

IADB's Share of voice on twitter from Colombia
By Wladimir Llanos
General Assembly DC -DAT5

Introduction

The concept of “listening to the conversation” which is at the core of most social media programs is really just a pseudonym for tracking content your “customers” create based on their satisfaction, or dissatisfaction. It's often described as “share of voice” - The percentage of all online content and conversations about your company, compared to your competitors.

There are many paid tools to calculate the share of voice but the complex it and high price is unnecessary for our organization. We want to know how many mentions do we have on twitter against other organizations and API twitter has some good options to collect data at the moment.

Data Mining

Data Mining is an analytic process designed to explore data (usually large amounts of data - typically business or market related - also known as "big data") in search of consistent patterns and/or systematic relationships between variables, and then to validate the findings by applying the detected patterns to new subsets of data.

Tasks

I've been collecting data from April 20th to May 28th using the twitter's API. Python has a lot of options to connect and work with twitter, Tweepy was the library that I've chose.

Keywords

The following words represent our organization and competitors:

@el_BID", 'el bid', 'Banco Interamericano de Desarrollo', 'Banco Mundial', '@BancoMundial', 'BancoMundial', 'el Banco Mundial', '@AgendaCAF', 'AgendaCAF', 'banco de desarrollo de América Latina', 'CAF-Banco de Desarrollo de América Latina

The function to collect data in realtime that I used to filter twitter streaming and bring you tweets in that matched with your keywords is:

```
twitterStream.filter(track=["keywords"])
```

It returns an xml with all tweet data:

```
"created_at": "Fri Jan 23 23:57:36 +0000 2015",
"id": 558775589612437504,
"id_str": "558775589612437504",
"text": "Thanks, polar vortex: Attendance dips at major Chicago |
"source": "\u003ca href=\"http://twitter.com\" rel=\"nofollow\"
"truncated": false,
"in_reply_to_status_id": null,
"in_reply_to_status_id_str": null,
"in_reply_to_user_id": null,
"in_reply_to_user_id_str": null,
"in_reply_to_screen_name": null,
"user": {
  "id": 7313362,
  "id_str": "7313362",
  "name": "Chicago Tribune",
  "screen_name": "chicagotribune",
  "location": "Chicago, IL",
  "url": "http://www.chicagotribune.com/",
  "description": "Chicago Tribune news, features and so much mor
  "protected": false,
  "verified": true,
  "followers_count": 321548,
  "friends_count": 523,
  "listed_count": 8074,
  "favourites_count": 34,
  "statuses_count": 47367,
  "created_at": "Sat Jul 07 14:10:07 +0000 2007",
  "utc_offset": -21600,
  "time_zone": "Central Time (US & Canada)",
  "geo_enabled": false,
  "lang": "en",
```

The next step was save the data in a csv file with 27 fields separated by "...".

```
tweetText, tweetId, tweetInReplyStatusId, tweetFavoriteCount, tweetRetweeted
, tweetCoordinates, tweetTimestamp_ms, tweetGeo, tweetLang, tweetDate
, tweetUserScreen_name, tweetUserProfile_image_url_https
, tweetUserFollowers_count, tweetUserListed_count, tweetUserFollowing_count
, tweetUserLocation, tweetUserGeo_enabled, tweetUserName, tweetUserLang
, tweetUserUrl, tweetUserCreated_at, tweetUserTime_zone, tweetUserIs_translator
```

Example:

1	1429560528.77::RT @SinEmbargoMX: Las pensiones de 80 millones de latinoamericanos se encuentran en peligro
2	1429560563.56::RT @SinEmbargoMX: Las pensiones de 80 millones de latinoamericanos se encuentran en peligro
3	1429560581.3::RT @SinEmbargoMX: Las pensiones de 80 millones de latinoamericanos se encuentran en peligro
4	1429560636.03::RT @BancoMundial: Reducir a cero la quema rutinaria de gas para 2030. Apoye esta iniciativa. http://t.co/9mz78aF8bv http://t.co/9mz78aF8bv
5	1429560638.12::RT @BancoMundial: Reducir a cero la quema rutinaria de gas para 2030. Apoye esta iniciativa. http://t.co/9mz78aF8bv http://t.co/9mz78aF8bv
6	1429560648.67::RT @BancoMundial: Acceso a una #educación sin calidad es lo mismo a no acceder a ella. Únase a la lucha para poner fin a la pobi
7	1429560649.71::RT @rubenromero685: @BancoMundial La caída de los precios del petróleo es inducida por los EEUU para frenar el avance de los p
8	1429560656.5::RT @rsg55: @BancoMundial i Saboteo !::590246364812091392::None::0::False::None::1429560674275::None::0::Mon Apr 20 .
9	1429560657.23::RT @BancoMundial: La dualidad de la caída de los precios del #petróleo ¿bonanza o #crisis? http://t.co/BjwmRZaY1e http://t.co/BjwmRZaY1e
10	1429560661.2::RT @BancoMundial: La dualidad de la caída de los precios del #petróleo ¿bonanza o #crisis? http://t.co/BjwmRZaY1e http://t.co/BjwmRZaY1e
11	1429560670.08::RT @NaranjaRA: 70 Bs el bidón de agua... todos los días compro al menos 2... Régimen es tan incapaz que ni el Agua puede asegi
12	1429560674.88::RT @SinEmbargoMX: Las pensiones de 80 millones de latinoamericanos se encuentran en peligro
13	1429560676.86::Banco Mundial destaca estrategia económica de México http://t.co/akpV75QDA6 ::590246449927041024::None::0::False::None::
14	1429560686.68::Las pensiones de 80 millones de latinoamericanos se encuentran en peligro

Data cleaning

Many lines in csv file had broken or corrupted data, I need to make a loop to check integrity and the results are below:

Row data	67,156
After cleaning	53,987

New fields:

Country

Every tweet data has informations but there is a problem with location, twitter allows the user to input text without specific format:

Bogot
Bogot
bogot
Bogot - Colombia
Bogot Colombia
Bogot Colombia
Bogot Colombia.
Bogot D.C
Bogot D.C - Colombia

Colombia

The best option to match the county was make a csv file with all Colombia's cities, states and counties. Regular expressions simplify the work. I've got great results

```

153 terms = pd.read_csv('Terms.csv', parse_dates=True)
154 #accent removed
155 termsCleaned = []
156 for term in terms['TERM']:
157     termsCleaned.append(term.encode("ascii", "ignore").strip())
158
159 #removing common places in other countris(the following list is irrelevant for our study)
160 #FLORENCIA, CORDOBA, AMAZONAS, LA PLATA, MADRID
161 termsCleaned.remove('FLORENCIA')
162 termsCleaned.remove('CRDOBA')
163 termsCleaned.remove('AMAZONAS')
164 termsCleaned.remove('LA PLATA')
165 termsCleaned.remove('FLORIDA')
166 termsCleaned.remove('MADRID')
167
168 #mapping twData.tweetUserLocation with terms cleaned
169
170 for location in tweetUserLocation:
171     flag = 0
172     op1 = '''\b'''
173     op2 = '''\b'''
174
175     for place in termsCleaned:
176         if not (re.search(r""+op1+place+op2, location, flags=re.IGNORECASE)) is None:
177             tweetisColombia.append("1")
178             flag = 1
179             break
180     if flag == 0:
181         tweetisColombia.append("0")

```

Brand

To tag the brand I've defined 4 names and proceeded to match every tweet with a brand:

BID	BANCO MUNDIAL	CAF	UNDEFINED
-----	---------------	-----	-----------

Example:

BID

RT @el_BID 5 maneras de usar la #tecnología para mejorar la #movilidad urbana <http://t.co/UXvUw703N6>
 @BID_Ciudades <http://t.co/VZB1bvpAEi>

RT @EFEnoticias: **EI BID** aconseja reformas para salvar la pensión de 80 millones de personas en A.Latina
<http://t.co/NKORj8qlzD>

BANCO MUNDIAL

Banco Mundial destaca estrategia económica de México <http://t.co/zVPds2xkUu::590246511759470592>

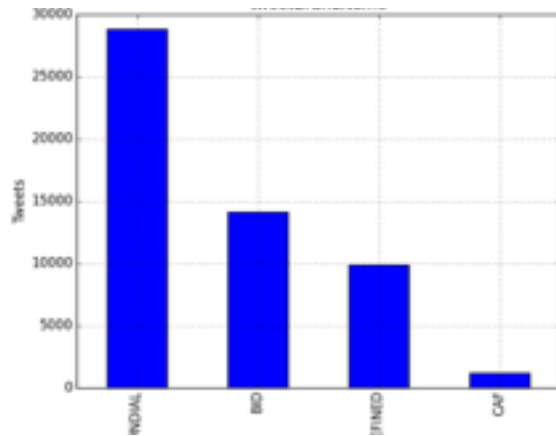
@BancoMundial Las razones en la caída de los precios del crudo son diversas

I've used regular expression again:

```
re.search(r"\b@el_BID\b",row.split(':')[1],flags=re.IGNORECASE)
```

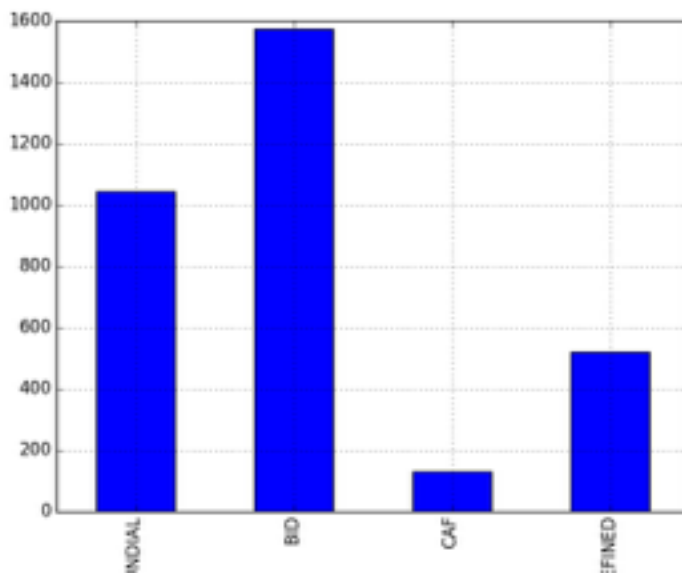
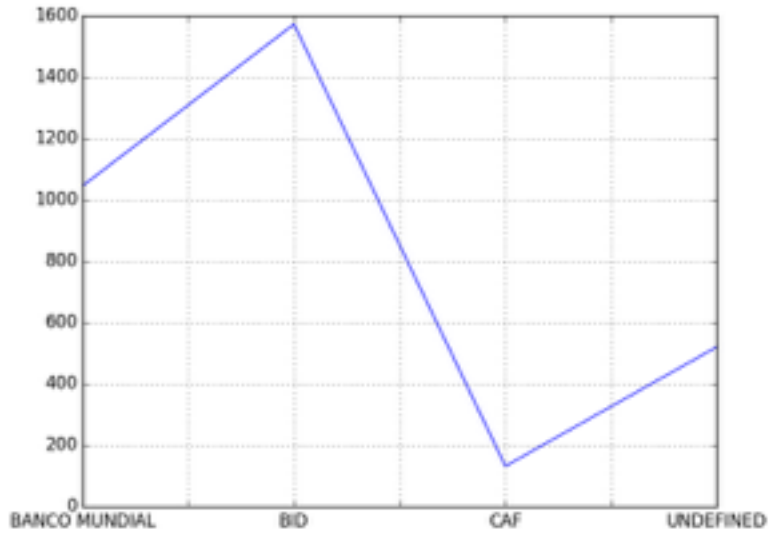
Data exploration

Not filtered by location



		TWEETS	%
1	BANCO MUNDIAL	28797	53
2	BID	14153	26
3	CAF	1190	2
	UNDEFINED	9847	18
	TOTAL:	53987	100

By country



		TWEETS	%
1	BID	1572	48
2	BANCO MUNDIAL	1045	32
3	CAF	131	4
	UNDEFINED	520	16
	TOTAL:	3268	100

Conclusion

Collecting data with python is the best and faster option. Performance and usability is easier than other languages like php, c#, c++.

The challenge is to build the data-frame. Once that is complete you only have to analyze the data, make connections and write the code to look for matches.