

# Certificate

This is to certify that the project entitled "**Twitter Data Analytics on Women Education and Employment**" is being submitted at IGDTUW, Delhi for the award of **Bachelor of Technology** in **Computer Science & Engineering** degree. It contains the record of bonafide work carried out by **NAME OF STUDENT** under my supervision and guidance. It is further certified that the work presented here has reached the standard of B.Tech and to the best of my knowledge has not been submitted anywhere else for the award of any other degree or diploma.

**Prof. Ela Kumar**

(Dean of Academic Welfare)

Computer Science Department

IGDTUW, Delhi

Date:

Place: New Delhi

# **Student's Declaration**

We, Group Id-1 hereby declared that the work being presented in this report entitled "**Twitter Data Analytics on Women Education and Employment**" submitted to INDIRA GANDHI DELHI TECHNICAL UNIVERSITY FOR WOMEN, New Delhi, for the award of degree of B.Tech CSE is an authentic record of our work carried out under the guidance of **Prof. Ela Kumar** (Dean of Academic Affairs), IGDTUW, Delhi.

The matter embodied in this report has not been submitted by us for the award of any other degree.

Dated:

## **GROUP MEMBERS :**

<b>Sonal</b>	<b>06713502711</b>
<b>Manali Verma</b>	<b>06813502711</b>
<b>Sonam Pal</b>	<b>07013502711</b>

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Under the Guidance of

**Prof. Ela Kumar**  
**(Dean of Academic Affairs)**  
**Computer Science Department**  
**IGDTUW, Delhi**

# Acknowledgement

We would like to express our deepest appreciation to all those who provided us the possibility to complete this report. First and foremost, we would like to express infinite thanks to our Mentor **Prof. Ela Kumar (Dean, Academic Affairs), IGDTUW, Delhi** for provision of expertise, and technical support in the implementation. Without their superior knowledge and experience, the Project would like in quality of outcomes, and thus their support has been essential. We would like to express our sincere thanks towards volunteer researchers who devoted their time and knowledge in the implementation of this project.

Nevertheless, we express our gratitude toward our families and colleagues for their kind co-operation and encouragement which help us in completion of this project.

Date:

Place: New Delhi

## **ABSTRACT**

---

This Document is submitted for the purpose of fulfillment of Degree of B. Tech in Computers Science and Engineering. The objective of this project is to develop an **Analysis on Twitter data on women education and employment**. In this project, our aim is to provide an enhancement in the women empowerment by performing extraction, transformation and standardization of the tweets from the Twitter. Twitter makes it easy to engage users and communicate directly with them, and in turn, users can provide word-of-mouth marketing for companies by discussing the products. The tweets are converted from JSON format into excel sheet. NodeXL is a free and open-source network analysis and visualization software package for Microsoft Excel 2007/2010. It is a popular package similar to other network visualization tools such as Pajek, UCINet, and Gephi. NodeXL imports UCINet and GraphML files, as well as Excel spreadsheets containing edge lists or adjacency matrices, into NodeXL workbooks. NodeXL also allows for quick collection of social media data via a set of import tools which can collect network data. NodeXL contains a library of commonly used graph metrics: centrality, clustering coefficient, diameter. NodeXL differentiates between directed and undirected networks. NodeXL implements a variety of community detection algorithms to allow the user to automatically discover clusters in their social networks.

This project is to provide an estimate and statistics of the education and employment of women all across the world especially in India. This project implements all the concepts learned in different subjects such as DBMS and ERP and provides a valuable experience to gain a sound understanding as well as utilize creativity to develop this solution. If successful, this project has the potential to be implemented within the organization as well.

# **CONTENTS**

---

<b>Certificate</b>	<b>1</b>
<b>Student's Declaration</b>	<b>2</b>
<b>Acknowledgement</b>	<b>3</b>
<b>Abstract</b>	<b>4</b>
	32
<b>Chapter 1: Introduction to Project</b>	<b>7</b>
1.1 Overview	7
1.2 Project Background	7
1.4 Thesis Objective and Scope	8
1.4 Thesis Outline	8
<b>Chapter 2: Literature Review</b>	<b>10</b>
2.1 General	10
2.2 Related Work	10
2.2.1 System Description and Workflow	13
<b>Chapter 3 Problem Identification</b>	<b>16</b>
3.1 Background	16
3.2 Problem in Big Data for Data Extraction, Transformation And Standardization	17
<b>Chapter 4 Proposed Approaches</b>	<b>19</b>
4.1 Motivation	19
4.2 Proposed Solution	19
4.2.1 Assumptions	20
4.2.2 Algorithm	20

4.3	Comparison of NodeXL and HootSuite	24
<b>Chapter 5</b>	<b>Implementation Details</b>	<b>25</b>
5.1	Getting the Twitter API Keys	25
5.2	Connecting to Twitter Streaming API and downloading data	26
5.3	NodeXL Installation to Import Tweets	28
5.4	Mining Tweets	29
<b>Chapter 6</b>	<b>Conclusion and Future work</b>	<b>48</b>
6.1	Conclusion	48
6.2	Future Work	48
	<b>References</b>	<b>49</b>

# **CHAPTER 1**

## **Introduction**

---

### **Overview**

The term Big Data is used almost anywhere these days. Big Data can be seen in the social media, finance and business where enormous amount of stock exchange, banking, online and onsite purchasing data flows through computerized systems every day and are then captured and stored for inventory monitoring, customer behavior, human relationship and market behavior. It can also be seen in the life sciences where big sets of data such as genome sequencing, clinical data and patient data are analyzed and used to advance breakthroughs in science in research.

### **Project Background**

Our aim for developing this project is to provide an estimate and statistics of the education and employment of women all across the world especially in India. In India, the women literacy rate is very low and the educated women are paid less than men. Worth of the work done or services rendered by women has not been recognized. We by this project are trying to give a hand in support for women empowerment. Since Social media has gained immense popularity with marketing teams, and Twitter is an effective tool for a company to get people excited about its products we are using Twitter's Data for our analytics. Twitter makes it easy to engage users and communicate directly with them, and in turn, users can provide word-of-mouth marketing for companies by discussing the products. Given limited resources, and knowing we may not be able to talk to everyone we want to target directly, marketing departments can be more efficient by being selective about whom we reach out to.

After extracting the number of tweets, we can perform analytics by using a query based interface tool and then can produce results

### **Thesis Outline**

This thesis report consists of seven chapters, References and Appendix. Rest of the thesis is organized as following:

**Chapter 2-** Various techniques and methodologies are using for extracting the data from social media. Here, we discuss the recent extraction details that are used for collecting, transforming and visualizing the data. Different organizations have different techniques to update the big data from social media.

**Chapter 3-**This chapter discusses the problem in extracting the huge amount of data from social media with seed words. We discuss the problem in transforming and standardize the data of the women' education and employment taken directly from twitter. This chapter also shows the position and implications of women in India.

**Chapter 4-** The proposed approach algorithm is described in this chapter followed by the system Requirements, Specifications, hardware architecture, software architecture, System integration and layered structure of the system. This chapter gives the details about the hardware components used for the system along with their specifications.

**Chapter 5-**The software implementation and schematic of the system and the integration of both hardware and software is explained in this section. The working of this project is explained with the help of a sequence diagram and the flow of software code is also explained in the form of a flowchart.

**Chapter 6-** Finally thesis ends with the conclusion and contributions of the system. The limitations of the system are also given in this chapter and at the end the scope for future work is given.

# CHAPTER 2

## Literature Review

To develop an **Analysis on Twitter data on women education and employment**. Our aim is to provide an enhancement in the women empowerment by performing extraction, transformation and standardization of the tweets from the Twitter.

### 2.1 General

Twitter is an online social networking service that enables users to send and read short 140-character messages called "tweets". Registered users can read and post tweets, but unregistered users can only read them. Users access Twitter through the website interface, SMS, or mobile device app. Twitter Inc. is based in San Francisco and has more than 25 offices around the world.

In order to retrieve tweets from Twitter in real time, querying Twitter data in a traditional RDBMS is inconvenient, since the Twitter Streaming API outputs tweets in a JSON format which can be arbitrarily complex. A Python library called Tweepy is used to connect to Twitter Streaming API and downloading the data. Tweepy is a Python 2.6 and 2.7 library for accessing Twitter. It provides access to all twitter RESTful API methods, including reading and posting of tweets. Tweepy supports OAuth authentication, as BasicAuth is no longer supported by the Twitter API.

**NodeXL** is a free and open-source network analysis and visualization software package for Microsoft Excel 2007/2010. It is a popular package similar to other network visualization tools such as Pajek, UCINet, and Gephi.

NodeXL is a set of prebuilt class libraries using a custom Windows Presentation Foundation control. Additional .NET assemblies can be developed as "plug-ins" to import data from outside data providers. Currently-implemented data providers for NodeXL include Facebook, Twitter, Wikipedia (theMediaWiki understructure), web hyperlinks, Microsoft Exchange Server.

NodeXL is intended for users with little or no programming experience to allow them to collect, analyze, and visualize a variety of networks. NodeXL integrates into Microsoft Excel 2007 and 2010 and opens as a workbook with a variety of worksheets containing the elements of a graph structure such as edges and nodes. NodeXL can also import a variety of graph formats such as edgelists, adjacency matrices, GraphML, UCINet .dl, and Pajek .net.

NodeXL imports UCINet and GraphML files, as well as Excel spreadsheets containing edge lists or adjacency matrices, into NodeXL workbooks. NodeXL also allows for quick collection of social media data via a set of import tools which can collect network data from e-mail, Twitter, YouTube, and Flickr. NodeXL requests the user's permission before collecting any personal data and focuses on the collection of publicly available data, such as Twitter statuses and follows relationships for users who have made their accounts public. These features allow NodeXL users to instantly get working on relevant social media data and integrate aspects of social media data collection and analysis into one tool.

NodeXL workbooks contain four worksheets: Edges, Vertices, Groups, and Overall Metrics. The relevant data about entities in the graph and relationships between them are located in the appropriate worksheet in row format. For example, the edges worksheet contains a minimum of two columns, and each row has a minimum of two elements corresponding to the two vertices that make up an edge in the graph. Graph metrics and edge and vertex visual properties appear as additional columns in the respective worksheets. This representation allows the user to leverage the Excel spreadsheet to quickly edit existing node properties and to generate new ones, for instance by applying Excel formulas to existing columns.

NodeXL contains a library of commonly used graph metrics: centrality, clustering coefficient, diameter. NodeXL differentiates between directed and undirected networks. NodeXL implements a variety of community detection algorithms to allow the user to automatically discover clusters in their social networks.

NodeXL generates an interactive canvas for visualizing graphs. The project allows users to pick from several well-known Force-directed graph drawing layout algorithms such as Fruchterman-Reingold and Harel-Koren. NodeXL allows the user to multi-select, drag and drop nodes on the canvas and to manually edit their visual properties (size, color, and opacity). In addition, NodeXL allows users to map the visual properties of nodes and edges to metrics it calculates, and in general to any column in the edges and vertices worksheet.

The software architecture comprises three extendable layers:

**Data Import Features:** NodeXL stores data in a pre-defined Excel template that contains the information needed for generating network charts. Data can be imported from existing Pajek files, other spreadsheets, comma separated value (CSV) files, or incidence matrices. NodeXL also extracts networks from a small but extensible set of data sources that includes email stored in the Windows Search Index and the Twitter micro-blogging network. Email reply-to information from personal e-mail messages is extracted from the Microsoft Windows Desktop Search index. Data can also be imported about which user subscribes to one another's updates in Twitter, a micro-blogging social network system. NodeXL has a modular architecture that allows for the integration of new components to extract and import network data from additional resources, services, and applications. The open source access to the NodeXL code allows for a community of programmers to extend the code and provide interfaces to data repositories, analysis libraries, and layout methods. Spreadsheets can then be used in a uniform way to exchange network data sets by a wide community of users.

**Network Analysis Module:** NodeXL represents a network in the form of edge lists, i.e., pairs of vertices which are also referred to as nodes. Each vertex is a representation of an entity in the network. Each edge, or link, connecting two vertices is a representation of a relationship that exists between them. This relationship may be directed or not. Some relationships are bidirectional (like marriage); others can be uni-directional (like lending money). the presence of a relationship. These lists can be extended with additional columns that can contain data about the relationship. NodeXL includes a number of software routines for calculating statistics about

individual vertices including in-degree, out-degree, clustering coefficient, and closeness, betweenness, and eigenvector centrality. Additional analyses features can be integrated by advanced users. The results of the network metric calculations are added to the spreadsheet as additional columns that can be further combined and reused in Excel formula during analysis and visualization. Spreadsheet features like data sorting, calculated formulae, and filters can be applied to network data sets directly.

**Graph Layout Engine:** NodeXL provides a canvas for displaying and manipulating network charts and data. Users can apply a range of controls to convert an edge list into a useful node-link chart. These include display options that specify the appearance of individual edges and nodes as well as the overall layout of the network. The lines between nodes that represent edges can have different thickness, color, and level of transparency depending on the attributes of the data or parameters specified by the user. Similarly, each node representing a vertex can be set to have a different location, size, color, transparency, or shape. Optionally, the user can specify images to replace the node shapes.

## 2.2.1 SYSTEM DESCRIPTION AND WORKFLOW

The core of NodeXL is a special Excel workbook template that structures data for network analysis and visualization. Six main worksheets currently form the template. There are worksheets for “Edges”, “Vertices”, and “Images” in addition to worksheets for “Clusters,” mappings of nodes to clusters (“Cluster Vertices”), and a global overview of the network’s metrics (“Overall Metrics”). NodeXL workflow typically moves from data import through processing, calculation and refinement before resulting in a network graph that tells a useful story. These steps include:

**Step 1: Import data.** Network data can be imported from one or more network data sources. Users may have data in files, e.g. in text format, separated by delimiters, or stored in relational databases. Wherever the data originates, it is entered into the NodeXL template in the “Edges” worksheet in the form of pairs of names, along with any additional attributes about relationship between them. Multiple edge lists can be stored in the same spreadsheet, expressing different relationships among a set of nodes or the same relationships at different times. Relationships can

be annotated with multiple additional columns, which can be used to set values for display attributes. Data about the “strength” of the relationship can be included, or edges can be annotated with the time slice or date range in which they occurred, allowing a single dataset to contain edges over multiple time periods.

**Step 2: Clean the data.** This typically involves eliminating duplicate edges when appropriate. Data sets can often be noisy and contain redundant data. In some cases network measures cannot be calculated correctly if multiple edges between the same pair of entities exist in a single data set. In these cases redundant edges may be aggregated into a single edge with a weighting that reflects the number of original instances.

**Step 3: Calculate graph metrics**-A range of measurements exist that capture the size and internal connectivity of a network as well as attributes of each node. NodeXL supports a minimal set of the most crucial network measures for individual nodes: in- and outdegree, clustering coefficient, and betweenness, closeness, and eigenvector centrality. This operation also populates the Vertex worksheet with a unique list of nodes and their network measures. In some cases multiple edges between nodes are part of what is interesting in a data set and should be retained despite the fact that some metrics will be inaccurate and should be ignored. NodeXL marks these network metrics with the Excel “Bad” format if asked to calculate some metrics with duplicate edges in the data set.

**Step 4: Create clusters**-Network nodes may share a variety of attributes. It is often useful to group and analyze them together. NodeXL has a clustering algorithm and allows users to create clusters and map nodes to them by editing the Clusters and Cluster Vertices worksheets. Each cluster can have its own display attributes with a distinctive shape, color, size, transparency or image. Users can toggle the display of clusters so that the display features for each node is replaced by the display features for its cluster (if any).

**Step 5: Create sub-graph images**-Whole graph images are often too large or dense to reveal details about individual nodes or clusters. Sub-graph images produce a local network that centers on each node at a time and encompasses the nodes to which it is immediately connected. These

extracted images can be useful representations of the range of variation in the local network structures of the population in the network.

**Step 6: Prepare edge lists-** Nodes can have a “Layout Order” value that governs the presentation of nodes in the graph display. Nodes and edges have attributes that can be used to order the data, for example, ordering nodes by their date of first appearance or rate of connection to other nodes. The value found in “Layout Order” governs the order in which nodes are laid out in the whole graph visualization.

**Step 7: Expand worksheet with graphing attributes-** Columns can be auto-filled to map data to display attributes. Graphical attributes of nodes and edges, their shape, color, opacity, size, label, and tooltip can be altered to convey additional information in the network visualization. Images listed in the “Images” worksheet can replace the shapes used to represent nodes. Users can insert additional numerical attributes about each node in adjacent columns. These attributes can be automatically scaled to display characteristics. For example, each node may have data about the income of the person it represents or the number of employees in an enterprise which could be mapped to the size, shape, or color of a node. Once set, these mappings apply to all networks created from that point on with NodeXL until reset “sticky” layout feature simplifies the creation of multiple networks while maintaining consistency of display mappings.

**Step 8: Show graph-**This opens or hides the graphical display pane in which NodeXL will render a visualization of the network. At each stage of the NodeXL workflow, the toolkit provides a number of options; the workflow is not rigidly prescribed. The user can iteratively refine any stage of the analysis and visualization. This may involve operations like: Read workbook. Load the current state of the network as stored in the spreadsheet and render it according to the selected layout. Adjust layout. Select among a number of automated layout options that govern where each node in the network will be located. Layouts include a force-directed Fruchterman-Reingold layout [10], that attempts to dynamically find a layout that clusters tightly connected nodes near one another as well as simple geometric layouts like circles or grids. Apply dynamic filters. Selectively hide edges and nodes, depending on the attributes of

the network. For example, a filter could hide all but the most connected nodes, or show only the edges that are ‘stronger’ than a selected threshold. Filtering may involve “trimming” parts of the network and then recalculating network metrics and layout based on the remaining population of nodes and edges. Re-render the graph. Redraw the network based on remaining nodes and edges and their changed display attributes. Finalize the network analysis

# CHAPTER 3

## Problem Identification

There are several technologies used for extracting the data from social media. Also several machine learning, open source, freeware and licensed software are available in the market. We have to select one among this wide range of components which appropriately suits the design and helps in achieving the goals and objectives set for this project. This chapter gives details about the problem arise in extracting the big data from social media directly and gives inefficient results.

### **3.1 Background**

Analyzing information involves examining it in ways that reveal the relationships, patterns, trends, etc. that can be found within it. That may mean subjecting it to statistical operations that can tell you not only what kinds of relationships seem to exist among variables, but also to what level you can trust the answers you're getting. It may mean comparing your information to that from other groups (a control or comparison group, statewide figures, etc.), to help draw some conclusions from the data. The point, in terms of your evaluation, is to get an accurate assessment in order to better understand your work and its effects on those you're concerned with, or in order to better understand the overall situation.

There are two kinds of data you're apt to be working with, although not all evaluations will necessarily include both. **Quantitative data** refer to the information that is collected as, or can be translated into, numbers, which can then be displayed and analyzed mathematically. **Qualitative data** are collected as descriptions, anecdotes, opinions, quotes, interpretations, etc., and are generally either not able to be reduced to numbers, or are considered more valuable or informative if left as narratives. As you might expect, quantitative and qualitative information needs to be analyzed differently.

The problem of taking the project on women's education and employment not only to show the impact on one side but also emphasis to each organization. There are some questions arise when we extract and transform the data sets on women's education and employment. To what extent legislative measures have been able to raise the status of women in India? Are women now feel empowered in the sense that they are being equally treated by men in all spheres of life and are able to express one's true feminine urges and energies? What initiatives have proven successful in helping women and girls transition from school/training to work? What about institutional efforts on gender equality and employment? What types of non-formal education and training do women and girls participate in? What about women employment after marriage? These are the important questions to be investigated with regard to women's empowerment in India.

### **3.2 Problem in Data Extraction, Transformation and Standardization**

The problems start right away during data acquisition, when the great volume of data flooded beyond the memory and requires us to make decisions. This continuously creates an problem to the data analyst, business intelligence manager and data scientist to take strategic decision to accommodate the petabytes of data. Currently there is an ad hoc manner, about what data to keep and what to discard, and how to store what we keep reliably with the right metadata. Another problem arises when the extracted data is not in the format that we use in our system. Then we manually manage the structured format of data. Naive user approximately analyzes the structured and standardized the data. But, if we have a petabytes of data in our data-sets, then there would be a problem to standardize the data. As we know that the data in the blogs, tweets, facebook, google hangout are not so much structured and have multiple format of texts, images audio and video for display and storage. So, we need to transform the data according to the user perspective into the standardized and structured way. The data should be platform independent and open source to everyone. This becomes a great challenge to manage the huge amount of data in the memory for collecting and transforming to conclude the effective results. Data should be integrated and credible to every data analyst for performing efficient analysis on seed words from social media. Today, mostly the data is directly extracted from the social media and stored for analysis, which gives inappropriate conclusion and percentage of growth. We have the opportunity and the challenge to facilitate the seed words to justify the queries of analyst. Data

analysis, organization, retrieval, and modeling are other foundational challenges. Algorithm design to get effectively retrieval of data is another increases the complexity at social media.

Even after data cleaning and error correction, some incompleteness and some errors in data are likely to remain. This incompleteness and these errors must be managed during data analysis.

# **CHAPTER 4**

## **Proposed Approach**

---

This chapter describes the proposed approach and conceptual approach of an application regarding social media. The approach related with the concepts of extracting the tweets from twitter for data analysts. The conceptual design of slow growth in extracting the huge amount of data is a concern for organizations. It comprises system design and methodology related with the requirement. It is an approach to clarify the techniques to extract the data effectively from social media.

### **4.1 Motivation**

Basically, we should take a stand to empower the women for education and employment. This should help the country to visualize the perspective related with women and girls. The subject of empowerment of women has becoming a burning issue all over the world including India since last few decades. They have demanded equality with men in matters of education, employment, inheritance, marriage, and politics and recently in the field of religion also to serve as cleric (in Hinduism and Islam). Women want to have for themselves the same strategies of change which men-folk have had over the centuries such as equal pay for equal work. Their quest for equality has given birth to the Big Data Analytics from Social Media on Women's education and employment.

### **4.2 Proposed Solution**

We designed the algorithm for keywords cleaning and eliminating the repeated tweets from a twitter user by the following mechanism:

1. Eliminating the non- English tweets.
2. Counting the retweets from the data sets, this refers to tweets from somebody and shared by any other person because this may conclude the enrichment in women's education and employment.
3. Elimination of URLs from the tweets.

4. Eliminating the hashtag (#) from the tweets because great volumes of data are attached with seed words, which tend to magnify the vocabulary size inadvertently.
5. Tackle the problem of repetitions of unnecessary letters in the tweets eg. Wwwwoommmeeeeennn to women.
6. Eliminating the format of @username from the tweets.
7. Condense each 2 repetitive letters into single letter if the number of such sequences is over a threshold. For example, ookk will be reduced to ok, while book remains intact.
8. Removing all the punctuation marks and stopword.
9. Eliminating the additional English grammar sentences.
10. Eliminating the non seed words. For example, Womania, educationaism, emplaymentnews and others.
11. Eliminating the name of the same twitterer who shared his/her tweets with others because they unnecessarily increase the volume of data.

#### **4.2.1 Assumption**

The algorithm of Fruchterman and Reingold added “even vertex distribution” to the earlier two criteria and treats vertices in the graph as “atomic particles or celestial bodies, exerting attractive and repulsive forces from one another.” The attractive and repulsive forces are redefined to

$$fa(d) = d^2/k \quad , \quad fr(d) = -k^2/d$$

in terms of the distance  $d$  between two vertices and the optimal distance between vertices  $k$  defined as

$$k = C (\text{area} / \text{number of vertices})^{1/2}$$

The algorithm of Fruchterman and Reingold adds the notion of “temperature” which could be used as follows: “the temperature could start at an initial value (say one tenth the width of the frame) and decay to 0 in an inverse linear fashion.” The temperature controls the displacement of vertices so that as the layout becomes better, the adjustments become smaller. The use of temperature here is a special case of a general technique called simulated annealing.

**Pseudo code :**

```

while temp > 0.5 and passes < maxIterations (500)
//calculate repulsive forces between each node
for v = 0 to numberOfNodes
  for u = v+1 to numberOfNodes
    calculateRepulsiveForce(v, u, temp)
    calculateAttractiveForce(v, u, temp)
    moveNode(v, u)
    updatePositions()
    decreaseTemp()
  
```

```

calculate the distance vector between the positions of v and u
calculate a displacement displaceVec = (distVec/|distVec|) * repulsion(|distVec|)
add displaceVec vector to v's displacement vector
subtract displaceVec from u's displacement vector
end
end
//calculate attractive forces
for e = 0 to numberOfEdges
get the nodes attached to the edge (v and u)
calculate the distance vector between the positions of v and u
calculate a displacement displaceVec = (distVec/|distVec|) * attraction(|distVec|)
subtract displaceVec vector from v's displacement vector
add displaceVec to u's displacement vector
end

```

calculate each nodes's displacement, but limit max displacement to temp

//decrease temperature parameter according to cooling schedule

if "Show updates" is true and this is an Nth pass, update the layout on screen  
at the end, go over all the nodes to find the max and min of the coords, rescale all coords so that  
network will fill the display end while

//repulsion function

repulsion(distance) = (optDist^2\*arcWeight^2) / layoutdistance

//attraction function attraction(distance) =  
layoutDistance^2 / (optDist\*arcWeight)

//cooling function

coolTemp(temp) = querys the CoolingSchedule for the the appropriate value for the iteration, as  
set by the user with the cooling function.

//optimal distance optimalDistance = layout parameter in apply layout settings

### **Asymptotic Notation and Complexity Calculation:**

Each iteration the basic algorithm computes  $O(|E|)$  attractive forces and  $O(|V|^2)$  repulsive forces.  
To reduce the quadratic complexity of the repulsive forces, Fruchterman and Reingold suggest  
using a grid variant of their basic algorithm, where the repulsive forces between distant vertices  
are ignored. For sparse graphs, and with uniform distribution of the vertices, this method allows

a  $O(|V|)$  time approximation to the repulsive forces calculation. This approach can be thought of as a special case of the multi-pole technique introduced in n-body simulations

### 4.2.2 Algorithm

Multi-level Graph algorithms are designed to make all the edges about the same length and to minimize line crossings, which can make for a more aesthetically pleasing and readable graph.

#### Notation and Definition:

Let  $G = G(V, E)$  be an undirected graph of vertices  $V$ , with edges  $E$  and which we will assume is connected. For any vertex  $v$  let  $\Gamma_v$  be the neighbourhood of, or set of vertices adjacent to,  $v$ , i.e.,  $\Gamma_v = \{u \in V : (u, v) \in E\}$ . We use the  $|.|$  operator to denote the size of a set so that  $|V|$  is the number of vertices in the graph and  $|\Gamma_v|$  is the number of vertices adjacent to  $v$  (the degree of  $v$ ). We also use  $|.|$  to denote the weight of a vertex; since weighted vertices in the coarsened graphs represent sets of vertices from the original graph, the weight of a coarsened vertex is just equivalent to the number of original vertices in the set it represents. We then use  $\|.\|$  to denote Euclidean distance in either 2D or 3D.

#### Algorithm:

```

{ initialisation }

function fr(x, w) := begin return -Cwk2/x end

function fa(x) := begin return x2/k end

t := t0;

Posn := NewPosn;

while (converged = 1) begin
    converged := 1;
    for v ∈ V begin
        OldPosn[v] = NewPosn[v]
    end
    for v ∈ V begin
        { initialise Θ, the vector of displacements of v }
        Θ := 0;
        { calculate (global) repulsive forces }
        for u ∈ V, u = v begin

```

```

 $\Delta := \text{Posn}[u] - \text{Posn}[v];$ 
 $\Theta := \Theta + (\Delta/\|\Delta\|) \cdot \text{fr}(\|\Delta\|, |u|);$ 
end
{ calculate (local) attractive/spring forces }
for  $u \in \Gamma v$  begin
 $\Delta := \text{Posn}[u] - \text{Posn}[v];$ 
 $\Theta := \Theta + (\Delta/\|\Delta\|) \cdot \text{fa}(\|\Delta\|);$ 
end
reposition  $v$ 
NewPosn[v] = NewPosn[v] + ( $\Theta/\|\Theta\|$ )  $\cdot \min(t, \|\Theta\|)$ ;
 $\Delta := \text{NewPosn}[v] - \text{OldPosn}[v];$ 
if ( $\|\Delta\| > k \cdot \text{tol}$ ) converged := 0;
end
{ reduce the temperature to reduce the maximum movement }
t := cool(t);
end

```

### **Complexity Calculation:**

If we set  $R$  to be the maximum distance over which repulsive forces will act we can then modify the algorithm by changing the global force calculation to:

```

function fr(x, w) :=
begin
  if ( $x \leq R$ ) return  $-Cw^2/x$ ;
  else return 0.0;
end

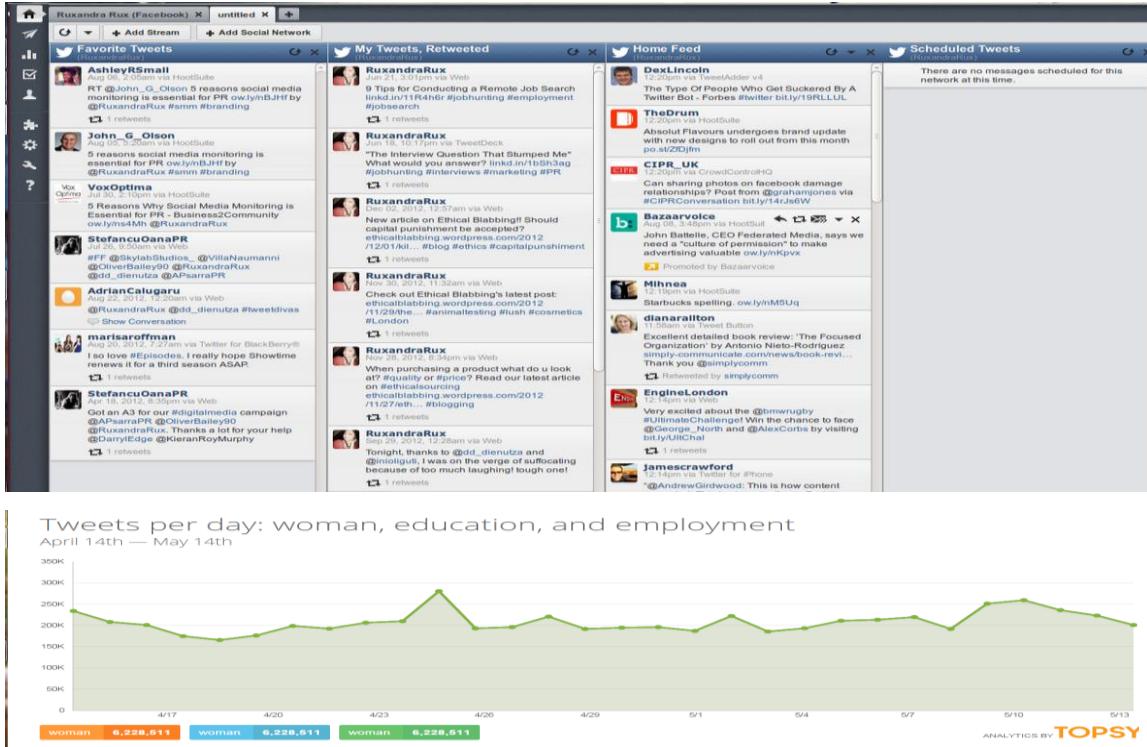
```

In itself this modification will do little or nothing to speed up the calculation as the complexity is still  $O(|V|)$ . However Fruchterman & Reingold, [7], showed that if the domain is divided up into regular square cells (or cube shaped cells in 3D) of size  $R^2$  (or  $R^3$  in 3D) then each vertex will only be affected by repulsive forces from vertices in its own and adjacent cells (including those diagonally adjacent). To implement this efficiently we simply visit every vertex at the start of each outer loop and add each to a linked list of vertices for the cell to which it belongs. Repulsive

forces can then be calculated for each vertex by using the linked lists of their own and adjacent cells.

## 4.3 Comparison of NodeXL and HootSuite

In this section we compare our analysis with HootSuite.



HootSuite has confusing interface because there is so much activity taking place within any one tab. It has longer learning curve. Hootsuite can incorporate analytics from Twitter, Facebook, LinkedIn, Google Analytics and its own Ow.ly shortened URL metrics. Costs can mount up. Although the Pro plan costs just \$9.99 per month, adding more than one team member increases the cost — adding a second team member doubles the price, to almost \$20 per month.

# CHAPTER 5

## IMPLEMENTATION

---

The people are responsible for extracting and transforming the huge amount of data to analyze various types of queries, so it starts by thinking about its structure. This describes the implementation details of the tweets from the twitter. What types of tweets are extracted from the twitter? What should each one do while transforming the data? And how should analyze the keywords of the twitter? Once code actually exists, they ask more questions. What does this class look like? What other classes is it related to? What's the sequence of calls from this method? All of these questions lend themselves to visual answers. In every case, creating diagrams that show what's going on can be the clearest path to understanding.

### 5.1 Getting the Twitter API Keys

1. Create account on Twitter Developer : <https://dev.twitter.com/apps> and log in with your twitter credentials.
2. Click "Create New App"
3. Fill out the form, agree to the terms, and click "Create your Twitter application"
4. In the next page, click on "API keys" tab, and copy your "API key" and "API secret".
5. Click "Create my access token", and copy your "Access token" and "Access token secret".

```
TwitterAgent.sources = Twitter
```

```
TwitterAgent.sinks = HDFS
```

```
TwitterAgent.sources.Twitter.type = com.cloudera.flume.source.TwitterSource
```

```
TwitterAgent.sources.Twitter.channels = MemChannel
```

```
TwitterAgent.sources.Twitter.accessToken = <accessToken>
```

```
TwitterAgent.sources.Twitter.keywords=women,education, employment
```

## Twitter67

Details    Settings    Keys and Access Tokens    Permissions

Data Analytics  
<https://github.com/sonal67>

**Organization**  
Information about the organization or company associated with your application. This information is optional.

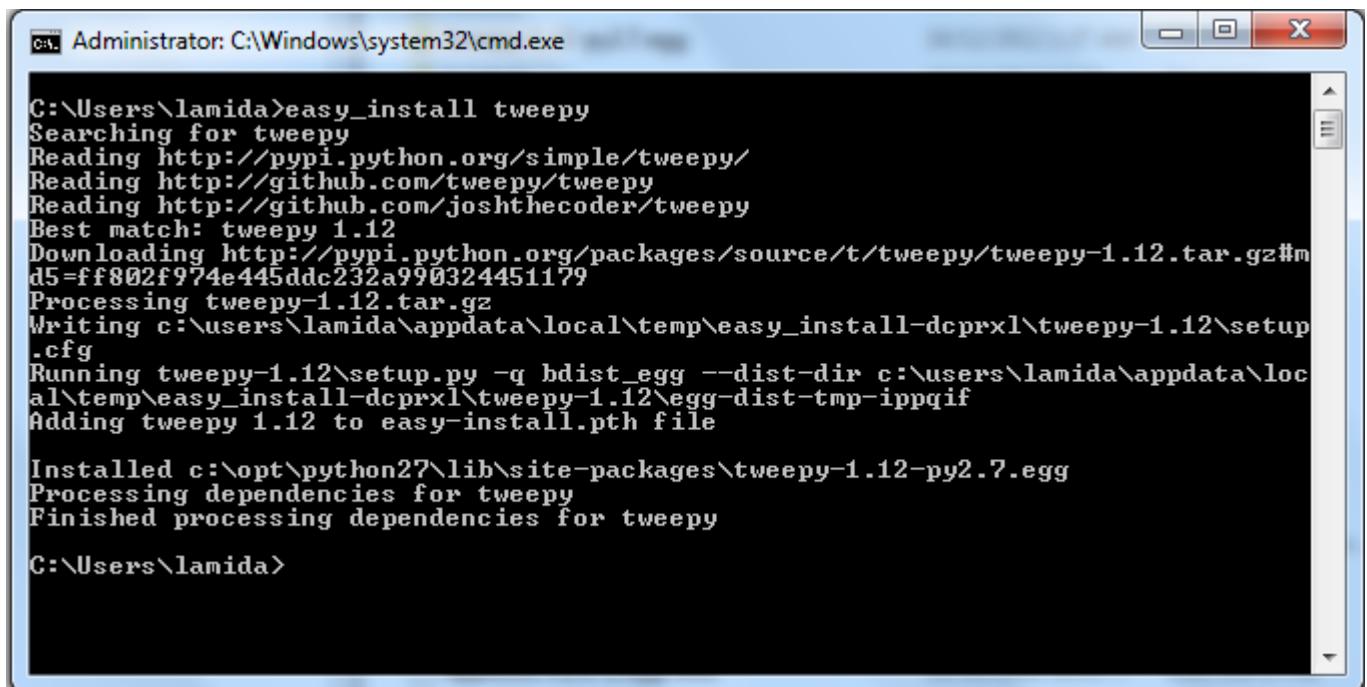
Organization	None
Organization website	None

**Application Settings**  
Your application's Consumer Key and Secret are used to authenticate requests to the Twitter Platform.

Access level	Read-only ( <a href="#">modify app permissions</a> )
Consumer Key (API Key)	NURjid2GRDeVhhz0HuulmsyAk ( <a href="#">manage keys and access tokens</a> )
Callback URL	None
Sign in with Twitter	No
App-only authentication	<a href="https://api.twitter.com/oauth2/token">https://api.twitter.com/oauth2/token</a>
Request token URL	<a href="https://api.twitter.com/oauth/request_token">https://api.twitter.com/oauth/request_token</a>

## 5.2 Connecting to Twitter Streaming API and downloading data

1. Install the tweepy python for importing the necessary methods from twitter library.

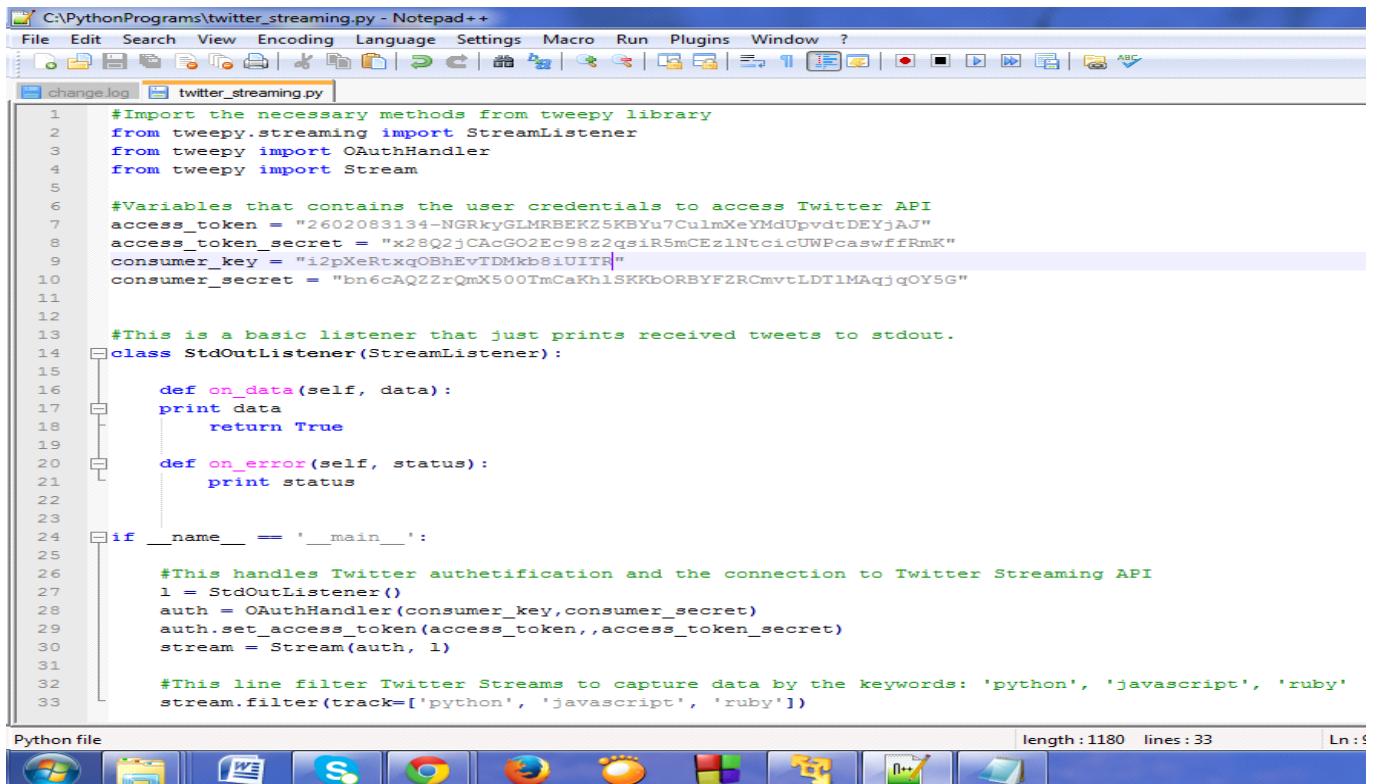


```
C:\Users\lamida>easy_install tweepy
Searching for tweepy
Reading http://pypi.python.org/simple/tweepy/
Reading http://github.com/tweepy/tweepy
Reading http://github.com/joshtechencoder/tweepy
Best match: tweepy 1.12
Downloading http://pypi.python.org/packages/source/t/tweepy/tweepy-1.12.tar.gz#md5=ff802f974e445ddc232a990324451179
Processing tweepy-1.12.tar.gz
Writing c:\users\lamida\appdata\local\temp\easy_install-dcpnx1\tweepy-1.12\setup.cfg
Running tweepy-1.12\setup.py -q bdist_egg --dist-dir c:\users\lamida\appdata\local\temp\easy_install-dcpnx1\tweepy-1.12\egg-dist-tmp-ipqqif
Adding tweepy 1.12 to easy-install.pth file

Installed c:\opt\python27\lib\site-packages\tweepy-1.12-py2.7.egg
Processing dependencies for tweepy
Finished processing dependencies for tweepy

C:\Users\lamida>
```

2. Implement the configuration keys on the notepad and save the file in the config folder:



```
C:\PythonPrograms\twitter_streaming.py - Notepad++
File Edit Search View Encoding Language Settings Macro Run Plugins Window ?
File change.log twitter_streaming.py
1 #Import the necessary methods from tweepy library
2 from tweepy.streaming import StreamListener
3 from tweepy import OAuthHandler
4 from tweepy import Stream
5
6 #Variables that contains the user credentials to access Twitter API
7 access_token = "2602083134-NGRkyGLMRBEKZ5KBYu7CulmKeYMdUpvdtDEYjAJ"
8 access_token_secret = "x28Q2jCACGO2Ec98z2qsiR5mCEzlNtcicUWPcaswffRmK"
9 consumer_key = "i2pXeRtxqOBhEvTDMkb8iUITR"
10 consumer_secret = "bn6cAQZZrQmX500TmCaKh1SKKbORBYFZRCmvtLDTlMAqj qOYSG"
11
12
13 #This is a basic listener that just prints received tweets to stdout.
14 class StdOutListener(StreamListener):
15
16     def on_data(self, data):
17         print data
18         return True
19
20     def on_error(self, status):
21         print status
22
23
24 if __name__ == '__main__':
25
26     #This handles Twitter authetification and the connection to Twitter Streaming API
27     l = StdOutListener()
28     auth = OAuthHandler(consumer_key,consumer_secret)
29     auth.set_access_token(access_token,access_token_secret)
30     stream = Stream(auth, l)
31
32     #This line filter Twitter Streams to capture data by the keywords: 'python', 'javascript', 'ruby'
33     stream.filter(track=['python', 'javascript', 'ruby'])

Python file length:1180 lines:33 Ln:5

```

3. Execute the following command on the terminal and saving the extracted tweets on data.txt :

**“ python twitterstreamingapi.py > data.txt ”**

4. Getting the extracted data in json format

```

cmd C:\Windows\system32\cmd.exe - python setup.py install
and 0.0. It is recommend to migrate to PEP 440 compatible versions.
PEP 440Warning: 'requests' <oauthlib-0.4.1> is being parsed as a legacy, non PEP 440, version. You
may find odd behavior and sort order. In particular it will be sorted as less th
an 0.0.1. It is recommend to migrate to PEP 440 compatible versions.
PEP 440Warning: 'requests' <oauthlib-0.4.2> is being parsed as a legacy, non PEP 440, version. You
may find odd behavior and sort order. In particular it will be sorted as less th
an 0.0.1. It is recommend to migrate to PEP 440 compatible versions.
PEP 440Warning:
Best match: requests-oauthlib 0.4.2
Downloading https://pypi.python.org/packages/source/r/requests-oauthlib/requests-
oauthlib-0.4.2.tar.gz#md5=930be3971f2118c67a8545d54661
Processing requests-oauthlib-0.4.2.tar.gz
Writing C:\Users\DELL\AppData\Local\Temp\easy_install-148wp6ug\requests-oauthlib
tar.gz
Running requests-oauthlib-0.4.2\setup.py -q bdist_egg --dist-dir C:\Users\DELL\AppData\Loca
al\Temp\easy_install-148wp6ug\requests-oauthlib-0.4.2\egg-dist-tmp-4ps
exqf3
creating c:\python34\lib\site-packages\requests_oauthlib-0.4.2-py3.4.egg
Extracting requests_oauthlib-0.4.2-py3.4.egg to c:\python34\lib\site-packages
Adding requests-oauthlib 0.4.2 to easy-install.pth file
Installed c:\python34\lib\site-packages\requests_oauthlib-0.4.2-py3.4.egg
Searching for requests>=2.4.3
Reading https://pypi.python.org/simple/requests/
Best match: requests 2.6.0
Downloading https://pypi.python.org/packages/source/r/requests/requests-2.6.0.t
ar.gz#md5=2522778fa306e207461112bb37656
Processing requests-2.6.0.tar.gz
Writing C:\Users\DELL\AppData\Local\Temp\easy_install-w_xggpw\requests-2.6.0\se
tup.cfg
Running requests-2.6.0\setup.py -q bdist_egg --dist-dir C:\Users\DELL\AppData\Lo
cal\Temp\easy_install-w_xggpw\requests-2.6.0\egg-dist-tmp-8aw9tn7
creating c:\python34\lib\site-packages\requests-2.6.0-py3.4.egg
Extracting requests-2.6.0-py3.4.egg to c:\python34\lib\site-packages
Adding requests-2.6.0 to easy-install.pth file
Installed c:\python34\lib\site-packages\requests-2.6.0-py3.4.egg
Searching for oauthlib<0.6.2
Reading https://pypi.python.org/simple/oauthlib/
Best match: oauthlib 0.7.2
Downloading https://pypi.python.org/packages/source/o/oauthlib/oauthlib-0.7.2.t
ar.gz#md5=9e029908dd48707b79a1cc92
Processing oauthlib-0.7.2.tar.gz
Writing C:\Users\DELL\AppData\Local\Temp\easy_install-jmx3dy01\oauthlib-0.7.2\se
tup.cfg
Running oauthlib-0.7.2\setup.py -q bdist_egg --dist-dir C:\Users\DELL\AppData\Lo
cal\Temp\easy_install-jmx3dy01\oauthlib-0.7.2\egg-dist-tmp-4lv3gv49

```

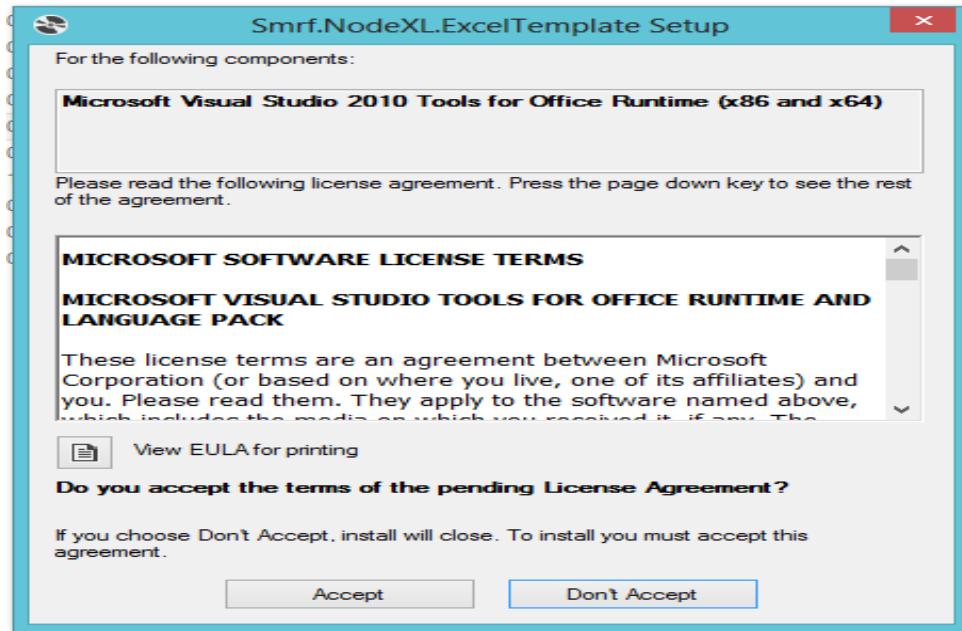
5. Run the program for 2 days (from 2015/01/15 till 2015/01/17) to get a meaningful data sample. This file size is 242 MB.

6. Standardize the extracted data in excel sheets format from [www.json-xls.com/](http://www.json-xls.com/).

### 5.3 NodeXL Installation to Import Tweets

NodeXL is a general purpose network analysis application that supports network overview, discovery and exploration. The tool enables the automation of a data flow that starts with the collection of network data and moves through multiple steps until final processed network visualizations and reports are generated.

1. Download the NodeXL setup from the
2. Run the NodeXL set up and accept the license.



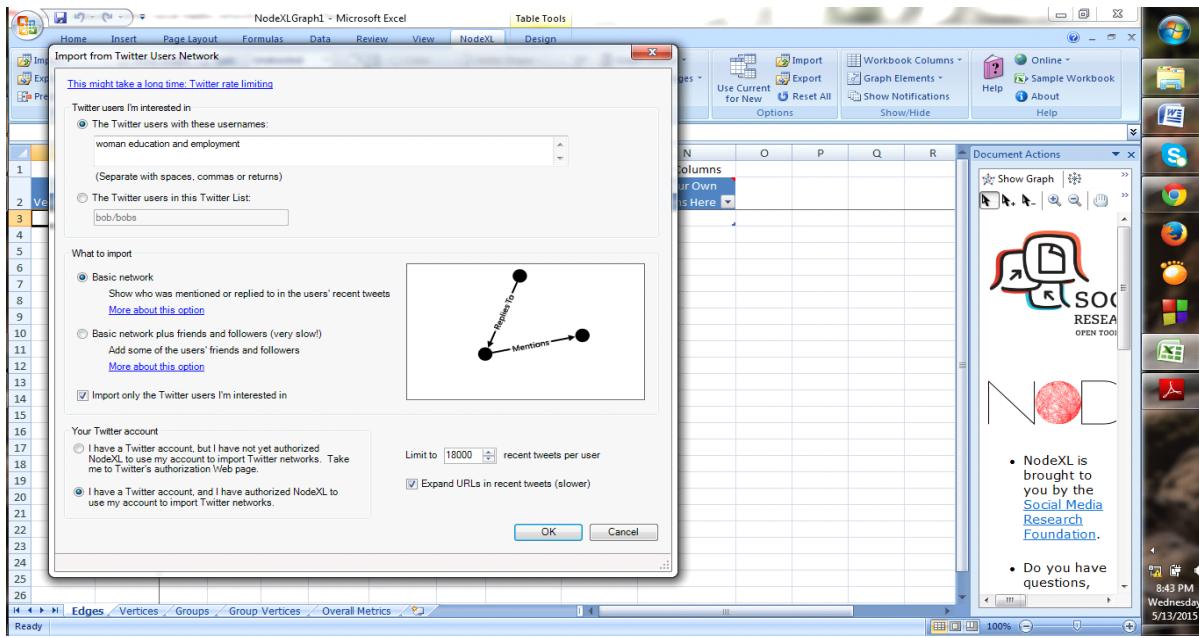
## 5.4 Mining Tweets in NodeXL

1. Open the NodeXL Sheets.

The screenshot shows the Microsoft Excel interface with the title bar 'NodeXLGraph1 - Microsoft Excel'. The ribbon tabs include Home, Insert, Page Layout, Formulas, Data, Review, View, NodeXL, and Design. The 'Table Tools' tab is selected. A tooltip is visible over cell A3, containing the text: 'Vertex 1 Name Enter the name of the edge's first vertex.' The right side of the screen features the 'Document Actions' pane, which includes icons for Show Graph, NodeXL Help, and various social media links. The bottom status bar shows the date and time: '8:27 PM Wednesday 5/13/2015'.

2. Import the standarized tweets in NodeXL excel sheet.

We used the NodeXL Twitter data import feature to extract networks. We configured our query using settings. We requested that NodeXL also add an “edge” which describes the connection between two twitter users that is formed when they follow, reply or mention one another. Data about each user along with the contents of their latest tweet were also selected to be added to the data set.



### 3. Getting tweets on NodeXL sheets.

	A	G	H	I	J	K	L	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM
1				Labels				Other Columns									
2	Vertex	File	Visibility	Label	Color	Position		Add Your Own Columns Here	Followed	Followers	Tweets	Favorites	Offset (Seconds)	Description	Time Zone UTC	Location	Web Zone
3	robertfrau	<a href="http://pbs.twimg.com/profile_images/598301624936">http://pbs.twimg.com/profile_images/598301624936</a>		robertfrau					699	727	49932	13236	-28800	Real Estate	Los Angeles, CA	Alask	
4	justonere	<a href="http://pbs.twimg.com/profile_images/598387539045">http://pbs.twimg.com/profile_images/598387539045</a>		justonere					103	122	4080	32	-25200	@Colorseatsbutt	▼	Tiju	
5	theregoki	<a href="http://pbs.twimg.com/profile_images/597089636361">http://pbs.twimg.com/profile_images/597089636361</a>		theregoki					3268	4106	68445	454	-18000	idc.		Quic	
6	fuckdeser	<a href="http://pbs.twimg.com/profile_images/96662723305">http://pbs.twimg.com/profile_images/96662723305</a>		fuckdeser					1950	2305	72103	50897	-18000		Englewood	Chicago, Centr	
7	drpatfarre	<a href="http://pbs.twimg.com/profile_images/341786372457">http://pbs.twimg.com/profile_images/341786372457</a>		drpatfarre					3749	4223	133763	27	-14400	Patricia Fan USA	<a href="http://t.co/...">http://t.co/...</a>	Easte	
8	drmerle	<a href="http://pbs.twimg.com/profile_images/533409622497">http://pbs.twimg.com/profile_images/533409622497</a>		drmerle@					1256	1264	20931	5028	-18000	Politics,Rigl Vaughan, O	<a href="http://t.co/...">http://t.co/...</a>	Centr	
9	larryasler	<a href="http://pbs.twimg.com/profile_images/582943445177">http://pbs.twimg.com/profile_images/582943445177</a>		larryasler					35	29	318	32			Libertarian, Toronto Marxist	WASTELA	
10	jikeriaaa	<a href="http://pbs.twimg.com/profile_images/589259290122">http://pbs.twimg.com/profile_images/589259290122</a>		jikeriaaa					4077	4796	78275	8375	-10800	I'm the rose that grew from conc	Atlanta		

Each “**edge**” represents a connection event between two people who tweeted within the data sample period. NodeXL constructs four different types of Twitter edges from the data it collects: follows, replies, mentions and tweet.

A “**follows**” edge is created if one author follows another who also tweeted in the sample dataset (the time stamp for a follows edge is the date of the query rather than the time when one user followed another user, which is information that is not available from Twitter).

A “**mentions**” edge is created when one user creates a tweet that contains the name of another user (indicated with a preceding “@” character, ex: “just spoke about social media with @marc\_smith”).

A “**reply**” relationship is a special form of “mention” that occurs when the user’s name is at the very start of a tweet (ex: “@itaih just spoke about social media”).

A tweet is a message that does not contain a reply or mention.

**4. Automate the NodeXL :** Using “Automate” NodeXL executes a series of user configured operations on the network without direct user control. The Automate dialog provides a good summary of the steps and operations applied to each network graph.

The screenshot shows a Microsoft Excel spreadsheet titled "NodeXLGraph1 - Microsoft Excel". The ribbon tabs are Home, Insert, Page Layout, Formulas, Data, Review, View, NodeXL, and Design. The NodeXL tab is selected. A sub-menu for "Automate" is open, displaying a list of tasks:

- Count and merge duplicate edges
- Show Graph
- Summary
- Graph Metrics
- Autofill columns
- Subgraph Images
- Show graph
- Save workbook to a new file if it has never been saved
- Save image to file
- Export graph to NodeXL Graph Gallery
- Export graph to email

Below the list are two radio buttons: "On this workbook" (selected) and "On every NodeXL workbook in this folder". There is also a "Browse..." button and a note "(Excludes open workbooks)". At the bottom are "Run" and "Cancel" buttons.

The main Excel window displays a table of network statistics. The columns are labeled K, L, AD, AE, AF, AG, AH, AI, and Other Columns. The "Other Columns" section includes Position, Tooltip, Ad, Followed, Followers, Tweets, Favorites, and Time Zone UTC. The "Time Zone UTC" column shows values such as -28800, -25200, -18000, -14400, -10800, and -10800. The "Followed" column shows counts like 699, 103, 3268, 1950, 3749, 1256, 35, 4077, etc. The "Followers" column shows counts like 727, 122, 4080, 2305, 4223, 1264, 29, 4796, etc. The "Tweets" column shows counts like 49932, 32, 454, 50897, 27, 5028, 32, 8375, etc. The "Favorites" column shows counts like 13236, -28800, -25200, -18000, -14400, -10800, -10800, etc. The "Time Zone UTC" column shows values like -28800, -25200, -18000, -14400, -10800, -10800, -10800, -10800, etc.

The screenshot shows a Microsoft Excel spreadsheet titled "NodeXLGraph1 - Microsoft Excel". The ribbon is visible at the top with tabs for Home, Insert, Page Layout, Formulas, Data, Review, View, NodeXL, and Design. The Design tab is currently selected. The main content area contains a table with several columns. The first column is labeled "Vertex" and contains names such as "robertfrau", "justonered", "theregokii", "fuckdeseh", "drpatfarre", "drmerle", "larryasler", and "jikeriaaa". The second column is labeled "Visual Properties" and includes dropdown menus for Color, Shape, Size, Opacity, and Visibility. The third column is labeled "Labels" and contains URLs for profile images. Subsequent columns represent network metrics: "Ad", "AE", "AF", "AG", "AH", "AI", "Followed", "Followers", "Tweets", "Favorites", and "Time Zone UTC Offset (Seconds)". Below the table, there are tabs for Edges, Vertices, Groups, Group Vertices, and Overall Metrics. The status bar at the bottom indicates "Count: 43" and the date "5/13/2015".

**5. Graph Metrics :** Each of the network metrics captures a different dimension of the size and shape of the graph as a whole and the location and connection properties of each person or entity in the network graph. We selected the creation all of the network metrics available through NodeXL.

Many of these metrics can be mapped to various network display attributes. For example, the size of a vertex representing a Twitter user can be scaled to represent the number of users who have chosen to Follow each user.

The screenshot shows a Microsoft Excel spreadsheet titled "NodeXLGraph1 - Microsoft Excel". The ribbon menu includes Home, Insert, Page Layout, Formulas, Data, Review, View, NodeXL, and Design. The NodeXL tab is selected.

The main area displays a table with columns for Visual Properties, Labels, and Graph Metrics. The Graph Metrics column contains a complex network diagram with various nodes and edges, color-coded by node properties.

	A	B	C	D	E	F	K	L	M	N	O	P	Q	R	
1	Visual Properties				Labels		Graph Metrics								
2	Group	Vertex Color	Vertex Shape	Visibility	Collapsed?	Label	Vertices	Unique Edges	Edges With Duplicates	Total Edges	Self-Loops	Reciprocated	Reciprocated Edge Ratio	Connected Components	Single-Value Components
3	G1	0, 12, 96	Disk					30	29	0	29	0	0.000		
4	G2	0, 136, 227	Disk					14	13	2	15	2	0.000		
5	G3	0, 100, 50	Disk					7	26	0	26	0	0.300		
6	G4	0, 176, 22	Disk					6	11	0	11	0	0.000		
7	G5	191, 0, 0	Disk					5	5	0	5	1	0.000		
8	G6	230, 120, 0	Disk					5	7	0	7	0	0.167		
9	G7	255, 191, 0	Disk					4	3	0	3	0	0.000		
10	G8	150, 200, 0	Disk					3	2	0	2	0	0.000		
11	G9	200, 0, 120	Disk					3	3	0	3	1	0.000		
12	G10	77, 0, 96	Disk					3	2	0	2	0	0.000		
13	G11	91, 0, 191	Disk					3	4	0	4	0	0.333		
14	G12	0, 98, 130	Disk					3	2	0	2	0	0.000		
15	G13	0, 12, 96	Solid Square					3	2	0	2	0	0.000		
16	G14	0, 136, 227	Solid Square					3	2	0	2	0	0.000		
17	G15	0, 100, 50	Solid Square					2	1	0	1	0	0.000		
18	G16	0, 176, 22	Solid Square					2	0	7	7	4	0.000		
19	G17	191, 0, 0	Solid Square					2	1	0	1	0	0.000		
20	G18	230, 120, 0	Solid Square					2	1	0	1	0	0.000		
21	G19	255, 191, 0	Solid Square					2	2	0	2	1	0.000		
22	G20	150, 200, 0	Solid Square					2	1	0	1	0	0.000		
23	G21	200, 0, 120	Solid Square					2	2	0	2	1	0.000		
24	G22	77, 0, 96	Solid Square					2	1	0	1	0	0.000		
25	G23	91, 0, 191	Solid Square					2	1	0	1	0	0.000		
26	G24	0, 98, 130	Solid Square					2	2	0	2	0	1.000		

Below the table, there are tabs for Vertices, Groups, Group Vertices, Overall Metrics, Group Edges, Twitter Search Ntwrk Top Items, and a few others that are partially visible.

### **6. Vertices:**

The Vertex worksheet displays attributes from Twitter along with a range of network metrics for each vertex in the network. Each row represents one Twitter user who appeared in the results of the initial search query. A set of network metrics is listed for each user. Network metrics capture a range of qualities about the location and connection pattern of each user within the larger network.

The screenshot shows a Microsoft Excel spreadsheet titled "NodeXLGraph1.xlsx - Microsoft Excel". The ribbon tabs include Home, Insert, Page Layout, Formulas, Data, Review, View, NodeXL, and Design. The active cell is A2, and the formula bar shows "Vertex".

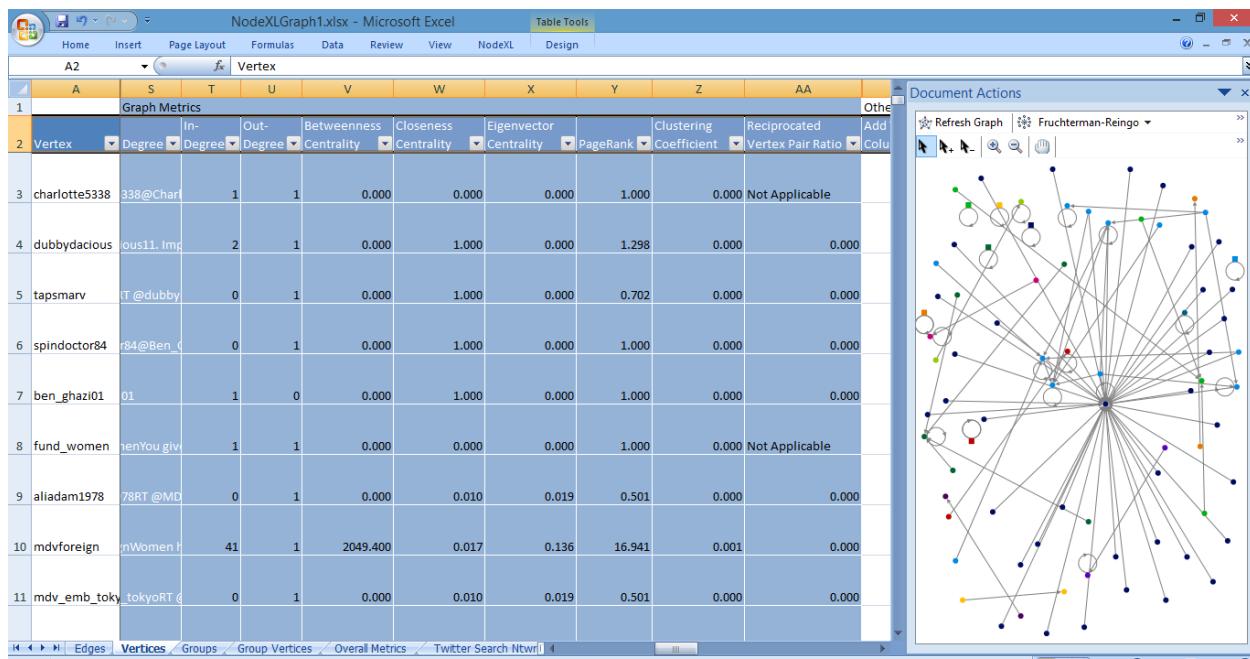
The table has three main sections:

- Visual Properties** (Columns A-F)
- Labels** (Columns G-L)
- Graph Metrics** (Columns M-T)

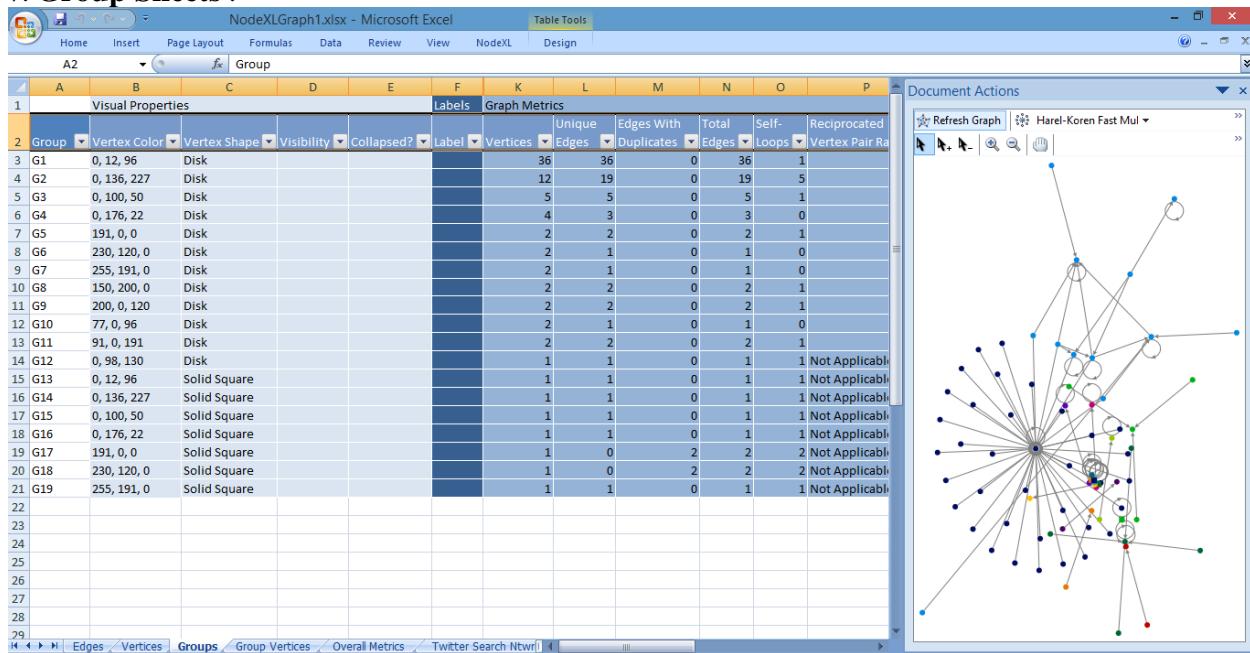
The data rows represent various Twitter users (vertices) and their connections (edges). The first row is a header, and the second row is a detailed example of a vertex entry.

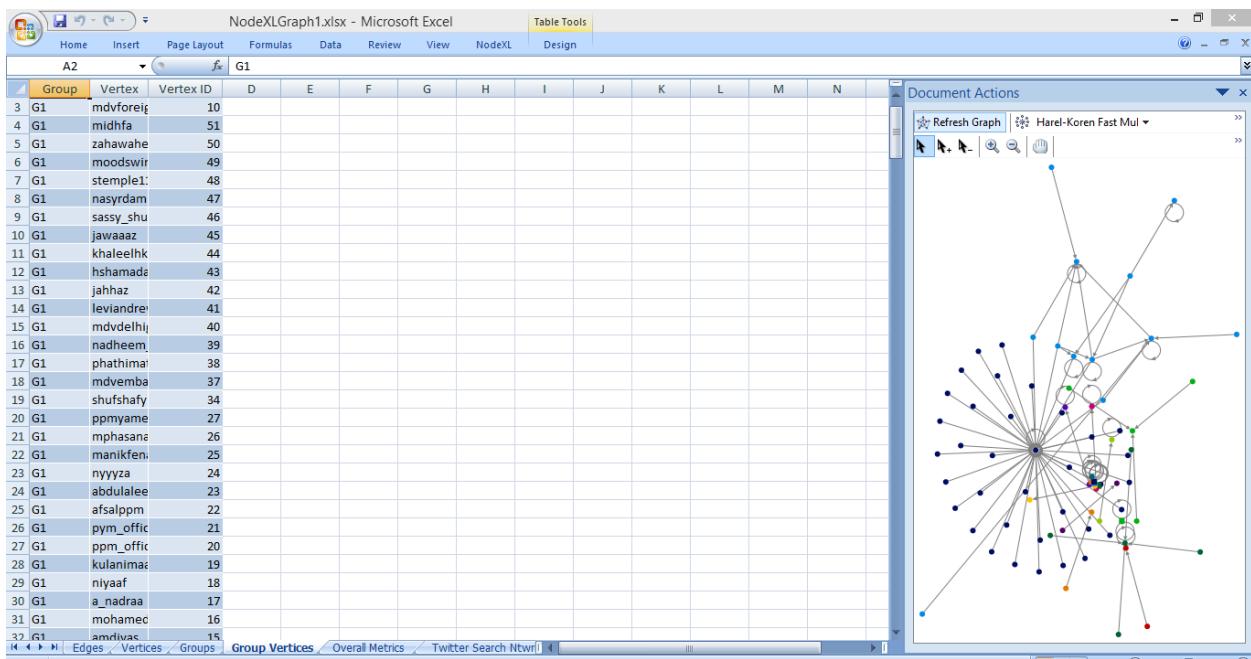
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
1		Visual Properties					Labels					Graph Metrics								
2	Vertex	Subgraph	Color	Shape	Size	Opacity	Image	File	Visibility	Label	Label Fill	Label Color	Position	Tooltip	Degree	In-Degree	Out-Degree			
3	charlotte5338	(circle)								http://pbs.twimg.com/profile_images/567781251256	charlotte5338@Char				1					
4	dubbydacious	(line)								http://pbs.twimg.com/profile_images/594514845507	dubbydacious11.Img				2					
5	tapsmarv	(line)								http://pbs.twimg.com/profile_images/591711062803	tapsmarvRT @duby				0					
6	spindoctor84	(line)								http://pbs.twimg.com/profile_images/483936785/un	spindoctor84@Ben				0					
7	ben_ghazi01	(line)								http://pbs.twimg.com/profile_images/595792820381	ben_ghazi01				1					
8	fund_women	(circle)								http://pbs.twimg.com/profile_images/578922986931	fund_womenYou give				1					
9	aliadam1978	(line)								http://pbs.twimg.com/profile_images/594728278419	aliadam1978RT @MD				0					
10	mdvforeign	(cluster)								http://pbs.twimg.com/profile_images/540754181573	mdvforeignWomen H				41					
11	mdv_emb_toky	(line)								http://pbs.twimg.com/profile_images/3280905712/a	mdv_emb_tokyRT @				0					

On the right side of the screen, there is a "Document Actions" panel with buttons for Refresh Graph and Fruchterman-Reingold, along with a zoom-in/out icon and a search icon. Below the table, there is a network graph visualization showing nodes as circles of varying sizes and colors, connected by lines representing edges.



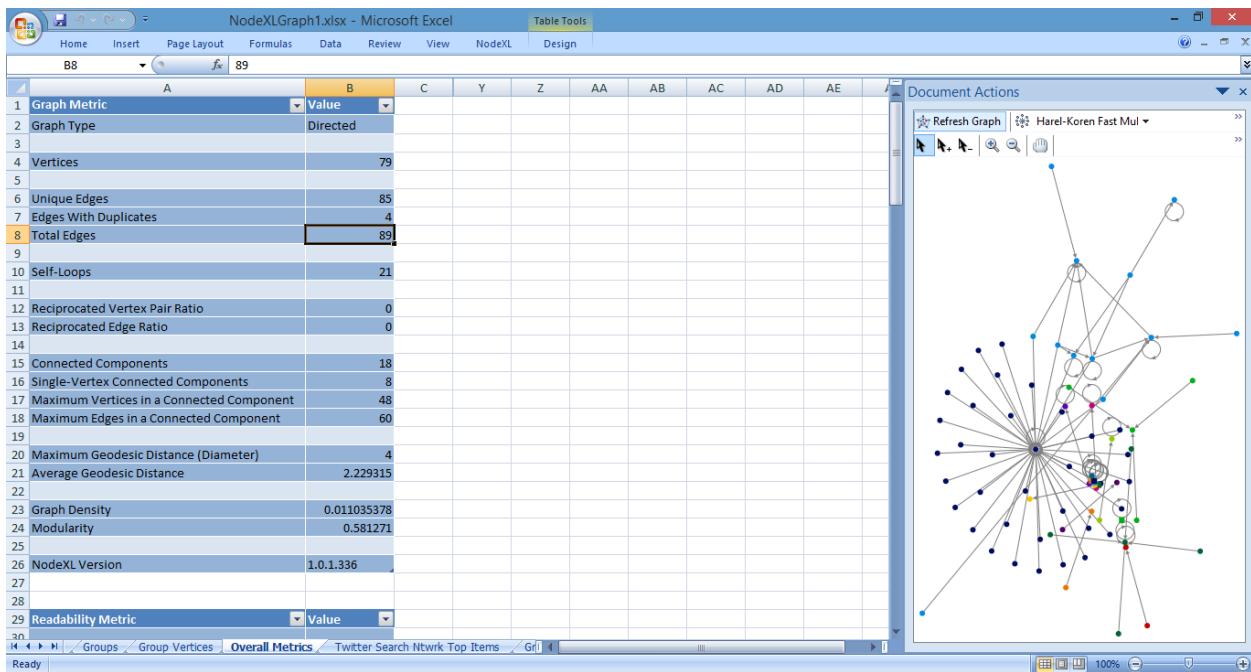
## 7. Group Sheets :

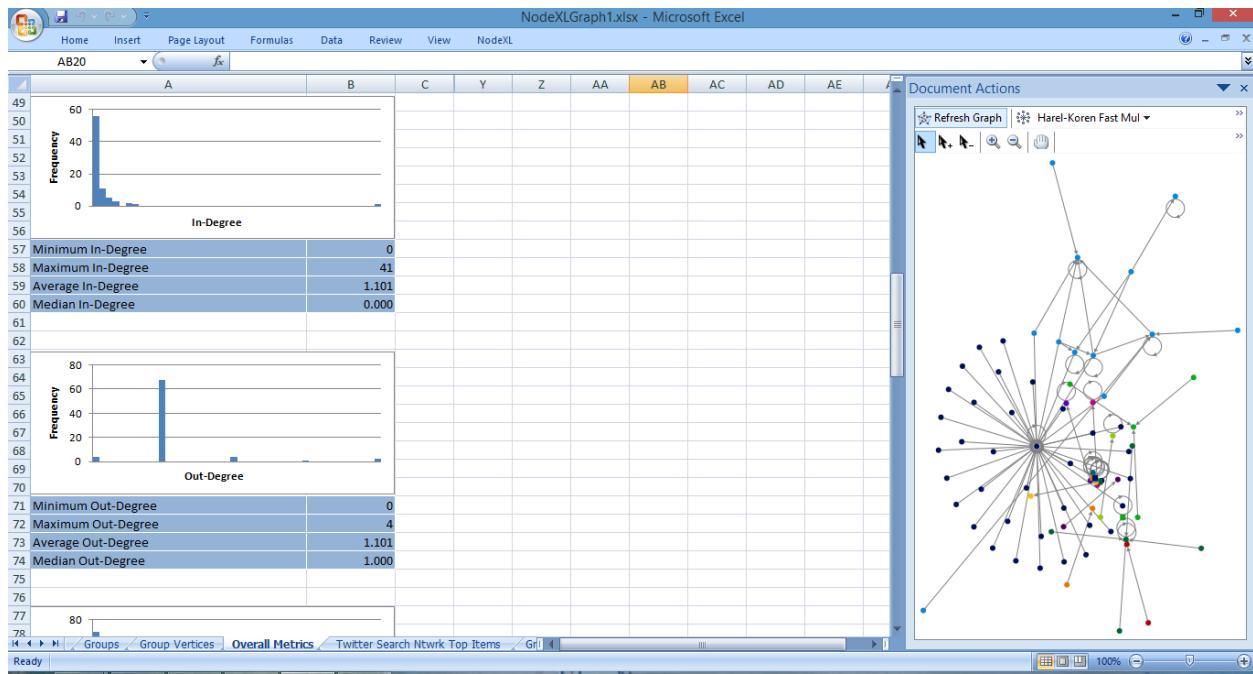




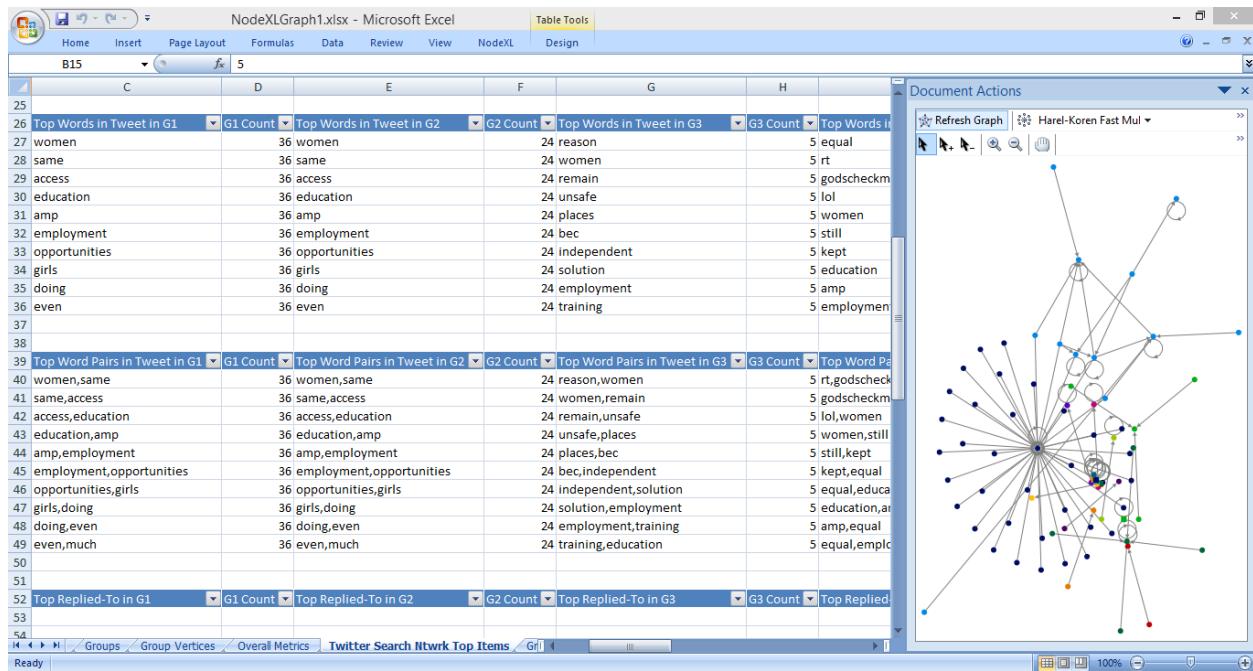
## 8. Overall Metrics Sheet:

These measures allow users to contrast networks to capture the ways networks in general and social media networks in particular compare to other networks and to themselves overtime.

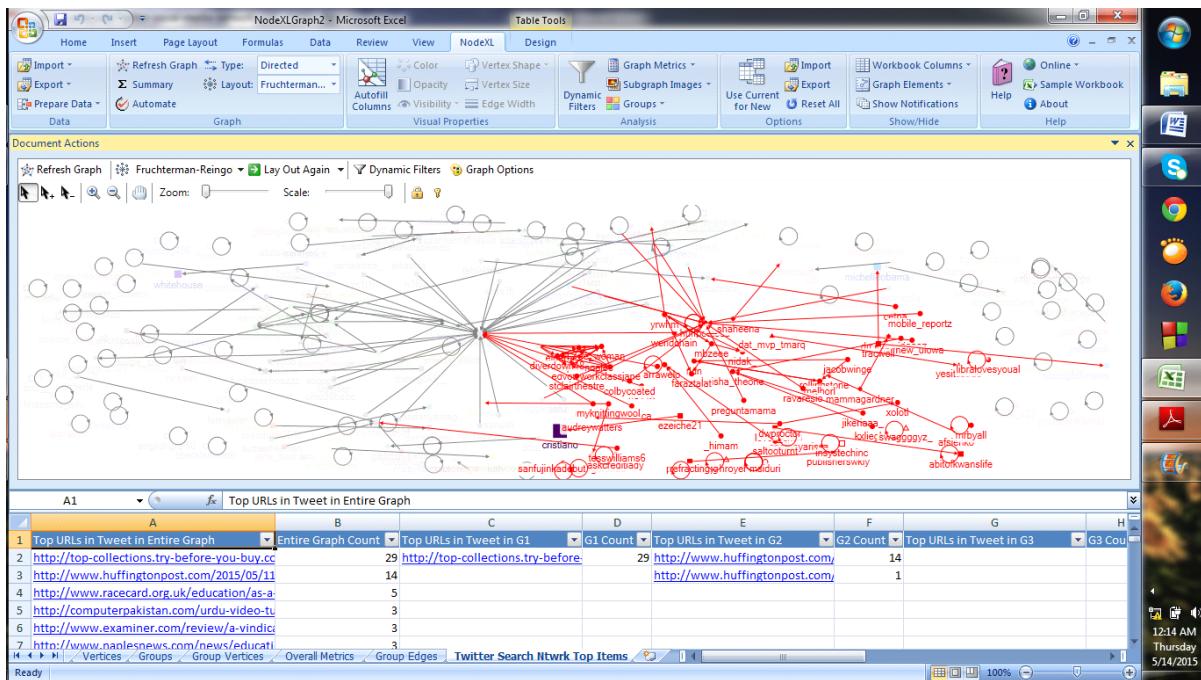




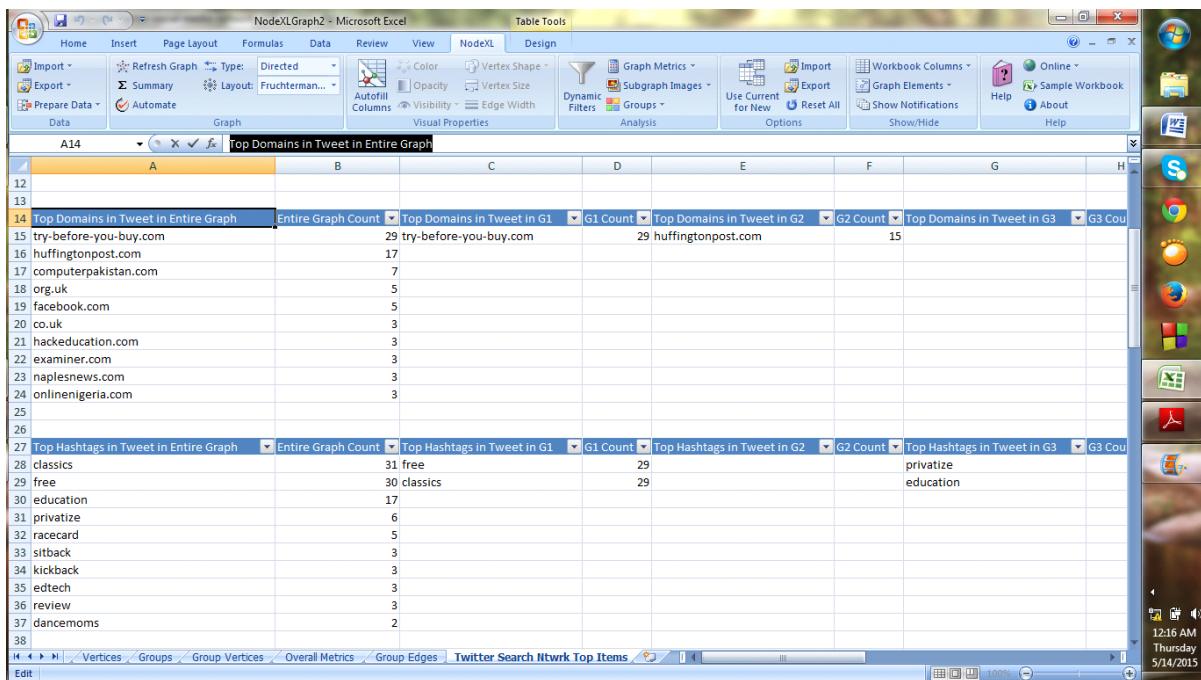
## 9. Twitter Search Network Top Items:



## 10. Top URLs in Tweet in Entire Graph

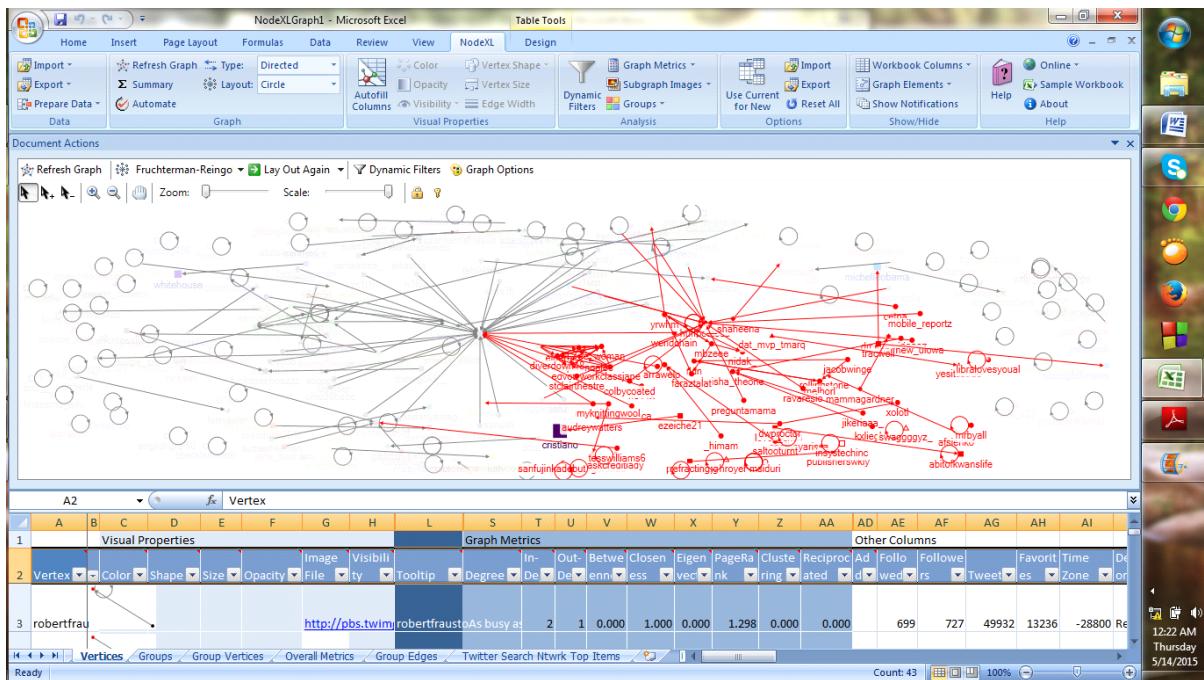


## 11. Top Domains and Hashtag in Tweet in Entire Graph



## 12 .Top Words in Tweets in Entire Graph

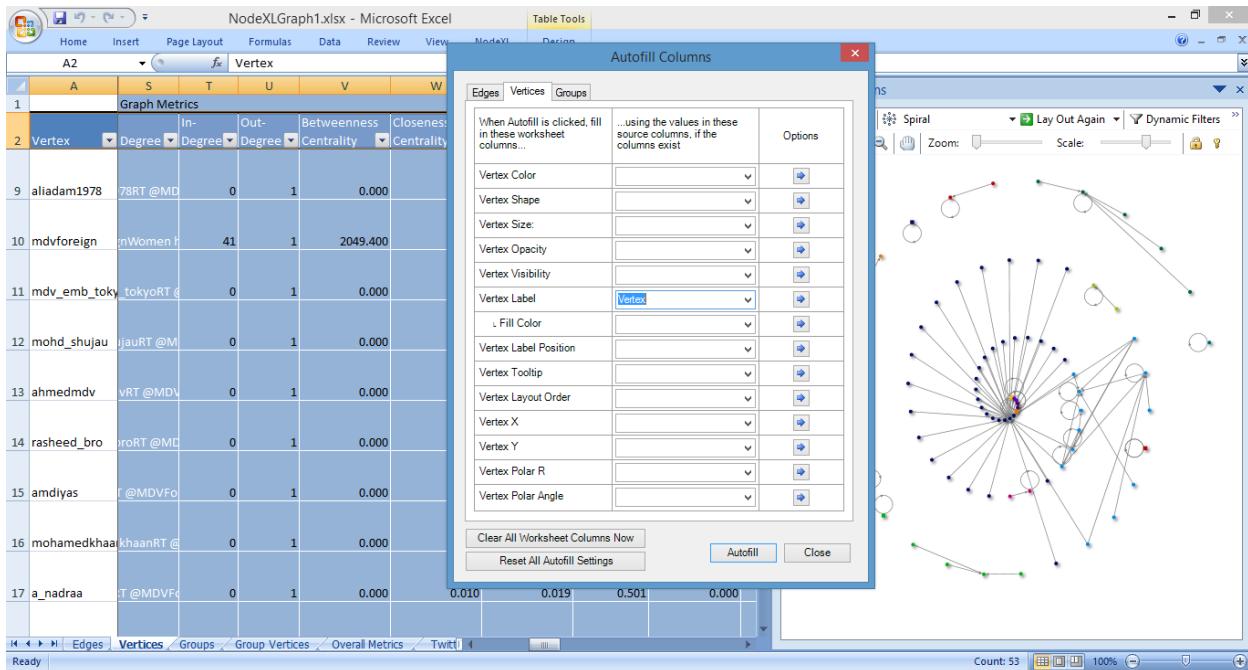
Tweets which are directly related to one common person



### 13. Count and Merge Duplicates

The screenshot shows a Microsoft Excel spreadsheet titled "NodeXLGraph2 - Microsoft Excel". The top ribbon has tabs for Home, Insert, Page Layout, Formulas, Data, Review, View, NodeXL, and Design. The NodeXL tab is selected. The main area displays a large table with many rows and columns. The columns include "Relationship", "Date (UTC)", "Tweet", "URLs in Tweet", "Domains in Tweet", "Hashtags in Tweet", "Twitter Page for Tweet", "Latitude", "Longitude", "Imported", "Tweet ID", and "Edge". The table contains numerous entries, mostly tweets from various users. At the bottom of the table, there are tabs for Edges, Vertices, Groups, Group Vertices, Overall Metrics, Group Edges, Twitter Search Ntwrk Top Items, and a status bar showing "Count: 2" and "100%".

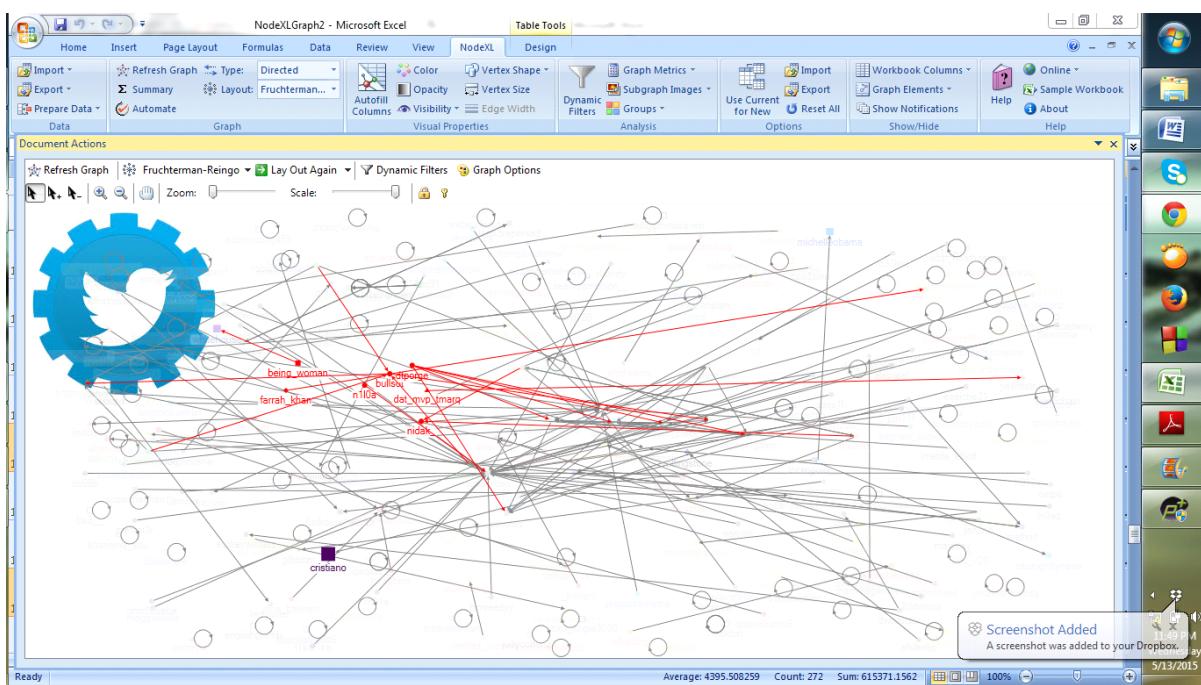
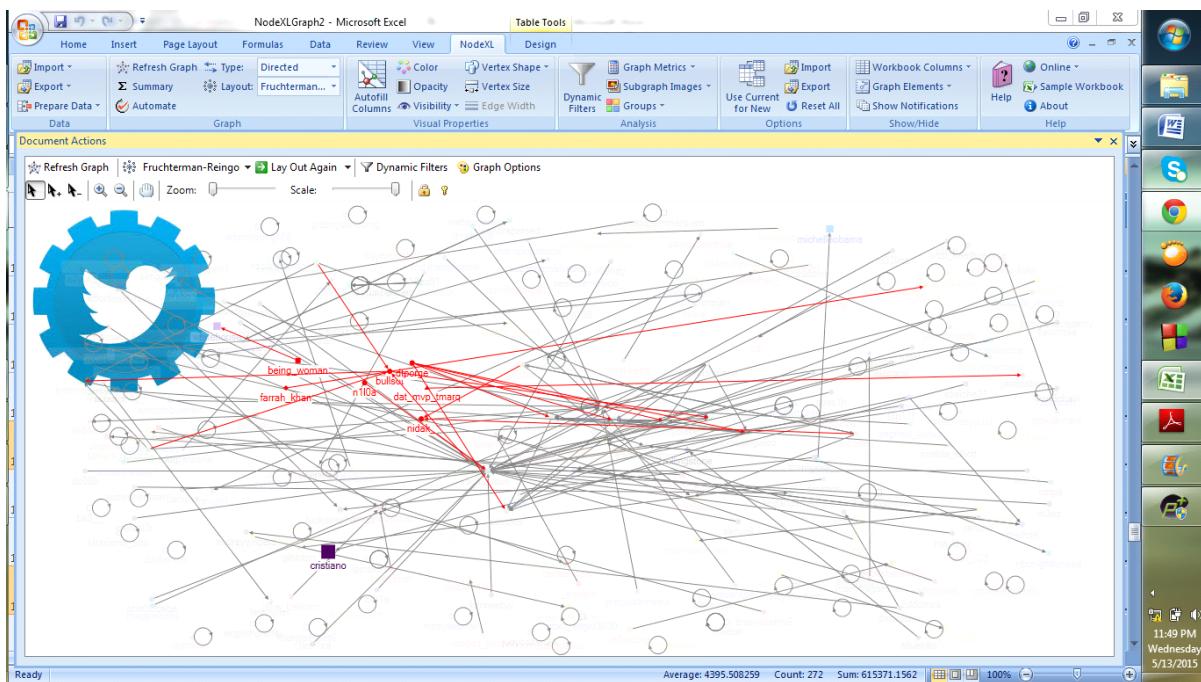
**14. AutoFill Columns :** NodeXL has a feature called “Autofill Columns” that makes it simple to pick attributes about each edge and vertex and map them to display attributes like the size, color, shape, or transparency of each vertex.



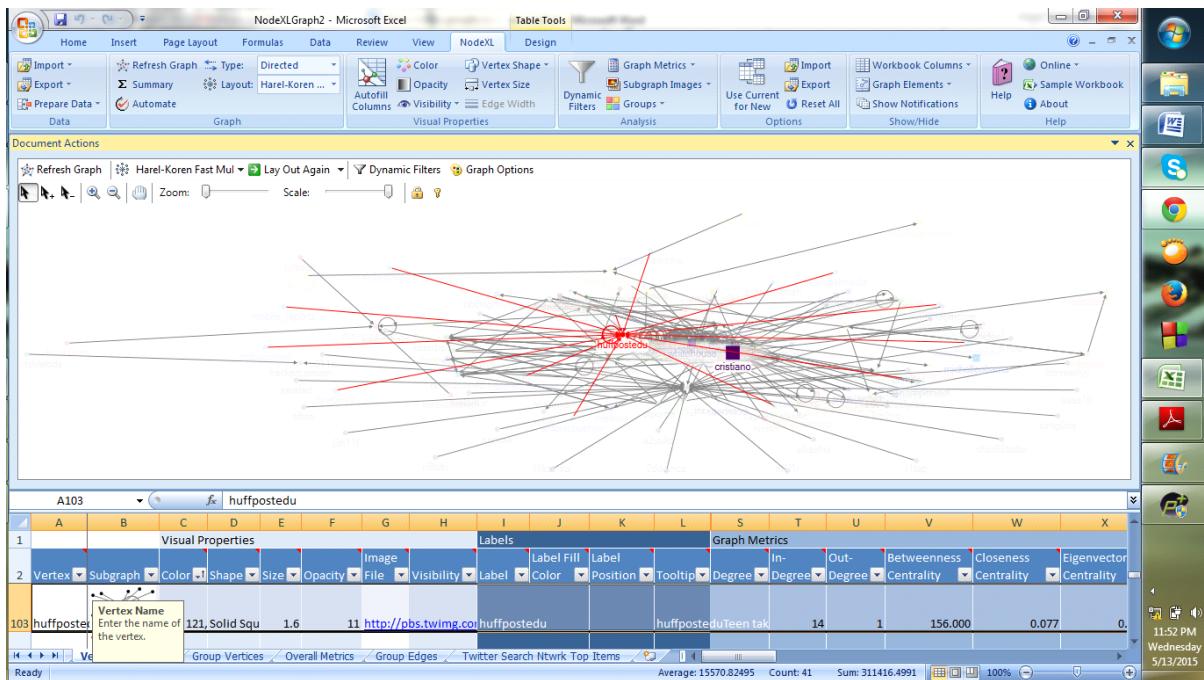
Using Autofill columns, we set the size of each vertex to be proportionate to the number of Followers the user had attracted. We also set the “opacity” of the vertex (which controls the transparency of the image) to be inverse to the number of Followers. This setting is controlled by the options which are accessed by the arrow on each row of the Autofill columns dialog. An inverse mapping for opacity means that the largest objects are somewhat transparent, allowing the other, smaller users to be seen through the images of the larger, more popular users.

The vertex label is set to the name of the user. The vertex label is drawn beneath the profile photo from Twitter that is used to represent each user. The order of the layout controls the way each vertex is placed on the screen. Setting it to equal the same value as the size value means that objects often line up in size order.

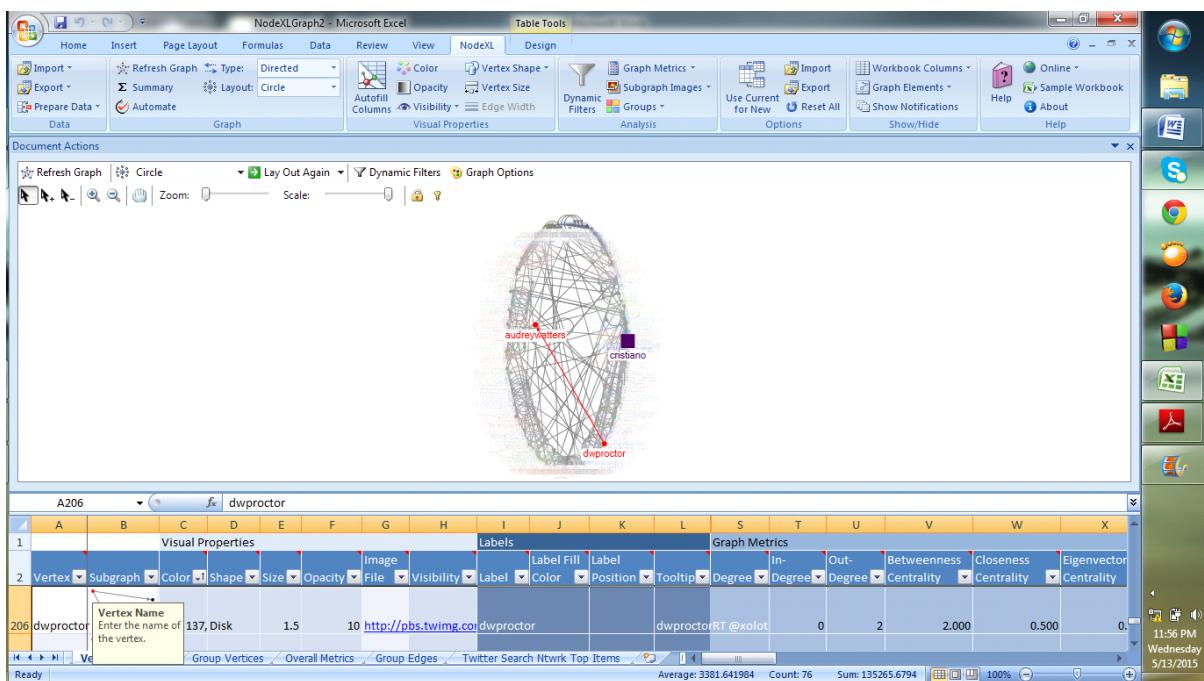
Frusta Graph

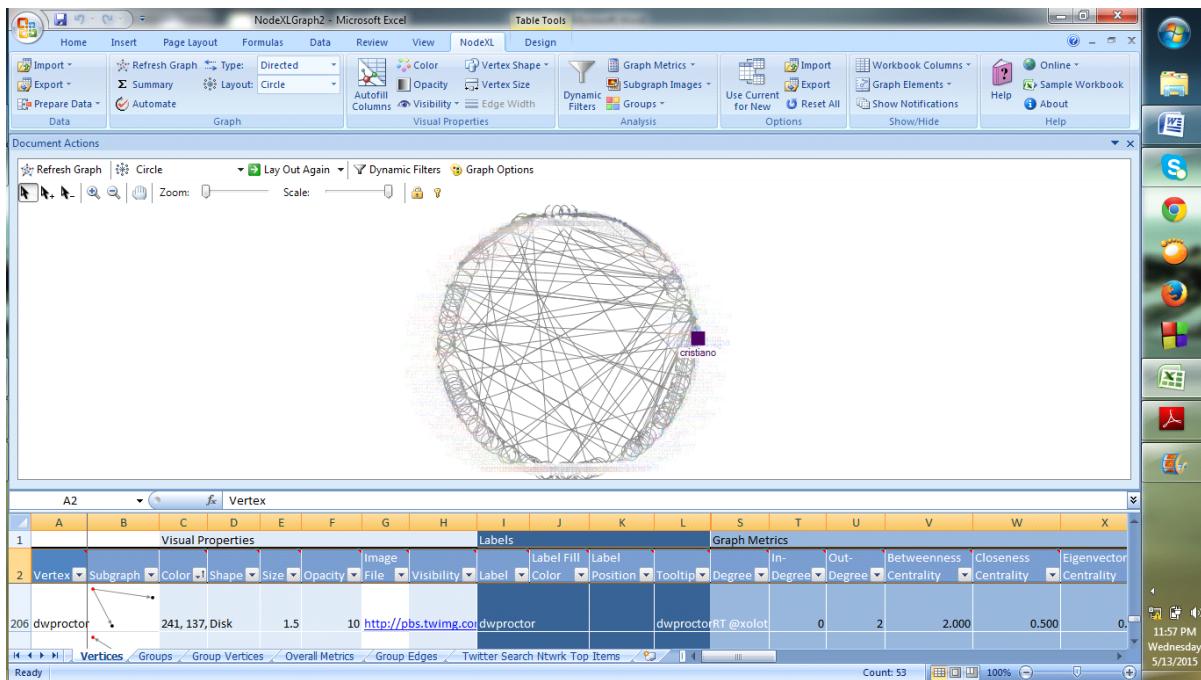


Harel Koren

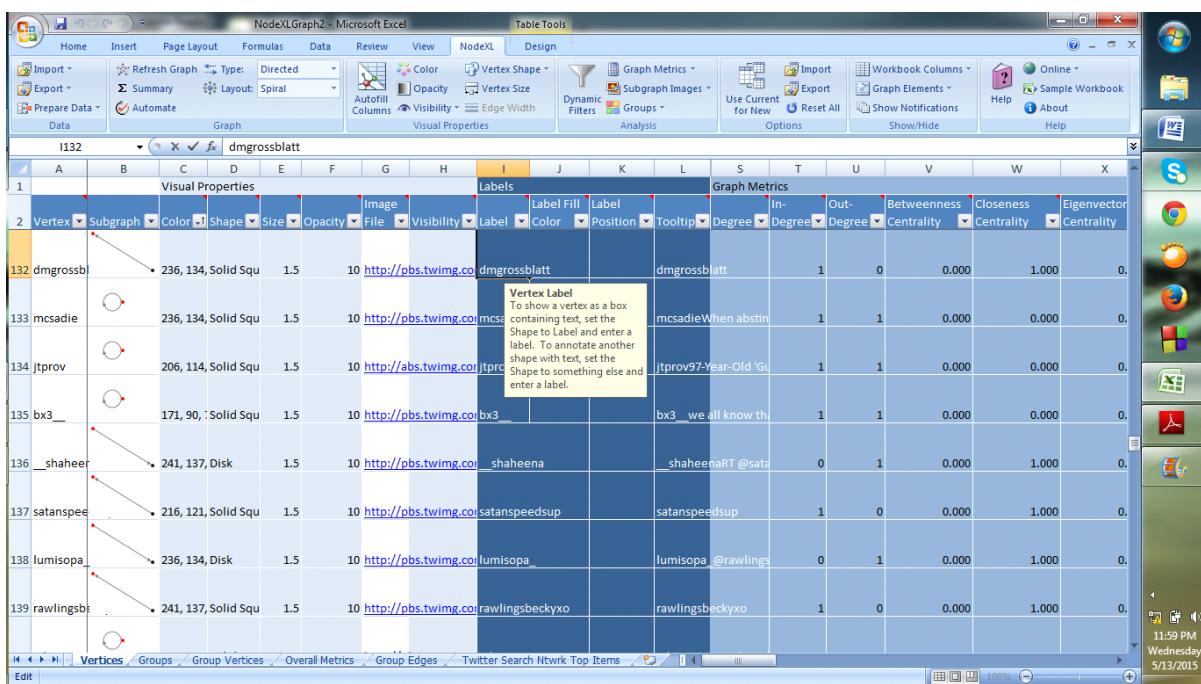


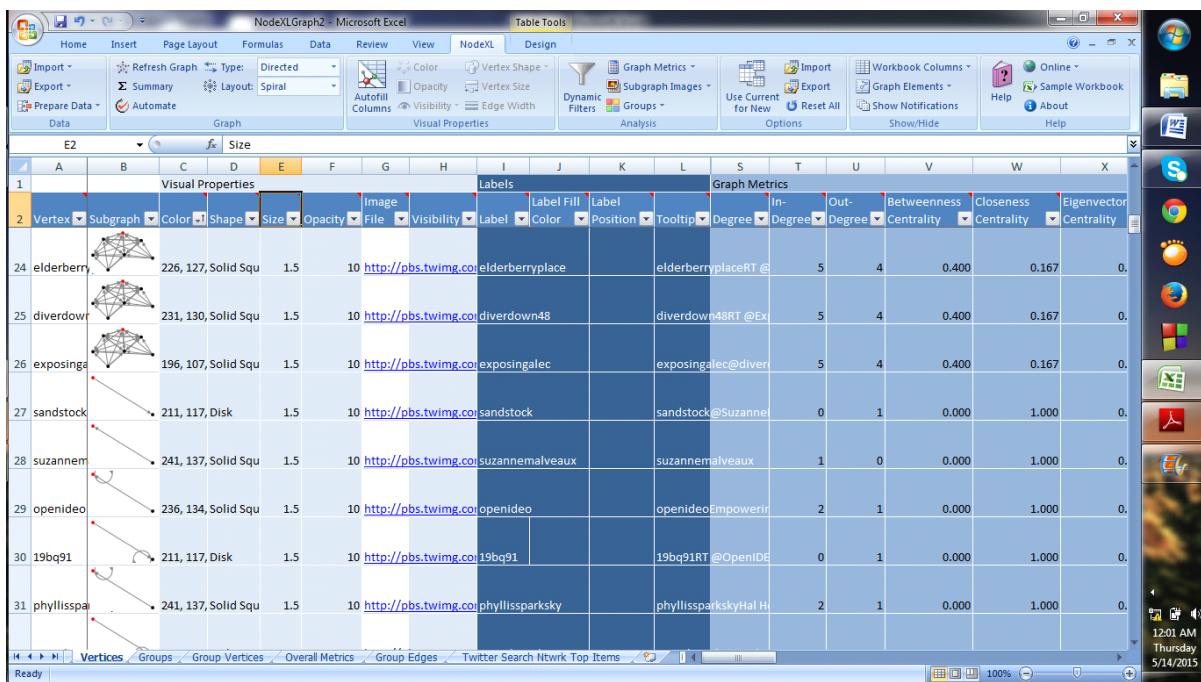
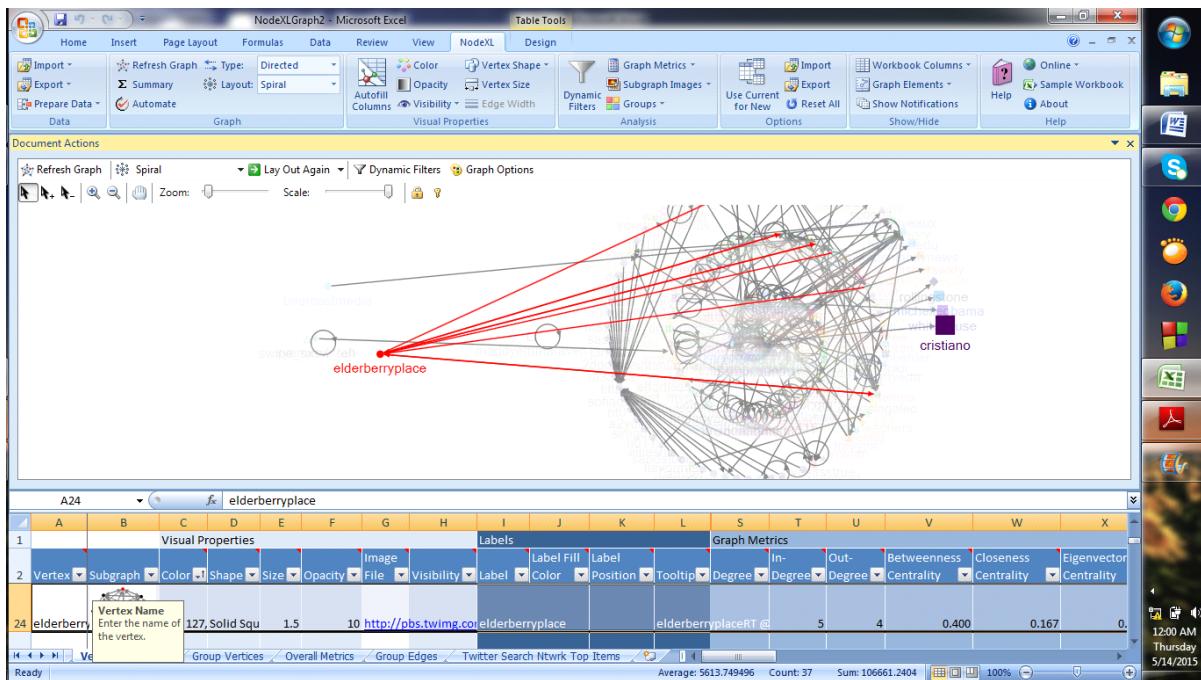
## Circle





## Spiral





Sugiyama

NodeXLGraph2 - Microsoft Excel

Home Insert Page Layout Formulas Data Review View NodeXL Format

Import Export Prepare Data Data Graph

Refresh Graph Type: Directed Layout: Sugiyama

Autofill Columns Visibility Edge Width

Color Vertex Shape Opacity Vertex Size

Dynamic Filters Graph Metrics Subgraph Images Groups Analysis

Use Current for New Import Export Workbook Columns Graph Elements Show Notifications Options Show/Hide Help

Document Actions

Refresh Graph Sugiyama Lay Out Again Dynamic Filters Graph Options

Zoom: Scale: Lock:

Subgraph-elderberry... fx

	A	B	C	D	E	F	G	H	I	J	K	L	S	T	U	V	W	X						
1	Visual Properties								Labels								Graph Metrics							
2	Vertex	Subgraph	Color	Shape	Size	Opacity	Image	Visibility	Label	Label Fill	Label	Position	Tooltip	Degree	In-Degree	Out-Degree	Betweenness	Closeness	Eigenvector					
24	elderberry	226, 127, Solid Squ	1.5	10	<a href="http://pbs.twimg.com">http://pbs.twimg.com</a>	elderberryplace			elderberryplaceRT @		5	4		0.400		0.167	0.							

Vertices Groups Group Vertices Overall Metrics Group Edges Twitter Search Ntwrk Top Items

Ready

NodeXLGraph2 - Microsoft Excel

Home Insert Page Layout Formulas Data Review View NodeXL Design

Import Export Prepare Data Data Graph

Refresh Graph Type: Directed Layout: Sugiyama

Autofill Columns Visibility Edge Width

Color Vertex Shape Opacity Vertex Size

Dynamic Filters Graph Metrics Subgraph Images Groups Analysis

Use Current for New Import Export Workbook Columns Graph Elements Show Notifications Options Show/Hide Help

Document Actions

Refresh Graph Sugiyama Lay Out Again Dynamic Filters Graph Options

Zoom: Scale: Lock:

A199 fx ravenonthedream

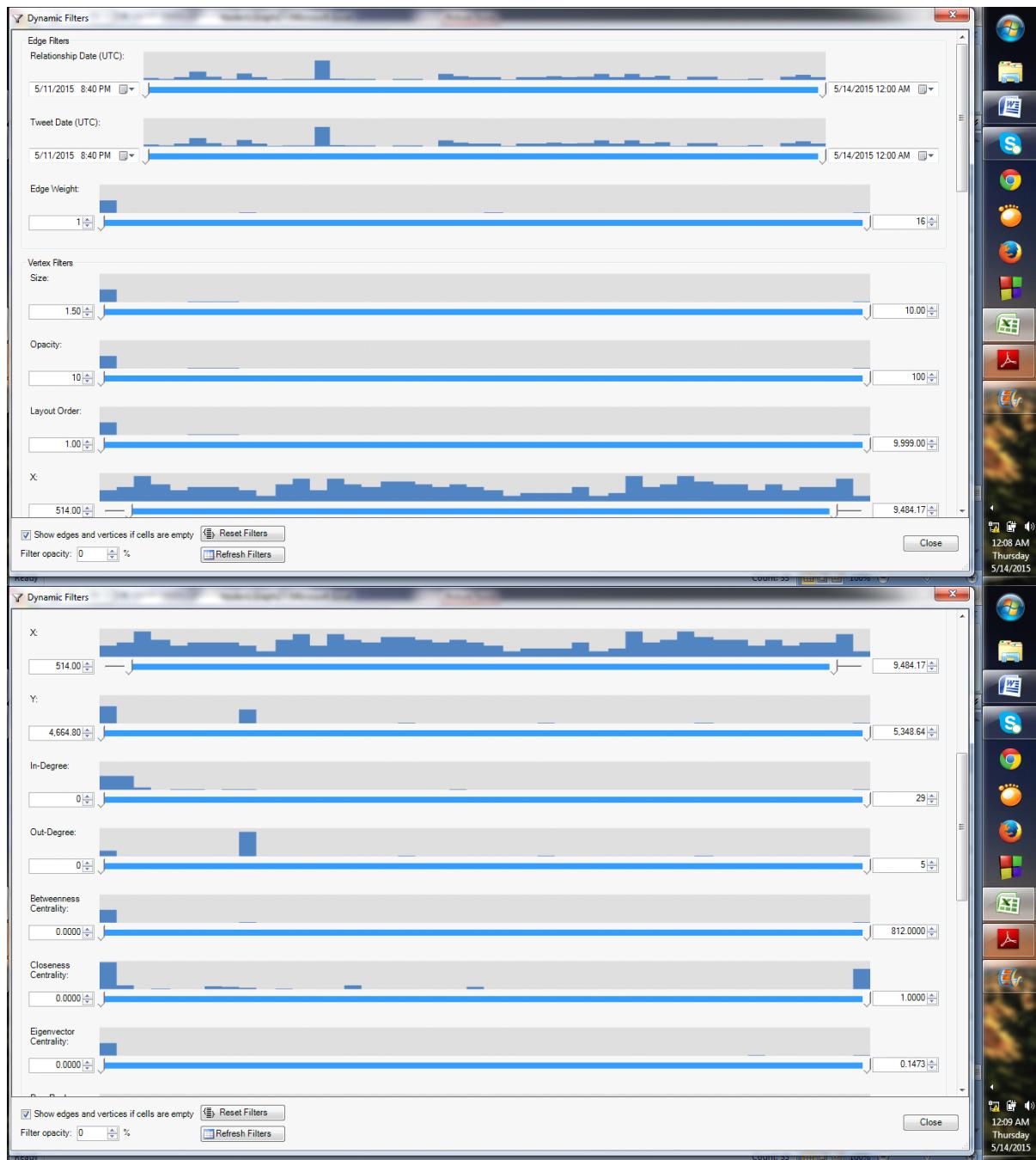
	A	B	C	D	E	F	G	H	I	J	K	L	S	T	U	V	W	X						
1	Visual Properties								Labels								Graph Metrics							
2	Vertex	Subgraph	Color	Shape	Size	Opacity	Image	Visibility	Label	Label Fill	Label	Position	Tooltip	Degree	In-Degree	Out-Degree	Betweenness	Closeness	Eigenvector					
199	ravenonthedream	Enter the name of the vertex.	93, (Solid Squ	1.5	10	<a href="http://pbs.twimg.com">http://pbs.twimg.com</a>	ravenonthedream		ravenonthedream		1	0		0.000	Screen Added	0.								

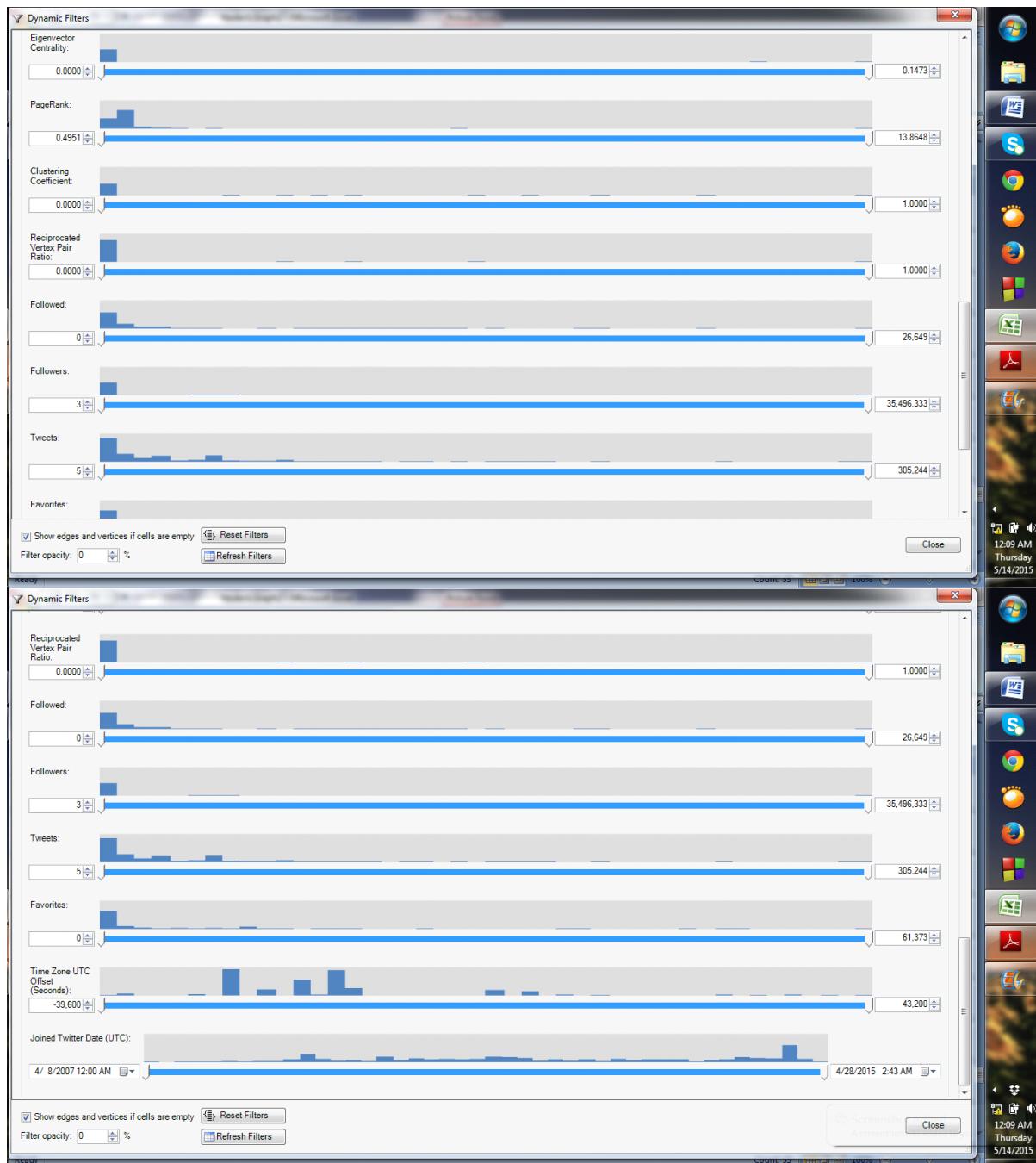
Group Vertices Overall Metrics Group Edges Twitter Search Ntwrk Top Items

Average: 32318.20558 Count: 3177 Sum: 52226220.21

Ready

## Dynamic Filter





# CHAPTER 6

## Conclusion and Future Work

---

This work was taken up with an objective to propose an algorithm to extract and transform the seed words from social media effectively and efficiently and eliminating the redundancy in the tweets. In this respect it was essential to first study the different technologies used till now for twitter data analytics which gives cost effective methodologies to collect the keywords from open source distributed platform.

### 6.1 Conclusion

This project helped me understand the intricate process of Data Extraction, Transformation and Standardization. We could strengthen our understanding about Data Analytics and how we conclude the various issues on women's education and employment. The following proposed algorithm gives the result effectively and efficiently. The proposed code acquainted us with various such scenarios where minor mistakes would result in erroneous conclusions. The project also helps to scan and analyze social media feeds so that we can permanently keep upto speed on the latest developments in the industry and its environment.

After immense toil and guidance, the algorithm was successfully realized and the extraction and statistical operations can now be smoothly executed within a matter of minutes. Further, we are going to implement this in real time environment.

### 6.2 Future Work

This includes

- User dependent threshold values.
- Updating the tweets every time from the twitter.
- Implementing with in-memory computing.
- Build it for another social media (Facebook, Youtube etc).
- Enable it for the real time query.
- To extend the limit of the number of tweets

## References

- [1] Smith and Shneiderman."Analyzing social media with nodexl." in University of Maryland . USA, 2011.
- [2] Stephen G. Kobourov." Spring Embedders and Force Directed Graph Drawing Algorithms." in *University of Arizona*, pp. 4-9, 2012.
- [3] Derek Hansen, Marc A. Smith, Ben Shneiderman." EventGraphs: Charting Collections of Conference Connections." in Connected Action Consulting Group, 2011.
- [4] Derek Hansen, Ben Shneiderman."Analyzing social media networks: learning by doing with Nodexl." in University of Maryland, 2009.
- [5] [Xindong Wu, Xingquan Zhu, Gong-Qing Wu, Wei Ding.]"Data Mining with Big Data," *IEEE Transactions on Knowledge and Data engineering.*, vol 26, no 1, January 2014.
- [6] Emilio Ferrara, Pasquale De Meo, Giacomo Fiumara, Robert Baumgartner." Web Data Extraction, Applications and Techniques: A Survey." july 2014.
- [7] Adilmoujaheed "An Introduction to Text Mining using Twitter Streaming API and Python", (adil moujahid) [online],2014.
- [8] Wang Ling, Luís Marujo." Crowdsourcing High-Quality Parallel Data Extraction from Twitter.",2013.
- [9] Shamanth Kumar, Fred Morstatter, Huan Liu."Twitter Data Analytics", in PROC Springer,2013