

Michael Lan  
Lingfei Zeng  
Xunjie Zhu

### Stage3 - The Implementation Stage.

This project was implemented using the Python programming language and the Linux operating system.

- Novelty

Users can already find tweets related to their topic by searching for keywords or hashtags related to their topic. However, this project will allow users to find tweets that are both semantically and structurally similar to their tweets, not just related to the topic of those tweets. Users will also be recommended hashtags, such as those of ‘intermediate popularity’, which users do not include in their tweets, but which are topically related to their tweets.

- Usefulness

This project will allow users to find which tweets and hashtags are similar to their own tweets. Thus, this project aims to bridge different Twitter communities by finding ‘hidden’ connections between them based on different measures of similarity. This will facilitate communication between previously isolated communities.

- Sample Small Data Snippet (This is a subset 3 tweets out of thousands of tweets used)

211 : RT @KCONTV: Fans dancing in front of the @HondaCenter before the #BTSinAnaheim concert! #TWTinAnaheim #BTS <https://t.co/QAVet0uqYn>

212 : RT @BTS\_V\_sana: 02line(01)ARMY all #ARMY #BTSRT #RT <https://t.co/170331>

213 : RT @bangtanitl: #BTS was mentioned on Kangnam Style, 170331 <https://t.co/uL8yrSGBJB>  
<https://t.co/qJ0IINaTrH>

- Sample Small Data Output

The word and paragraph vectors, along with the document-term frequency matrix, are huge, so we did not include them in Phase 3 write-up. Additionally, they are only intermediate outputs. The dendrogram of the clustering results is the final output, and a visual and description of it is given below the pipeline description in this section.

- Pipeline Description

Since it would be quite lengthy to include all the code used, only the hierarchical clustering algorithm code has been given at the end of this section. Also, the bipartite graph and Simrank code was completed but are not part of the main pipeline yet, so only a histogram of their output is given; the code is not included.

The rest of the code can be given upon request, and may be included in the final report. The pipeline is briefly described:

The data crawler uses a listener to find new tweets that contain certain keywords or hashtags, collectively called a ‘filter’. To find tweets that are related to certain topics but do not necessarily contain the hashtags of those topics, we first gathered tweets containing initial keywords or hashtags (say 5). New keywords, hashtags, and topics were found from in this first iteration of tweets, and were added to the filter. Then, we ran the data crawler again using this new filter, and repeated this multiple times. Essentially, we used previous iterations of keywords or hashtags to find new topics, and we only ran clustering on the tweets found in the final iteration. Each iteration is called a ‘generation’.

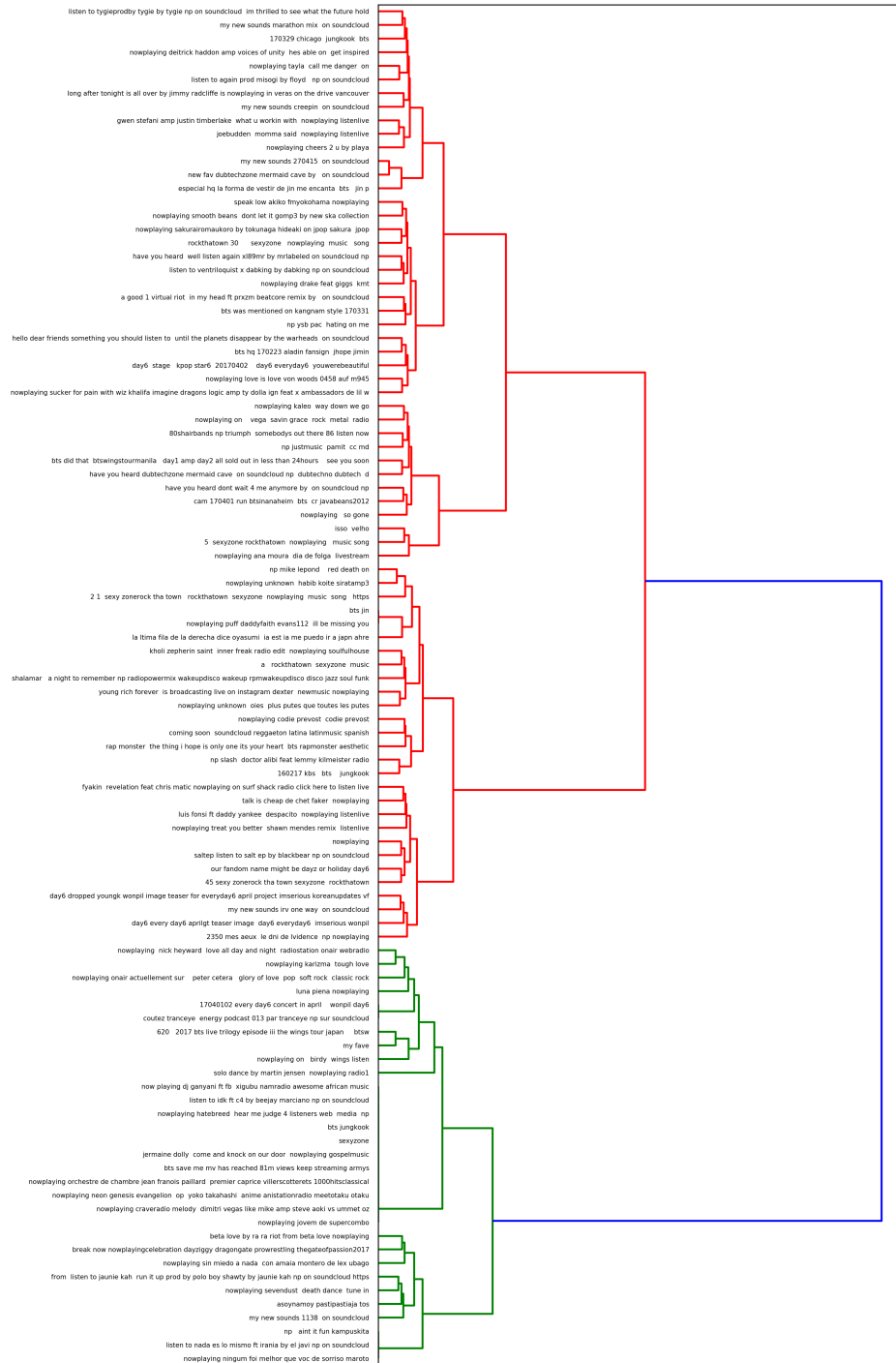
There were two separate methods we used to find new topics: 1) We took the 10 most frequent new, non-trivial keywords hashtags found in this first iteration of tweets and added them to the filter 2) We used LDA to find new topics. IE) Filter in Gen 1: music, concerts Filter in Gen 2: music, concerts, classical, metal Filter in Gen 3: music, concerts, classical, metal, bach

Later, the crawler and topic finding scripts will be deployed to a linux server to be run periodically. This will be implemented by crontab.

The pipeline is described as follows:

- `Twitter_streaming.py`: Starts up a listener of incoming tweets. Uses the Tweepy library to connect to Twitter Streaming API. Outputs a JSON text file of data
- `Tweet_json_to_plain.py`: Reads a JSON text file and outputs a file of tweets
- `Toptopics.py`: Use term frequency and latent dirichlet allocation respectively to find new topics from existing tweets
- `IOUtil.py`: Reads file of tweets and uses regular expressions to do some preprocessing
- `Word2vec.py`: Trains word vector and tf-idf weighted sum of word vector to form document vector. Computes similarity between tweets. Outputs word and paragraph vectors
- `HierarchicalCluster.py`: Clusters the tweets into topics and outputs a dendrogram
- Interesting findings in Dendrogram

Since there were thousands of tweets that were gathered, we created a dendrogram for 100 tweets. There appear to be two main groups, and within the red group, there were 4 main subgroups. This sample shows that our algorithms work. However, since we did not input in meaningful data yet, and since this dendrogram only shows 100 of thousands of tweets, we do not expect this sample to give meaningful results. Later, we will gather better data and perform an analysis.



- Working Code
- 

```
class HierarchicalClusterer(object):
    def __init__(
        self,
        distance_matrix,
        linkage_function,
        minium_clusters=2):
        self.distance_matrix = distance_matrix
        self.linkage_function = linkage_function
        self.minium_clusters = minium_clusters

    def build_clusters(self):
        print 'a'

    def do_link_clustering(self, cluster_id_list, cluster_nodes):

        def small_big(first, second):
            if first < second else (second, first)

        def merge(cluster1,
                   cluster2,
                   distance,
                   cluster_ids,
                   cluster_nodes):
            cluster1, cluster2 = small_big(cluster1, cluster2)
            cluster_ids[cluster1].extend(cluster_ids[cluster2])
            cluster_ids[cluster2] = []
            node = Node()

            if cluster_nodes[cluster1]:
                node.left = cluster_nodes[cluster1]
                cluster_nodes[cluster1].parent = node
            else:
                node.left_instance = cluster1

            if cluster_nodes[cluster2]:
                node.right = cluster_nodes[cluster2]
                cluster_nodes[cluster2].parent = node
            else:
                node.right_instance = cluster2

            node.set_height(distance, distance)
            cluster_nodes[cluster1] = node
```

```

length = self.distance_matrix.shape[0]
cluster_number = self.distance_matrix.shape[0]
distance_list = [
    ClusterPair(first,
                second,
                self.distance_matrix[first][second], 1, 1)
    for first in range(length) for second in range(first +

distance_queue = MyHeap(
    initial=distance_list,
    key=lambda x: x.distance)

while cluster_number - self.minium_clusters > 0:
    pair = distance_queue.pop()
    if not pair:
        break

    if not valid_pair(pair, cluster_id_list):
        continue

    merge(pair.cluster1, pair.cluster2,
          pair.distance, cluster_id_list, cluster_nodes)

    for index in range(length):
        if index != pair.cluster1 \
            and cluster_id_list[index]:
            cluster1, cluster2 = \
                small_big(cluster1, cluster2)

            new_distance = self.linkage_function(
                self.distance_matrix,
                cluster_id_list[smaller],
                cluster_id_list[bigger])

            distance_queue.push(ClusterPair(
                smaller,
                bigger,
                new_distance,
                len(cluster_id_list[smaller]),
                len(cluster_id_list[bigger])))

cluster_number -= 1

```

---

- SimRank output (to be integrated into main pipeline later)

This outputs structural similarity scores between tweets. As described in the

‘pruning’ section of the original SimRank paper (cited in Phase 2), we did not compute distances between tweet that were too far away from one another; thus, score of 0 are not presented in the histogram. Over a hundred thousand tweets were read.

Later, we may use Edit Distance to compute similarity scores and compare the clustering performance on using Edit Distance vs. using Document Vector.

- `simrank_tweets.py`: contains functions for running SimRank on bipartite graph
- `twitter_bipgraph.py`: create bipartite graph from tweets and run SimRank on it.

