# Final Project
# (ID2221)

Chuting Zhu           Suhas Sheshadri
chuting@kth.se         suhass@kth.se

November 4, 2018

## 1    Topic

We implement sentiment analysis on "Sweden's football journey during the FIFA World Cup-2018" based on the tweets collected from Twitter. Our project shows the sentiments of people during the whole world cup vs when Sweden lost matches. And also, the Term Frequency (tf) and Document Frequency (df) for all the words form the tweets collected.

## 2    Dataset

We collect data from Twitter for a given search query and date range. Search queries would be relevant to Sweden and Football and we use "Sweden football" as the query to get related tweets. The date range would be the duration of the world cup tournament, which was from 14 June to 15 July 2018.

We also collected data separately for the days of 23 June and 7 July when Sweden lost the matches on these two days and combine them into a single file.

There are two files, one is $AllTweets.txt$ corresponding to the tweets relevant to "Sweden football" during the whole world cup, the other is $LossTweets.txt$ which are tweets related when Sweden lost.

## 3    Methods

- Collect relevant data from Twitter.

  Since TwitteR API allows us to collect the tweets only from last 7 days, we had to use a different scrapping method. We have used the GetOldTweets implemented by Jefferson-Henrique for scrapping old tweets. We have modified it to suit our need. [2]

  The above obtained data are in .CSV format. We clean the data and convert the files to .TXT format.

- Calculate the sentiment score of the collected tweets. The sentiment score of a tweet is equivalent to the sum of the sentiment scores for each term in the tweet. And we check all the tweets for words which are present in $AFINN - 111$ file. [1]

  AFINN is a list of English words rated for valence with an integer between minus five (negative) and plus five (positive). The words have been manually labeled by Finn rup Nielsen in 2009-2011.

  If the word is present in this file, we will add the word's score obtained from the file to the sentiment score of the corresponding tweet. Moreover, a positive sentiment has a positive score and vice versa.

- Utilize MapReduce Framework to calculate the sentiment score for each tweet.

- Utilize MapReduce Framework to calculate the term frequency and document frequency of all the words from the collected tweets.

# 4 Run the application

1. Analyze the sentiment of all tweets during the whole world cup.

```
python tweets_sentiment.py AFINN-111.txt AllTweets.txt
```

2. Analyze the TF-DF for AllTweets.txt file

```
python tfdf.py AllTweets.txt
```

3. Analyze the sentiment when Sweden lost.

```
python tweets_sentiment.py AFINN-111.txt LossTweets.txt
```

4. Analyze the TF-DF for LossTweets.txt file

```
python tfdf.py LossTweets.txt
```

# 5 Results

The Results are in the format:

Sentiment Analysis: (Tweet number, Sentiment Score)

For reference:
[1, 2.0]  Sentiment score of tweet 1 is 2
[5, 0]  Sentiment score of tweet 5 is 0

Tfdf: (Word, Number of occurrences, [(Tweet number, Frequency in this Tweet)])

For reference:
["hello", 4, [[1, 1], [2, 3], [6, 1], [9, 2]]]  This shows that term hello has df 4 and occurs in tweets 1, 2, 6, 9 with tf 1, 3, 1, 2 respectively.
["world", 1, [[3, 1]]]  This shows that term world has df 1 and occurs in tweet 3 with tf 1.



Figure 1: Sentiment Analysis for AllTweets.txt

Figure 2: TF-DF for AllTweets.txt
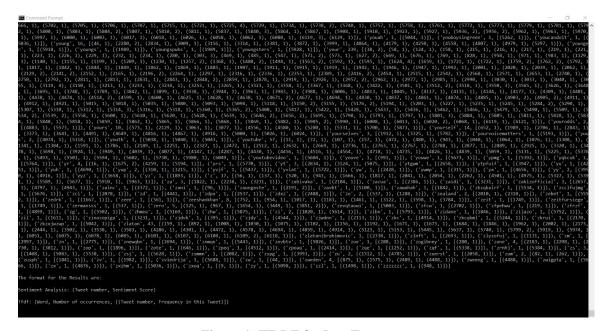


Figure 3: Sentiment Analysis for LossTweets.txt

Figure 4: TF-DF for LossTweets.txt

# References

[1] Afinn. http://www2.imm.dtu.dk/pubdb/views/publication$_details.php?id = 6010$.

[2] Github. https://github.com/Jefferson-Henrique/GetOldTweets-python.