



Software Solutions

Twitter Data Analysis

Manual Técnico

Manejo de Datos no Estructurados

Dr. Francisco J. Moo Mena

Integrantes:

Daniel Fernando Baas Colonia

Rodrigo Castilla López

Luis Gerardo Castillo Pinkus

Rafael Rodríguez Guzmán

Instalaciones necesarias

- Python: el proyecto fue realizado con la versión 3.6 de Python
 - <https://www.python.org/downloads/release/python-367/>
- pip3: se usó el manejador de paquetes pip3 para instalar las bibliotecas de Python necesarias.
 - `sudo apt install python3-pip`
- Tweepy: se usó la biblioteca de tweepy para realizar la conexión a la api de Twitter.
 - `pip3 install tweepy`
- Matplotlib: biblioteca que permite generar gráficas.
 - `pip3 install matplotlib`
- Pymongo: driver para Python usado para manejar la interacción con la base de datos de MongoDB.
 - `pip3 install pymongo`
- WordCloud: biblioteca que permite generar nubes de palabras con Python.
 - `pip3 install wordcloud`
- Numpy: biblioteca que permite realizar cómputo científico con Python.
 - `pip3 install numpy`
- MongoDB: base de datos no relacional. Se usa para almacenar la información de los tweets recolectados.
 - <https://www.mongodb.com/cloud/atlas>

MongoDB como base de datos

Esta aplicación hace uso de MongoDB como base de datos no estructurada. Más específicamente, se usa el clúster de Atlas en su versión gratuita.

Cada tweet representa un documento, y la información se almacena es la siguiente:

```
_id: ObjectId("5bdaf8bf622ac410a401fe64")
created_at: 2018-11-01 12:59:37.000
text: "Not gonna be able to watch this Smash Bros Nintendo Direct live so I'l..."
user: Object
  name: "Joe Dillon"
  screen_name: "joffocakes"
  location: "Glasgow, Scotland"
  url: "http://joffocakes.tumblr.com/"
  description: "Mutton dressed as lambrini. I like Street Fighter and I like eating t..."
  source: "Twitter for Android"
  extended_tweet: ""
  retweeted: false
  lang: "en"
place: Object
  country: "United Kingdom"
  full_name: "Glasgow, Scotland"
  name: "Glasgow"
truncated: true
```

Cabe destacar que no todos los tweets contienen la información completa para todos los campos. Por ejemplo, sólo 361 de los 72595 tweets totales contienen información en el campo *place*.

Documentación del código

En esta aplicación está compuesta de dos partes: la recolección de tweets y el análisis de ellos. La colección se realiza ejecutando el script *get_tweets.py*, mientras que el análisis de los datos se realiza ejecutando el script *procesing.py*.

procesing.py

Lo primero que hace al ejecutar este archivo es crear una conexión al servidor Atlas de MongoDB. En este servidor se encuentra la base de datos que contiene todos los tweets. Una vez encontrada la base de datos y la colección, se muestra en pantalla un mensaje al usuario indicando la conexión exitosa. Algo importante es que `client = MongoClient()` establece una conexión con MongoDB. Dado que usamos el clúster de Atlas de MongoDB, esa función lleva una cadena como parámetro. Dicha cadena se tiene que generar desde la página de Atlas una vez hayas creado tu servidor.

```
15 # Creamos una conexión con el servidor
16 client = MongoClient(
17     "mongodb://rafael:S7BNKu4WAj2z3ca6@nonstructureddatamanagementclass-shard-00-00-xa2lb.mongodb.net:27017,"
18     "nonstructureddatamanagementclass-shard-00-01-xa2lb.mongodb.net:27017,"
19     "nonstructureddatamanagementclass-shard-00-02-xa2lb.mongodb.net:27017/twitter_db?ssl=true&replicaSet"
20     "=NonStructuredDataManagementClass-shard-00&authSource=admin&retryWrites=true")
21
22 # Nuestra base de datos se llamará "nintendo"
23 db_name = "nintendo"
24
25 # Intentamos realizar una conexión con la base de datos. Si la BD no existe, la crea.
26 db = client[db_name]
27
28 # Recuperamos el nombre de todas las bases de datos en el servidor
29 db_list = client.list_database_names()
30
31 # Revisamos si la base de datos existe
32 if db_name in db_list:
33     print("\nLa base de datos \"" + db_name + "\" existe!")
34
35 # Una colección es un grupo de documentos almacenados en MongoDB
36 # Nuestra colección se llamara "tweets"
37 collection_name = "tweets"
38
39 # Crea una nueva colección o recupera el contenido de una existente
40 # Si la colección se está creando, MongoDB no la va a crear hasta que se le inserte documentos
41 collection = db[collection_name]
42
43 # Recuperamos el nombre de todas las colecciones
44 collection_list = db.list_collection_names()
45
46 # Revisamos si la colección existe. Si la colección creada no tiene elementos, MongoDB no la va a crear hasta
47 # insertar un documento
48 if collection_name in collection_list:
49     print("La coleccion \"" + collection_name + "\" existe!\n")
```

is_valid_tweet(tweet): recibe una cadena de texto que contiene el tweet y revisa si contiene alguna palabra correspondiente a una lista de términos que no son de nuestro interés. Si la cadena contiene alguna palabra se retorna falso, sino, se retorna verdadero, indicando que el tweet si es de nuestro interés.

```

77 # Verifica que el tweet no sea uno relacionado con términos que no son relevantes para nosotros
78 def is_valid_tweet(tweet):
79     # Dados los tags usados para recopilar los tweets, no todos resultaron tener información útil. Así que se tiene
80     # una lista de tweets cuyo tema no es de nuestro interés
81     filter_list = [
82         "64",
83         "classic",
84         "n64",
85         "#lookingfor",
86         "$239.99",
87         "#splatoon2",
88         "wii"
89     ]
90
91     for fw in filter_list:
92         if fw in tweet:
93             return False
94
95     return True
96

```

generate_word_cloud(tweets): esta función crea una imagen de una nube de palabras. Primero recibe una lista de tweets. El contenido de texto de cada tweet se concatena en una variable para poder pasar la cadena de texto completa a la biblioteca de wordcloud. Una vez que todos los textos de los tweets están unidos generamos una lista de palabras que el wordcloud va ignorar y no las va a colocar al final. Después que se genere la lista se carga una imagen con una forma que tomará la nube de palabras conocida como máscara (mask). En nuestro caso la máscara tiene forma de Mario (Super Mario Bros) pues va con la temática de Nintendo. Finalmente se usa la función WordCloud para generar la nube de palabras y la guarda en una imagen en la ruta `/plots/wordcloud.png`

```

109 # Genera una imagen .png con una nube de palabras de los términos más usados en los tweets
110 def generate_wordcloud(tweets):
111     print("\nGenerando nube de palabras...\n")
112
113     tweet_list = []
114     for tweet in tweets:
115         tweet_list.append(tweet["text"])
116
117     # Unimos el texto de todos los tweets en una sola variable
118     text = "\n".join(tweet_list)
119
120     # Cargamos la lista de stopwords y añadimos más que consideramos no relevantes para mostrar en la nube de palabras
121     # Las stopwords son palabras que serán ignoradas
122     stopwords = set(STOPWORDS)
123     stopwords.add("nintendodirect")
124     stopwords.add("smash")
125     stopwords.add("bros")
126     stopwords.add("smashbros")
127     stopwords.add("supersmashbros")
128     stopwords.add("ultimate")
129     stopwords.add("smashbrosultimate")
130     stopwords.add("supersmashbrosultimate")
131     stopwords.add("co")
132     stopwords.add("nintendo")
133
134     # Cargamos la máscara que usaremos para darle forma a la nube
135     mask = np.array(Image.open("sources/mario.png"))
136
137     # Genera una nube de palabras. Las stopwords serán ignoradas y no aparecerán en la gráfica
138     wordcloud = WordCloud(width=1079, height=1623, max_words=10000, relative_scaling=1, stopwords=stopwords,
139         mask=mask, contour_color="white").generate(text)
140
141     # Guardamos la nube de palabras generada en una imagen en la misma carpeta del script
142     wordcloud.to_file("plots/wordcloud.png")

```

`generate_app_plot(tweets)`: esta función genera una gráfica de pastel que muestra las siete aplicaciones más usadas para realizar los tweets. Recibe como parámetro una lista de tweets, y genera una lista de objetos de tipo {"device": aplicación, "count": número de tweets para esa aplicación}, para ir contabilizando cuántos tweets se realizaron desde cada aplicación. Dada la enorme cantidad de aplicaciones, sólo se seleccionaron las siete más usadas. Al finalizar, se genera la gráfica de pastel y la guarda en una imagen en la ruta `/plots/applications.png`.

```
145 # Genera una gráfica de pastel con las aplicaciones desde donde se realizaron la mayoría de los tweets
146 def generate_app_plot(tweets):
147     print("\nGenerando gráfica de dispositivos...\n")
148
149     # Obtenemos una lista de objetos con el nombre de la aplicación y el número de tweets
150     source_list = generate_source_list(tweets)
151
152     # Dado que hay una enorme cantidad de aplicaciones, sólo vamos a mostrar las 7 más usadas, y las demás se contarán
153     # dentro de "Otras" aplicaciones
154
155     src_count = 0
156     max_src = 7
157     src_plot = []
158
159     # Revisamos la lista de aplicaciones y generamos una con las más usadas para hacer los tweets
160     for src in source_list:
161         if src_count < max_src:
162             src_plot.append(src)
163             src_count = src_count + 1
164
165         if src_count == max_src:
166             src_plot.append({
167                 "source": "Others",
168                 "count": 1
169             })
170
171             src_count = src_count + 1
172
173         if src_count > max_src:
174             src_plot[max_src]["count"] = src_plot[max_src]["count"] + 1
175
176     # Nombre de las aplicaciones
177     labels = []
178     # Número de tweets por aplicación
179     sizes = []
180
181     # Añadimos la información para las etiquetas en la gráfica
182     for src in src_plot:
183         labels.append(src["source"])
184         sizes.append(src["count"])
185
186
187     # Ajusta el tamaño de la gráfica para que se adapte al tamaño del contenido
188     plt.tight_layout()
189     # Ajustamos el tamaño de la letra
190     plt.rcParams["font.size"] = 7.0
191     # Añadimos un título
192     plt.title("Most used applications")
193     # Generamos una gráfica de pastel
194     plt.pie(sizes, labels=labels, shadow=True, radius=10, autopct="%0.2f%%")
195     plt.axis("equal")
196
197     # Guardamos la gráfica en una imagen en la misma carpeta que este script
198     plt.savefig("plots/applications.png", dpi=300)
199
200     plt.show(block=False)
201     plt.close()
```


`generate_country_plot(tweets)`: esta función genera una gráfica de barras que muestra los diez países que realizaron más tweets. Recibe como parámetro una lista de tweets, y genera una lista de objetos de tipo {"country": país, "count": número de tweets para ese país}, para ir contabilizando cuántos tweets se realizaron desde cada cada país. Dada la enorme cantidad de países, sólo se seleccionaron diez que más tweets realizaron. Al finalizar, se genera la gráfica de barras y la guarda en una imagen en la ruta `/plots/countries.png`. Como nota se puede decir que, como no todos los tweets contienen información sobre localización, la gráfica representa alrededor de 300-400 tweets, que son los que si contienen información de localización.

```
204 # Genera una gráfica de barras con los países que más tweets realizaron.
205 # Dado que no todos los tweets tienen información de ubicación, esta lista de tweets tiene alrededor de 300+ tweets
206 def generate_country_plot(tweets):
207     print("\nGenerando gráfica de países...\n")
208
209     # Obtenemos una lista de objetos que contiene el nombre del país con el número de tweets
210     country_list = generate_country_list(tweets)
211
212     # Dado que la lista de países es bastante grande, sólo mostraremos los primeros 10 con más tweets
213
214     country_count = 0
215     max_country = 10
216     country_chart = []
217
218     # Generamos una lista con los 10 países que más tweetearon
219     for country in country_list:
220         if country_count < max_country:
221             country_chart.append(country)
222             country_count = country_count + 1
223
224     # Nombre de los países
225     labels = []
226     # Número de tweets
227     sizes = []
228
229     # Añadimos la información a las etiquetas de la gráfica de barras
230     for country in country_chart:
231         labels.append(country["country"])
232         sizes.append(country["count"])
233
234     # Prepara la información necesaria para generar la gráfica de barras
235     y_pos = np.arange(len(labels))
236     plt.yticks(y_pos, labels)
```

```

238 # Añadimos una etiqueta al eje vertical
239 plt.ylabel("Country")
240 # Añadimos una etiqueta al eje horizontal
241 plt.xlabel("# of tweets")
242 # Añadimos un título
243 plt.title("Top 10 countries with the most tweets")
244 # Ajusta el tamaño de la gráfica para que se adapte al tamaño del contenido
245 plt.tight_layout()
246
247 # Añade una etiqueta que muestra el número de tweets a lado de cada barra
248 for i, v in enumerate(sizes):
249     plt.text(v + 5, i - 0.1, str(v), color="#1F76B3", fontweight="bold")
250
251 # Genera la gráfica
252 plt.barh(y_pos, sizes)
253
254 # Guarda la gráfica como imagen en la misma carpeta que este script
255 plt.savefig("plots/countries.png", dpi=300)
256
257 plt.show(block=False)
258 plt.close()

```

`generate_language_plot(tweets)`: esta función genera una gráfica de pastel que muestra los siete lenguajes más usados para realizar los tweets. Recibe como parámetro una lista de tweets, y genera una lista de objetos de tipo {"lang": lenguaje, "count": número de tweets para ese lenguaje}, para ir contabilizando cuántos tweets se realizaron por cada lenguaje. Dada la enorme cantidad de lenguajes, sólo se seleccionaron los siete más usados. Al finalizar, se genera la gráfica de pastel y se guarda en una imagen en la ruta `/plots/languages.png`.


```

261 # Genera una gráfica de pastel que representa los idiomas más usados en los tweets
262 # Dado que hay muchos idiomas solo te van a tomar como máximo los definidos en la variable max_lang
263 def generate_language_plot(tweets):
264     print("\nGenerando gráfica de idiomas...\n")
265
266     # Obtenemos una lista de objetos con las siglas del idioma y el número de tweets en ese idioma
267     language_list = generate_language_list(tweets)
268
269     # Dado que hay una enorme cantidad de idiomas, sólo tomaremos un número limitado, definido por la variable max_lang
270     lang_count = 0
271     max_lang = 7
272     lang_plot = []
273
274     # Revisamos la lista de aplicaciones y generamos una con las más usadas para hacer los tweets
275     for lang in language_list:
276         if lang_count < max_lang:
277             lang_plot.append(lang)
278             lang_count = lang_count + 1
279
280         if lang_count == max_lang:
281             lang_plot.append({
282                 "language": "Others",
283                 "count": 1
284             })
285
286             lang_count = lang_count + 1
287
288         if lang_count > max_lang:
289             lang_plot[max_lang]["count"] = lang_plot[max_lang]["count"] + 1
290
291     # Listado de idiomas a usar como etiquetas en la gráfica
292     labels = []
293     # Número de tweets por idioma
294     sizes = []

```

```

296     # Añadimos la información para las etiquetas en la gráfica
297     for lang in lang_plot:
298         labels.append(lang["language"])
299         sizes.append(lang["count"])
300
301     # Ajusta el tamaño de la gráfica para que se adapte al tamaño del contenido
302     plt.tight_layout()
303     # Ajustamos el tamaño de la letra
304     plt.rcParams["font.size"] = 7.0
305     # Añadimos un título
306     plt.title("Most used languages")
307     # Generamos una gráfica de pastel
308     plt.pie(sizes, labels=labels, shadow=True, radius=10, autopct="%0.2f%%")
309     plt.axis("equal")
310
311     # Guardamos la gráfica en una imagen en la misma carpeta que este script
312     plt.savefig("plots/languages.png", dpi=300)
313
314     plt.show(block=False)
315     plt.close()

```

`print(tweet)_single_tweet` y `print_tweets(tweets)`: la primera función imprime en pantalla el contenido completo de un solo tweet, mientras que la segunda itera sobre la lista completa de tweets, llamando a `print_single_tweet(tweet)` para la impresión en pantalla por cada elemento.

```

52 # Imprime el contenido del tweet recuperado de MongoDB
53 def print_single_tweet(tweet):
54     print("\ncreated_at: " + str(tweet["created_at"]))
55     print("text: " + tweet["text"])
56     print("user.name: " + tweet["user"]["name"])
57     print("user.screen_name: " + tweet["user"]["screen_name"])
58     print("user.location: " + str(tweet["user"]["location"]))
59     print("user.url: " + str(tweet["user"]["url"]))
60     print("user.description: " + str(tweet["user"]["description"]))
61     print("source: " + tweet["source"])
62     print("extended_tweet: " + tweet["extended_tweet"])
63     print("retweeted: " + str(tweet["retweeted"]))
64     print("lang: " + tweet["lang"])
65     print("place.country: " + tweet["place"]["country"])
66     print("place.full_name: " + tweet["place"]["full_name"])
67     print("place.name: " + tweet["place"]["name"])
68     print("truncated: " + str(tweet["truncated"]))
69
70
71 # Imprimimos todos los tweets de un conjunto dado
72 def print_tweets(tweets):
73     for tweet in tweets:
74         print_single_tweet(tweet)

```

generate_plot_list.py

Este script permite generar listas de objetos dados ciertos campos de los tweets. Se usa para generar la lista de aplicaciones, la de países y la de lenguajes. Estas listas tiene una forma:

{“field”: campo que se utiliza para contabilizar, “count”: número de ocurrencias}

“field” puede tomar el valor de “device”, “country” o “lang”, dependiendo de si la lista a generar es una de aplicaciones, países o lenguajes. “count” siempre representa el número de tweets que se han contabilizado para dicho campo de referencia.

Al final de la ejecución, retorna una lista de objetos basados en la estructura mencionada anteriormente y el campo sobre el cuál se contabilizaron los tweets.

is_in_list(element, field): verifica si el elemento ya se encuentra en la lista.

add_element(element, field): añade un nuevo elemento a la lista.

```

4 # Lista de todas las aplicaciones desde las que se realizaron los tweets con sus conteos
5 element_list = []
6
7 |
8 # Verifica si la aplicación ya se está contabilizando o es una nueva
9 def is_in_list(element, field):
10     # La lista está vacía
11     if len(element_list) == 0:
12         return False
13
14     # Verifica que la aplicación se encuentre en la lista
15     for elmt in element_list:
16         if element == elmt[field]:
17             return True
18
19     return False

```

```

36 # Añadimos una nueva aplicación a la lista de fuentes de tweets
37 def add_element(element, field):
38     new_element = {
39         field: element,
40         "count": 1
41     }
42
43     element_list.append(new_element)
44
45
46 # Imprime la lista de dispositivos con la cantidad de tweets de cada uno
47 def print_elements(elmt_list, field):
48     for elmt in elmt_list:
49         print(field + ": " + elmt[field])
50         print("total: " + str(elmt["count"]))

```

generate_source_list(tweets): genera una lista que contabiliza el número de tweets por cada aplicación bajo la cual fueron realizados.

```

53 # Organiza las aplicaciones desde las que se realizaron los tweets, aumentando el contenido de tweet por aplicación o
54 # añadiendo nuevas aplicaciones al conteo
55 def generate_source_list(tweets):
56     element_list.clear()
57
58     field = "source"
59
60     for tweet in tweets:
61         if is_in_list(tweet[field], field):
62             increment_element(tweet[field], field)
63         else:
64             add_element(tweet[field], field)
65
66     new_list = sort_list(element_list, "count")
67
68     #print_elements(new_list, field)
69
70     return new_list

```

generate_country_list(tweets): genera una lista que contabiliza el número de tweets por cada país del cual se realizó el tweet.

```

75 def generate_country_list(tweets):
76     element_list.clear()
77
78     field = "country"
79
80     for tweet in tweets:
81         if is_in_list(tweet["place"]["country"], field):
82             increment_element(tweet["place"]["country"], field)
83         else:
84             add_element(tweet["place"]["country"], field)
85
86     new_list = sort_list(element_list, "count")
87
88     #print_elements(new_list, field)
89
90     return new_list

```

generate_language_list(tweets): genera una lista que contabiliza el número de tweets por cada lenguaje en el que se realizó cada tweet.

```

92 def generate_language_list(tweets):
93     element_list.clear()
94
95     field = "language"
96
97     for tweet in tweets:
98         if is_in_list(tweet["lang"], field):
99             increment_element(tweet["lang"], field)
100        else:
101            add_element(tweet["lang"], field)
102
103     new_list = sort_list(element_list, "count")
104
105     #print_elements(new_list, field)
106
107     return new_list

```

get_tweets.py

Este script realiza una conexión al servidor para acceder a la base de datos, y de no existir, crea una. Una vez que ya está conectado a la base de datos accede

a la colección donde se almacenarán los tweets que se recuperen. El siguiente paso es realizar una conexión a la api de Twitter para empezar un stream de tweets que se están publicando en este mismo instante. Los tweets que se encontrarán corresponden a una serie de términos clave proporcionados. El stream permanecerá activo mientras el script esté en ejecución y haya conexión al clúster de Atlas. Para terminar su ejecución es necesario acabar con el proceso de manera manual.

`MyStreamListener(tweepy.StreamListener)`: es una clase que permite el stream de tweets. Dentro de ella obtenemos cada tweet por el objeto `"self"`. Debemos acceder a los campos que sean de nuestro interés del tweets, y almacenarlos en una variable para poder guardar dicha información en MongoDB. `collection_insert_one(tweet)` inserta el tweet que generamos a la base de datos.

```
# Sirve para realizar el stream de tweets
class MyStreamListener(tweepy.StreamListener):

    # Para cada tweets obtenido, recuperamos la información que se muestra
    def on_status(self, status):
        tweet = {
            "created_at": status.created_at,
            "text": status.text,
            "user": {
                "name": status.user.name,
                "screen_name": status.user.screen_name,
                "location": status.user.location,
                "url": status.user.url,
                "description": status.user.description
            },
            "source": status.source,
            "extended_tweet": "",
            "retweeted": status.retweeted,
            "lang": status.lang,
            "place": {
                "country": "",
                "full_name": "",
                "name": ""
            },
            "truncated": status.truncated
        }
```

```
# Ciertos atributos de los tweets pueden o no pueden estar, así que debemos verificar que si se encuentren
# antes de acceder a ellos u obtendremos errores por null
if hasattr(status, "extended_tweet"):
    if (hasattr(status.extended_tweet, "full_text")):
        tweet["extended_tweet"] = status.extended_tweet.full_text
if hasattr(status, "place"):
    if hasattr(status.place, "country"):
        tweet["place"]["country"] = status.place.country
    if hasattr(status.place, "full_name"):
        tweet["place"]["full_name"] = status.place.full_name
    if hasattr(status.place, "name"):
        tweet["place"]["name"] = status.place.name

# Imprimimos los tweets para verificar que el stream sigue activo
print(status.text + "\n")

# Insertamos el tweet a la colección
inserted = collection.insert_one(tweet)
```

consumer_key, *consumer_secret*, *access_token* y *access_token_secret* es información proporcionada por Twitter. Para obtener dichas claves es necesario registrarse en Twitter como desarrollador (<https://developer.twitter.com/>).

Una vez obtenidas las llaves, es necesario realizar una autenticación con ellas e iniciar el stream de tweets. Se establecieron términos que filtran qué tweets son los que queremos obtener.

```
# Datos de acceso a la api
consumer_key = '0zLVRnLCdP0JuUxjsM1B5P6id'
consumer_secret = 'UE0vI6l0QFgaWxIzee3LJM0t5C59M5YhL9yh40d0q0ZGBXs85n'
access_token = '940271662959906816-CEzHk0ji1IYR1Z1VZeXdSYDrNetgRVT'
access_token_secret = 'wXev8vS56LhX06adG0gDX8HFTTrAcjhwSztQGBGlzQePVA'

# Autenticación y validación con la api
auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_token_secret)
api = tweepy.API(auth)

# Iniciamos un nuevo stream de tweets
myStreamListener = MyStreamListener()
myStream = tweepy.Stream(auth=api.auth, listener=myStreamListener, tweet_mode='extended')

# Establecemos los temas sobre los que queremos obtener información
myStream.filter(
    track=['Nintendo', 'NintendoDirect', 'NintendoSwitch', 'SmashBros', 'SuperSmashBros', 'SuperSmashBrosUltimate',
          'SmashBrothers', 'SuperSmashBrothers', 'SuperSmashBrothersUltimate', 'SmashUltimate'])
```