

## ▼ Import modules

```
1 import json
2 import tweepy
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import missingno as msno
7 import time
8 import seaborn as sns
9
10 from datetime import datetime
```

## ▼ Reading data

```
1 data = pd.read_csv('Full_data.csv')
```

## ▼ Look for dimension, types, missing values and duplicates

```
1 print(data.shape)
2 print(data.columns)
3 data.head(3)
```



```
(10137, 23)
Index(['status_id', 'user_id', 'created_at', 'screen_name', 'text', 'source',
      'reply_to_status_id', 'reply_to_user_id', 'reply_to_screen_name',
      'is_quote', 'is_retweet', 'favourites_count', 'retweet_count',
      'country_code', 'place_full_name', 'place_type', 'followers_count',
      'friends_count', 'account_lang', 'account_created_at', 'verified',
      'lang', 'source_tweet'],
```

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10137 entries, 0 to 10136
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  -
0   status_id             10137 non-null  int64
1   user_id               10137 non-null  int64
2   created_at            10137 non-null  object
3   screen_name           10137 non-null  object
4   text                  10137 non-null  object
5   source                10137 non-null  object
6   reply_to_status_id    1451 non-null   float64
7   reply_to_user_id      1794 non-null   float64
8   reply_to_screen_name  1794 non-null   object
9   is_quote              10137 non-null  bool
10  is_retweet             10137 non-null  bool
11  favourites_count       10137 non-null  int64
12  retweet_count          10137 non-null  int64
13  country_code           9825 non-null   object
14  place_full_name        9825 non-null   object
15  place_type             9825 non-null   object
16  followers_count        10137 non-null  int64
17  friends_count          10137 non-null  int64
18  account_lang           0 non-null      float64
19  account_created_at     10137 non-null  object
20  verified               10137 non-null  bool
21  lang                   10137 non-null  object
22  source_tweet           10137 non-null  object
dtypes: bool(3), float64(3), int64(6), object(11)
memory usage: 1.6+ MB
```

```
1 msno.matrix(data, figsize=(13, 4))
```

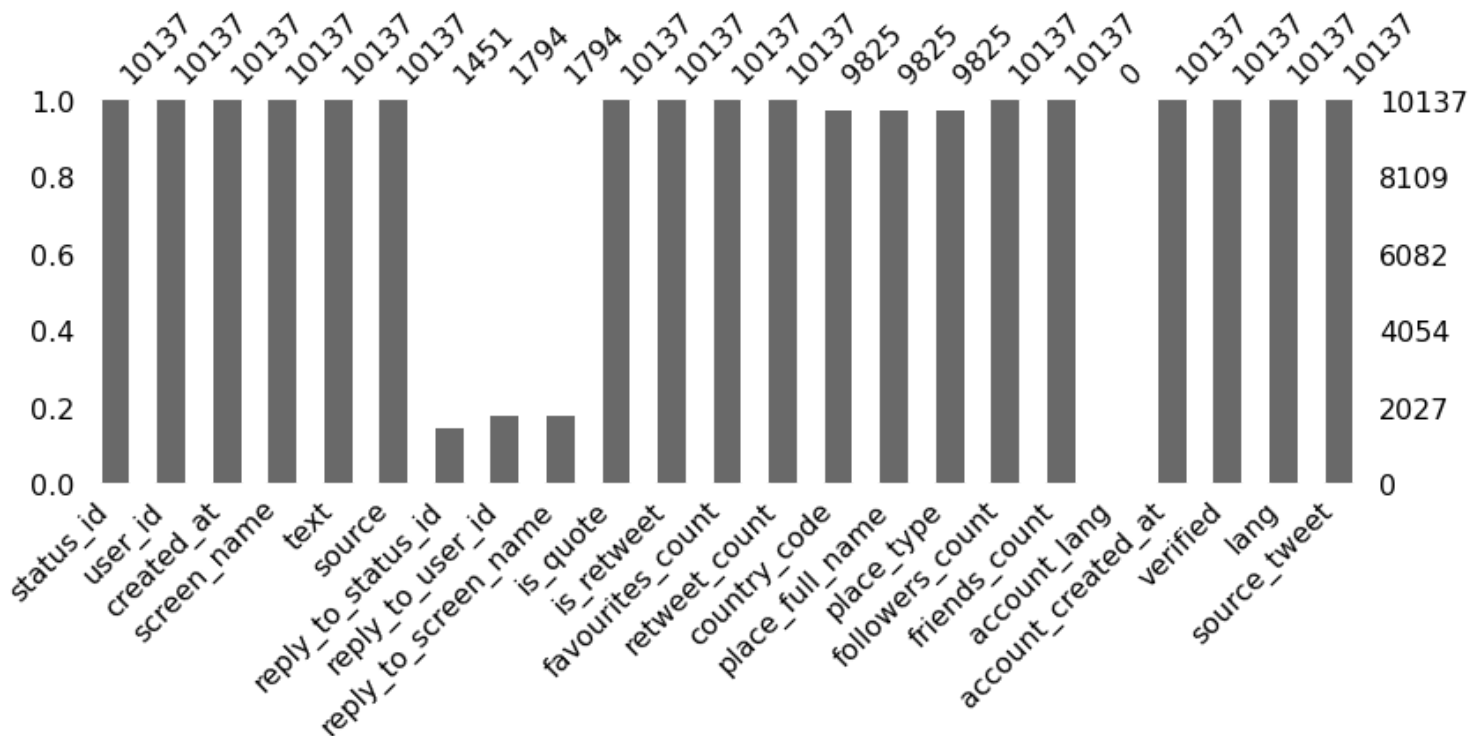
```
<=>
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5e1f2fdd8>
```

```
1
status_id
user_id
created_at
screen_name
text
source
reply_to_status_id
reply_to_user_id
reply_to_screen_name
is_quote
is_retweet
favourites_count
retweet_count
country_code
place_full_name
place_type
followers_count
friends_count
account_lang
account_created_at
verified
lang
source_tweet
```

```
1 msno.bar(data, figsize=(13, 4))
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5df6406d8>
```



▼ Since the columns `account_lang` doesn't have any value, then we proceed to drop it

```
1 data_ = data.drop(['account_lang'], axis=1).copy()
```

▼ Here we notice that there are about 400 tweet repeat it, so we proceed to delete them.

```
1 print("Unique values:")
2 data_.nunique()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5df6406d8>
```

```
Unique values:
status_id      9703
user_id        3387
created_at     9682
screen_name    3393
text           9702
source         11
reply_to_status_id 1009
reply_to_user_id 852
reply_to_screen_name 854
is_quote       2
is_retweet     1
favourites_count 4977
retweet_count  113
country_code   2
place_full_name 432
place_type     4
followers_count 3634
friends_count  2631
```

```
1 data_.groupby('status_id').count().sort_values(by='user_id', ascending=False).head()
```



	user_id	created_at	screen_name	text	source	reply_to_status_id	reply_to
status_id							
<b>1246926034873827328</b>	12	12	12	12	12	12	
<b>1244712607069745152</b>	9	9	9	9	9	9	
<b>1243245611102732288</b>	8	8	8	8	8	8	
<b>1244291628204019712</b>	8	8	8	8	8	8	
<b>1244357680673759232</b>	8	8	8	8	8	8	

```
1 data_[data_['status_id'] == 1246926034873827328]['text']
```



```
9688 @CoordinaValle La marcha del 8M y el evento de...
9689 @CoordinaValle La marcha del 8M y el evento de...
9735 @CoordinaValle La marcha del 8M y el evento de...
9736 @CoordinaValle La marcha del 8M y el evento de...
9797 @CoordinaValle La marcha del 8M y el evento de...
9798 @CoordinaValle La marcha del 8M y el evento de...
9876 @CoordinaValle La marcha del 8M y el evento de...
9877 @CoordinaValle La marcha del 8M y el evento de...
9966 @CoordinaValle La marcha del 8M y el evento de...
9967 @CoordinaValle La marcha del 8M y el evento de...
10066 @CoordinaValle La marcha del 8M y el evento de...
10067 @CoordinaValle La marcha del 8M y el evento de...
Name: text, dtype: object
```

```
1 data1 = data_.drop_duplicates('status_id', keep='first').copy()
```

```
1 print(data1.shape)
```

```
2 data1.head(3)
```

```
data1.head(3)
```

```
(9703, 22)
```

		status_id	user_id	created_at	screen_name	text	source	repl
0	1244051801516711938	803282972317204480		2020-03-29 00:00:37	redcomunitariat	Este lunes estaremos hablando sobre la situaci...	Twitter for iPhone	
1	1244052036511051778		2476348920	2020-03-29 00:01:33	SebasCamposCol	Aquí con frío 🥶 viendo cómo pasa la cuarentena 🧐...	Twitter for Android	
2	1244052338412847104		239176842	2020-03-29 00:02:45	Jonathan_518	Hoy es #sábado, apenas es hora de bañarme y or...	Twitter for iPhone	

There's something weird, there are tweets with the same text, but they seem to have different ids, and be publish in differents dates

```
1 for i in data1[data1['text'].duplicated()].index[:3]:
2     print(f"[ Id: {data1.loc[i][data1.columns[0]]},\n"
3           f" user id: {data1.loc[i][data1.columns[1]]},\n"
4           f" date: {data1.loc[i][data1.columns[2]]},\n"
5           f" text: {data1.loc[i][data1.columns[4]]}\n\n")
```

```
[ Id: 1245834468935340033,
user id: 376451765,
date: 2020-04-02 22:04:18,
text: #FelizJueves #QuedateEnCasa #EstaEnTusManos #LaVidaEsSagrada #COVID19 #HaciendoCiudad #Con
]

[ Id: 1245881819741990914,
user id: 376451765,
date: 2020-04-03 01:12:27,
text: #FelizJueves #QuedateEnCasa #EstaEnTusManos #LaVidaEsSagrada #COVID19 #HaciendoCiudad #Con
]

[ Id: 1247637085756035075,
user id: 376451765,
date: 2020-04-07 21:27:15,
text: #FelizMartes #QuedateEnCasa #EstaEnTusManos #LaVidaEsSagrada #COVID19 #HaciendoCiudad #Con
]
```

▼ Categorical variables

- screen\_name
- source
- reply\_to\_screen\_name
- is\_quote
- place\_full name
- place\_type
- lang
- country\_code

```
1 print(f"Number of screen_names: {data1['screen_name'].shape[0]}")
2 print(f"Number of unique screen_names: {len(set(data1['screen_name'].unique()))}")
3 first_twenty = data1['screen_name'].value_counts()[:25].index
4
5 plt.figure(figsize=(10, 8))
6 sns.countplot(y='screen_name',
7               data = data1[data1['screen_name'].apply(lambda x: x in first_twenty)],
8               order = data1[data1['screen_name'].apply(lambda x: x in first_twenty)]
9               ['screen_name'].value_counts().sort_values(ascending=False).index)
10
11 plt.title("25 most-frequent screennames for tweets", fontsize=20, verticalalignment='bottom')
```



Number of screen\_names: 9703

Number of unique screen\_names: 3393

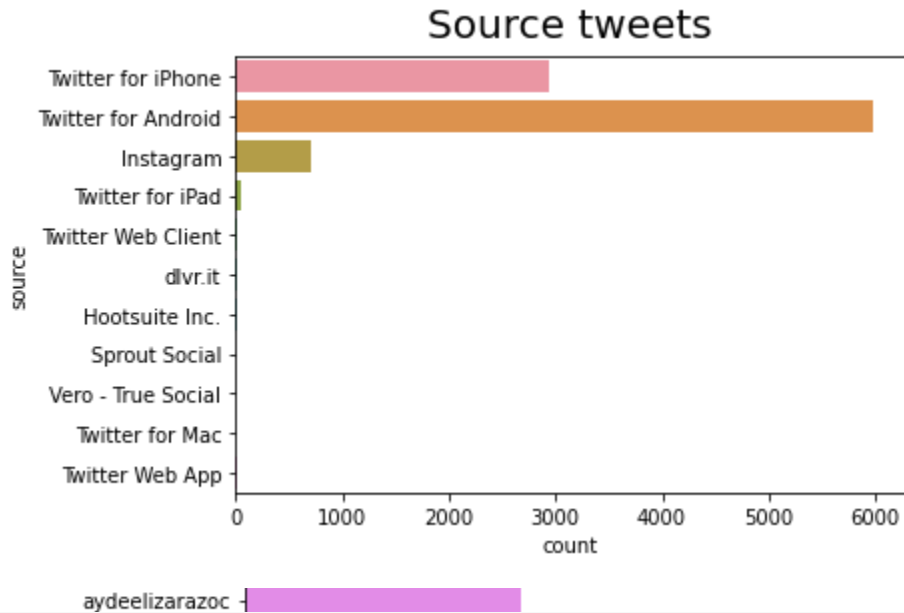
Text(0.5, 1.0, '25 most-frequent screennames for tweets')

## 25 most-frequent screennames for tweets

```
1 sns.countplot(y='source', data = data1)
```

```
2 plt.title("Source tweets", fontsize=20, verticalalignment='bottom')
```

```
3 Text(0.5, 1.0, 'Source tweets')
```



```
1 print(f"Number of replies screen_names: {data1['reply_to_screen_name'].shape[0]}")
```

```
2 print(f"Number of unique replies screen_names: {len(set(data1['reply_to_screen_name'].unique()))}")
```

```
3 first_twenty = data1['reply_to_screen_name'].value_counts()[:25].index
```

```
4
```

```
5 plt.figure(figsize=(10, 8))
```

```
6 sns.countplot(y='reply_to_screen_name',
```

```
7             data = data1[data1['reply_to_screen_name'].apply(lambda x: x in first_twenty)],
```

```
8             order = data1[data1['reply_to_screen_name'].apply(lambda x: x in first_twenty)]
```

```
9             ['reply_to_screen_name'].value_counts().sort_values(ascending=False).index)
```

```
10
```

```
11 plt.title("25 most-frequent screennames for replies tweets", fontsize=20, verticalalignment='bottom')
```

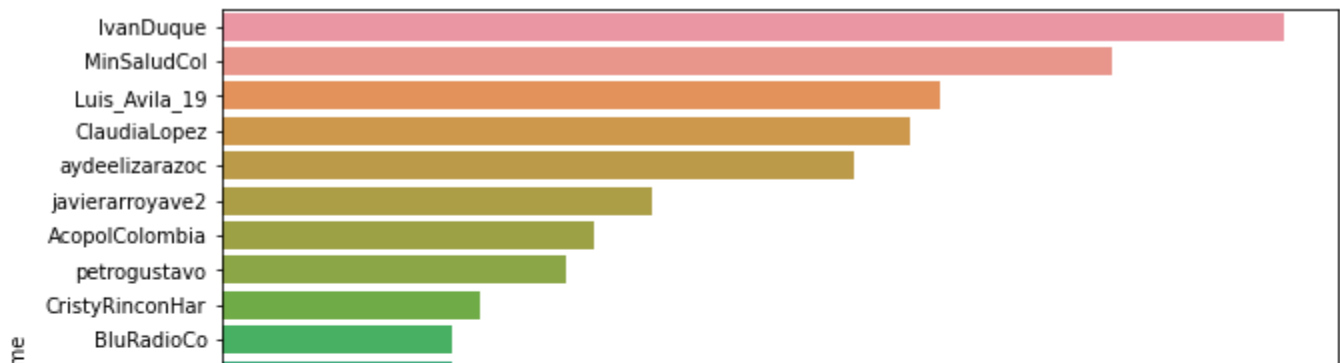
```
12
```

Number of replies screen\_names: 9703

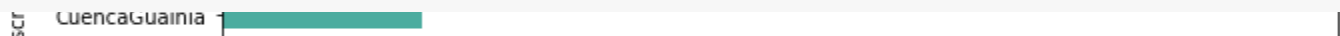
Number of unique replies screen\_names: 855

Text(0.5, 1.0, '25 most-frequent screennames for replies tweets')

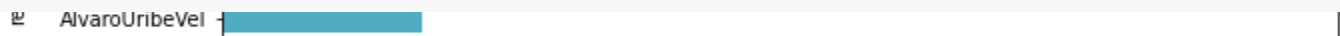
## 25 most-frequent screennames for replies tweets



```
1 df_temp = data1[~data1['reply_to_screen_name'].isnull()].copy()
```



```
1 # print("\n\n".join(df_temp[df_temp['reply_to_screen_name'] == 'MinSaludCol']['text'].apply(lambda x
```



Here we notice that the tweets that has more replies are the ones related to a colombian politician or a colombian entity.

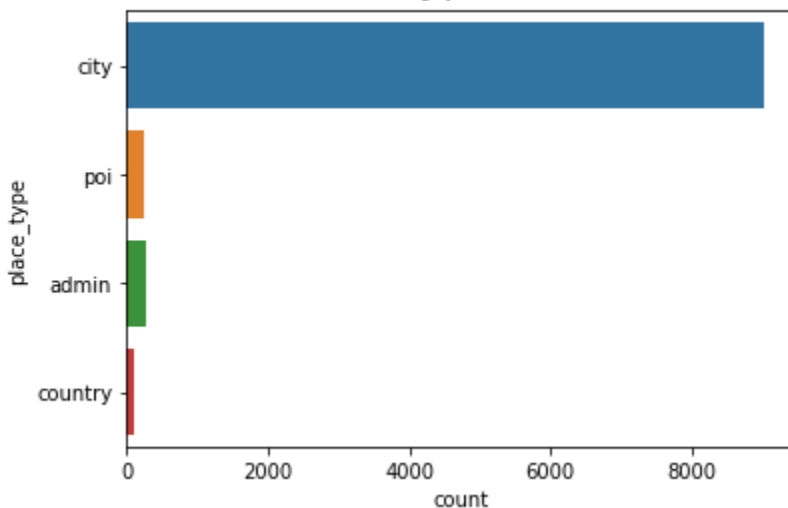


```
1 sns.countplot(y='place_type', data = data1)
```

```
2 plt.title("Place type tweets", fontsize=20, verticalalignment='bottom')
```

Text(0.5, 1.0, 'Place type tweets')

## Place type tweets



```
1 df = data1[data1['place_type'] == 'city'].copy()
```

```
2
```

```
3 print(f"Number of cities: {df['place_full_name'].shape[0]}")
```

```
4 print(f"Number of unique cities: {len(set(df['place_full_name'].unique()))}")
```

```
5 first_twenty = df['place_full_name'].value_counts()[:10].index
```

```
6
```

```
7 plt.figure(figsize=(10, 5))
```

```
8 sns.countplot(y='place_full_name',
```

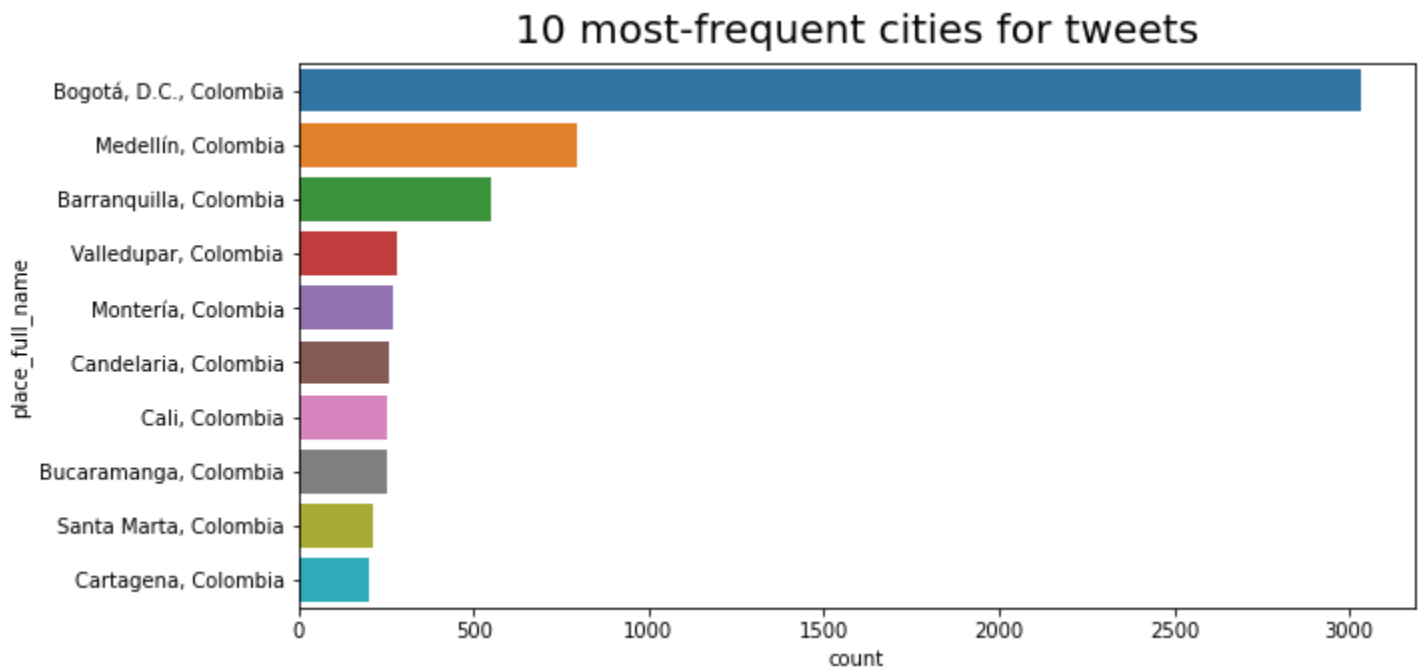


```

9         data = df[df['place_full_name'].apply(lambda x: x in first_twenty)],
10         order = df[df['place_full_name'].apply(lambda x: x in first_twenty)]
11             ['place_full_name'].value_counts().sort_values(ascending=False).index)
12
13 plt.title("10 most-frequent cities for tweets", fontsize=20, verticalalignment='bottom')
14

```

➞ Number of cities: 9023  
Number of unique cities: 258  
Text(0.5, 1.0, '10 most-frequent cities for tweets')



```

1 df = data1[data1['place_type'] == 'poi'].copy()
2
3 print(f"Number of poi: {df['place_full_name'].shape[0]}")
4 print(f"Number of unique poi: {len(set(df['place_full_name'].unique()))}")
5 first_twenty = df['place_full_name'].value_counts()[:10].index
6
7 plt.figure(figsize=(10, 5))
8 sns.countplot(y='place_full_name',
9             data = df[df['place_full_name'].apply(lambda x: x in first_twenty)],
10             order = df[df['place_full_name'].apply(lambda x: x in first_twenty)]
11                 ['place_full_name'].value_counts().sort_values(ascending=False).index)
12
13 plt.title("10 most-frequent poi for tweets", fontsize=20, verticalalignment='bottom')
14
15

```

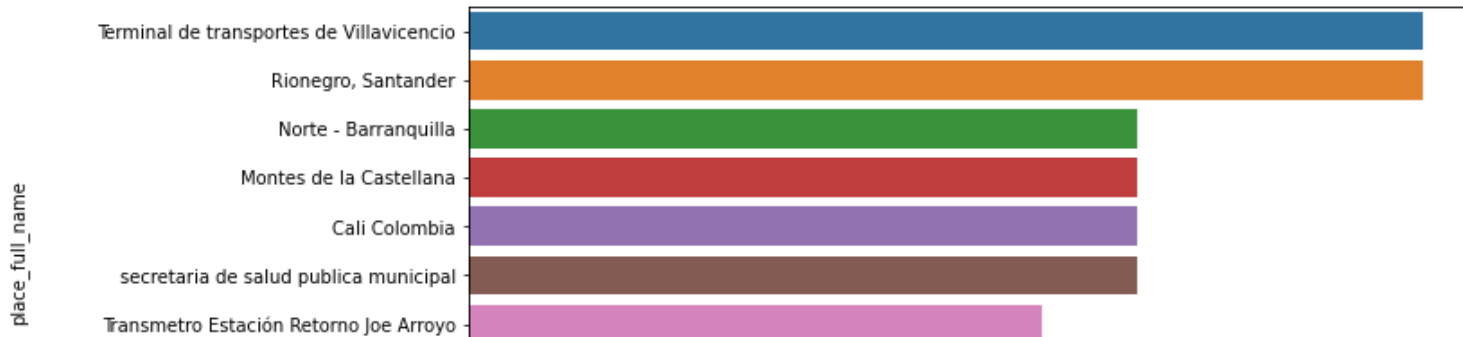
➞

Number of poi: 239

Number of unique poi: 154

Text(0.5, 1.0, '10 most-frequent poi for tweets')

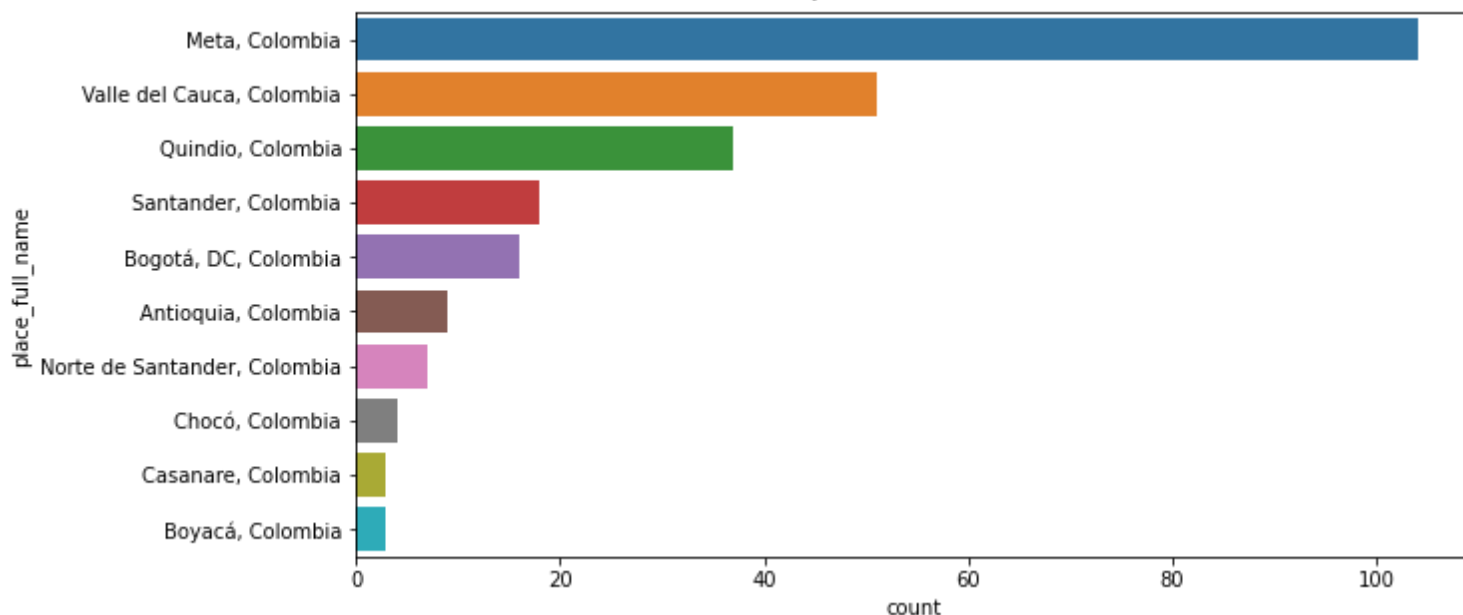
10 most-frequent poi for tweets



```
1 df = data1[data1['place_type'] == 'admin'].copy()
2
3 print(f"Number of admin: {df['place_full_name'].shape[0]}")
4 print(f"Number of unique admin: {len(set(df['place_full_name'].unique()))}")
5 first_twenty = df['place_full_name'].value_counts()[:10].index
6
7 plt.figure(figsize=(10, 5))
8 sns.countplot(y='place_full_name',
9               data = df[df['place_full_name'].apply(lambda x: x in first_twenty)],
10              order = df[df['place_full_name'].apply(lambda x: x in first_twenty)]
11                    ['place_full_name'].value_counts().sort_values(ascending=False).index)
12
13 plt.title("10 most-frequent admin for tweets", fontsize=20, verticalalignment='bottom')
14
```

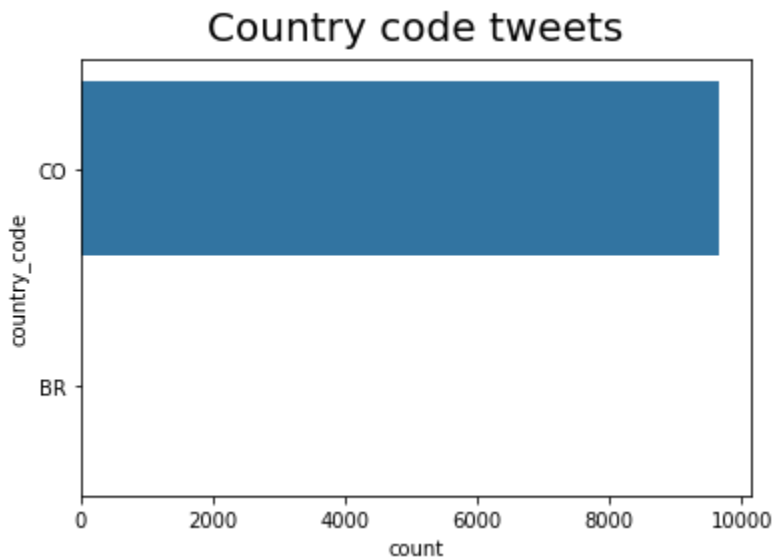
➡ Number of admin: 265  
Number of unique admin: 19  
Text(0.5, 1.0, '10 most-frequent admin for tweets')

10 most-frequent admin for tweets



```
1 sns.countplot(y='country_code', data = data1)
2 plt.title("Country code tweets", fontsize=20, verticalalignment='bottom')
```

```
Text(0.5, 1.0, 'Country code tweets')
```



▼ Since all tweets are from Colombia is not need it to have the feature Country code

```
1 data2 = data1[data1['country_code'] == 'CO'].drop(['country_code'], axis=1).copy()
2 data2['place_full_name'] = data2['place_full_name'].apply(lambda x: x.split(',')[0])
```

```
1 cities_population = {'Bogotá': 7413000,
2                       'Medellín': 2529403,
3                       'Barranquilla': 1274250,
4                       'Valledupar': 490075,
5                       'Montería': 460223,
6                       'Candelaria': 23985,
7                       'Cali': 2471474,
8                       'Bucaramanga': 581130,
9                       'Santa Marta': 515556,
10                      'Cartagena': 1036412
11                      }
```

```
1 print(f"Number of places: {data2['place_full_name'].shape[0]}")
2 print(f"Number of unique places: {len(set(data2['place_full_name'].unique()))}")
3 first_twenty = data2['place_full_name'].value_counts()[:10].index
4
5 plt.figure(figsize=(10, 8))
6 sns.countplot(y='place_full_name',
7               data = data2[data2['place_full_name'].apply(lambda x: x in first_twenty)],
8               order = data2[data2['place_full_name'].apply(lambda x: x in first_twenty)]
9                     ['place_full_name'].value_counts().sort_values(ascending=False).index)
10
11 plt.title("10 most-frequent place_full_name for tweets", fontsize=20, verticalalignment='bottom')
```

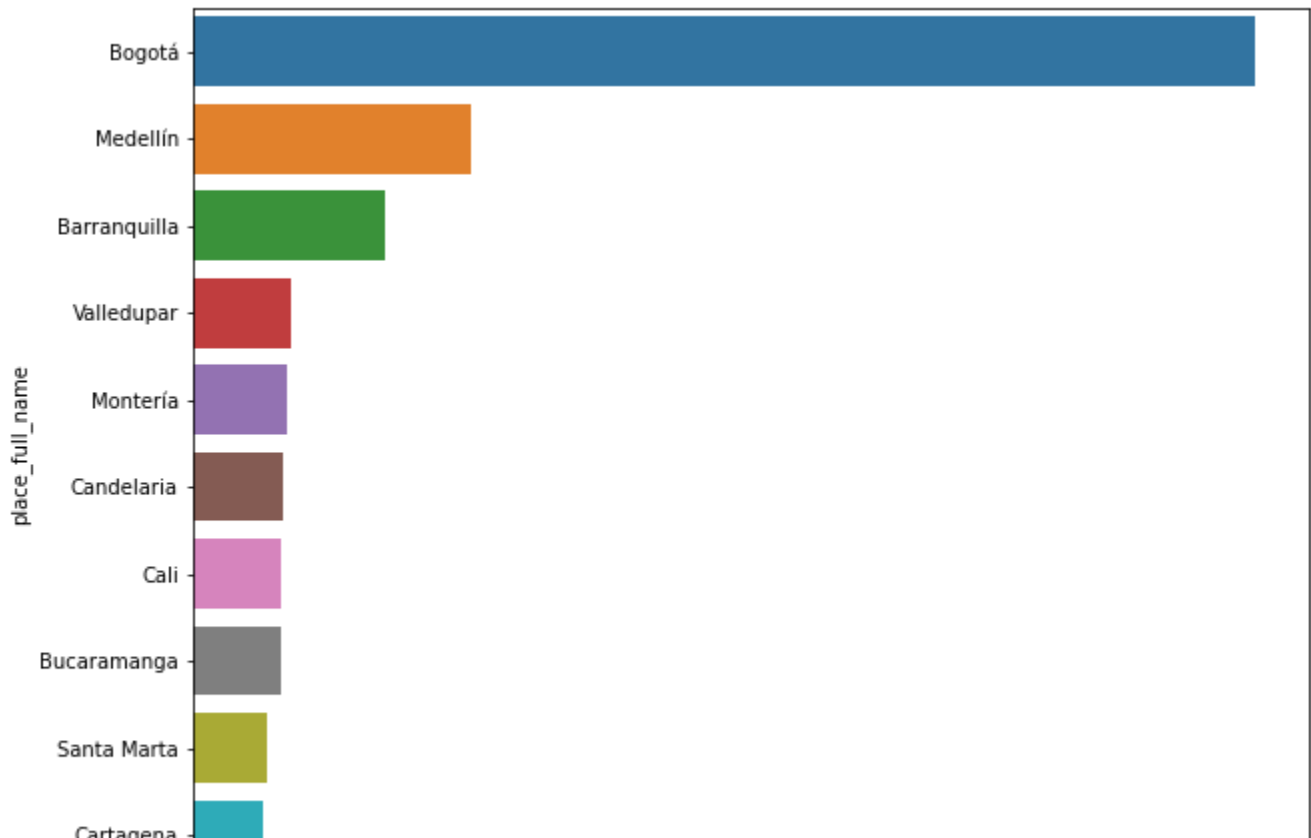
```
Text(0.5, 1.0, 'Country code tweets')
```

Number of places: 9642

Number of unique places: 418

Text(0.5, 1.0, '10 most-frequent place\_full\_name for tweets')

## 10 most-frequent place\_full\_name for tweets

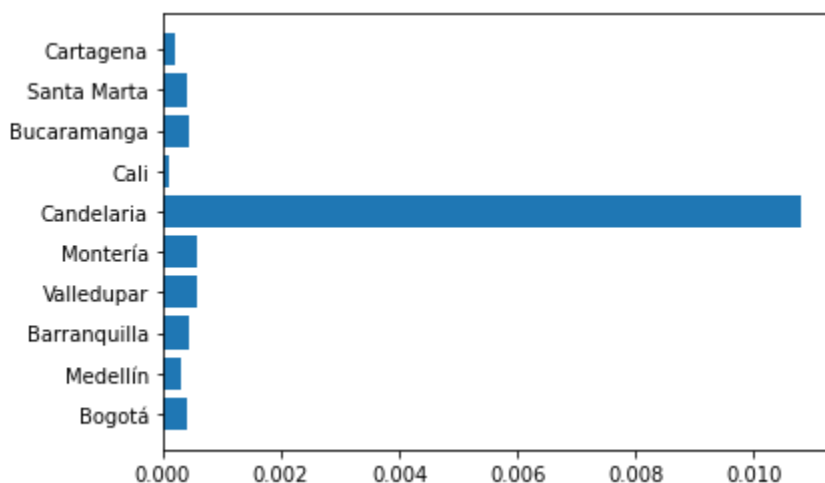


We notice that most of the tweets are from Bogotá, however there are more people in Bogotá than in other cities, so we're going to normalize data with population per city

```
1 cities_names = data2['place_full_name'].value_counts()[:10].index.tolist()
2 cities_count = data2['place_full_name'].value_counts()[:10].tolist()
3 cities_count_ = [cities_count[i]/cities_population[cities_names[i]] for i in range(len(cities_count_))]
```

```
1 plt.barh(cities_names, cities_count_)
```

<BarContainer object of 10 artists>



Another interesting thing is that there are plenty of tweets of a Bogotá neighborhood (Candelaria). So, we're going to add it to Bogotá

```
1 cities_names_ = cities_names.copy()
2 cities_count_ = cities_count.copy()
3
```

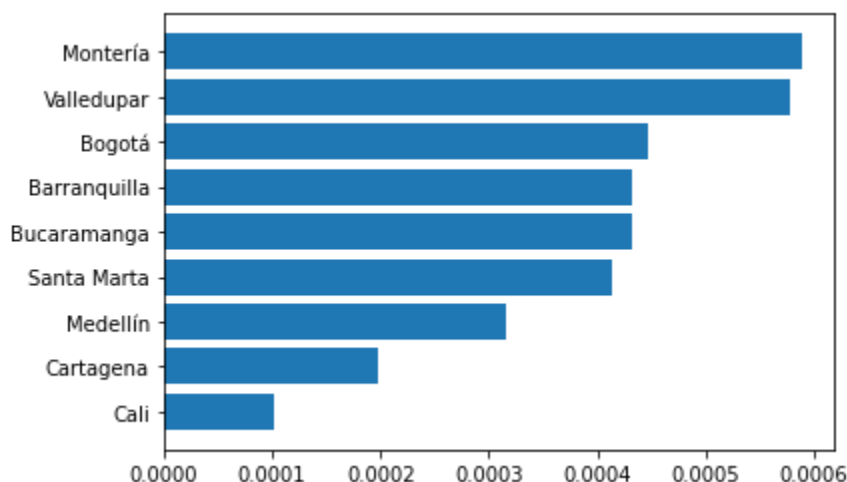
```
1 cities_names_.pop(5)
2 cities_count_[0] += cities_count_.pop(5)
3 cities_count_ = [cities_count_[i]/cities_population[cities_names_[i]] for i in range(len(cities_count_))]
```

```
1 cities_most_freq = [(cities_names_[i], cities_count_[i]) for i in range(len(cities_count_))]
2 cities_most_freq.sort(key=lambda x:x[1])
3 cities_most_freq
```

```
[('Cali', 0.00010236806051773153),
 ('Cartagena', 0.00019683291972690399),
 ('Medellín', 0.00031628016571499284),
 ('Santa Marta', 0.0004131461955636245),
 ('Bucaramanga', 0.00043191712697675217),
 ('Barranquilla', 0.00043241122228762017),
 ('Bogotá', 0.00044637798462161067),
 ('Valledupar', 0.0005774626332704178),
 ('Montería', 0.0005888449729804898)]
```

```
1 plt.barh([i[0] for i in cities_most_freq], [i[1] for i in cities_most_freq])
```

```
[<BarContainer object of 9 artists>
```

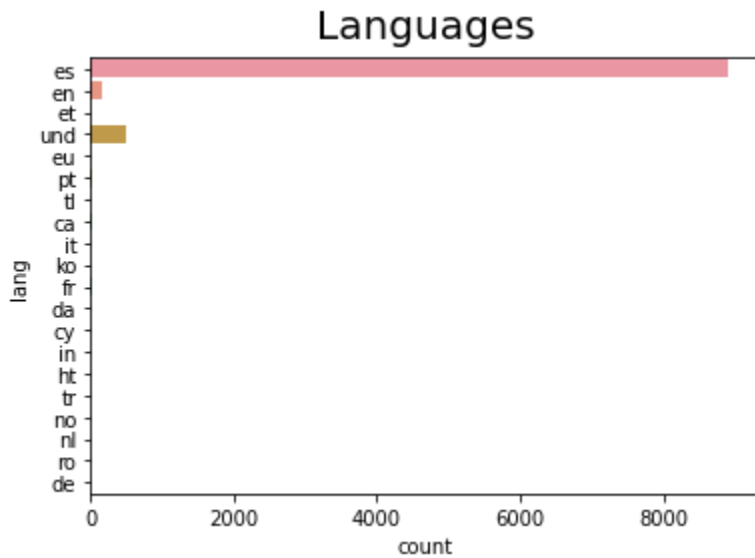


Now, we see that actually the city with a bigger percentage of tweets is Monteria, however Bogotá doesn't stand behind, it is in third place.

```
1 sns.countplot(y='lang', data = data2)
2 plt.title("Languages" , fontsize=20 , verticalalignment='bottom')
```

```
2 plt.title('Languages', fontsize=20, verticalalignment='bottom')
```

```
↳ Text(0.5, 1.0, 'Languages')
```



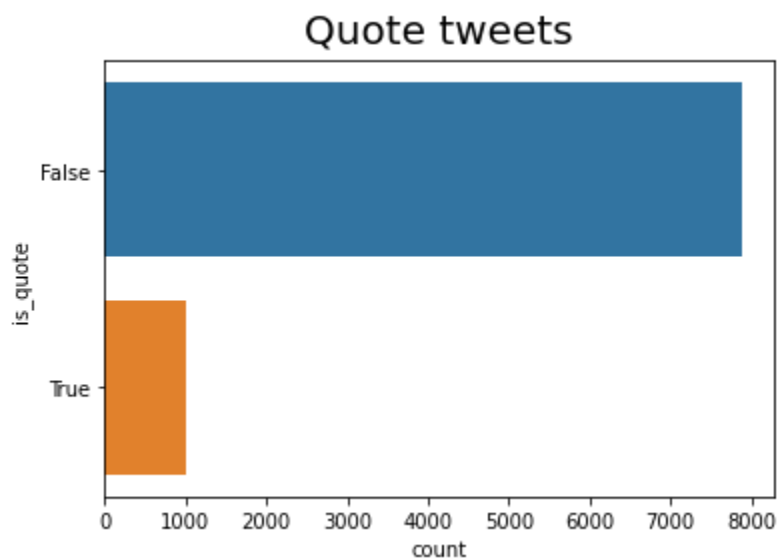
There are plenty of languages, since we are going to use NLP for spanish text, then we're going to drop tweets with other languages

```
1 data3 = data2[data2['lang'] == 'es'].copy()  
2 print(data3.shape)
```

```
↳ (8909, 21)
```

```
1 sns.countplot(y='is_quote', data = data3)  
2 plt.title("Quote tweets", fontsize=20, verticalalignment='bottom')
```

```
↳ Text(0.5, 1.0, 'Quote tweets')
```



Variables numericas

- favorite\_count
- retweet\_count
- followers\_count
- friends\_count

To see the behavior of these features over time, we're gonna plot them all tweets, and then split the data into replies and original tweets to do the same plots

```
1 data3.columns
```

```
Index(['status_id', 'user_id', 'created_at', 'screen_name', 'text', 'source',
      'reply_to_status_id', 'reply_to_user_id', 'reply_to_screen_name',
      'is_quote', 'is_retweet', 'favourites_count', 'retweet_count',
      'place_full_name', 'place_type', 'followers_count', 'friends_count',
      'account_created_at', 'verified', 'lang', 'source_tweet'],
      dtype='object')
```

```
1 data3['created_at'] = pd.to_datetime(data3['created_at'])
2 data3['account_created_at'] = pd.to_datetime(data3['account_created_at'])
```

```
1 data3['day'] = data3['created_at'].dt.to_period('d')
```

```
1 data3[['favourites_count', 'retweet_count', 'followers_count', 'friends_count']].describe()
```

```

favourites_count  retweet_count  followers_count  friends_count
count            8909.000000      8909.000000      8.909000e+03      8909.000000
mean             8267.765518         3.173196      8.079546e+03      1508.830508
std             17778.678265        21.844543      5.086454e+04      4239.547439
min              0.000000         0.000000      0.000000e+00         0.000000
25%              533.000000         0.000000      2.350000e+02      283.000000
50%             2315.000000         0.000000      8.230000e+02      687.000000
75%             7631.000000         2.000000      3.525000e+03     1644.000000
max            272462.000000       1163.000000      1.857053e+06     104637.000000

```

```

1 def plot_time_series(df):
2     fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(15, 8))
3     cols = ['favourites_count', 'retweet_count', 'followers_count', 'friends_count']
4     for i in range(2):
5         for j in range(2):
6             col = cols[i+j] if i==0 else cols[i+j+1]
7             axes[i][j].set_title(f'{col} over time', fontsize=15)
8             df.groupby('day')[col].sum().plot(ax=axes[i][j])
9             df.groupby('day')[col].sum().rolling(window=3).mean().plot(ax=axes[i][j])

```

```

9         df.groupby('day')[col].sum().rolling(window=5).mean().plot(ax=axs[1][1])
10

```

```

1 print(f"First day: {min(data3['created_at'])}")
2 print(f>Last day: {max(data3['created_at'])}")

```

```

➤ First day: 2020-03-26 18:37:06
  Last day: 2020-04-30 23:59:12

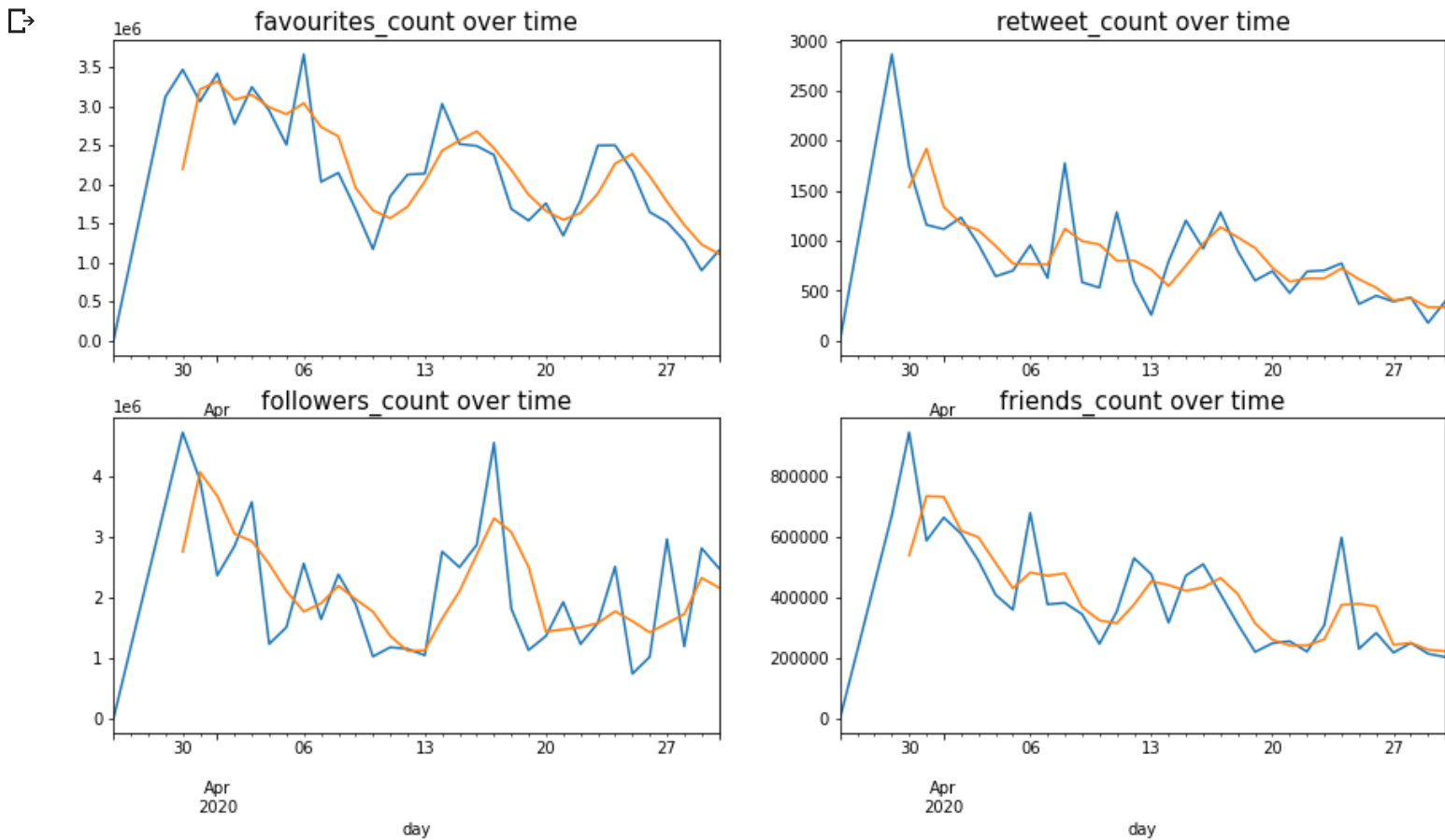
```

## ▼ All tweets

```

1 plot_time_series(data3)

```



## ▼ Original tweets

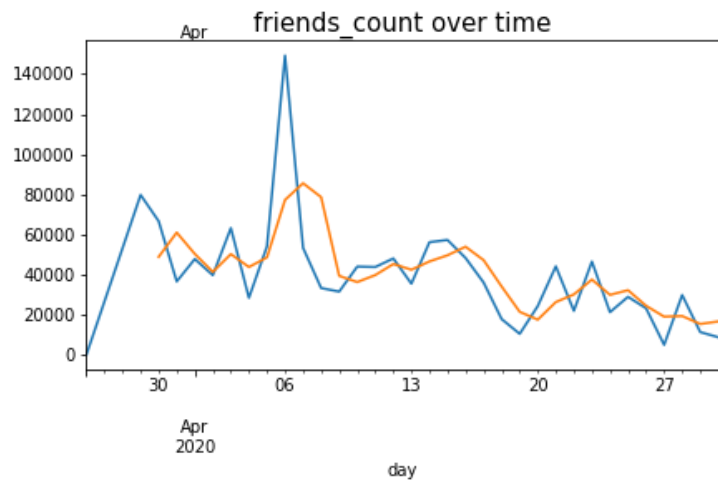
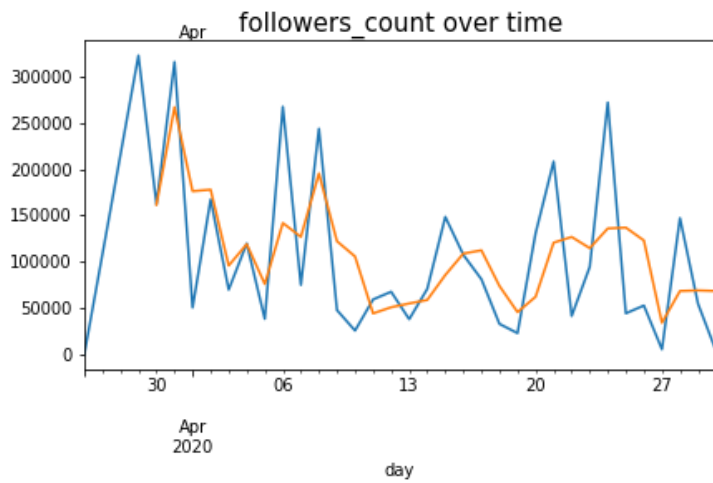
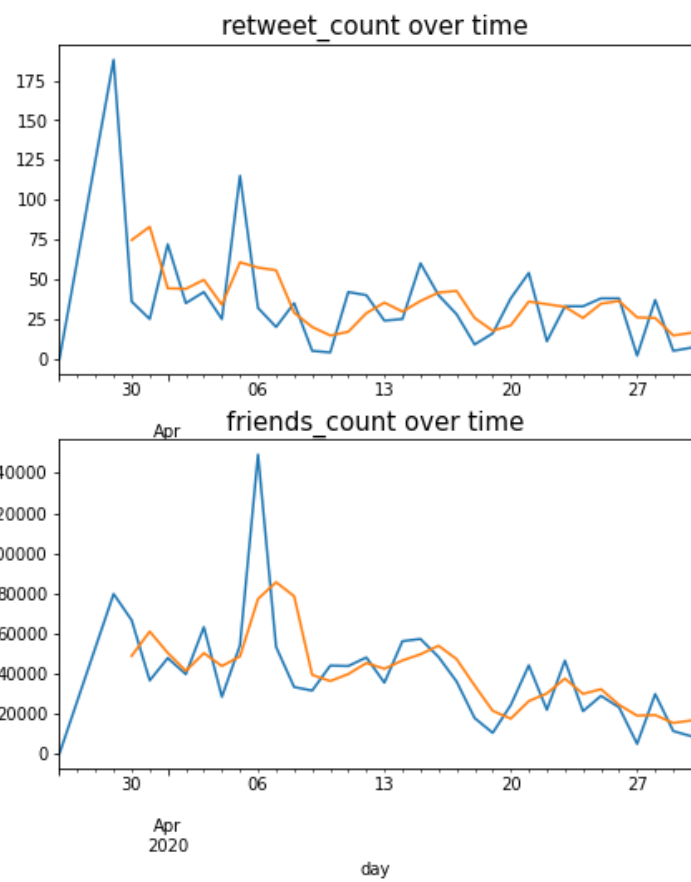
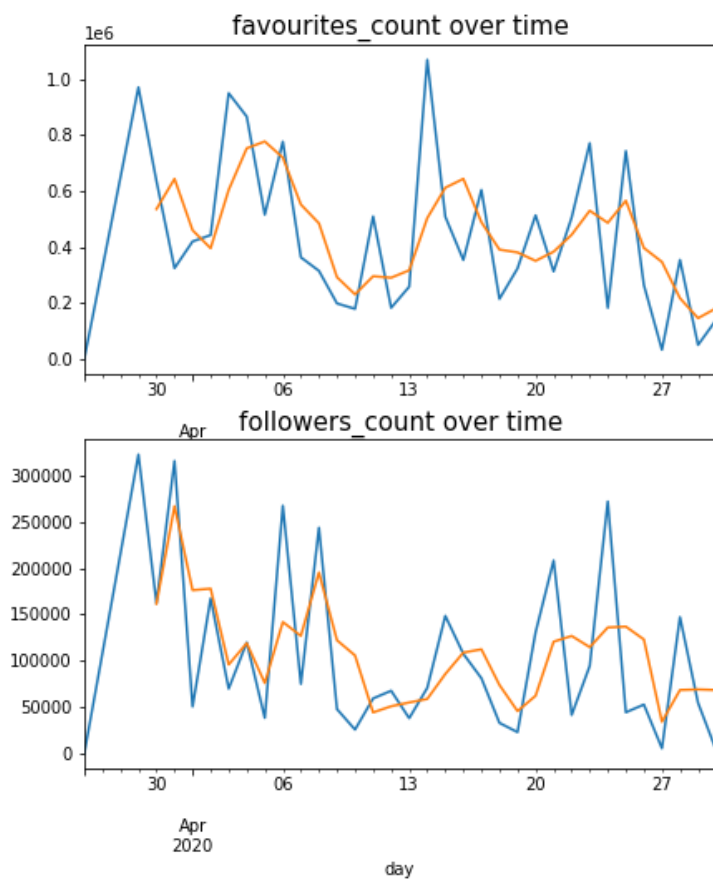
```

1 original_tweets = data3[~data3['reply_to_status_id'].isnull() ]
2 plot_time_series(original_tweets)

```

➤

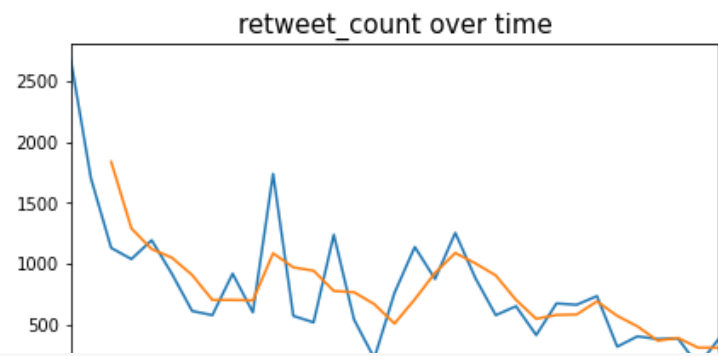
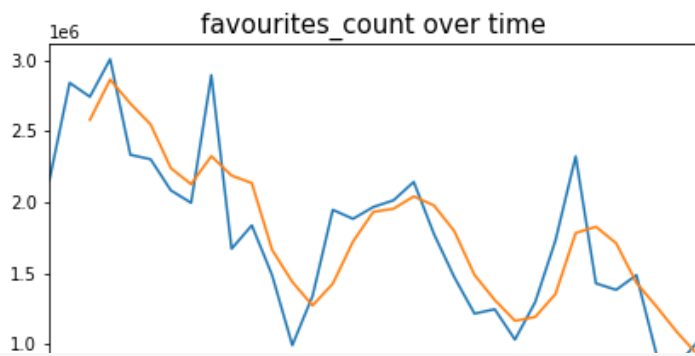




## ▼ Replies tweets

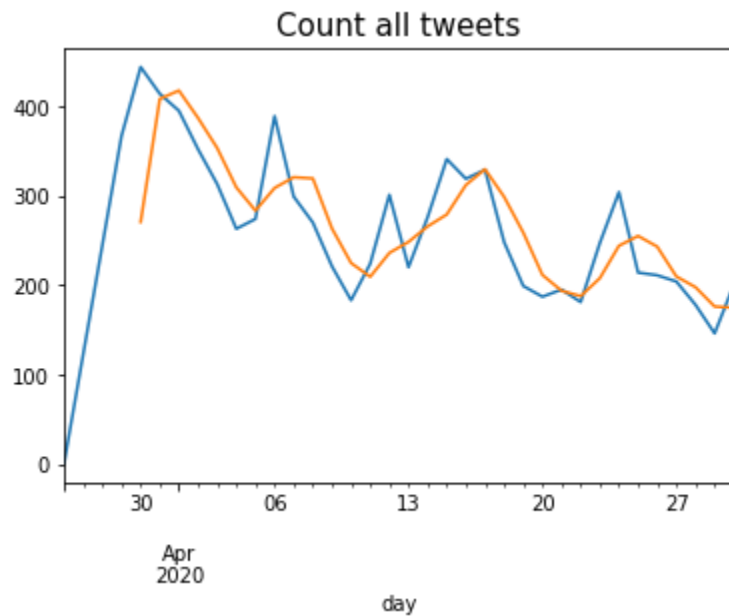
```
1 replies_tweets = data3[data3['reply_to_status_id'].isnull() ]
2 plot_time_series(replies_tweets)
```





```
1 plt.title('Count all tweets', fontsize=15)
2 data3.groupby('day')['status_id'].count().plot()
3 data3.groupby('day')['status_id'].count().rolling(window=3).mean().plot()
4
```

☞ <matplotlib.axes.\_subplots.AxesSubplot at 0x7fb5db69a7b8>



```
1 plt.title('Count originals tweets', fontsize=15)
2
3 original_tweets.groupby('day')['status_id'].count().plot()
4 original_tweets.groupby('day')['status_id'].count().rolling(window=3).mean().plot()
```

☞

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5db595ef0>
```

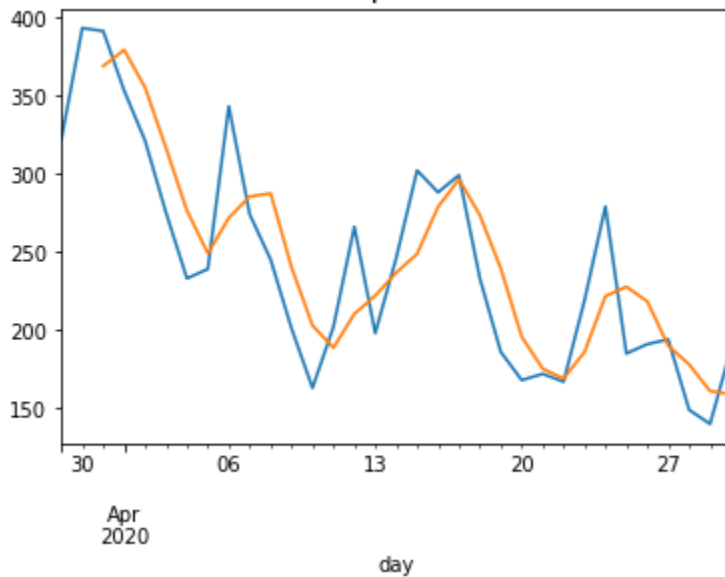
Count originals tweets



```
1 plt.title('Count replies tweets', fontsize=15)
2
3 replies_tweets.groupby('day')['status_id'].count().plot()
4 replies_tweets.groupby('day')['status_id'].count().rolling(window=3).mean().plot()
```

```
↳ <matplotlib.axes._subplots.AxesSubplot at 0x7fb5db586470>
```

Count replies tweets



```
1
```

## ▼ Natural Language Proccesing

```
1 # !pip install wordcloud
```

```
1 import nltk
2 import contractions
3 import re
4 import emoji
5 import wordcloud as wc
6 import string
7
8
9 from collections import Counter
10
11 nltk.download('punkt')
12 nltk.download('wordnet')
13 nltk.download('stopwords')
14 nltk.download('averaged_perceptron_tagger')
15 stop_words = nltk.corpus.stopwords.words('spanish')
```

```
➤ [nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /root/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
```

1

```
1 def create_wordcloud(text, plt_):
2     wordcloud = wc.WordCloud(max_font_size=50,
3                               max_words=100,
4                               scale=10,
5                               background_color="white")\
6                               .generate(text)
7     plt.figure(figsize=(10,5))
8     plt_.imshow(wordcloud, interpolation="bilinear")
9     plt_.axis("off")
10 #     plt.show()
11
12 def remove_links(text):
13     text_ = re.sub(r'https?:\//.*[\r\n]*', '', text, flags=re.MULTILINE)
14     return text_
15
16
17 def remove_space(text):
18     text_ = re.sub(r'[\r|\n|\r\n|]+', '\n', text)
19     return text_
20
21
22 def extract_emojis(text):
23     return ''.join(c for c in text if c in emoji.UNICODE_EMOJI)
24
25 def get_hashtags(text):
26     patter = re.compile(r"# (\w+)")
27     all_hashtags = [f"#{i}" for i in patter.findall(text)]
28     return all_hashtags
29
30 def get_mentions(text):
31     patter = re.compile(r"@ (\w+)")
32     all_hashtags = [f"@{i}" for i in patter.findall(text)]
33     return all_hashtags
34
35 def remove_special_characters(text):
36     punctuation = string.punctuation+";"
37     text = re.sub(f"[{punctuation}]", '', text)
38     return text
39
40 def remove_repeated_word(text):
```

```

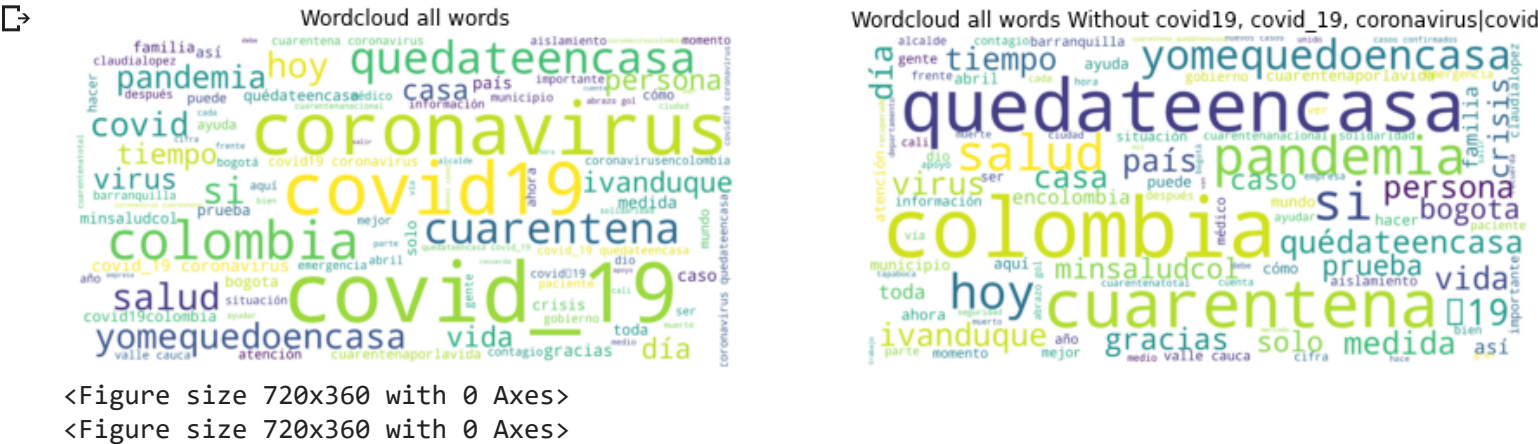
41     tokens = nltk.word_tokenize(" ".join(data3['text']))
42
43
44 def remove_stopwords(text):
45     tokens = nltk.word_tokenize(text)
46     tokens = [token.strip() for token in tokens]
47     filtered_tokens = [token for token in tokens if token.lower() not in stop_words]
48     filtered_text = " ".join(filtered_tokens)
49     return filtered_text
50
51 def clean_text(text):
52     text1 = remove_links(text)
53     text2 = remove_space(text1)
54     text3 = remove_stopwords(text2)
55     return text3

```

- ▼ Wordcloud of all words

```
1 text = clean_text(" ".join(data3['text']))
```

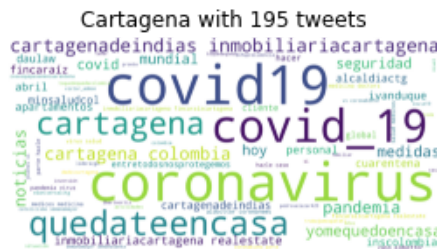
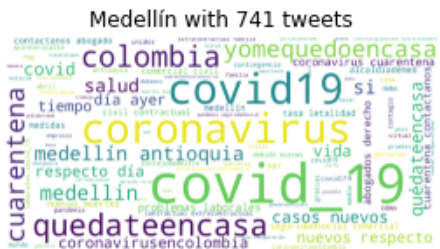
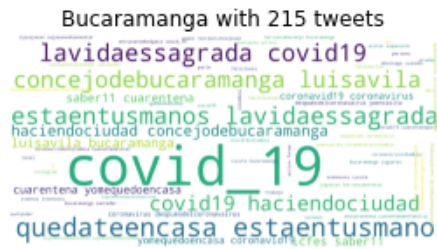
```
1 fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(15, 8))
2
3 axes[0].set_title('Wordcloud all words')
4 axes[1].set_title('Wordcloud all words Without covid19, covid_19, coronavirus|covid')
5
6 create_wordcloud(text.lower(), axes[0])
7 create_wordcloud(re.sub(r"covid19|covid_19|coronavirus|covid", "", text.lower()), axes[1])
```



```
1 create_wordcloud(re.sub('(@\w*)|(#\w*)|(covid19)|(covid_19)|(coronavirus)|(covid)', '', text.lower(
```







<Figure size 720x360 with 0 Axes>

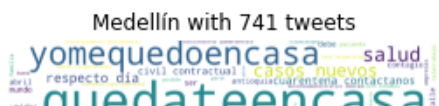
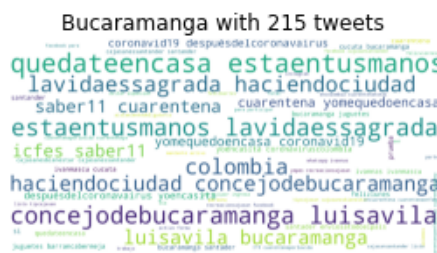
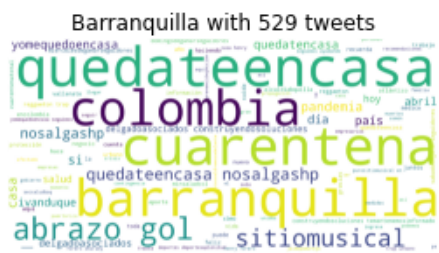
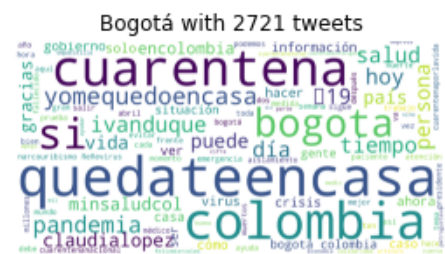
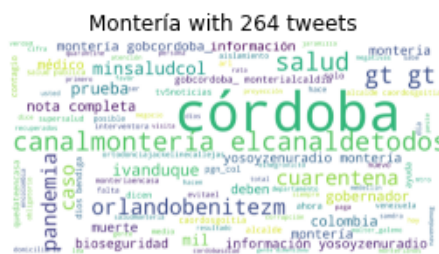
## Without [covid19|covid\_19|coronavirus|covid]

<Figure size 720x360 with 0 Axes>

```
1 fig, axes = plt.subplots(nrows=3, ncols=3, figsize=(15, 8))
2 c = 0
3 for city in most_frequent_cities:
4     temp_data = data3[data3['place_full_name'] == city].copy()
5     # print(f"{city} has {temp_data.shape} tweets")
6     temp_text = clean_text(" ".join(temp_data['text']))
7     axes[c//3][c%3].set_title(f"{city} with {temp_data.shape[0]} tweets")
8     create_wordcloud(re.sub(r"covid19|covid_19|coronavirus|covid", "", temp_text.lower()), axes[c//3][c%3])
9     c+=1
10
```



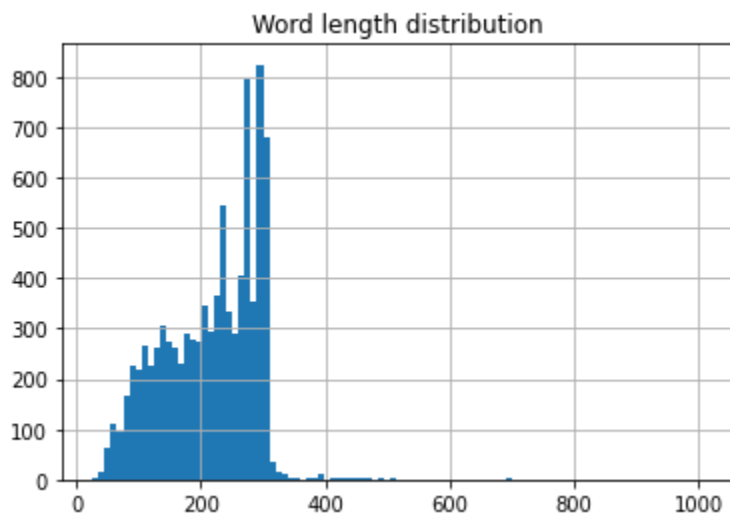




```
1 data3['len_text'] = data3['text'].apply(lambda x: len(x))
```

```
1 print(data3.shape)
2 plt.title('Word length distribution')
3 data3['len_text'].sort_values().hist(bins=100)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5d2373eb8>
```



- Also we have plenty of hashtags, mentions and emojis. So, let's see what we got
- Emojis

```
1 import regex
```

```
1 def get_emojis(text):
2     emoji_list = []
```



```

3 data = regex.findall(r'\X', text)
4 for word in data:
5     if any(char in emoji.UNICODE_EMOJI for char in word):
6         emoji_list.append(word)
7
8 return emoji_list

```

```
1 emojis_ = extract_emojis(text)
2 emojis_[:200]
```


'  \U0001f9a0  \U0001f9ee                                                  

```
1 # Counter(emojis_).most_common(50)
```

```
1 emojis2 = Counter(get_emojis_(text)).most_common()
2 x = [i[0] for i in emojis2]
3 y = [i[1] for i in emojis2]
4 print([i[0] for i in emojis2[1:250]])
```

☞ [🤪, 🖐️, \U0001f9a0', ⬇️, ⬇️, 👉, 👈, 👉, 🖐️, 👈, 🔑, 💎, ➡️, ✅, ]

```
1 plt.plot(x[:10], y[:10])
```

☞ Figure 10.10 illustrates the relationship between the  $\beta$  coefficient and the correlation coefficient  $r$ .

```
[<matplotlib.lines.Line2D at 0x7fb5d34c7be0>]/usr/local/lib/python3.6/dist-packages/matplotlib/ba
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Gl
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Gl
font.set_text(s, 0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Gl
font.set_text(s, 0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Gl
font.set_text(s, 0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Gl
font.set_text(s, 0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Gl
```

## ▼ Hashtags

```
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Gl
```

```
1 hashtags = get_hashtags(text)
2 print(len(hashtags))
3 print(hashtags[:50])
```

```
➞ 12769
['#Covid_19', '#Coronavirus', '#Covid_19', '#QueVainabella', '#Covid_19', '#SabadoDeGanarSeguidor
```

```
1 hashtags2 = Counter(hashtags).most_common(30)
2 x = [i[0] for i in hashtags2]
3 y = [i[1] for i in hashtags2]
4 print(hashtags2)
```

```
➞ [('Covid_19', 1165), ('coronavirus', 807), ('COVID19', 545), ('QuedateEnCasa', 369), ('Coron
```

```
1 # plt.figure(figsize=(10, 7))
2 # plt.barh(x[::-1], y[::-1])
3 x[:10]
```

```
➞
```

```
['#Covid_19',
 '#coronavirus',
 '#COVID19',
 '#QuedateEnCasa',
 '#Coronavirus']
```

## ▼ Daiy use of hashtags

```
1 df_count_hashtags = data3[['day']].drop_duplicates().copy()
```

```
1 ten_most_frequent_hashtags = x[:15]
2 dates = data3['day'].unique()
3
4 for hashtag in ten_most_frequent_hashtags:
5     daily_use = []
6     for date_ in dates:
7         temp_df = data3[data3['day'] == date_].copy()
8         daily_use.append(sum(temp_df['text'].apply(lambda x: 1 if hashtag in x else 0)))
9     df_count_hashtags[hashtag] = daily_use
10
11 #     print(daily_use)
```

```
1 df_count_hashtags.head()
```



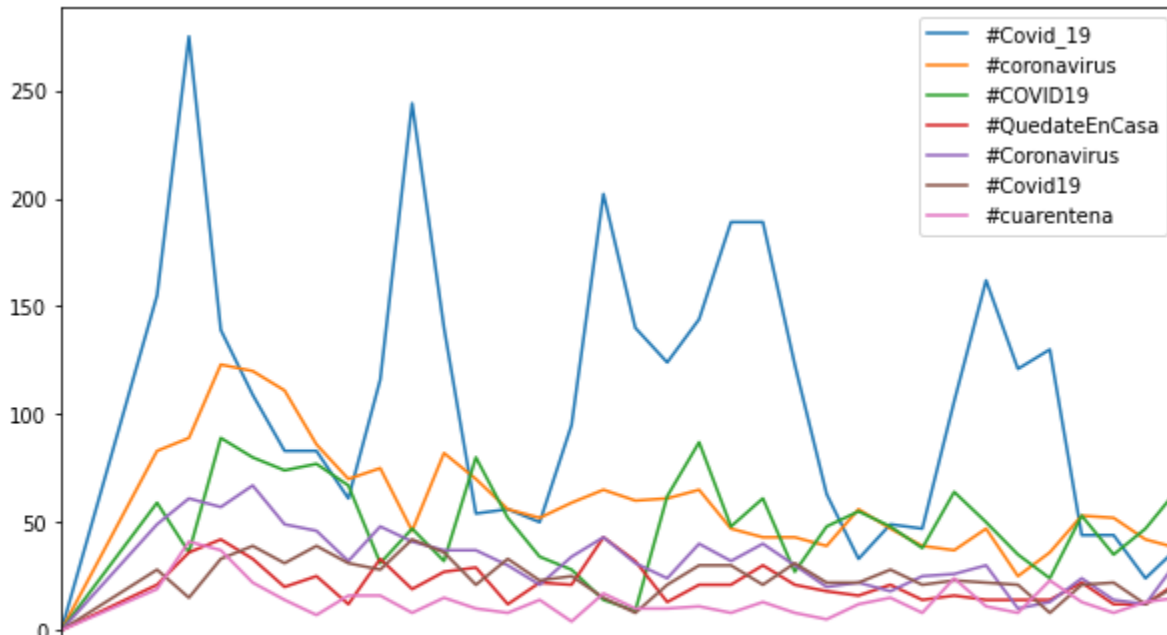
	day	#Covid_19	#coronavirus	#COVID19	#QuedateEnCasa	#Coronavirus	#Covid19	#cuarente
0	2020-03-29	155	83	59	21	49	28	
386	2020-03-30	275	89	36	36	61	15	
858	2020-03-31	139	123	89	42	57	33	
1305	2020-04-01	109	120	80	33	67	39	
1726	2020-04-02	83	111	74	20	49	31	

```
1 df_count_hashtags.index = df_count_hashtags['day']
```

```
1 df_count_hashtags[df_count_hashtags.columns[1:8]].sort_index().plot(figsize=(10, 6))
```



<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb5d236ee10>



Abr

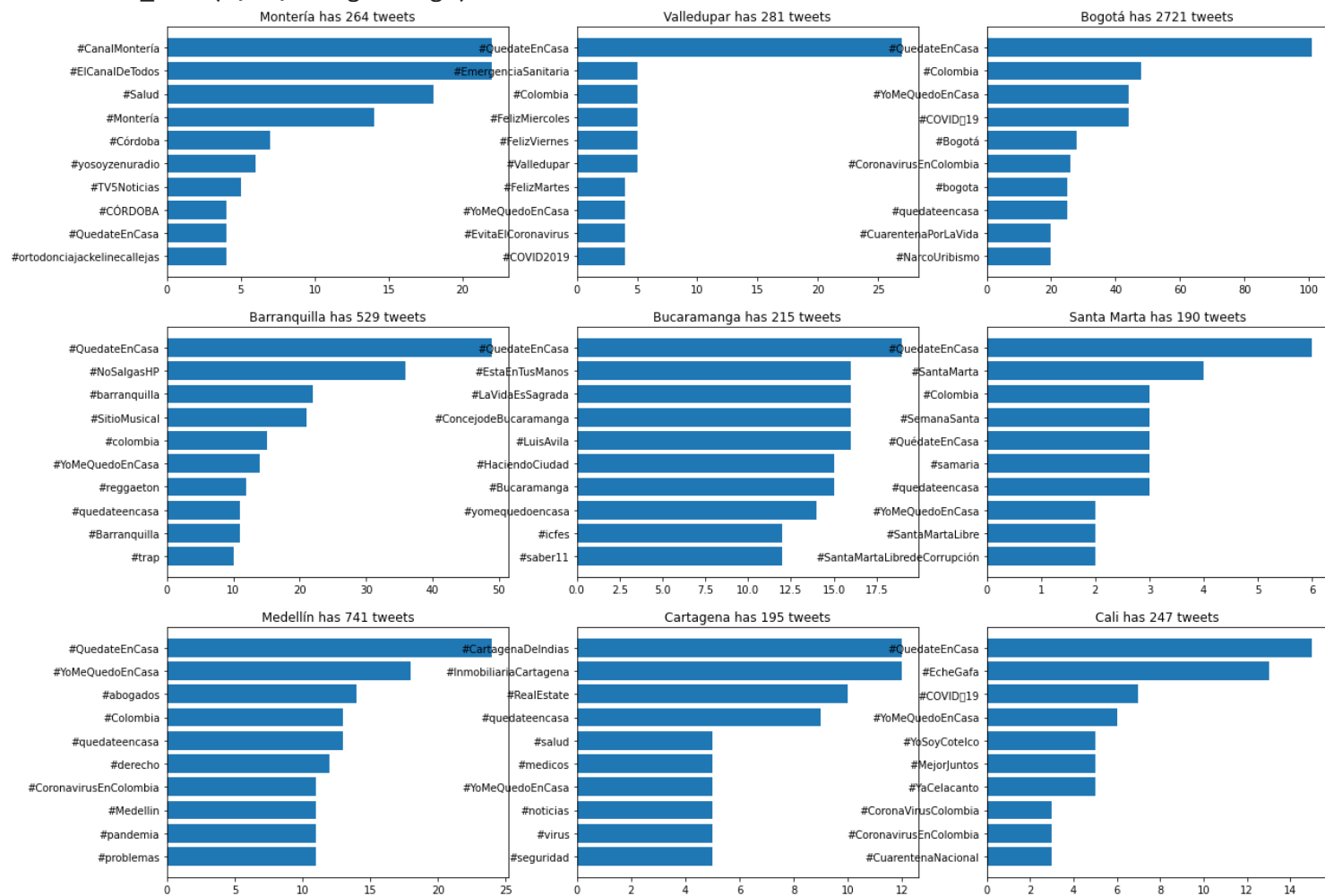
```
1 fig, axes = plt.subplots(nrows=3, ncols=3, figsize=(20, 15))
2 c=0
3 for city in most_frequent_cities:
4     temp_data = data3[data3['place_full_name'] == city].copy()
5     # print(f"{city} has {temp_data.shape} tweets")
6     temp_text = clean_text(" ".join(temp_data['text']))
7     temp_hashtags = get_hashtags(temp_text)
8     temp_hashtags2 = Counter(temp_hashtags).most_common()
9     del_covid_words = ["#covid", "#covid19", "#coronavirus", "#cuarentena", "#covid19colombia", "#c"]
10    temp_hashtags2 = [i for i in temp_hashtags2 if not i[0].lower() in del_covid_words][:10]
11    x = [i[0] for i in temp_hashtags2]
12    y = [i[1] for i in temp_hashtags2]
13    axes[c//3][c%3].set_title(f"{city} has {temp_data.shape[0]} tweets")
14    axes[c//3][c%3].barh(x[::-1], y[::-1])
15    c+=1
```



```

/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:214: RuntimeWarning: Glyph not found for set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.6/dist-packages/matplotlib/backends/backend_agg.py:183: RuntimeWarning: Glyph not found for set_text(s, 0, flags=flags)

```



## ▼ Mentions

```

1 mentions = get_mentions(text)
2 print(len(mentions))
3 # mentions[:15]

```

📄 2805

```

1 mentions2 = Counter(mentions).most_common(25)
2 x = [i[0] for i in mentions2]
3 y = [i[1] for i in mentions2]
4 mentions2[:10]

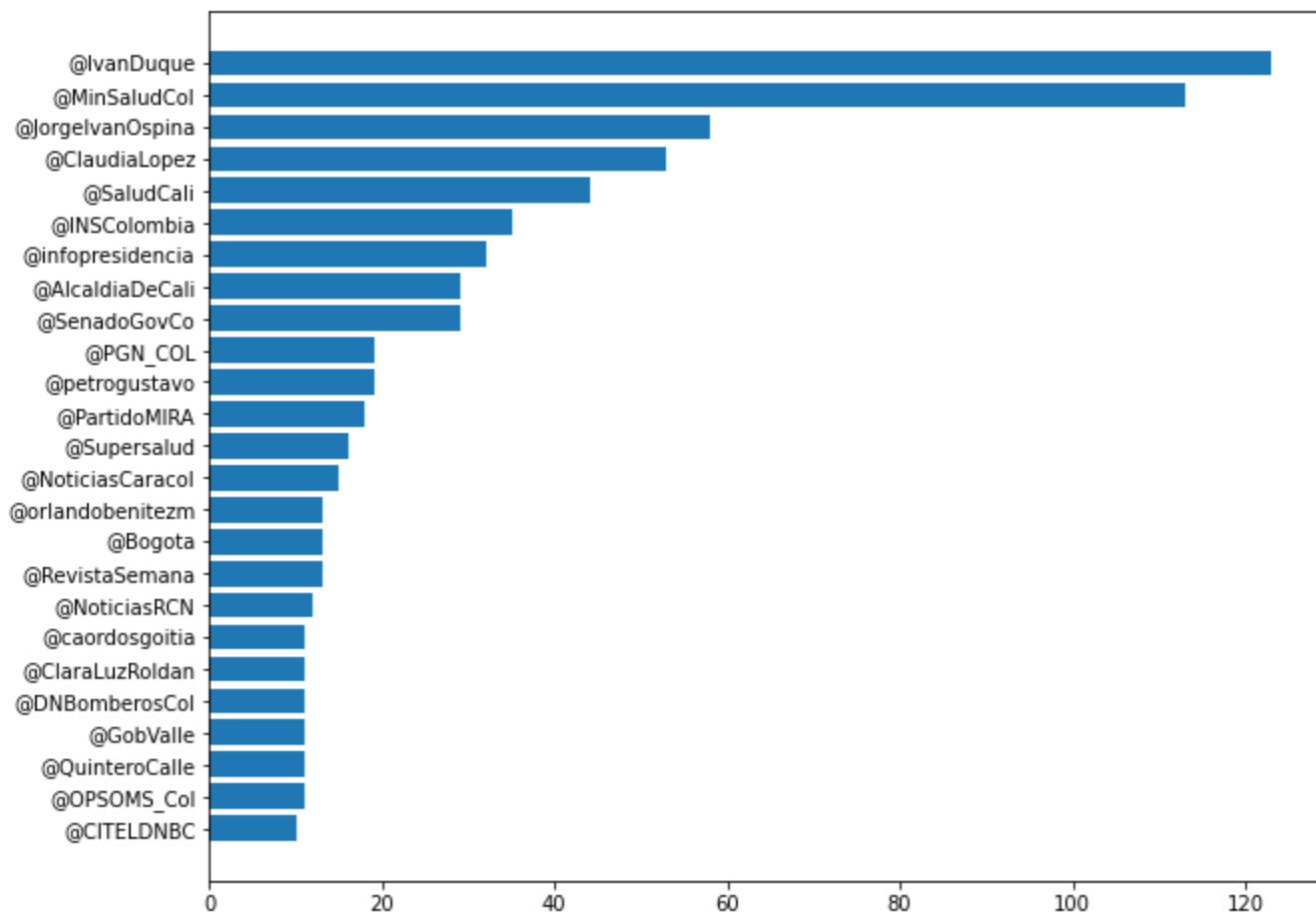
```

📄

```
[('@IvanDuque', 123),
 ('@MinSaludCol', 113),
 ('@JorgeIvanOspina', 58),
 ('@ClaudiaLopez', 53),
 ('@SaludCali', 44),
 ('@INSColombia', 35),
 ...]
```

```
1 plt.figure(figsize=(10, 8))
2 plt.barh(x[::-1], y[::-1])
```

↗ <BarContainer object of 25 artists>



## ▼ Daily use of mentions

```
1 data3.head()
```

↗

	status_id	user_id	created_at	screen_name	text	source	repl
0	1244051801516711938	803282972317204480	2020-03-29 00:00:37	redcomunitariat	Este lunes estaremos hablando sobre la situaci...	Twitter for iPhone	
1	1244052036511051778	2476348920	2020-03-29 00:01:33	SebasCamposCol	Aquí con frío 🥶 viendo cómo pasa la cuarentena 📺...	Twitter for Android	
2	1244052338412847104	239176842	2020-03-29 00:02:45	Jonathan_518	Hoy es #sábado, apenas es hora de bañarme y or...	Twitter for iPhone	
			2020-03-29		En el unico PAIS donde la	Twitter	

```

1 def get_hashtags_(text):
2     patter = re.compile(r"#(\w+)")
3     all_hashtags = [f"#{i}" for i in patter.findall(text)]
4     return all_hashtags

```

```

1 df_count_mentions = data3[['day']].drop_duplicates().copy()

```

```

1 data3['hashtags'] = data3['text'].apply(lambda x: get_hashtags_(x))

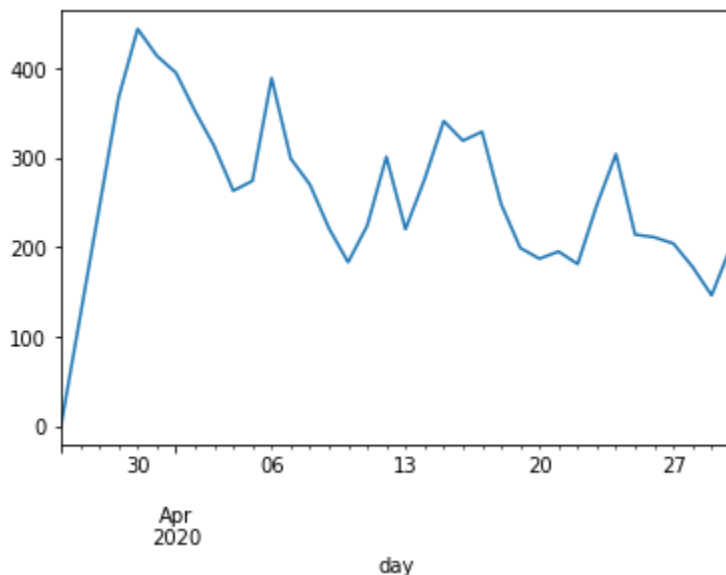
```

```

1 data3.groupby('day')['hashtags'].count().plot()

```

☞ <matplotlib.axes.\_subplots.AxesSubplot at 0x7fb5d2b827f0>



```

1 ten_most_frequent_mentions = x[:15]
2 dates = data3['day'].unique()
3
4 for mention in ten_most_frequent_mentions:
5     daily_use = []
6     for date_ in dates:
7         temp_df = data3[data3['day'] == date_].copy()
8         daily_use.append(sum(temp_df['text'].apply(lambda x: 1 if mention in x else 0)))
9     df_count_mentions[mention] = daily_use
10
11 #     print(daily_use)

```

```

1 df_count_mentions.head()

```

	day	@IvanDuque	@MinSaludCol	@JorgeIvanOspina	@ClaudiaLopez	@SaludCali	@INSColombia
0	2020-03-29	13	17	4	9	1	7
386	2020-03-30	20	15	1	5	0	7
858	2020-03-31	26	16	4	15	3	4
1305	2020-04-01	27	17	3	11	3	6
1726	2020-04-02	16	17	5	8	5	6

```

1 df_count_mentions.index=df_count_mentions['day']

```

```

1 # df_count_mentions.groupby('day').count().plot()

```

```

1 df_count_mentions[df_count_mentions.columns[1:]].sort_index().plot(figsize=(10,6))

```

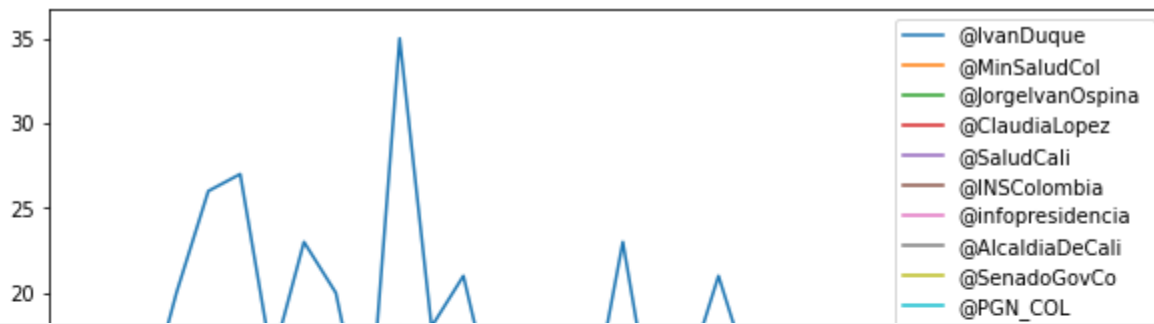
```

↳

```

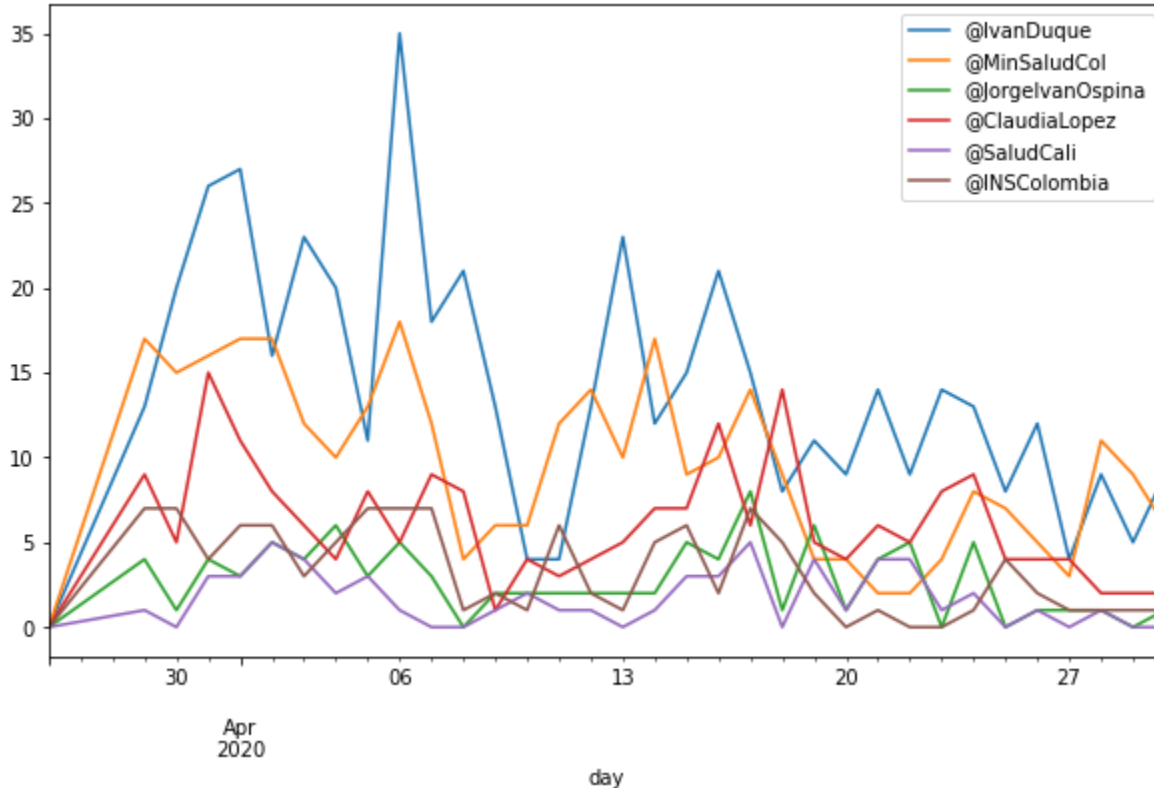


```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5d2ed2908>
```



```
1 df_count_mentions[df_count_mentions.columns[1:7]].sort_index().plot(figsize=(10,6))
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb5d304b048>
```



## ▼ Mentions separated by region

```
1 fig, axes = plt.subplots(nrows=3, ncols=3, figsize=(20, 15))
2 c=0
3 for city in most_frequent_cities:
4     temp_data = data3[data3['place_full_name'] == city].copy()
5     # print(f"{city} has {temp_data.shape} tweets")
6     temp_text = clean_text(" ".join(temp_data['text']))
7     temp_mentions = get_mentions(temp_text)
8     temp_mentions2 = Counter(temp_mentions).most_common()
9     del_covid_words = ["#covid", "#covid19", "#coronavirus", "#cuarentena", "#covid19colombia", "#c"]
10    temp_mentions2 = [i for i in temp_mentions2 if not i[0].lower() in del_covid_words][:10]
11    x = [i[0] for i in temp_mentions2]
12    y = [i[1] for i in temp_mentions2]
13    axes[c//3][c%3].set_title(f"{city} has {temp_data.shape[0]} tweets")
14    axes[c//3][c%3].barh(x[:-1], y[:-1])
```

