

Analyzing Racist Tweets on Soccer Players

Report 1

Aravind Padmanabhan
Computer Science
Binghamton University
Binghamton, New York
apadman1@binghamton.edu

Chethana Gopinath
Computer Science
Binghamton University
Binghamton, New York
cgopinal@binghamton.edu

Ridhuvarshan Natarajan
Computer Science
Binghamton University
Binghamton, New York
rnatar2@binghamton.edu

INTRODUCTION

Our project, as mentioned in the proposal is focused on analyzing racist tweets made on soccer players. We used the “Tweepy” library for python and streamlined the live tweets with the help of Twitter APIs.

CHALLENGES

Our project, as mentioned in the proposal is focused on analyzing racist tweets made on soccer players. Our initial methodology of data collection was to collect historical tweets using the `tweepy.cursor()` method. But the dataset collected as a result of this, failed to produce enough racist content to perform any reasonable analysis. Many of the collected tweets didn’t hold much relevance to our objective of collecting racist data. This inference was mainly drawn due to Twitter’s moderation policy of removing abusive tweets. Due to the irrelevancy and sparseness of the collected historical data, the consensus was to start collecting live tweets from Twitter as and when a match happens. As the data was being collected, the trend observed was that many tweets just had headlines of articles without actual proof of racist verbal attacks happening on any particular player from any club. Without actual proof of racism, there wouldn’t be a metric which decided if the information in the tweet was veritable or not. Finally, we ended up collecting all the tweets based on hashtags of clubs and tournaments who we wanted to follow. This was done during the period of the match from Twitter. This would give a wider dataset for performing sentiment analysis as next steps. This dataset could be refined and cleaned further to derive the results that are required.

METHODOLOGY

In order to collect tweets pertaining to the requirement, a certain number of hashtags were isolated and based on them, tweets were collected. Using the real-time tweet streamlining API, we were able to gather around 400,000 tweets that matched the hashtags which we used as filters for the data collection. The gathering of data was done on match days - before, during and after the match. Our choice of hashtags for the tweets collection was based on tournaments, clubs and keywords that we wanted to closely follow. We chose

three tournaments namely English Premier League, UEFA Champions League and UEFA Europa league and monitored tweets that were posted during the matches in which these clubs were participating. The hashtags used for the soccer matches were of the format #SHUARS (for a match between Sheffield United and Arsenal), #UEL (for Europa league games) etc. We also collected tweets that contained the hashtag #NoRoomForRacism which is an official hashtag released by the English Premier League to call out racist behavior. Though tweets containing this hashtag may not be racist, they may be helpful for deriving some valuable observations such as analyzing racist incidents by fans. This in turn, can provide further information to collect even more tweets relating to those particular incidents. In addition to this, we used our Professor’s suggestion of expanding our scope to include NBA games as well. As the new season started last week, collecting these tweets would enlarge our dataset adding new data which can be analyzed along with our existing data. For the NBA games, we used hashtags such as #NBA, #LAKERS etc.

IMPLEMENTATION

At first, we ran our code on one of our local machines and collected data there. This was done in python using the Jupyter Notebook. We used the “Tweepy” Library for streamlining the live tweets and stored the collected data in a json file. This was later moved to the mongo database that was installed in the virtual machine.

A) First, we copied the consumer keys and access keys to a json file so that it is easy to refer the keys and ping the API. The code snippet below is used to ping Twitter API every time we run the program.

The code snippet below simply receives the incoming stream of tweets and saves them in a json file.

A

B

[illegible]

<http://docs.tweepy.org/en/latest/>