

# Twitter Hashtag Prediction

---

**Mukund Manikarnike**  
**ASU ID: 1208597425**

---

**Social Media Mining**  
CSE 472/598 – Phase – 2  
Fall 2015

---

# Contents

---

Contents .....	ii
Tables and Figures .....	0
<b>Abstract</b> .....	1
<b>Introduction</b> .....	1
<b>Requirements Specification</b> .....	2
<b>Design Choices</b> .....	2
<b>Implementation</b> .....	3
<b>Results</b> .....	1
<b>Conclusion</b> .....	3
<b>Bibliography</b> .....	4

# Tables and Figures

---

## Figures

---

Figure 1 Similarity Matching - NLP .....3  
Figure 2 Pre-processed Data .....1  
Figure 3 Similarity Matching.....1  
Figure 4 Predicted Outputs.....2

## Abstract

---

As of today, there is a lot of data in the form of blogs, posts, and microblogs from users from varied communities, locations and coming from different backgrounds. Given the varied background of people that use all of these platforms and create such varied data sources, there is often a need to retrieve this data. One such mechanism that Twitter (a microblogging site) uses is a concept called Hashtags which are an ever growing set of indices for retrieval. The way it varies from the typical indices is that these indices are ever-growing, user created and similar such hashtags could be named differently although in reality, they represent the same thing. The problem this project addresses is of achieving prediction of hashtags for a tweet that doesn't have any hashtags by learning from a tweet dataset that contains hashtags. The basic principle that this project explores is learning based on the similarity of words in the tweets.

### Keywords

#Twitter #HashtagPrediction

## Introduction

---

### Goal Description

The goals that this project sets out to achieve are as follows

1. Creation of an algorithm that is able to learn about similarities between tweets in the dataset that is provided.
2. Creation of an algorithm that can predict new hashtags for the other tweets in the dataset.

### Assumptions

The project makes a few assumptions which are listed below.

1. The project assumes that the dataset is a static dataset that doesn't grow with time as stated in the abstract.
2. The project doesn't use information about the user, location from where the tweet was sent out to perform the classification. It is purely based on the words in the tweets and their similarities.

### Dataset Description

The dataset used for this project was provided as part of the course. It consists of the following information

1. User Graph
  - a. This is a table of all the user IDs and the user IDs that they follow.
2. Tweet Data
  - a. This is a table of all the user IDs, the 140 character tweet that they've sent out along with latitude and longitude from which the tweet was sent out.

### Dataset Statistics

The dataset contains

1. About 60,000 unique users
2. 677,732 Tweets

This project used 1000 tweets as the training set and the remaining 1000 tweets as the test set. The similarity between users was also calculated only for the users who were involved in these 1000 tweets. The reason for this is stated in further sections.

### Overview

The solution to the problem is described in the following few sections as described below

1. The Requirements Specification section describes detailed requirements that the project aims to meet.

2. The Design Choices section describes details of the algorithm chosen and the reasons for the same.
3. The Implementation section describes certain choices made during implementation and how each of the algorithms were implemented, issues that came up during the same.
4. The Results section describes the tests that were carried out using the implementation and the outputs that were obtained.
3. The Conclusion section describes how the approaches compare against existing approaches and future enhancements that can be made.

## Requirements Specification

---

This section introduces each requirement that the project aims to achieve.

- [REQ\_1] It is required to use the dataset provided as a basis for any decisions that are made.
- [REQ\_2] All decisions that are made using the data will have to be carried out using social media mining techniques like similarity calculation, classification and so on that were discussed as part of the course.
- [REQ\_3] In order to carry out any of the algorithms designed for this project, none of the machine learning libraries shall be used. Libraries to carry out mathematical or linear algebraic tasks are allowed.
- [REQ\_4] The algorithms chosen shall use the dataset provided to make any decisions required and shall not make assumptions outside of the dataset about the general trend in twitter datasets.
- [REQ\_5] All ideas behind the design choices made shall be clearly stated along with results obtained for those design choices.

## Design Choices

---

- [DSN\_CH\_1] This design choices made for this project were mainly inspired by the paper referenced in [1]
- [DSN\_CH\_2] The ideas that were referenced by the paper were
- a. To treat tweets as a set of function which varies according to the words that are present in the tweet.
  - b. Usage of the natural language processing library NLTK to compute the similarity of words in order to find how the tweets rank against each other.
  - c. Prediction of tags using the dominant tag approach.
- [DSN\_CH\_3] The design provides two kinds of similarity matching techniques
- a. Similarity matching using word similarities in tweets
  - b. Similarity matching using user based similarities.
- [DSN\_CH\_4] In order to compute user based similarities, cosine similarity was used.
- [DSN\_CH\_5] The procedure chosen to perform matching using similarities in words is detailed in the next sub-section

### Similarity Matching using NLP

- [DSN\_CH\_6] The approach used here is to find similarities between nouns in each of the tweets.
- [DSN\_CH\_7] Natural Language processing enables us to find similarities between any 2 words. By using the word similarity measure, the similarities between tweets is computed as a function of frequencies and semantic similarity as shown below.

$$I_{p,q} = S_{i,j} * (f_p - f_q)$$

$I_{p,q} \rightarrow$  importance of word  $p$  over  $q$

$S_{i,j} \rightarrow$  Semantic similarity obtained using NLP

$f_p, f_q \rightarrow$  frequency of words  $p$  and  $q$

[DSN\_CH\_8] By using this metric, the importance of the noun is obtained. A summation of each such importance would give the rank of the noun with respect to the other nouns.

[DSN\_CH\_9] By using this metric, one would get the importance of each noun in the tweet. A summation of all the importance of each noun would give a measure of the importance of each tweet. The logic is that frequency and similarity put together would give an idea of what the tweets are talking about.

[DSN\_CH\_10] If the same procedure is carried out on a tweet for which hashtags need to be predicted, we would get an importance value for the tweet of interest as well.

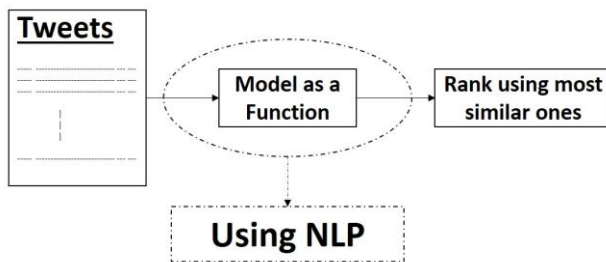


Figure 1 Similarity Matching - NLP

## Prediction Using Dominant Tags

[DSN\_CH\_11] After the ranking of tweets is performed, the prediction is done using the 10 most similar tweets.

[DSN\_CH\_12] Out of the 10 most similar tweets, the average number of tags for all tweets is obtained. The final number of suggested tags would be as many as this number.

[DSN\_CH\_13] In order to actually predict the tags, the average number of dominant tags in the entire set of tags from the top 10 tweets in their vicinity. A summary of the design is as shown in **Figure 1**

# Implementation

## Overview

[IMP\_CH\_1] The implementation was carried out using Python 2.7. The libraries that were used were

- a. NLTK
  - i. To extract nouns from tweets and similarities between them.

## Data Pre-processing

[IMP\_CH\_2] The given data contained tweet texts which contained “RT”, “@mentions”, special characters, numbers, URLs, #Tags and all of them were in different cases.

[IMP\_CH\_3] In order for the data to be easily processed by the NLP library, all of these special words had to be removed as part of pre-processing.

## Finding Similarities

[IMP\_CH\_4] Similarities were found using the NLTK library by extracting nouns and computing semantic similarities and frequencies as highlighted in the design section.

[IMP\_CH\_5] Cosine Similarity between the users who were involved in the tweet dataset chosen was calculated using the graph of users provided as an adjacency list.

## Prediction using Dominant tags

[IMP\_CH\_6] This was done by finding the most dominant tags according to tweet similarity as well as user-cosine similarity as highlighted in the design section.

## Results

This section will briefly talk about the output file formats generated by this project and the results that were seen for each step.

There were several intermediate files that were generated and have been submitted as part of the package which will be talked about in the readme and not here.

### Data Pre-processing

```

0      111648679    all teargas in one place      ['teargas', 'place']  ['lstaïd4', 'egypt', 'syria', 'libya', 'yemen', 'bahrain',
1      53229950    officials cia report blamed militants not mob within hours of attack  ['report', 'attack']  ['benghazi']
2      111648679    arabic ['arabic']  ['lstaïd4', 'egypt', 'bahrain', 'syria', 'libya', '25jan', 'yemen', 'firstaid', 'tunisia']
3      456193667    bahraini rioters attacks the citizens on the road to kill them and burns cars http  ['bahraini', 'road']  ['h
4      354708537    alqaeda seeks total control over with mountain bases to harass amp libyan coast arabafrican neighbors  ['alqaeda',
5      821403404    which one of the bedwetters at nbc finally came out with that obvious statement  ['statement']  ['obamafail
6      272187450    video if obama knew it was an act of terror how does he justify what followed  ['video', 'act', 'terror']  ['be
7      268263316    we r in under seige since yesterday after a police vehicle crashed down inside the village  ['seige', 'yesterda
8      268263316    moi claimed that a police man was killed bcuz of what happened to the vehicle  ['moi', 'police', 'man', 'bcuz', 've
9      268263316    moi claims r just lies to justify the raids on houses and attaks on peaceful ppl in  ['moi', 'r', 'ppl']  ['el
10     198075393    ambassador stevens repeatedly conveyed worries to obama admin about lawlessness amp violence in  ['ambassador

```

Figure 2 Pre-processed Data

As indicated in **Figure 2**, the pre-processed data contains the Tweet ID, User ID, Tweets stripped of the data that was mentioned in the design section, nouns in the tweet, hashtags in the tweet. The tweet ID is an ID that was assigned by the pre-processing step for easy identification.

As shown in **Figure 3**, the similarities between each user ID involved in the Tweets that are being processed is provided in cosine Similarity. As part of the Similarity using NLP, the importance of each of the words as rated according to the other words in the whole dataset is shown.

### Cosine Similarity

```

200134659 222089225 0.0776150525706
200134659 25448459 0.0192910588216
200134659 525821964 0.0835236861855
200134659 400067888 0.0101477094258
200134659 252569616 0.0
200134659 247963674 0.0450902709032
200134659 339034131 0.0031196123536
200134659 156501196 0.0020319747644
200134659 15962135 0.00323677144175
200134659 27277336 0.068569189614
200134659 36540441 0.0164646389985
200134659 834852890 0.0
200134659 733620253 0.180159417934
200134659 727161668 0.0600603308735
200134659 351205408 0.19508658889
200134659 17274913 0.0412725861943
200134659 83453988 0.0345192628084
200134659 90918950 0.124460972047

```

### Similarity Using NLP

```

gadhafi 0
protest -424.576866842
jihad -411.703272789
sleep -663.157189706
mirage -370.543936129
battleground -227.86231185
integrity -224.592018276
voter -23.0200786529
tweet 139.216186836
tenacity -148.500555615
lord -475.918757631
sorry 0
risk -590.883075258
rise -723.893384393
school 150.613766789
waldo 0
solution 232.538080955
cholesterol -304.523297042
triumph -607.978934627
force 176.877511378
street -190.106166464

```

Figure 3 Similarity Matching

## Prediction using NLP

```
D:\Courses\3-FL-15\CSE472-598\Projects\20151017-Phase-2\Project\Code\TweetS
Enter TweetId between 1001 and 1999 for which you want to obtain tags:1005
Noun Importances Found
Suggested Tags are
['lebanon\n']
User used Tags are
['benghazi', 'optimal', 'obama?\n']
```

```
D:\Courses\3-FL-15\CSE472-598\Projects\20151017-Phase-2\Project\Code\TweetSi
Enter TweetId between 1001 and 1999 for which you want to obtain tags:1376
Noun Importances Found
Suggested Tags are
['tcot', 'sgp']
User used Tags are
['ekerseige', 'ekerseigelies', 'bahrain', 'bbc', 'cnn', 'usa', 'bahrain\n']
```

```
D:\Courses\3-FL-15\CSE472-598\Projects\20151017-Phase-2\Project\Code\TweetS
Enter TweetId between 1001 and 1999 for which you want to obtain tags:1792
Noun Importances Found
Suggested Tags are
['benghazi', 'israel']
User used Tags are
['tripoli\n']
```

```
D:\Courses\3-FL-15\CSE472-598\Projects\20151017-Phase-2\Project\Code\TweetS
Enter TweetId between 1001 and 1999 for which you want to obtain tags:1222
Noun Importances Found
Suggested Tags are
['bahrain', 'israel', 'iran']
User used Tags are
['istaid4', 'egypt', 'libya', 'yemen', 'bahrain', 'firstaid', '\n']
```

Figure 4 Predicted Outputs

As shown in **Figure 4**, the predicted tags and their actual counterparts are indicated for 4 different input tweet IDs. A detailed analysis of the accuracy of the results aren't available. But, a theoretical analysis is provided as part of the conclusion.



# Conclusion

---

This section talks about the analysis of the model that was used, the things that were learned as part of this project and further enhancements that could be made.

## Theoretical Analysis

### Prediction using Cosine Similarity

This method considers the neighborhood of the users and how similar they are. Using such a measure for prediction of tweets wouldn't be a very good measure because of the following reason reasons

Similar neighborhoods don't indicate similar users because different kinds of people could have similar neighborhoods around them.

However if they were similar, given that the dataset is centered around a particular topic, i.e. the tweets are all talking about something very specific, it might be very wise to suggest that similar users following similar people might want to tweet about a particular thing.

### Prediction using NLP

Given that we've just established that the data in the tweets is centered on a topic, this is a very good mechanism because it factors in the frequencies of words and the semantic similarities between them to establish similarity between tweets.

However, if the dataset were to have material covering different topics, it would be better to combine the User similarity and NLP technique to predict suitable hashtags.

## The Learning Curve

The learning curve for this project was very steep since it involved

1. Understanding usage of Natural Language Processing libraries
2. Reading up on existing techniques to perform Hashtag prediction
3. Carrying out complex algorithms on large datasets and ensuring that they run fast.

## Future Enhancements

As on 5<sup>th</sup> of December, 2015, the following things weren't implemented which were originally part of the plan.

1. Prediction Using Cosine Similarity
  - a. Cosine similarity was calculated. The last part where matching using the similarities has to be done wasn't carried out.
2. Testing the existing implementation on the entire dataset.

As on 5<sup>th</sup> of December, 2015, the following things could be done to make the project better

1. Changing the Cosine Similarity to a better measure giving a better user similarity metric
2. Correlating the 2 similarities to factor both in.
3. Enhancing the speed of the algorithm by using a distributed setup like Hadoop.

**Notes to the Grader**

1. Since no participation from both of the other group members (Harshil Shah and Diego Henrique Ferreira), their names have been omitted from this submission.
2. Due to this, several things that were planned weren't completed. However, 90% of the project has been implemented and tested on a small dataset.

**Bibliography**

---

[1]	Tianxi Li, Yu Wu, Department of Computer Science Stanford University and Yu Zhang, Department of Computer Science Trinity University , <i>Twitter Hash Tag Prediction Algorithm</i>
[2]	Python NLTK Documentation, <a href="http://www.nltk.org/index.html">http://www.nltk.org/index.html</a>