

BCSC 4227 DATA SCIENCE

BITC01/0054/2019

KIMUTAI GRIFFINS

Project Scope: Analyzing tweets collected using the query term “World Cup Finals”

Project Report

Text mining, also known as text data mining, is the process of transforming unstructured text into a structured format to identify meaningful patterns and new insights.

Data Collection

The first step of text mining is gathering the data. With this, the python tweepy library came in handy to facilitate the pulling of data from Twitter.

The data pulled from the API is unstructured data- (This data does not have a predefined data format. It can include text from sources, like social media or product reviews, or rich media formats like, video and audio files.)

The data collected has to be cleaned and transformed into usable format. It involves the use of techniques such as language identification, tokenization and syntax parsing.

```
In [2]: # Authentication
API_KEY = "tPQsvRTg0Xj6LdwWEE6jU5Z7A"
API_SECRET_KEY = "eb00BgMIrZwJ4BxTW6bTWbAhm7bTzaRdhZNxE1UDnYcUtofs4"
ACCESS_TOKEN = "770261890257686528-yV49uHpcD3HN3FZBbiiJFqw8IkjxV0b"
ACCESS_SECRET_TOKEN = "YNItLRZclGYcidN5pGDYBtprFQnKnuyjD0oMM963JZzQM"

auth = tweepy.OAuthHandler(API_KEY, API_SECRET_KEY)
auth.set_access_token(ACCESS_TOKEN, ACCESS_SECRET_TOKEN)
api = tweepy.API(auth)
print(api)

<tweepy.api.API object at 0x7f0df1866470>
```

Information retrieval

Returning relevant information based on a pre-defined set of queries.

Tasks performed include :

Tokenization - breaking out long-form text into sentences

```
In [17]: # Tokenization of words
tweetWords = word_tokenize(cleanTweets)
print(tweetWords)
```

```
[['!', '[', ']', ' ', 'FIFAWorldCup', 'quarter-finals', 'start', 'today', '!', 'It', '"', 's', 'Brazil', 'vs', 'Croatia',
'starting', 'at', '8:30', 'PM', 'IST', 'which', 'will', 'be', 'followed', 'by', 'Argentina', 'vs', 'N.', '"', 's', 'B
0.0', '0.1', ' ', 'FIFAWorldCup', 'quarter-finals', 'start', 'today', '!', 'It', '"', 's', 'B
razil', 'vs', 'Croatia', 'starting', 'at', '8:30', 'PM', 'IST', 'which', 'will', 'be', 'followed', 'by', 'Argentina
', 'vs', 'N.', '"', 's', '500', 'Today', 'the', 'FIFA', 'World', 'Cup', 'quarter-finals', 'begin',
', '!', '🇧🇷', '\\n\\nWho', 's', 'going', 'to', 'the', 'semi-finals', '?', '🇺🇦', '\\n\\nCroatia', 'vs', '🇸🇰',
'Brazil\\n\\nNetherlands...', '"', '0.0', '0.0', ']', '500', 'Today', 'the', 'FIFA', 'World',
Cup', 'quarter-finals', 'begin', '!', '🇧🇷', '\\n\\nWho', 's', 'going', 'to', 'the', 'semi-finals', '?', '🇺🇦',
'\\n\\nCroatia', 'vs', '🇸🇰', 'Brazil\\n\\nNetherlands...', '"', '[', 'The', 'last', 'time', 'these',
'two', 'met', 'at', 'the', 'quarter-finals', 'of', 'the', 'World', 'Cup', 'was', 'in', '1998', 'Dennis', 'Bergk
amp', 'scored', 'a', 'clutch', 'goal', 'on', 'th...', '0.0', '0.06666666666666667', ']',
'The', 'last', 'time', 'these', 'two', 'met', 'at', 'the', 'quarter-finals', 'of', 'the', 'World', 'Cup', 'was', 'in
1998', 'Dennis', 'Bergkamp', 'scored', 'a', 'clutch', 'goal', 'on', 'th...', '[', '\\n\\n', '🇸🇰',
'2022', 'World', 'Cup', 'quarter-finals', 'games', 'schedule', 'for', 'today', ':', '\\n\\n\\n', 'Croatia', '🇸🇰', 'vs
🇸🇰', 'Brazil', '(', '4pm', 'CET', ')', '\\n\\n\\n', 'Netherlands', '🇸🇰', 'vs', '🇺🇦', 'Ar...', '"', '0.0',
'0.0', ']', '\\n\\n', '"', '2022', 'World', 'Cup', 'quarter-finals', 'games', 'schedule', 'for', 'today',
':', '\\n\\n\\n', 'Croatia', '🇸🇰', 'vs', '🇸🇰', 'Brazil', '(', '4pm', 'CET', ')', '\\n\\n\\n', 'Netherlands', '🇸🇰', 'vs',
'🇺🇦', 'Ar...', '[', 'On', 'the', 'eve', 'of', 'the', 'WorldCup2022', 'Quarter', 'Finals', 'the', 'w
e', '"', 're', 'giving', 'away', 'a', '2022', 'Football', 'World', 'Cup', 'Shirt', 'of', 'your', 'choice', '!', '🇺🇦', '\\n\\n\\n']]
```

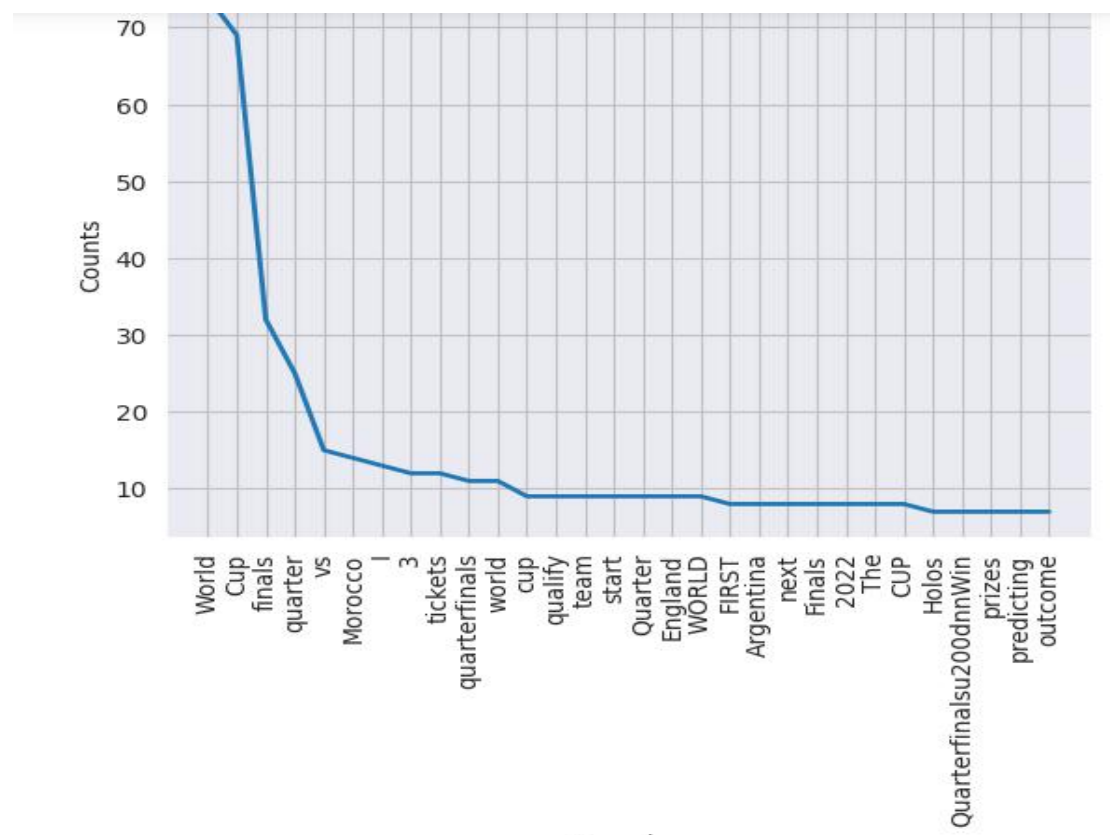
Tokenization of sentences

Sentiment analysis: This task detects positive or negative sentiment from data sources, allowing you to track changes in customer attitudes over time.

Summarization : provides a synopsis of long pieces of text to create a concise, coherent summary of a document's main points.

Information extraction

Word frequency -counts the *frequency* of each and every *word* in a text, helping you understand the keyword density for SEO or the rate of *word* repetition



or locations. For example, NER identifies “California” as a location and “Mary” as a woman’s name.

Importance of the performed analysis

Sentiment analysis-

provides information about perceptions of brands, products, and services.

These insights can propel businesses to connect with customers and improve processes and user experiences.

Providing a mechanism for companies to prioritize key pain points for their customers, allowing businesses to respond to urgent issues in real-time and increase customer satisfaction.

Risk management

Give insights into industry trends and financial markets by tracking shifts in sentiment and pulling information from analyst reports and whitepapers This is especially beneficial to financial organizations since it gives them more confidence when contemplating company investments in diverse areas.

Text mining tools analyze documents to identify entities and extract relationships between them, unlocking hidden information to help researchers:

Challenges encountered during the process

Short Informal texts - one of the difficulties in sentiment analysis is short informal text. They are restricted in length, usually spanning one or less sentences. They frequently use slang words, misspellings, and truncated word forms.

Many of the collected tweets are stripped of context or too short to serve as real carriers of meaning. Most of the tweets have the meaning hidden in images or links to websites.

Twitter API has a request limit that limits the amount of data to collect and analyse. (TooManyRequests: 429 Too Many Requests 88 - Rate limit exceeded)