date_res_m = date.resample('MS').count() # resample por mes
date_res_m.head()
```

```
[28]:      Freq
date
2012-01-01    331
2012-02-01  2693
2012-03-01  2880
2012-04-01  2769
2012-05-01  2899
```

```
[29]: plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.plot(date_res_m, color = '#1b1b1e')
plt.xlabel('Año', fontsize = 12)
plt.ylabel('Cantidad publicaciones', fontsize = 12)
plt.title('Cantidad de publicaciones por año remuestreo mensual', fontsize = 18)
# plt.ylim([100, 700])
plt.show()
```



Frecuencias de la variable date por mes y por día

```
[30]: # buscar frecuencia por mes y por día
date_t = pd.DataFrame(df['date'])
```

```
[31]: date_t['Mes'] = pd.DatetimeIndex(date_t['date']).month_name()
date_t.head()
```

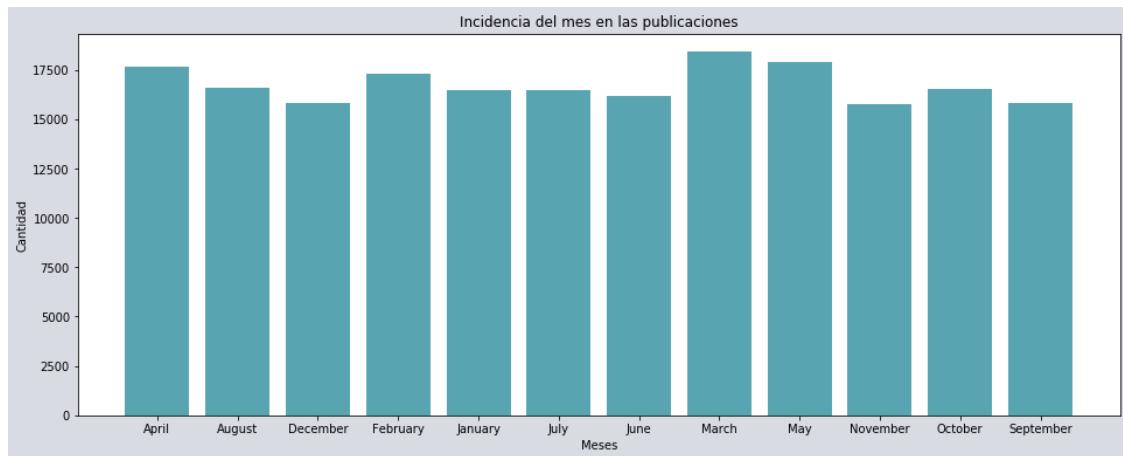
```
[31]:      date   Mes
0  2018-05-26  May
1  2018-05-26  May
2  2018-05-26  May
3  2018-05-26  May
4  2018-05-26  May
```

```
[32]: date_t_group = date_t
date_t_group = date_t_group.groupby(['Mes']).count()
date_t_group['Meses'] = date_t_group.index
date_t_group = date_t_group.reset_index(drop = True)
date_t_group
```

```
[32]:      date      Meses
0    17635     April
1    16576    August
2    15816 December
3    17286 February
4    16493 January
5    16470     July
6    16183     June
7    18418     March
8    17917     May
9    15756 November
10   16505 October
```

```
11 15798 September
```

```
[33]: plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.bar(date_t_group.Meses, date_t_group.date, color = '#58a4b0')
plt.xlabel('Meses')
plt.ylabel('Cantidad')
plt.title('Incidencia del mes en las publicaciones')
plt.show()
```



```
[34]: mayor = date_t_group.date.max()
menor = date_t_group.date.min()
rango = mayor - menor
print("El mayor valor de publicaciones sumando los meses es: {} y el menor es :{} para un rango de {}".
      format(mayor, menor, rango))
```

El mayor valor de publicaciones sumando los meses es: 18418 y el menor es : 15756, para un rango de 2662

Con respecto a esta variable no se puede ver una estacionalidad clara ni un mes en el que sea mayor la incidencia de publicaciones el rango es de 2662

1.1.5 Variables authors & category

```
[35]: top_authors_nw.head(10)
```

```
[35]:   Frecuencia      authors
1        2423      Lee Moran
2        1913    Ron Dicker
3        1562  Reuters, Reuters
4        1322       Ed Mazza
5        1140     Cole Delbyck
6        1068  Andy McDonald
```

```

7      1059 Julia Brucculieri
8      1054 Carly Ledbetter
9      1020 Curtis M. Wong
10     974 Mary Papenfuss

```

[36]: aut_cat = pd.merge(top_authors_nw, df) # joint de los datos en funcion de el
 ↪top 10 de los autores

[37]: top = len(aut_cat)
all_pub = len(df)
aut_porcent = np.round(((top/all_pub) * 100), 2)
print("El top 10 de autores producen el {}% del total de los articulos".
 ↪format(aut_porcent))

El top 10 de autores producen el 6.74% del total de los articulos

[38]: aut_cat.Frecuencia = 1
aut_cat.head()

[38]: Frecuencia authors category date \
0 1 Lee Moran COMEDY 2018-05-24
1 1 Lee Moran COMEDY 2018-05-24
2 1 Lee Moran POLITICS 2018-05-24
3 1 Lee Moran COMEDY 2018-05-24
4 1 Lee Moran MEDIA 2018-05-24

 headline \
0 'Late Night' Writer's Breathless Royal Wedding...
1 Seth Meyers Gives Donald Trump Some Valuable M...
2 Chrissy Teigen Taunts Donald Trump Over Twitte...
3 Samantha Bee Torchies ICE: 'Let's Shut It The F...
4 Jake Tapper Shreds Donald Trump With A Long Li...

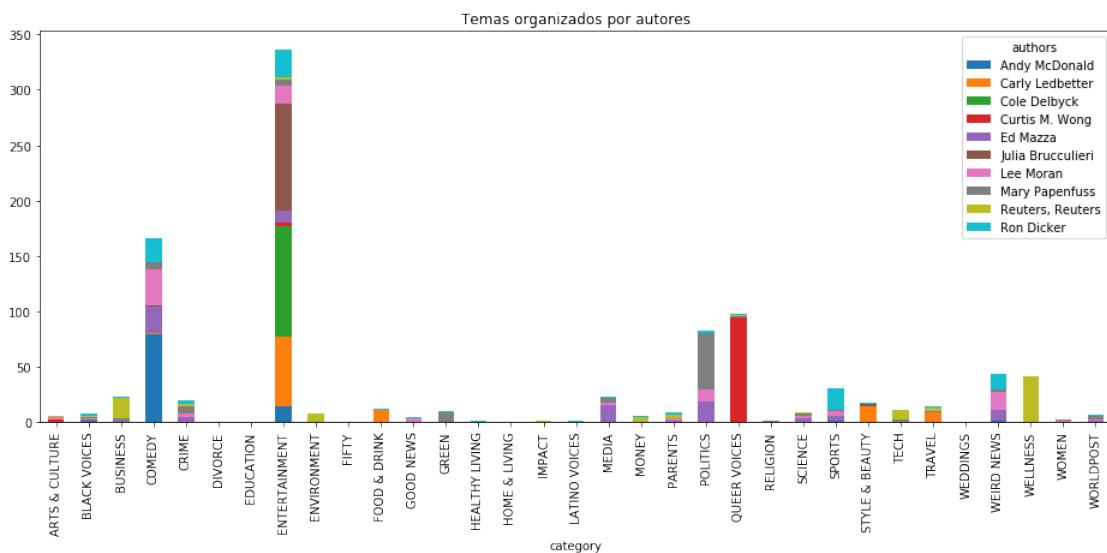
 link \
0 https://www.huffingtonpost.com/entry/royal-wed...
1 https://www.huffingtonpost.com/entry/seth-meye...
2 https://www.huffingtonpost.com/entry/chrissy-t...
3 https://www.huffingtonpost.com/entry/samantha-...
4 https://www.huffingtonpost.com/entry/jake-tapp...

 short_description
0 Then he bites his lip, like, Hmm, girl.
1 Stop trying to make Fetch happen."
2 Other tweeters, including Mexico's former pres...
3 "Im serious. Its awful. And we dont actuall...
4 "I could go on, but this is just an hour show."

Con respecto a estas dos variables se puede ver que el top de los 10 primeros autores en publicar artículos solo cubren el 6.74% de todas las publicaciones, el peso de las publicaciones que no

tienen autores es del 18.23% con 36620 publicaciones con la observación de la variable authors en blanco, esto hace que el 75.03% de las publicaciones tengan una variabilidad muy alta en cuanto a autores lo cual dificulta analizar esta variable

```
[39]: plot = pd.crosstab(index = aut_cat['category'],
                       columns = aut_cat['authors'] # tabla de frecuencias por categoria
                      ).apply(lambda r: r/r.sum() * 100, # promedio ponderado
                             axis=0).plot(kind='bar', stacked=True,
                                         figsize = (16, 6),
                                         title = "Temas organizados por autores"
                                         )
```



Se puede observar en los datos que los temas donde se publican mas noticias no son los mismo de los autores que mas publican, el top 10 de autores en publicar publican en la categoria ENTERTAINMENT, esto hace mas desbalanceados los datos y al publicar un articulo en con estos datos la mayor probabilidad es que sea de politica pero no de un autor del top 10

1.1.6 Variables authors, category & date

```
[40]: aut_cat['Año'] = pd.DatetimeIndex(aut_cat['date']).year
```

```
[41]: aut_cat.head()
```

```
[41]:   Frecuencia    authors category      date \
0          1  Lee Moran    COMEDY 2018-05-24
1          1  Lee Moran    COMEDY 2018-05-24
2          1  Lee Moran  POLITICS 2018-05-24
3          1  Lee Moran    COMEDY 2018-05-24
4          1  Lee Moran     MEDIA 2018-05-24
```

```

headline \
0 'Late Night' Writer's Breathless Royal Wedding...
1 Seth Meyers Gives Donald Trump Some Valuable M...
2 Chrissy Teigen Taunts Donald Trump Over Twitte...
3 Samantha Bee Torches ICE: 'Let's Shut It The F...
4 Jake Tapper Shreds Donald Trump With A Long Li...

link \
0 https://www.huffingtonpost.com/entry/royal-wed...
1 https://www.huffingtonpost.com/entry/seth-meye...
2 https://www.huffingtonpost.com/entry/chrissy-t...
3 https://www.huffingtonpost.com/entry/samantha-...
4 https://www.huffingtonpost.com/entry/jake-tapp...

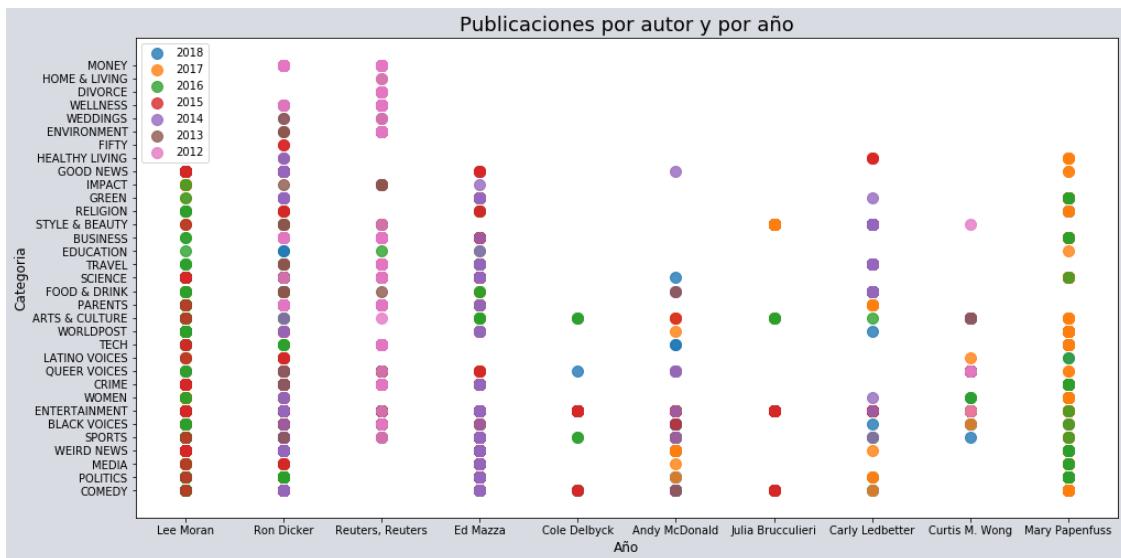
short_description   Año
0      Then he bites his lip, like, Hmm, girl.  2018
1      Stop trying to make Fetch happen."  2018
2 Other tweeters, including Mexico's former pres...  2018
3 "Im serious. Its awful. And we dont actuall...  2018
4 "I could go on, but this is just an hour show."  2018

```

[42]: Años_pub = aut_cat.Año.unique()
Años_pub

[42]: array([2018, 2017, 2016, 2015, 2014, 2013, 2012])

[43]: plt.figure(figsize=(16,8), facecolor = '#d8dbe2') # tamaño del plot
for i in Años_pub:
f = aut_cat[aut_cat['Año'] == i]
plt.scatter(f.authors, f.category, label = i, s = 100, alpha = 0.8)
plt.legend(loc = 'upper left')
plt.title("Publicaciones por autor y por año", fontsize = 18)
plt.xlabel("Año", fontsize = 12)
plt.ylabel("Categoria", fontsize = 12)



Con respecto a estas tres variables se puede observar como los autores que mas publican no son tan especializados y algunos no son tan activos en los diferentes años

1.2 Procesamiento de lenguaje natural de los datos

1.2.1 Alistamiento de los datos para determinar categoría según descripción y los titulares

```
[44]: df_class = df[['category', 'headline', 'short_description']]
df_class.head()
```

```
[44]: category                                              headline \
0      CRIME      There Were 2 Mass Shootings In Texas Last Week...
1  ENTERTAINMENT  Will Smith Joins Diplo And Nicky Jam For The 2...
2  ENTERTAINMENT      Hugh Grant Marries For The First Time At Age 57
3  ENTERTAINMENT  Jim Carrey Blasts 'Castrato' Adam Schiff And D...
4  ENTERTAINMENT  Julianna Margulies Uses Donald Trump Poop Bags...

                                                short_description
0  She left her husband. He killed their children...
1                      Of course it has a song.
2  The actor and his longtime girlfriend Anna Ebe...
3  The actor gives Dems an ass-kicking for not fi...
4  The "Dietland" actress said using the bags is ...
```

```
[192]: df_class_joint = df_class
```

```
[193]: p = df_class.headline + ' ' + df_class.short_description
# p = df_class.headline
df_class_joint['text'] = p
df_class_joint = df_class_joint.drop(columns=['headline', 'short_description'])
df_class_joint.head()
```

```
/home/oscar/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:3:  
SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

This is separate from the ipykernel package so we can avoid doing imports until

```
[193]:      category          text  
0        CRIME  There Were 2 Mass Shootings In Texas Last Week...  
1 ENTERTAINMENT  Will Smith Joins Diplo And Nicky Jam For The 2...  
2 ENTERTAINMENT  Hugh Grant Marries For The First Time At Age 5...  
3 ENTERTAINMENT  Jim Carrey Blasts 'Castrato' Adam Schiff And D...  
4 ENTERTAINMENT  Julianna Margulies Uses Donald Trump Poop Bags...
```

```
[194]: len(df_class_joint)
```

```
[194]: 200853
```

```
[195]: df_class_joint.head()
```

```
[195]:      category          text  
0        CRIME  There Were 2 Mass Shootings In Texas Last Week...  
1 ENTERTAINMENT  Will Smith Joins Diplo And Nicky Jam For The 2...  
2 ENTERTAINMENT  Hugh Grant Marries For The First Time At Age 5...  
3 ENTERTAINMENT  Jim Carrey Blasts 'Castrato' Adam Schiff And D...  
4 ENTERTAINMENT  Julianna Margulies Uses Donald Trump Poop Bags...
```

```
[196]: #limpieza de los datos  
df_class_joint.text = [x.lower() for x in df_class_joint.text] # Convertir los  
→textos a minúscula  
df_class_joint.text = [''.join(c for c in x if c not in string.punctuation) for  
→x in texts] # Eliminamos signos de puntuación  
df_class_joint.text = [''.join(c for c in x if c not in '0123456789') for x in  
→texts] # Eliminamos los números  
df_class_joint.text = [' '.join(x.split()) for x in texts] # Eliminar espacios  
→en blanco y separadores extras
```

```
[197]: df_class_joint.head()
```

```
[197]:      category          text  
0        CRIME  there were mass shootings in texas last week b...  
1 ENTERTAINMENT  will smith joins diplo and nicky jam for the w...  
2 ENTERTAINMENT  hugh grant marries for the first time at age t...  
3 ENTERTAINMENT  jim carrey blasts castrato adam schiff and dem...  
4 ENTERTAINMENT  julianna margulies uses donald trump poop bags...
```

1.2.2 Modelo Naive Bayes

```
[198]: from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.naive_bayes import ComplementNB

[199]: X_train, X_test, y_train, y_test = train_test_split(df_class_joint['text'],
                                                       df_class_joint['category'],
                                                       random_state = 137,
                                                       train_size = 0.8)
print("X_train {} y X_test {}".format(len(X_train), len(X_test)))
```

```
/home/oscar/anaconda3/lib/python3.7/site-
packages/sklearn/model_selection/_split.py:2179: FutureWarning: From version
0.21, test_size will always complement train_size unless both are specified.
FutureWarning)
```

```
X_train 160682 y X_test 40171
```

```
[200]: %%time
count_vect = CountVectorizer()
X_train_counts = count_vect.fit_transform(X_train)
tfidf_transformer = TfidfTransformer()
X_train_tfidf = tfidf_transformer.fit_transform(X_train_counts)
clf = ComplementNB(alpha = 0.8).fit(X_train_tfidf, y_train)
```

```
CPU times: user 17.6 s, sys: 110 ms, total: 17.7 s
Wall time: 17.7 s
```

```
[201]: probe = pd.DataFrame(df_class_joint.text[X_test])
```

```
[202]: clf.get_params
```

```
[202]: <bound method BaseEstimator.get_params of ComplementNB(alpha=0.8,
   class_prior=None, fit_prior=True, norm=False)>
```

```
[203]: pos = randint(0, len(probe))
probe.index[pos]
```

```
[203]: 'lessons from the fast lane i was in the mood to try something entirely new and
different something that would shake up my routine and give me a fresh
perspective on life zipping around a racecourse at more than miles an hour
seemed like just the thing provided i didnt kill myself'
```

```
[204]: pred = clf.predict(count_vect.transform([probe.index[pos]]))
pred[0]
```

```
[204]: 'WELLNESS'
```

```
[205]: df_class_joint[df_class_joint['text'] == probe.index[pos]]
```

```
[205]: category text  
168182 WELLNESS lessons from the fast lane i was in the mood t...
```

```
[206]: %%time  
pred_model = []  
for i in range(len(X_test)):  
    p = clf.predict(count_vect.transform([X_test.iloc[i]]))  
    pred_model.append(p[0])  
    if(i)%1000 == 0:  
        print("Iteración {}".format(i))
```

```
Iteración 0  
Iteración 1000  
Iteración 2000  
Iteración 3000  
Iteración 4000  
Iteración 5000  
Iteración 6000  
Iteración 7000  
Iteración 8000  
Iteración 9000  
Iteración 10000  
Iteración 11000  
Iteración 12000  
Iteración 13000  
Iteración 14000  
Iteración 15000  
Iteración 16000  
Iteración 17000  
Iteración 18000  
Iteración 19000  
Iteración 20000  
Iteración 21000  
Iteración 22000  
Iteración 23000  
Iteración 24000  
Iteración 25000  
Iteración 26000  
Iteración 27000  
Iteración 28000  
Iteración 29000  
Iteración 30000  
Iteración 31000  
Iteración 32000  
Iteración 33000  
Iteración 34000  
Iteración 35000  
Iteración 36000
```

```

Iteración 37000
Iteración 38000
Iteración 39000
Iteración 40000
CPU times: user 12min 18s, sys: 1.96 s, total: 12min 20s
Wall time: 12min 24s

```

[207]: `len(pred_model)`

[207]: 40171

[208]: `Acc_NB = pd.DataFrame()`

[209]: `Acc_NB['Real'] = y_test`
`Acc_NB['Predicho'] = pred_model`

[210]: `Acc_NB.head()`

	Real	Predicho
198274	ENTERTAINMENT	ENTERTAINMENT
9281	POLITICS	POLITICS
85726	WORLDPOST	POLITICS
156506	WELLNESS	WELLNESS
72313	LATINO VOICES	STYLE & BEAUTY

[211]: `Acierto = []`
`for i in range(len(Acc_NB)):`
 `if Acc_NB.Real.iloc[i] == Acc_NB.Predicho.iloc[i]:`
 `n = 1`
 `else:`
 `n = 0`
 `Acierto.append(n)`

[212]: `Acc_NB['Acierto'] = Acierto`

[213]: `Acc_NB.tail()`

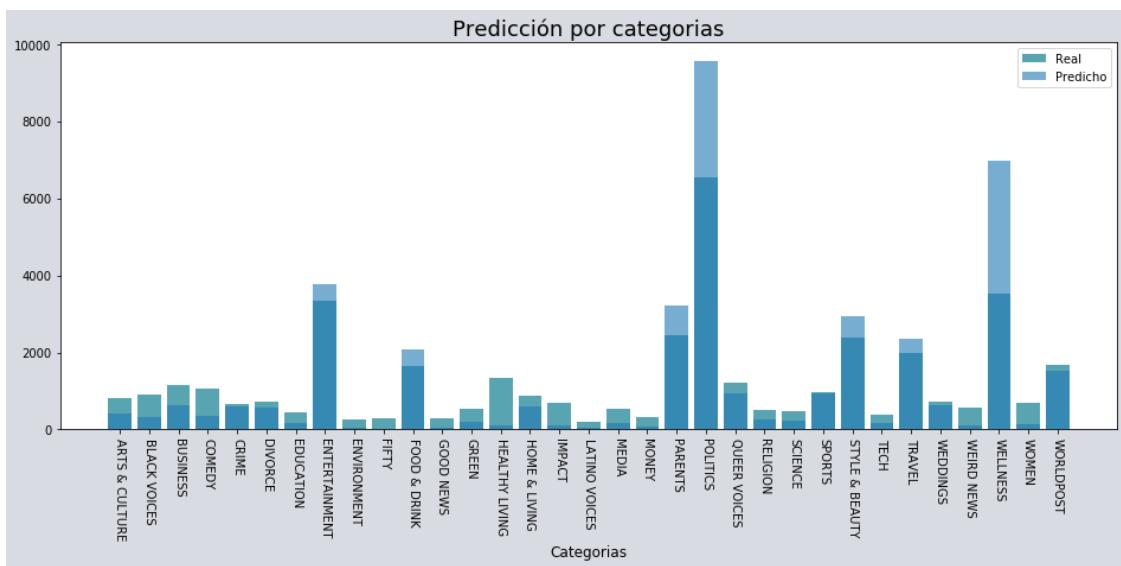
	Real	Predicho	Acierto
73859	POLITICS	POLITICS	1
6393	IMPACT	POLITICS	0
92933	SCIENCE	SCIENCE	1
4603	WORLDPOST	POLITICS	0
168889	WELLNESS	WELLNESS	1

[215]: `Acc_mod_NB = np.mean(Acc_NB.Acierto)`
`print("La efectividad del modelo es de {}%".format(np.round(Acc_mod_NB * 100, 2)))`

La efectividad del modelo es de 61.32%

[257]: `Acc_NB_Plot_1 = Acc_NB.groupby(['Real']).count()`
`Acc_NB_Plot_2 = Acc_NB.groupby(['Predicho']).count()`

```
[262]: plt.figure(figsize=(16,6), facecolor = '#d8dbe2') # tamaño del plot
plt.bar(Acc_NB_Plot_1.index, Acc_NB_Plot_1.Acierto, label = 'Real', color = '#58a4b0')
plt.bar(Acc_NB_Plot_2.index, Acc_NB_Plot_2.Acierto, label = 'Predicho', alpha = 0.6)
plt.xlabel('Categorías', fontsize = 12)
plt.title('Predicción por categorías', fontsize = 18)
plt.xticks(rotation=270)
plt.legend()
plt.show()
```



Se puede observar como el modelo sobre clasifica algunas categorias y otras no, para próximos análisis es bueno en base a esta información ajustar los pesos, de la categoria del modelo

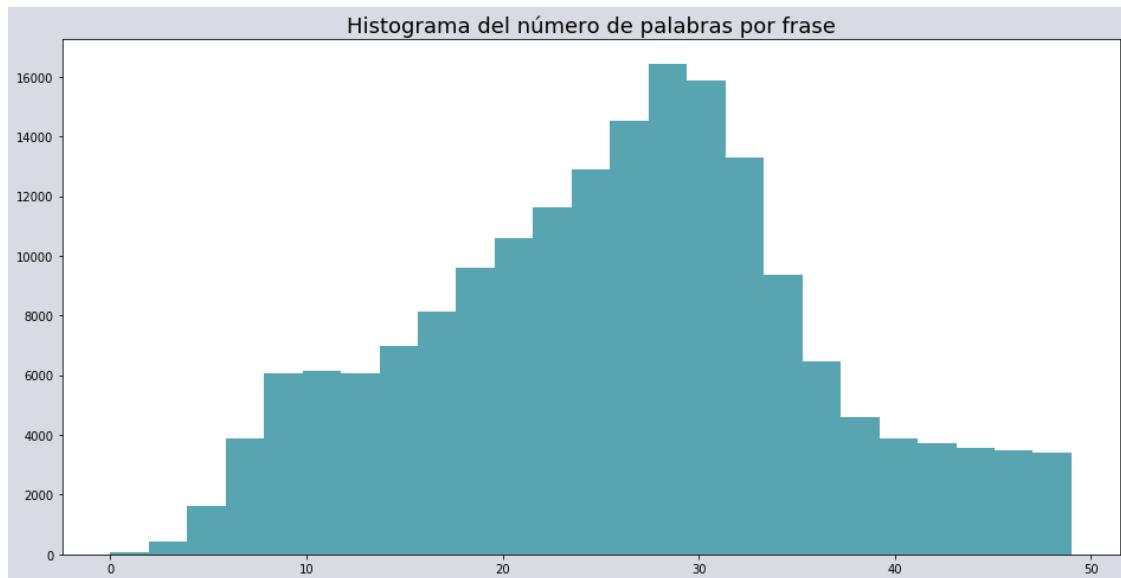
1.2.3 Modelo Linear SVC

```
[65]: texts = df_class_joint.text
df_class_joint.head()
```

```
[65]:      category          text
0        CRIME  There Were 2 Mass Shootings In Texas Last Week...
1  ENTERTAINMENT  Will Smith Joins Diplo And Nicky Jam For The 2...
2  ENTERTAINMENT  Hugh Grant Marries For The First Time At Age 5...
3  ENTERTAINMENT  Jim Carrey Blasts 'Castrato' Adam Schiff And D...
4  ENTERTAINMENT  Julianna Margulies Uses Donald Trump Poop Bags...
```

```
[67]: plt.figure(figsize=(16,8), facecolor = '#d8dbe2') # tamaño del plot
texts_lengths = [len(x.split()) for x in texts]
texts_lengths = [x for x in texts_lengths if x < 50]
plt.hist(texts_lengths, bins=25, color = '#58a4b0' )
```

```
plt.title("Histograma del número de palabras por frase", fontsize = 18)
plt.show()
```



```
[68]: sentence_size = 48
min_word_freq = 3
```

```
[69]: vocab_processor = learn.preprocessing.VocabularyProcessor(sentence_size,
                                                               →min_frequency=min_word_freq)
vocab_processor.fit_transform(texts)
transformed_texts = np.array([x for x in vocab_processor.transform(texts)])
embedding_size = len(np.unique(transformed_texts))
```

```
WARNING:tensorflow:From <ipython-input-69-1c35fd779a57>:1:
VocabularyProcessor.__init__ (from
tensorflow.contrib.learn.python.learn.preprocessing.text) is deprecated and will
be removed in a future version.
Instructions for updating:
Please use tensorflow/transform or tf.data.
WARNING:tensorflow:From /home/oscar/anaconda3/lib/python3.7/site-
packages/tensorflow/contrib/learn/python/learn/preprocessing/text.py:154:
CategoricalVocabulary.__init__ (from
tensorflow.contrib.learn.python.learn.preprocessing.categorical_vocabulary) is
deprecated and will be removed in a future version.
Instructions for updating:
Please use tensorflow/transform or tf.data.
WARNING:tensorflow:From /home/oscar/anaconda3/lib/python3.7/site-
packages/tensorflow/contrib/learn/python/learn/preprocessing/text.py:170:
tokenizer (from tensorflow.contrib.learn.python.learn.preprocessing.text) is
deprecated and will be removed in a future version.
```

Instructions for updating:
Please use tensorflow/transform or tf.data.

```
[70]: transformed_texts # texto como bag of words
```

```
[70]: array([[ 70,    73, 1281, ...,     0,     0,     0],  
           [ 42, 1570, 2859, ...,     0,     0,     0],  
           [ 5744, 5103, 7760, ...,     0,     0,     0],  
           ...,  
           [ 4009,    75, 4942, ...,     0,     0,     0],  
           [28704, 1570, 1046, ...,     0,     0,     0],  
           [15537, 4538, 3574, ...,     0,     0,     0]])
```

```
[71]: transformed_texts.shape
```

```
[71]: (200853, 48)
```

```
[72]: transformed_texts.size
```

```
[72]: 9640944
```

```
[73]: category_text = np.array(df_class_joint.category)  
category_text
```

```
[73]: array(['CRIME', 'ENTERTAINMENT', 'ENTERTAINMENT', ..., 'SPORTS', 'SPORTS',  
         'SPORTS'], dtype=object)
```

```
[74]: from sklearn.preprocessing import OneHotEncoder  
from sklearn.preprocessing import LabelEncoder  
  
# encoder entero  
label_encoder = LabelEncoder()  
integer_encoded = label_encoder.fit_transform(category_text)  
  
# encoder binario  
onehot_encoder = OneHotEncoder(sparse=False)  
integer_encoded = integer_encoded.reshape(len(integer_encoded), 1)  
target = onehot_encoder.fit_transform(integer_encoded)
```

```
/home/oscar/anaconda3/lib/python3.7/site-  
packages/sklearn/preprocessing/_encoders.py:368: FutureWarning: The handling of  
integer data will change in version 0.22. Currently, the categories are  
determined based on the range [0, max(values)], while in the future they will be  
determined based on the unique values.  
If you want the future behaviour and silence this warning, you can specify  
"categories='auto'".  
In case you used a LabelEncoder before this OneHotEncoder to convert the  
categories to integers, then you can now use the OneHotEncoder directly.  
warnings.warn(msg, FutureWarning)
```

```
[75]: integer_encoded[0]
```

```
[75]: array([4])
[76]: label_encoder.classes_
[76]: array(['ARTS & CULTURE', 'BLACK VOICES', 'BUSINESS', 'COMEDY', 'CRIME',
   'DIVORCE', 'EDUCATION', 'ENTERTAINMENT', 'ENVIRONMENT', 'FIFTY',
   'FOOD & DRINK', 'GOOD NEWS', 'GREEN', 'HEALTHY LIVING',
   'HOME & LIVING', 'IMPACT', 'LATINO VOICES', 'MEDIA', 'MONEY',
   'PARENTS', 'POLITICS', 'QUEER VOICES', 'RELIGION', 'SCIENCE',
   'SPORTS', 'STYLE & BEAUTY', 'TECH', 'TRAVEL', 'WEDDINGS',
   'WEIRD NEWS', 'WELLNESS', 'WOMEN', 'WORLDPOST'], dtype=object)
[77]: target[0] # categorias como one hot
[77]: array([0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
   0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.])
[78]: target.shape
[78]: (200853, 33)
```

1.2.4 Modelo Linear SVC

```
[131]: from sklearn.model_selection import train_test_split
from sklearn.svm import LinearSVC

model = LinearSVC(verbose = 1, max_iter = 300)

[132]: X_train, X_test, y_train, y_test = train_test_split(transformed_texts,
                                                       df_class_joint.category,
                                                       test_size = 0.2,
                                                       random_state = 137)

[133]: %%time
model.fit(X_train, y_train)
y_pred = model.predict(X_test)

[LibLinear]CPU times: user 20min 43s, sys: 2.73 s, total: 20min 45s
Wall time: 20min 52s

/home/oscar/anaconda3/lib/python3.7/site-packages/sklearn/svm/base.py:922:
ConvergenceWarning: Liblinear failed to converge, increase the number of
iterations.
"the number of iterations.", ConvergenceWarning)
```

```
[134]: y_pred.shape
[134]: (40171,)
[135]: y_pred
[135]: array(['WORLDPOST', 'STYLE & BEAUTY', 'POLITICS', ..., 'COMEDY',
   'ENTERTAINMENT', 'PARENTS'], dtype=object)
```

```
[136]: Acc_SVC = pd.DataFrame()
[137]: Acc_SVC['Real'] = y_test
Acc_SVC['Predicho'] = y_pred
[138]: Acc_SVC.head()
[138]:
          Real      Predicho
198274  ENTERTAINMENT    WORLDPOST
9281      POLITICS  STYLE & BEAUTY
85726     WORLDPOST      POLITICS
156506     WELLNESS      POLITICS
72313   LATINO VOICES    WORLDPOST
```

```
[139]: Acierto_SVC = []
for i in range(len(Acc_SVC)):
    if Acc_SVC.Real.iloc[i] == Acc_SVC.Predicho.iloc[i]:
        n = 1
    else:
        n = 0
    Acierto_SVC.append(n)
[140]: Acc_SVC['Acierto'] = Acierto_SVC
[141]: Acc_SVC.head()
[141]:
          Real      Predicho  Acierto
198274  ENTERTAINMENT    WORLDPOST      0
9281      POLITICS  STYLE & BEAUTY      0
85726     WORLDPOST      POLITICS      0
156506     WELLNESS      POLITICS      0
72313   LATINO VOICES    WORLDPOST      0
```

```
[142]: Acc_mod_SVC = np.mean(Acc_SVC.Acierto)
print("La efectividad del modelo es de {}%".format(np.round(Acc_mod_SVC * 100, 2)))
La efectividad del modelo es de 10.99%
```

1.2.5 Red neuronal

```
[117]: import tflearn # libreria simplificada por encima de tensorflow
[118]: X_train_t, X_test_t, y_train_t, y_test_t = train_test_split(transformed_texts,
                                                               target,
                                                               test_size = 0.2,
                                                               random_state = 137)
[119]: print(len(X_train_t))
print(len(X_test_t))
print(len(y_train_t))
print(len(y_test_t))
```

```
160682  
40171  
160682  
40171
```

```
[120]: X_train_t.shape
```

```
[120]: (160682, 48)
```

```
[121]: entradas = 48 # 48 por el largo del vector  
capa1 = 60 # las neuronas son un valor que se debe ajustar dependiendo del problema  
clases = 33 # estas son las clases para la clasificación al final del entrenamiento de la red neuronal  
learning_Rate = 0.001
```

```
[122]: def crear_modelo():  
    tf.reset_default_graph() # este comando hace un reset en caso de tener algo en memoria  
    red = tflearn.input_data([None, entradas]) # esta es la entrada de los datos de la red neuronal, None para enviar la información por lotes  
    red = tflearn.fully_connected(red, capa1, activation='ReLU') # son las conexiones de las capas de la red neuronal  
    red = tflearn.fully_connected(red, clases, activation='softmax') # capa de salida en este caso por ser clasificación la activación es "softmax"  
    red = tflearn.regression(red, optimizer = 'adam', # optimizador de gradiente descendente estocástico  
                             learning_rate = learning_Rate, loss = 'categorical_crossentropy') # la tasa de aprendizaje es un valor pequeño  
    modelo = tflearn.DNN(red, tensorboard_verbose = 3) # la variable de modelo cargar toda la estructura red DNN red neuronal profunda tensorboard_verbose = 3 para mostrar graficas  
    return modelo
```

```
[123]: modelo = crear_modelo() # para aplicar la función de crear_modelo con los parámetros anteriores
```

```
[124]: %%time  
modelo.fit(X_train_t, y_train_t, # variables de entrenamiento  
            validation_set = (X_test_t, y_test_t), # % de datos que dejo para la validación del modelo, no confundir con los datos de validación, los datos de validación no los debe ver el modelo por temas de over fitting  
            show_metric = True, # muestra las metricas  
            batch_size = 2000, # cantidad de batch de datos que voy a enviar por iteración al modelo  
            n_epoch = 200) # es que tantas veces se va a barrer el conjunto completo de datos
```

```
Training Step: 16199 | total loss: 2.98672 | time: 1.855s
```

```

| Adam | epoch: 200 | loss: 2.98672 - acc: 0.1937 -- iter: 160000/160682
Training Step: 16200 | total loss: 2.98483 | time: 2.881s
| Adam | epoch: 200 | loss: 2.98483 - acc: 0.1946 | val_loss: 2.99892 - val_acc:
0.1928 -- iter: 160682/160682
--
CPU times: user 11min 48s, sys: 46 s, total: 12min 34s
Wall time: 10min 26s

```

1.2.6 Estilos de escritura

```
[125]: from PIL import Image
from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator

[126]: label_encoder.classes_

[126]: array(['ARTS & CULTURE', 'BLACK VOICES', 'BUSINESS', 'COMEDY', 'CRIME',
       'DIVORCE', 'EDUCATION', 'ENTERTAINMENT', 'ENVIRONMENT', 'FIFTY',
       'FOOD & DRINK', 'GOOD NEWS', 'GREEN', 'HEALTHY LIVING',
       'HOME & LIVING', 'IMPACT', 'LATINO VOICES', 'MEDIA', 'MONEY',
       'PARENTS', 'POLITICS', 'QUEER VOICES', 'RELIGION', 'SCIENCE',
       'SPORTS', 'STYLE & BEAUTY', 'TECH', 'TRAVEL', 'WEDDINGS',
       'WEIRD NEWS', 'WELLNESS', 'WOMEN', 'WORLDPOST'], dtype=object)

[127]: len(label_encoder.classes_)

[127]: 33

[128]: p = df[['category', 'headline']]
p.head()

[128]:      category                      headline
0        CRIME  There Were 2 Mass Shootings In Texas Last Week...
1  ENTERTAINMENT  Will Smith Joins Diplo And Nicky Jam For The 2...
2  ENTERTAINMENT    Hugh Grant Marries For The First Time At Age 57
3  ENTERTAINMENT  Jim Carrey Blasts 'Castrato' Adam Schiff And D...
4  ENTERTAINMENT  Julianna Margulies Uses Donald Trump Poop Bags...

[129]: def my_cloud_wors(category_cloud):
    a = 'category'
    b = category_cloud
    c = p[p[a] == b]

    comment_words = ''
    stopwords = set(STOPWORDS)

    # iterate through the csv file
    for val in c.headline:

        # typecaste each val to string
        val = str(val)
```

```

# split the value
tokens = val.split()

# Converts each token into lowercase
for i in range(len(tokens)):
    tokens[i] = tokens[i].lower()

for words in tokens:
    comment_words = comment_words + words + ' '

wordcloud = WordCloud(width = 800, height = 800,
                      max_words = 60,
                      background_color ='white',
                      stopwords = stopwords,
                      min_font_size = 10).generate(comment_words)

# plot the WordCloud image
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud, interpolation='bilinear')
plt.title(category_cloud, fontsize = 18)
plt.axis("off")
plt.tight_layout(pad = 0)

plt.show()

```

```
[130]: for i in label_encoder.classes_:
    my_cloud_wors(i)
```

ARTS & CULTURE

BLACK VOICES

BUSINESS

COMEDY

know say samantha bee christmas
stephen colbert ad
will comedy thing reveal perfect jimmy kimmel
obama twitter
donald trump
president make look spoof
want trevor noah real
one time gets seth meyer
comedian - take back trump snl dog
really SNL first back jon stewart tweet
comedian jimmy fallon watch
day every best guy parody new
people tell john oliver photo
give way talk bill maher love late night
hilarious america show james corden

CRIME

california killing year old
toddler mom police officer father
girl teen dead shot new york
video victim dead inmate
police alleged prison fire year
suspect family home mother
former student die
woman dog sex
found arrested face
officer son
death shooter
charged new
accused florida
man driver say
allegedly missing attack
shooting charge
trial gun case

DIVORCE

split need couple may life
divorced blended family reason love single mom
husband friend sign
bad women getting video
thing star way infidelity
one tell affair make kid
divorce wife post split children holiday parenting celebrity
cheating single photo new say
spouse time woman men best sex step
ex better man relationship
breakup married advice dating will
post divorce

EDUCATION

student loan
learning
career
high school
grad
rape class
charter school
public
plan job
american thing
parent

university
make
classroom
know future
day top
education part say
professor
life want change state
campus year
search education america

best president
california one
take
now
back help
new first frat
college student graduation
will
way time
trump

student
student
teacher
college
education
school
sexual assault

ENTERTAINMENT

tv
best
netflix
fan
photo
make
will
way
dead
call
time
people
first
death
day
actor
snl'
singer
taylor swift
kanye west
beyoncé
love
video
game
thrones'
now
take
film
dead
call
will
way
get
think
oscar
back
one
life
song
watch
year
want
look
star wars'
one life
director
reveal
hollywood
trailer
trump
kim kardashian
women
star
box office
donald trump
twitter
set
watch
year
want
look
star wars'

ENVIRONMENT

climate change
food pet water fracking
free tip first
storm ocean year tornado time
dog show energy
will way obama video animal
extreme weather big photo
panda weather
photos week snow hurricane sandy
global warming sea make act
climate day animal day
home oil zoo species photo life
national cat whale
animal photo
keystone xl new baby green
shark may elephant say arctic find
pictures capture

FIFTY

alzheimer
advice
back
family
help
kid
people one
women
boomer
new
old
older
thing
make
will
best
first
say
mother never
take
time
age
day
every
year old woman
learned
home
sex
age
love
world
father
better
good
year
need
want
way
aging
reason
living
job
know
find
mom
life
live
simple
parent
really
feel
step
summer
tip
retirement

FOOD & DRINK

cooking summer restaurant salad eat breakfast chicken pizza day food drink beer cheese cake pie soup way ice cream chocolate favorite bacon dinner delicious need make recipe day best time perfect love thanksgiving healthy wine burger one holiday dessert coffee cookie cook thing good egg will taste test tip chef video top taste easy huffpost tastemakers cocktail kitchen

GOOD NEWS
teen saved
photo
school
find take puppy video world
car turn pit bull
homeless
rescued make
pup give use
year old
student home
animal cancer
best
woman
baby
help now
need christmas
love
kitten buy
friend year
time
free
one
watch cop gets
cat
life
man
boyday
dog
he
kid
Save

GREEN

big us environmental elephant
one time global future wild state
earth day animal pipeline pet take
say year watch help
look rescued need kill way end
epa oil storm
zoo food photo make first world green may
america coal will ocean home
california wildfire
climate change
baby save love people trump wildfire
live life

HEALTHY LIVING

HOME & LIVING

home finds style us look top sprout home show
make idea remove diy idea
best idea space furniture
holiday photos craft tiny home
clean new gift ideas sale
home color one roundup ebay
craft day garden new use
vintage home finds photos house tour
photo use table decor party worth
kitchen tip perfect
house weekly roundup apartment

IMPACT

global refugee crisis take want ebola
women need us
make veteran good thing war change cancer
world job american giving health work hope
education million children
one day fight leader
back end people home kid live video
time american child food service
will give helping poverty
aid homeless school way future
help american free
homeless

LATINO VOICES

A word cloud visualization where the size and color of words represent their frequency and association within the dataset. The most prominent words are 'puerto rico', 'trump', and 'latino'. Other significant words include 'immigration', 'people', 'star', 'say', 'mexican', 'american', 'donald trump', 'latina', 'immigrant', and 'day'.

The word cloud includes many smaller, colored words that provide context and associations:

- Associated with Puerto Rico:** life, anti, immigration, first, white, make, help, people, death, dreamer, undocumented, teen, award, show, want, mexico, america, shakira, hispanic, hispanic, president, latin, women, ferrera.
- Associated with Donald Trump:** call, timetake, daca, kid, love, way, vote, family, star, photo, explain, spanish, mexican, student, eva longoria, new, american, video, school, obama, manuel, manuel miranda, one.
- Associated with Latino/a:** latina, immigrant, hispanic, hispanic, president, latin, us, watch, day, women, ferrera.

MEDIA

question show want newspaper megyn kelly
president anderson cooper take msnbc gawker debate
report editor one vice
new magazine obama sean Hannity election
call media watch press
york times new york
make year first network american
donald trump back
claim bill o'reilly story
cnn video fox new fox slam
reporter interview tv war say
twitter journalism time washington post

MONEY

financial consumer may time save best
job car irs five
will taxes free credit card
tip find study percent plan
year new day make report avoid
tax bank pay people
college mortgage video saving
american way bad
price retirement
student cash loan
know student loan bill
million spending fee
thing cost credit buy debt
customer black friday
home sale
credit score
man photo
bill
spending
fee
credit score
home
sale
black friday

PARENTS

A word cloud centered around the word "mom" in large green letters. Other prominent words include "parent" (blue), "children" (green), "kid" (purple), "video" (yellow), and "baby" (green). Various family-related terms like "father", "son", "daughter", "toddler", "pregnancy", and "toy" are scattered throughout in different colors.

POLITICS

president white house
campaign stop win senate
donald trump
american new plan want
report war clinton first
year democratamerica
obama women vote take
republican election state
health care bill call say
ted cruz will people big
supreme congress support paul ryan
trump now still law
gop make end help may
attack bernie sanders
think hillary clinton
fight one back morning email

QUEER VOICES

queer
lgbt
trans
gay
transgender

open
pride
family
man
take
back son
year
activist
right
marriage equality

new
star
say
video
one
come first
community
america

look
talk
two
black
hiv
drag race'

love
trump
people
show
life
gay men
son
now
rupaul drag
kid
watch
men make
anti gay
sex
photo
mom
couple
time
gender

RELIGION

SCIENCE

universe discovered animal
see image people

nasas life earth

astronaut study suggests astronomer pluto

face mars

astronauts say new

say found

new study

year space make brain

first day take women suggest secret

watch mission reveal photo

show ancient star look world

way dinosaur help

light asteroid climate change sun

one big moon scientist planet

will

Video

SPORTS

STYLE & BEAUTY

A word cloud graphic centered around fashion and style. The most prominent words are "fashion" (purple), "style" (green), "photo" (yellow), and "video" (green). Other significant words include "new york" (blue), "hair" (brown), "model" (teal), "dress" (green), "look" (green), "red carpet" (green), "kate middleton" (teal), "men" (brown), "woman" (brown), "obama" (brown), "kardashian" (brown), "makeup" (blue), "perfect" (blue), "street style" (brown), "year" (brown), "love" (brown), "us" (brown), "girl" (brown), "star" (brown), "ad" (brown), "first" (brown), "will" (brown), "fall" (brown), "time" (brown), "outfit" (brown), "shoe" (brown), "thing" (brown), "want" (brown), and "style evolution" (brown).

TECH

The word cloud illustrates various tech-related topics and companies. The most prominent words are 'iphone' (blue), 'facebook' (green), 'apple' (green), 'willgoogle' (green), 'twitter' (green), 'internet' (purple), 'new' (purple), 'tech' (blue), 'data' (blue), 'watch' (green), 'app' (teal), 'youtube' (yellow), 'video' (teal), 'help' (green), 'still' (green), 'photo' (green), 'watch' (purple), 'top' (purple), 'thing' (green), 'day' (yellow), 'report' (green), 'way' (green), 'take' (green), 'make' (green), 'say' (yellow), 'phone' (green), 'may' (green), 'uber' (purple), 'millions' (green), 'fbi' (green), 'samsung' (green), 'privacy' (green), 'people' (blue), 'videos' (blue), 'week' (blue), 'technology' (blue), 'apple' (green), 'rumors' (purple), 'future' (green), 'billion' (yellow), 'year' (green), 'top' (purple), 'youtube' (green), 'deal' (yellow), and 'internet' (purple).

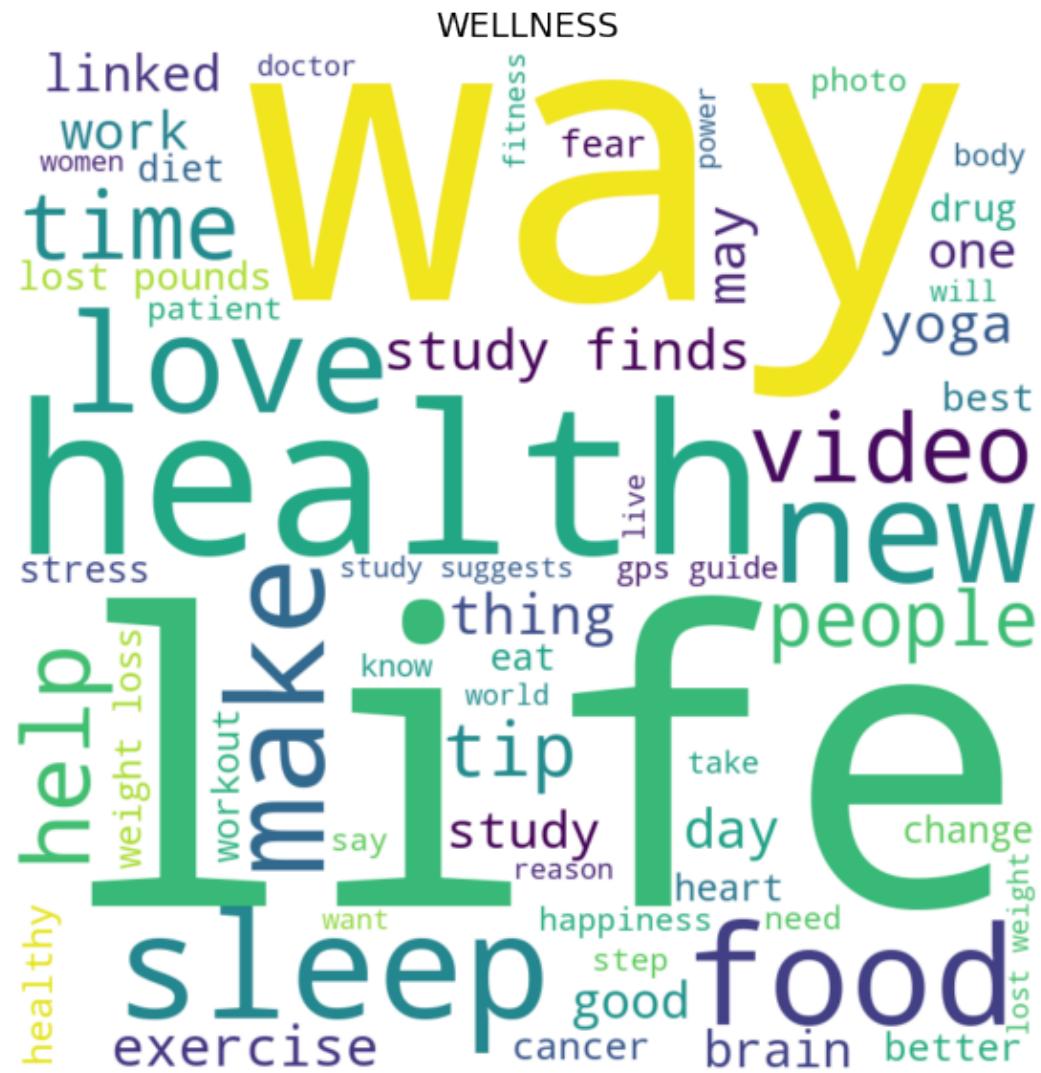
TRAVEL

travel around world flight beautiful
destination love looney front
great road trip london reason
trip paris see day go island adventure
way time best tsa tourist beach city
make trip one summer winter guide life
traveling new york europe town tour
vacation video american america take
food passenger

WEDDINGS

WEIRD NEWS

A word cloud visualization where the size of each word indicates its frequency. The most prominent words are 'man' (large blue), 'woman' (large blue), 'watch' (large blue), 'people' (large green), 'cat' (large teal), 'dog' (large purple), 'guy' (medium green), 'video' (medium green), 'police' (medium green), 'car' (medium purple), 'cop' (medium green), 'shark' (medium purple), 'fire' (medium yellow), 'weird' (medium yellow), 'new' (medium yellow), 'fark' (medium yellow), 'allegedly' (medium yellow), 'turn' (medium yellow), 'bear' (medium purple), 'one' (medium purple), 'take' (medium yellow), 'found' (medium yellow), 'donald' (medium purple), 'trump' (medium purple), 'head' (medium green), 'world' (medium green), 'way' (medium green), 'time' (medium green), 'sex' (medium green), 'want' (medium green), 'make' (medium green), 'know' (medium green), 'show' (medium green), 'christmas' (medium green), 'news' (medium green), 'quiz' (medium green), 'florida' (medium green), 'internet' (medium green), 'fark' (medium green), 'year old' (medium purple), 'baby' (medium purple), 'day' (medium purple), 'now accused' (medium purple), 'accused' (medium purple), 'home' (medium green), 'real' (medium green), 'truck' (medium green), 'turn' (medium green), 'thing' (medium green), 'year' (medium green), 'look' (medium purple), 'will' (medium purple), 'take' (medium purple), 'see' (medium purple), 'say' (medium green), and 'life' (medium green).



A word cloud centered around the word "WOMEN". The words are arranged in a cluster, with the size of each word indicating its frequency or importance. The colors of the words vary, including shades of blue, green, yellow, and red. Some of the prominent words include "funniest", "woman", "tweets", "women", "men", "thing", "woman", "week", "trump", "say", "abortion", "feminism", "work", "body", "think", "call", "real", "powerful", "one", "love", "girl", "feminist", "man", "look", "time", "still", "stop", "sex", "will", "make", "right", "sexist", "back", "show", "want", "way", "sexual", "assault", "world", "know", "really", "day", "new", "made", "take", "rape", "clinton", "hillary", "never", "people", "life", "watch", "first", "now", "single", "donald", "trump", "planned", "parenthood", and "year".

WORLDPOST