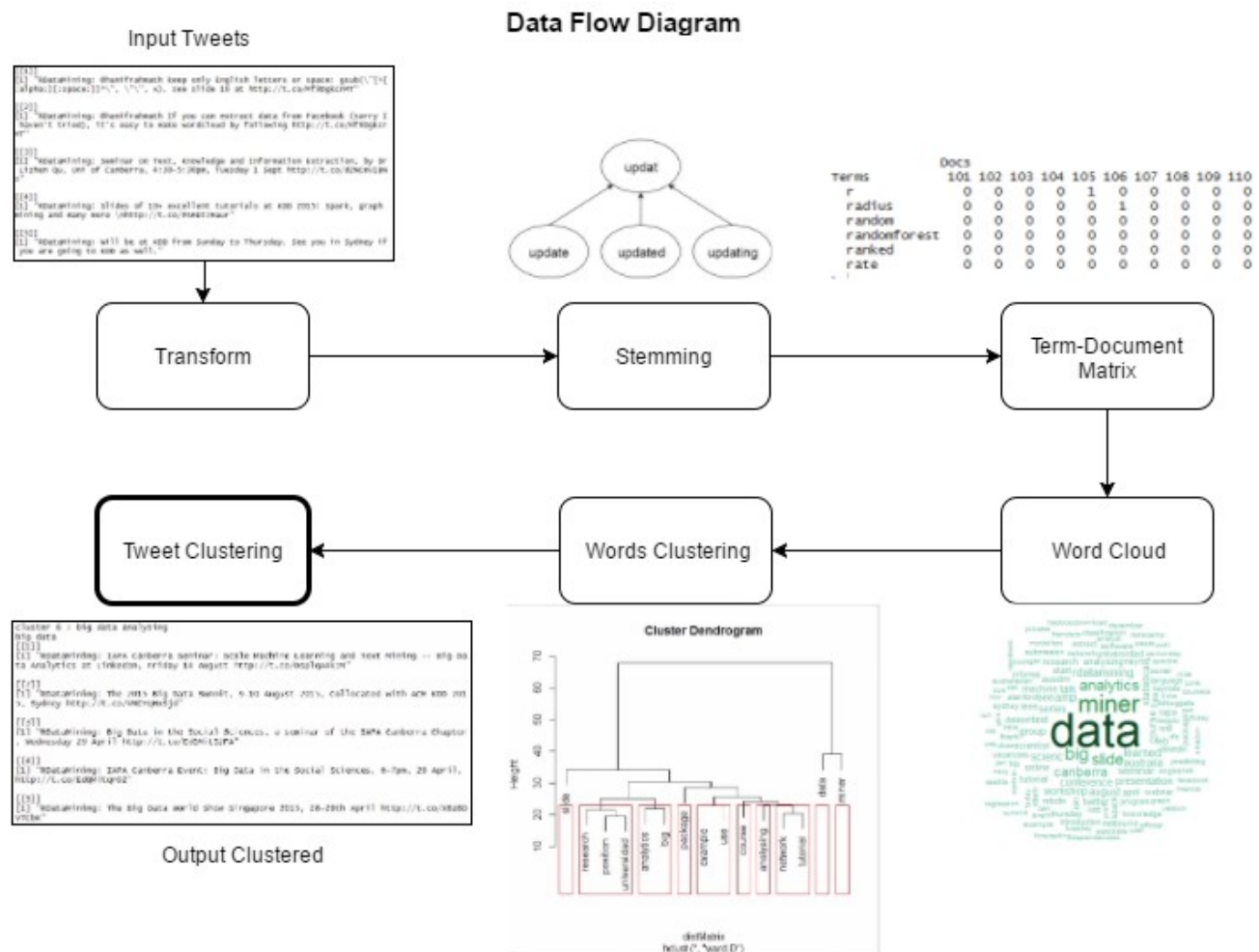# Twitter Text Mining Using Clustered Search

Anjali Malav(UG201310006)

Muttineni Navya(UG201310020)

**Problem Statement :** Given a search word, the objective is to return results in this case tweets in clustered manner from a large database of tweets.

**Introduction :** Twitter is a major microblogging service. When a user searched these tweets with a search term, thousands of tweets are displayed. The user has to go through every tweet in order to see if that tweet is of his relevance or not. This burden would be reduced if the tweets were clustered first and then returned to the user topic wise.



Data Flow Diagram

**Methodology :** The methodology is as given in the above diagram. For

clustering of tweets we have used two algorithms, K-means and K-medoids.
<u>Comparison between K-means and K-medoids</u>

1. Number of clusters needs to be specified in advance where as in k-medoids is not required.
2.  k-medoids is based on centroids (or medoids) calculating by minimizing the absolute distance between the points and the selected centroid, rather than minimizing the square distance. As a result, it's more robust.
3.  K-means is faster than K-medoids.
4. Implementation of k-means is easier than k-medoids which is complex.

**Applications :**

| Application | Benefit | Example |
|---|---|---|
| Search result clustering | more effective information presentation to user | |
| Scatter-Gather | alternative user interface: search without typing | |

**Conclusions :** The input data i.e the number of tweets being retrieved is small (200-600). Scalability is not handled as working on large data sets needs a lot of computation which would take a lot of time on 1 CPU. This can be handled by computing on more than 1 CPUs. This can be one of the future scopes of this project with others being sentiment analysis and social network analysis.