

CS 6375.003 Machine Learning Assignment 4 Coding Part 2 Report

SHIVAM GUPTA (NET ID: SXG190040)

BHAVYA SREE BOMBAY (NET ID- BXB180036)

We have developed k-means clustering from scratch. In the **Pre-processing steps**, we have done the following:

- Removed the tweet id and timestamp
- Removed any word that starts with the symbol @ e.g. @AnnaMedaris.
- Removed any hashtag symbols e.g. convert #depression to depression.
- Removed any URLs.
- Converted every word to lowercase.

- The Tweet dataset being used for the following table:
foxnewshealth.txt – The total number of tweets in this dataset: **2000**
- We have also tested on different tweets text files for eg: msnhealthnews.txt, usnewshealth.txt, bbchealth.txt, etc
- The Code can be tested on other datasets by uncommenting line number 114, 115, etc. and change the value of K in line number 117.
- We have implemented the convergence criteria when all the centroids remain constant.

Table showing the Sum of squared error w.r.t Value of K and the size of Cluster

<u>Value of K</u>	<u>Sum of Squared Error (SSE)</u>	<u>Size of Each Cluster</u> <u>(Cluster Number : Number of tweets)</u>
<u>10</u>	<u>1447.7925723000096</u>	<u>1 : 47</u> <u>2 : 329</u> <u>3 : 382</u> <u>4 : 500</u> <u>5 : 97</u> <u>6 : 41</u> <u>7 : 41</u> <u>8 : 312</u> <u>9 : 90</u> <u>10 : 161</u>
<u>8</u>	<u>1157.3646461599678</u>	<u>1 : 431</u> <u>2 : 71</u> <u>3 : 88</u> <u>4 : 216</u> <u>5 : 491</u> <u>6 : 358</u> <u>7 : 240</u> <u>8 : 105</u>
<u>6</u>	<u>847.2669163699676</u>	<u>1 : 257</u> <u>2 : 297</u> <u>3 : 837</u> <u>4 : 42</u> <u>5 : 502</u> <u>6 : 65</u>
<u>5</u>	<u>701.1716522199871</u>	<u>1 : 495</u> <u>2 : 365</u> <u>3 : 85</u> <u>4 : 928</u> <u>5 : 127</u>
<u>4</u>	<u>563.2013628100026</u>	<u>1 : 107</u> <u>2 : 571</u> <u>3 : 354</u> <u>4 : 968</u>
<u>2</u>	<u>275.94057882000055</u>	<u>1 : 1305</u> <u>2 : 695</u>

Some of the screenshots of the results for the ushealth.txt tweet dataset:

```
Anaconda Prompt
[59, 237, 243]
Kmeans clustering converged after iteration no. = 3 for the value of k = 3
(base) F:\University of Texas at Dallas(UTD) MS-CS Fall 2019\CS-6375.003 Machine
Learning\Assignmnet_4>python Tweets_kmeans_clustering.py
Updated Centroids after iteration no. 1
[413, 355, 898, 237, 1303]
Cluster Number : Number of Tweets
1 : 491
2 : 137
3 : 196
4 : 369
5 : 207
Sum of square Error after iteration no. 1
652.8012201000088
Updated Centroids after iteration no. 2
[413, 1136, 598, 237, 59]
Cluster Number : Number of Tweets
1 : 202
2 : 264
3 : 249
4 : 149
5 : 536
Sum of square Error after iteration no. 2
617.3948241200093
Updated Centroids after iteration no. 3
[413, 243, 598, 237, 59]
Cluster Number : Number of Tweets
1 : 287
2 : 625
3 : 74
4 : 152
5 : 262
Sum of square Error after iteration no. 3
614.2651944900092
Updated Centroids after iteration no. 4
[413, 243, 598, 237, 59]
Kmeans clustering converged after iteration no. = 4 for the value of k = 5
(base) F:\University of Texas at Dallas(UTD) MS-CS Fall 2019\CS-6375.003 Machine
Learning\Assignmnet_4>
```

Anaconda Prompt

```

(base) F:\University of Texas at Dallas(UTD) MS-CS Fall 2019\CS-6375.003 Machine
Learning\Assignmnet_4>python Tweets_kmeans_clustering.py
Updated Centroids after iteration no. 1
[237, 598, 898, 243, 413, 268, 242, 59]
Cluster Number : Number of Tweets
1 : 317
2 : 240
3 : 129
4 : 294
5 : 133
6 : 70
7 : 41
8 : 176
Sum of square Error after iteration no. 1
1061.6291287100164
Updated Centroids after iteration no. 2
[237, 598, 898, 243, 413, 268, 1365, 59]
Cluster Number : Number of Tweets
1 : 216
2 : 70
3 : 460
4 : 171
5 : 101
6 : 98
7 : 182
8 : 102
Sum of square Error after iteration no. 2
1019.5997830300249
Updated Centroids after iteration no. 3
[237, 598, 898, 243, 413, 268, 1136, 59]
Cluster Number : Number of Tweets
1 : 201
2 : 68
3 : 392
4 : 169
5 : 138
6 : 104
7 : 198
8 : 130
Sum of square Error after iteration no. 3
1003.9866305300237
Updated Centroids after iteration no. 4
[237, 598, 898, 243, 413, 268, 1136, 59]
Kmeans clustering converged after iteration no. = 4 for the value of k = 8

(base) F:\University of Texas at Dallas(UTD) MS-CS Fall 2019\CS-6375.003 Machine
Learning\Assignmnet_4>_

```

```
Anaconda Prompt
(base) F:\University of Texas at Dallas(UTD) MS-CS Fall 2019\CS-6375.003 Machine Learning\Assignmnet_4>python Tweets_kmeans_clustering.py
Updated Centroids after iteration no. 1
[237, 59, 243, 898, 413, 956, 372, 598, 707, 67]
Cluster Number : Number of Tweets
1 : 233
2 : 125
3 : 215
4 : 165
5 : 142
6 : 37
7 : 188
8 : 53
9 : 104
10 : 138
Sum of square Error after iteration no. 1
1335.8627341099855
Updated Centroids after iteration no. 2
[237, 59, 243, 898, 413, 956, 372, 598, 1136, 67]
Cluster Number : Number of Tweets
1 : 150
2 : 143
3 : 128
4 : 261
5 : 46
6 : 28
7 : 88
8 : 55
9 : 158
10 : 343
Sum of square Error after iteration no. 2
1281.448766729977
Updated Centroids after iteration no. 3
[237, 59, 243, 898, 413, 956, 372, 598, 1136, 67]
Kmeans clustering converged after iteration no. = 3 for the value of k = 10
(base) F:\University of Texas at Dallas(UTD) MS-CS Fall 2019\CS-6375.003 Machine Learning\Assignmnet_4>
```

Results and Conclusions:

- We saw that for different values of k we are getting different values of sum of squared error.
- The table above shows the sum of squared errors (SSE) for different values of K and also the size of tweets in each cluster.
- In this tweet dataset (foxnewshealth.txt) we observed that the sum of squared error got decreased by decreasing the value of K .
- There is no specific trend in increment and decrement of the SSE. This trend is going to vary from dataset to dataset and it also depends on the initial seeds (centroids) and we have taken it randomly.