# Half-Day Tutorial (dg.o 2011 Conference)

## Collecting, Analyzing and Visualizing Tweets using Open Source Tools

Seungwon Yang and Andrea L. Kavanaugh
Department of Computer Science, Virginia Tech
{seungwon, kavan}@vt.edu

6/12/2011

This tutorial will teach participants how to collect, analyze and visualize results from twitter data. We will demonstrate several different free, open-source web-based tools that participants can use to collect twitter data (e.g., Archivist, 140kit.com, TwapperKeeper), and show them a few different methods, tools or programs they can use to analyze the data in a given collection. Finally, we will show participants visualization tools and programs they can use to present the analyses, such as tag clouds, graphs and other data clustering techniques. As much as possible this will be a hands-on tutorial, so participants can learn by making their own twitter data collection, analysis and visualization as part of the tutorial.

## Table of Contents

# Part 1. Web-Based Tools

## 1. Collecting Tweets

Twitter provides tweets through their REST & Search API and through their Streaming API[1]. The REST & Search API is used to find relevant tweets that are already archived in Twitter's servers. Tweets as old as 7 days usually can be collected from this API. The Streaming API, also called a 'fire hose' API, provides current tweets that are posted in real time. Most tweet collection tools use both of these APIs to archive existing older tweets as well as incoming new tweets.

### Terminology

*Hashtag*
They are a community-driven convention for adding additional context and metadata to your tweets.   It is added inline to your post.  You create a hashtag simply by prefixing a word with a hash symbol (e.g., #twitter, #japanearthquake, etc.).  For more information, please see http://twitter.pbworks.com/w/page/1779812/Hashtags

*Retweet (RT)*
It is similar to forwarding email to another person.   It helps quickly share tweets with all of your followers.  The retweeted tweets look like "RT @VerifiedQuotes: I'm not addicted to #twitter. I only tweet…", with 'RT' in front of the text.  For more information, please see http://support.twitter.com/entries/77606-what-is-retweet-rt

### 1.1. The Archivist

**Web version**

*The (Online) Archivist*[2] developed by Mix Online provides quick and easy creation of tweet visualizations.  However, due to Twitter's API Terms and Service, collections created by users reside in the company's servers.  In addition, only three collections per account can be created.  Export/download of the collected tweets is not allowed. However, it provides six visualizations for the basic analysis of archived tweets:
- Tweet volume over time: the number of tweets posted during a period of time
- Top users: the user IDs of people who tweeted much
- Tweet vs. Retweet:  the ratio of original tweets and retweets
- Top words: a list of frequent words in tweets
- Top URLs: a list of frequent URLs found in tweets
- Source: the name of an online service that is used to post tweets

---

[1] Twitter API FAQ. http://dev.twitter.com/pages/api_faq
[2] http://archivist.visitmix.com/

**Instructions to create a new tweet archive and view the visualization of the basic analyses in The Archivist (Web version)**

Step 1. You can sign-in using your Twitter account.  Click 'Sign In To Twitter' link on the top right corner of the main page.



Step 2. Type in your Twitter account info, and click 'Authorize app' button.

Step 3. Enter a key phrase or a hashtag of your interest. Then click 'Start analysis' button right next to the text box. After some seconds, The Archivist will show six different visualizations, each of which has a link to its full-page view.



Step 4. To keep archiving the tweets, you should save the archive in your profile page by clicking 'Save this archive'.

Step 5. To remove an archive, simply click the 'Remove archive' link on your profile page.



Note: to have more than 3 archives, you can sign out from Twitter and create an archive. Then, click 'Save this archive' link.  It will direct you to the Twitter login page.  Once you login, the newly created archive is added to your profile page even if your total number of archive has already exceeded the limit of 3.   (May be there is a bug in the software's logic.)

(You need a Twitter account to do this exercise. Please make one if you don't have it yet.)

*Exercise*
1. Participants get together and have a short brainstorming session.

2. Discuss current issues of interest.
3. Each participant selects an event/issue that he/she would like to collect tweets about.
4. Develop a keyword or a key phrase that represents the event/issue.
5. Use the developed key word or a key phrase as a query to a search engine. See if it returns a good result.
6. Use the same keyword or a key phrase as a search term in Twitter. See if it returns relevant tweets.
7. Modify the keyword or a key phrase.
8. Create an archive using The Archivist Web version using your keyword or key phrase.
9. From six different visualizations, which one(s) do you like the most? Which one will be most useful for your research?

\* The same search keyword / key phrase will be used to create archives in other tools in this tutorial session to compare the collection results.
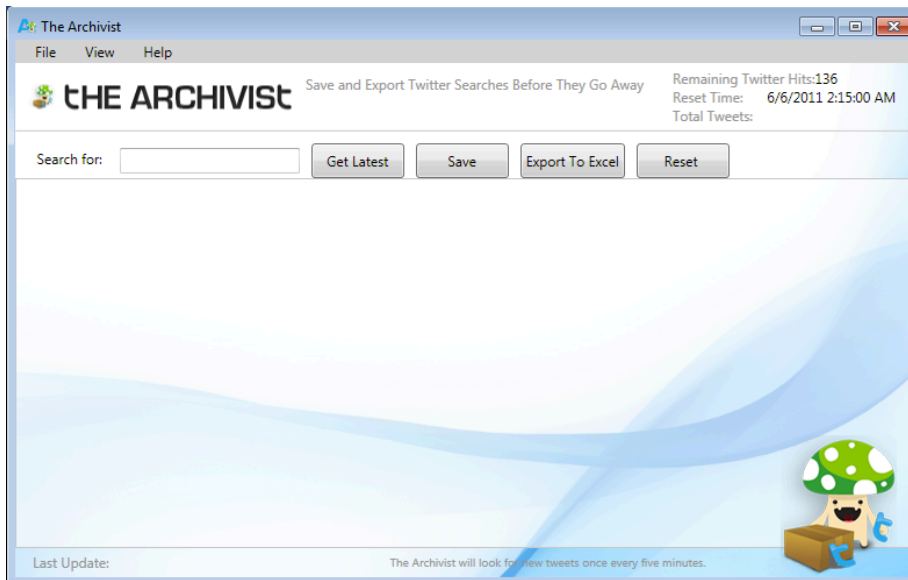

**Stand-alone version (Windows PC only)**

*The Archivist Desktop*[3] can run on a user's own **Windows** machine and continuously collect tweets. Tweets can be exported into an XML file or tab-delimited text file for later processing using Excel. Its pie chart visualization shows tweet volumes per Twitter ID. Tweet volume is visualized as a line graph. *The Archivist Desktop* runs continuously in the background and takes up much computing resources so it might result in performance degradation when multiple archives are created in a single machine.
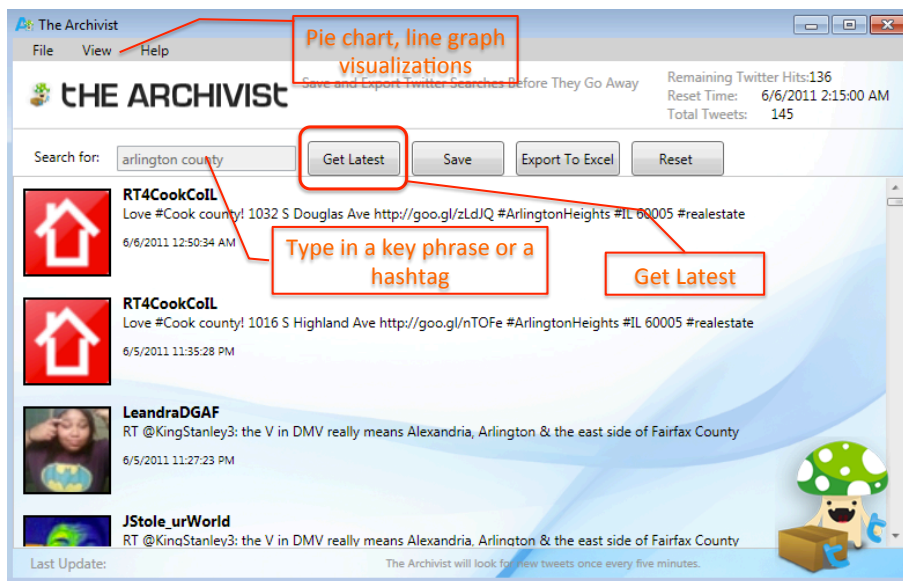
Step 1. Installing the tool
The software can be downloaded from http://visitmix.com/labs/archivist-desktop/
Once installed by saving and double-clicking the setup.exe file, run the tool. You should see its user interface that looks like the image below:

_____

[3] http://visitmix.com/labs/archivist-desktop/
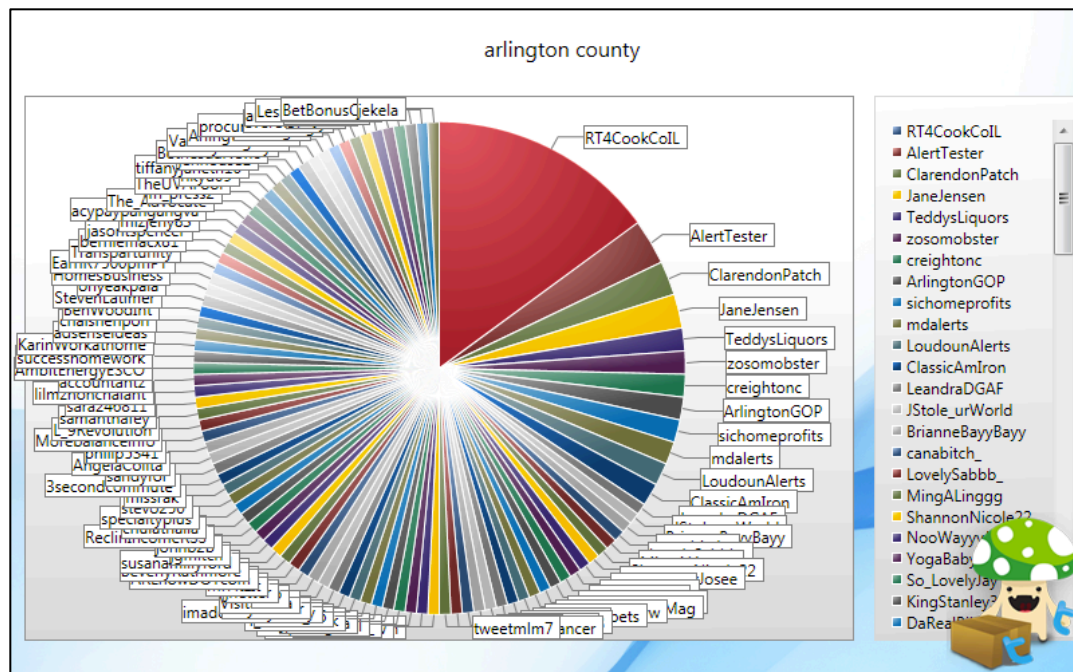
Step 2. Collecting tweets
Type a key phrase or a hashtag into the box, which is located next to 'Search for'. Then click 'Get Latest' button (the left-most one among control buttons) to begin archiving tweets.
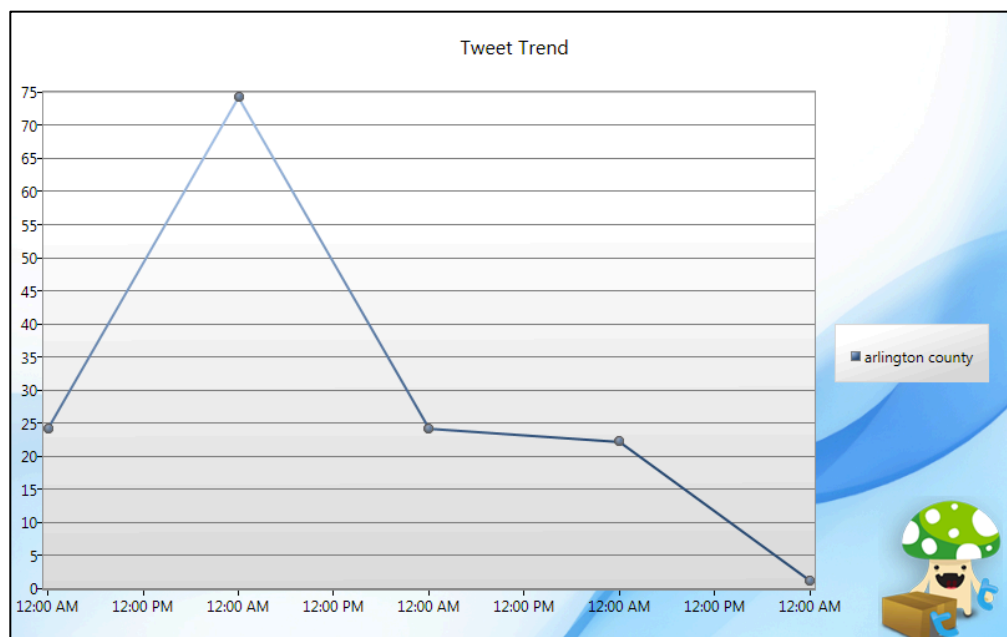


Step 3. Two visualizations.

*Pie chart visualization of tweets per tweet IDs:*
Go to 'View' on the task bar and select 'View Pie'. The pie chart shows the proportion of tweets from each person.
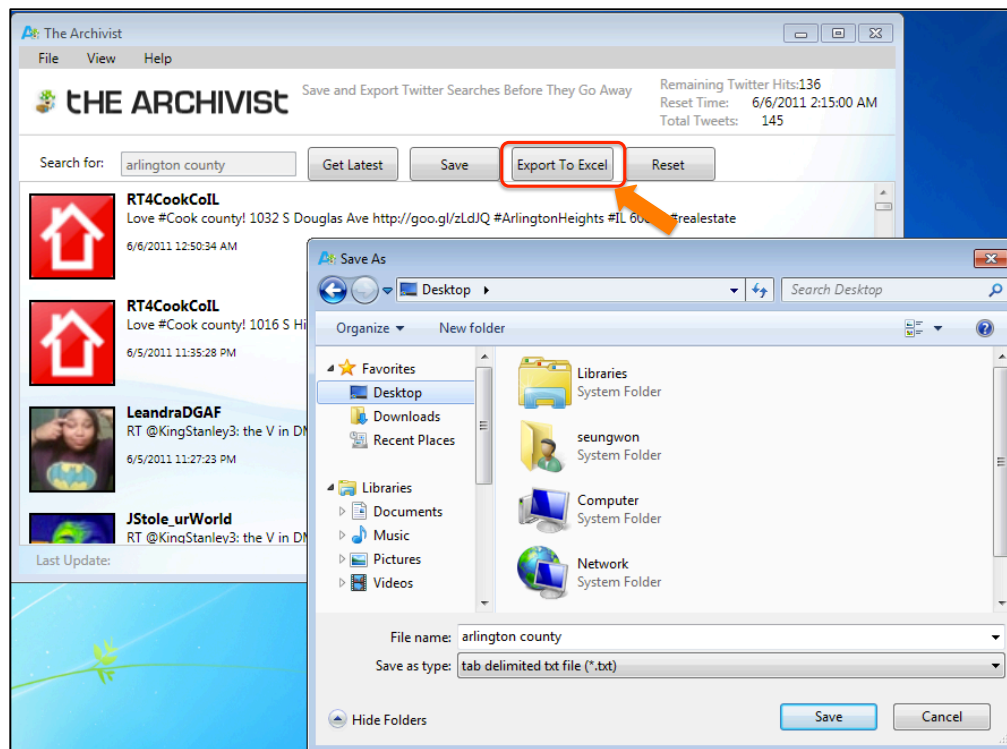
*Tweet volume visualization overtime:*
Go to 'View' on the task bar and select 'View Chart'.  The graph shows the number of tweets over a period of time.



Step 4. Exporting tweets.
There are two options for exporting collected tweets.   First option is to export as a tab-delimited text file, which can be imported into Excel for analysis. For this, click 'Export to Excel' button on the software interface as shown in Figure 16.  It will open a window, where you can select the folder to store the exported file.
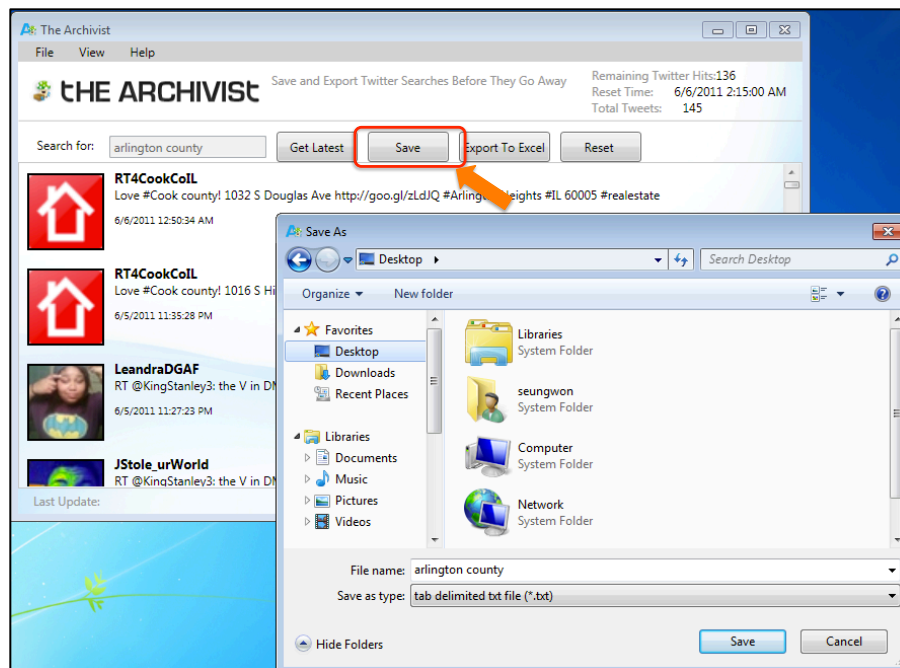
(Tab-delimited text file example)



Note: After exporting tweets as a tab-delimited text file, please change the filename to avoid the file content becoming an xml file.

The other option is to export as an XML file that can be parsed with scripts written in Python or PHP. For this, click the 'Save' button.
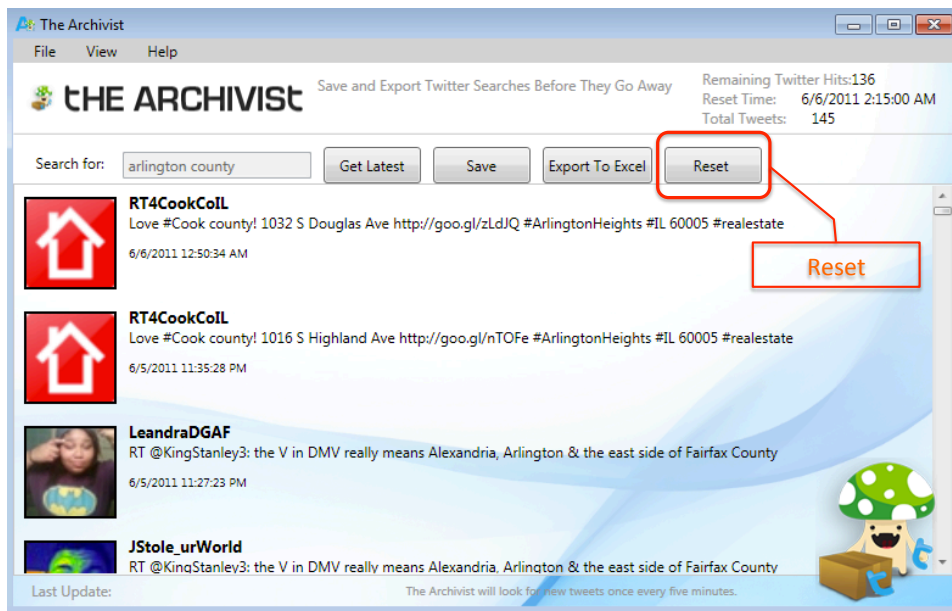
(XML export file example)



Step 5. Stop collecting tweets.

After running the tool for a period of time to archive tweets, you may want to stop archiving.   Clicking the 'Reset' button or closing the tool's interface by clicking 'X' on the top right corner will stop archiving the tweets.   Please don't forget to save your tweets.

Note: Since *The Archivist Desktop* instances are running in the computer's main memory and not storing the collected tweets into a database. Therefore, all the tweets will be lost if the tool stops running (e.g., due to automatic rebooting of the computer, or by the user's decision to stop archiving tweets.). If the computer turns into a sleep mode, the collection process stops, and starts again once the computer returns to a normal mode.

### 1.2. yourTwapperKeeper

To avoid violation of Twitter's API Terms of Service, *TwapperKeeper*[4], the public web service for archiving tweets, had discontinued its export and download features, which are essential for researchers to analyze tweet content.

*YourTwapperKeeper*[5], the open source version of *TwapperKeeper*, can run on a user's own machine to archive tweets with certain hashtags (e.g., #libya) and key words (e.g., japan earthquake). Users can export and download archived tweets into various formats such as an Excel file, RSS, JSON and HTML. Because *YourTwapperKeeper* resides on a non-public machine, exporting tweets is not a violation of any terms of service.

**Instructions to develop an archive and export its content**

*Creating an archive*
Step 1. Click the blue button on top of the user interface.

---

[4] http://twapperkeeper.com/index.php
[5] http://your.twapperkeeper.com/

Step 2. You let *yourTwapperKeeper* access your Twitter account info. Then the control is passed back to the application.



Step 3. Type in a key phrase/hashtag, description of your archive and optional tags.  Then, click 'Create Archive' button.

Step 4. Once the archive is successfully created, it is displayed on the list of archives as shown below.  Only the creator of an archive can remove the archive or edit its description.  Users have access to all archives.  They can view and export tweets from archives that were created by anyone.



Clicking a binocular icon on the right side of an archive opens a tab on a browser, where the raw tweets can be exported in various formats such as HTML, RSS, Excel, Simple Table, or JSON API.

*Exporting an archive as an Excel file*
Step 5. Select settings for tweet export.  Users can select tweets within a certain time period, the maximum number of tweets to export, tweets from a certain user, tweets that has a certain text, and whether to include retweets or not.

   After selecting and typing in all the settings, click 'Query' button to generate download links for various export formats.

Step 6. Exported tweets in an Excel file.
Columns include the tweet text, to-user ID, from-user, tweet ID, language setting, geo location coordinates, tweet time, etc.



*Exercise*
1. Create an archive using yourTwapperKeeper tool at http://virginia.cc.vt.edu/
2. Use your keyword or key phrase as the search term. You may add a hash '#' to make your keyword a hashtag (Search Twitter or The Archivist Web version with your hashtag and see if it returns a reasonable result).
3. Refresh the browser to update the number of tweets archived.
4. Open a tweet export page by clicking the binocular icon in one of the archives.
5. Adjust the export setting and export the tweets following the instructions above.
6. Download the excel file and examine its content.

### 1.3. 140kit

140kit[6] is online services to collect, analyze, and visualize tweets. Collection of tweets can continue for a maximum of seven days, and then potentially be extended for further archiving. Once the collection is completed, the tool provides basic analyses and graph visualizations of tweets. CSV files of the basic analyses can be downloaded from the site.

Users can search existing collections and visualize the tweets by using the features provided. 140kit used to allow exporting of the raw tweet archives, but this feature is no longer available due to the updated Twitter API Terms of Service at
http://dev.twitter.com/pages/api_terms

*Especially, under 4.A:*

    *"You may export or extract non-programmatic, GUI-driven Twitter Content as a PDF or spreadsheet by using "save as" or similar functionality. Exporting Twitter Content to a data store as a service or other cloud based service, however, is not permitted."*

**Instructions on creating a new tweet archive and viewing the pie charts of the basic analysis**

Step 1. After creating an account at the sign up page, please login by typing in your User Name and Password and then click the 'Log In' button.



Step 2. Click the 'Add a new collection' link on your main page.

---

[6] http://140kit.com/

15

Step 3. Click 'Search/Term scrape' link to collect tweets using a hashtag or terms.



Step 4. Enter a search term or a hashtag (e.g., "#japanearthquake") to start archiving tweets, which contain them. Click 'Create' button on the bottom of the page to begin archiving tweets.

Step 5. Once archiving the tweet is completed, users can view pie charts of basic analysis results by clicking links under 'Histograms' and 'Word Frequencies'. Also the retweet network graph is provided. Please see the sections enclosed by rounded rectangles in the image below.



Note: it usually takes some time to see the results from the 140kit tool.

(You need to make an account at 140kit.com to do this exercise.)

*Exercise*
1. Create an archive using 140kit tool. Use your keyword or key phrase as the search term. You may add a hash '#' to make your keyword a hashtag (Search Twitter or The Archivist Web version with your hashtag and see if it returns a reasonable result).
2. Leave the tool to archive tweets.
3. Select 'Collections' link on top of the 140kit homepage.
4. Select one of the collections and click the name.
5. Examine various histograms and word frequency graphs, as well as the retweet graph.


## 2. Visualizing Tweets

Word cloud visualizations show frequently appearing words from input text. More frequent words appear as bigger fonts. Colors are used to visually distinguish words and make word clouds visually appealing.


### 2.1. Wordle.net

The image above shows a word cloud created by the Wordle[7] Web service using 100 tweets about the Japan earthquake disaster that occurred in March, 2011.  Users can create word clouds and refine them by using the features:

- Remove uninteresting and common words directly in the GUI
- Limit the maximum words
- Change font, upper/lower cases, color schemes, layouts
- Assign weights on words, change background color

**Instructions to create a word cloud**

Step 1. Click the 'Create' button on the task bar to start the process.



---

[7] http://www.wordle.net/

Step 2. Copy and paste tweet texts from an excel file into the text box.  Or type in a URL of the data file or a Web page.  You can make multiple word clouds by each date or each hour of a specific date to see the different topical words emerge.



Step 3. Adjust layout and remove commonly repeated words.
Usually, 'Horizontal' layout is comfortable to read words in the cloud.  To change the organization of the words, select 'Re-layout with current settings' under 'Layout'.

The key phrase that was used to collect tweets can be removed since we know that the tweets are about that key phrase. For example, if a keyword 'victims' was used to collect tweets, it is shown with a bigger font size due to its frequency. By removing this dominating word, we can examine other interesting words closely.

Right click on the word to remove and select 'Remove "<word>"' option to delete it from the visualization canvas.



Step 4. Adjust the number of words to visualize. If there are too many words to visualize, the font of the less frequent words becomes too small to read. For this, please select 'Maximum words' under 'Layout' menu in the tool bar.

Step 5. Change fonts. Teen and Coolvetica fonts usually make a nice word cloud. For this, select 'Font' menu in the tool bar.

(Teen font)



(Coolvetica font)

Step 6. Change color settings.
Try different color settings to find the one you like. Sometimes dark backgrounds make
the words in the cloud more visible.

Step 7. Capturing the created word cloud.
You can use 'Print' feature to make a PDF file of the cloud or capture the screen.

Step 8. Sharing the cloud with others.
You can make your word cloud public so that others can view it, too. Please click 'Save
to public gallery…' on the bottom right corner of the interface. Then, enter the title,
username and comments, and hit 'OK' button.

Step 9. Word counts.

To get a list of word frequency, please go to 'Language' and select 'Show word counts…'.  You can reorder the words in either ascending or descending order by clicking 'Frequency' on the Word Counts pop up window.



(You can use your own tweet archive developed using yourTwapperKeeper at http://virginia.cc.vt.edu/ or any archives at http://mule.dlib.vt.edu/ )

*Exercise*

1.  Export tweets in the archive as an excel file (You can either include or exclude retweets when exporting).

2. Group tweets in the excel file by dates (if you use an existing archive with lots of tweets collected for many days), by the hour, or by 10 minutes.
3. Create at least 4 word clouds using the grouped tweet data in the previous step.
4. Be sure to remove the key phrase, which was used to archive tweets.
5. Compare the word clouds.
6. Do you see the differences? What do you think the word clouds as a visualization tool for tweets?

## 3. Analyzing Tweets Using Web-Based Tools

The basic algorithm for text analysis in word cloud creation is to count word frequencies in input texts. For more meaningful content analysis, terms and phrases can be extracted using web services and Natural Language Toolkit.

### 3.1. Terminology Extraction

*Translated Labs*
The terminology extraction web service from The Translated Labs[8] identifies terms from texts. The basic idea is to compare the frequency of words in an input text with their frequency in the language. Their assumption is that the words, which appear very frequently in the document but rarely in the language, are probably terms.
Extracted terms are Google-searched when they are clicked.



---

*AlchemyAPI*

AlchemyAPI provides more detailed analyses of the input text such as named entity extraction, concept tagging, keyword/term extraction, sentiment analysis, etc. Its demo site provides Web interface to accept raw texts and then process them.



(Please use the grouped tweet data, which was used in the previous exercise.)

*Exercise*
1. Enter each of the grouped data into the text box in AlchemyAPI demo site.
2. Copy and Paste the analysis results to a separate sheet in the tweet excel file.
3. You will have four groups of results for each tweet group.
4. Compare the terms and concepts extracted, with the corresponding word clouds you previously created.
5. Do they match well?  Which one summarizes the group of tweets better for human understanding?

# Part 2. Tools with Scripting

## 1. WordCram with Processing Language

Processing[9] is an open source programming language and environment that is gaining wide acceptance from people in various fields. By using the WordCram library[10] with Processing, users can develop dynamic word clouds. For example, new tweets for the last 10 minutes in the database about the Japan earthquake disaster can be accessed and then converted into a word cloud. The codes can be exported as an applet to be uploaded to a server for online access. Users might be able to monitor current events based on this dynamic word cloud.

Example dynamic word clouds can be accessed at:

(Left)  Japan Earthquake Disaster: http://mule.dlib.vt.edu/~seungwon/japan.html
(Right) Libya Revolution: http://mule.dlib.vt.edu/~seungwon/libya.html



### 1.1. Procedure

---

[9] Processing programming language and environment. http://processing.org
[10] http://code.google.com/p/wordcram/

*Tasks running in the background*
1. YourTwapperKeeper archives tweets into MySQL database tables.
2. For each time period (e.g., every 5 minutes), a script accesses one of database tables and fetches text in tweets.  Then it writes the raw text into a file.
3. Following execution of the script at step 2, another script processes the raw text file to remove stop-words (e.g., 'the', 'a', 'about', 'it', etc.) and symbol characters (e.g., '$', '!', '?', etc.), and creates a processed text file.

*Tasks for creating word clouds*
4. Processing sketch (i.e., code) with WordCram library accesses the processed text file developed at step 3.
5. Methods to change the layouts, color settings, canvas size, maximum and minimum font sizes, etc. are applied to make the words on the canvas comfortable to read and visually appealing.
6. Export Processing sketch to an applet, which can be inserted in a Web site.

## 1.2. Installing Processing and WordCram Library

*Processing Installation on a Windows PC*
Step 1. Go to http://processing.org/download. Download and save a zipped file for your operating system.  For Windows version, file size is 85.77 MB.  In most cases, you might want to download 'Windows' version and not the 'Windows (Without Java)' version unless you know that JDK (Java Development Kit) is already installed on your computer.

By downloading the software from this page, you agree to the specified terms.

**1.5.1 | 15 May 2011**

↓ Linux          ↓ Windows
↓ Mac OSX        ↓ Windows (Without Java)*

**The** list of revisions **covers the differences between releases in detail. Please read the** changes **if you're new to the 1.0 series.** Also check the known problems for this release.

*\* The Windows version without Java is for users who can take care of installing a JDK (not JRE) themselves. It should only be downloaded by advanced users who are familiar with Java.*

**Resources**

» Tutorials
» Examples
» FAQ
» Troubleshooting
» Supported Platforms
» Processing Wiki
» Processing Forum
» Report a bug
» Download Source

**Announcements**

Email address

Submit

If you are interested in receiving updates about Processing, submit your email through this form. *Your email will only be used to send infrequent updates about Processing. It will not be sold or shared.*

Step 2. Double-click (total twice) the zipped file to see its content. You will see folders such as 'java', 'lib', 'modes', and 'tools' as well as files such as 'processing' and 'revisions'. Double-click 'processing' and follow the instruction to extract the content onto your computer.

Step 3. After extraction is done, make a shortcut icon in the Desktop screen for easy access.

Step 4. Run Processing by double-clicking the shortcut icon. Its interface window appears.



Step 5. Test with an example sketch. Go to 'File' → 'Examples'. It will open a 'Standard Examples' window. Select an example sketch by double-clicking one, and then run it by clicking the play button below the task bar, or by going to 'Sketch' → 'Run'. The image below shows an example sketch, 'Animator'.

This completes an installation of Processing on a Windows PC.

*Processing Installation on a Mac Machine*
Step 1. Go to http://processing.org/download. Click 'Mac OSX' link to download a zipped file. Processing-1.5.1-macosx.zip is 50.5 MB in size.



Step 2. Double-click the file to unzip it. Then an application with the Processing logo appears in Downloads folder.

Step 3. Move the Processing application from Downloads folder to Applications folder using drag-and-drop.



Step 4. Click the Processing icon in Applications folder to run it. The processing interface window appears.



Step 5. Test with an example sketch.   Go to 'File' → 'Examples'.  It will open a 'Standard Examples' window.  Select an example sketch by double-clicking one, and then run it by clicking the play button below the task bar, or by going to 'Sketch' → 'Run'.  The image below shows an example sketch, 'RadialGradient2' in 'Color' folder.

This completes an installation of Processing on a Mac machine.

*WordCram Installation (after installing Processing)*

Step 1. Go to WordCram installation page at
http://code.google.com/p/wordcram/wiki/Installing.  Download a zipped file
'wordcram.0.5.0.zip' by following the 'Download' link under 'Details'.  Once
downloaded, double-click it to unzip it.  The unzipped folder name is 'wordcram.0.5.0'.



Step 2. Find out your Processing Sketchbook folder.
Go to 'File' → 'Preferences'.  It will open a window.  Processing Sketchbook folder
location is specified below 'Sketchbook location'.   The location of the Sketchbook folder
in the example below is at 'C:\Users\seungwon\Documents\Processing'.

For a Mac machine, the Sketchbook location looks, for example, "/Users/seungwon/Documents/Processing".



Step 3. Open a file browser and go to that Sketchbook folder. Then, create a folder named, 'libraries' if it does not exist there.

(Windows screen shot)



(Mac screen shot)



Step 4. Copy the folder named 'WordCram', which is located inside the unzipped folder, 'wordcram.0.5.0'. (For a Mac machine, 'WordCram' folder is created in the same folder, where wordcram.0.5.0.zip is located, once this zip file is unzipped.)

Then paste the entire content of the 'WordCram' directory into the 'libraries' folder that you created in Step 3.  Restart Processing to make changes to be effective.

(Windows screen shot)

(Mac screen shot)



Step 5.  Testing an example of the WordCram library.
Go to 'File' and select 'Examples'.  'Standard Examples' window will open.  Expand 'WordCram' folder under 'Contributed Libraries' folder on the bottom of the window. Expand 'Other Examples' folder and double-click 'usConstitution' sketch to open.



Step 6. Click the play button under the task bar to visualize a word cloud of the US constitution.  If you see a beautiful word cloud opened up in a window, then your WordCram library is successfully installed.

At this point, both Processing and WordCram have been installed successfully. The next is to connect them to a dynamically updated tweet text file on a server.

**1.3. Connecting WordCram/Processing Sketch with Processed Tweet Text File**

Step 1. Using usConstitution example as a template.
Go to 'File' and select 'Save as' to save the 'usConstitution' example as another name, for example, 'myCloud'. You can save it under the folder, where usConstitution is saved.



Step 2. Accessing a text file in a server.
Comment out '.fromTextFile("../../usconst.txt")' line by adding two slashes in front of it, and add '.fromWebPage("<URL to the text file>")' below the commented line. It allows the Processing sketch to use the text file in a server as an input. You can experiment with a URL to your homepage and see how it looks.

```
WordCram wordCram;

void setup() {
  size(800, 600);
  background(255);
  colorMode(HSB);

  initWordCram();
}

void initWordCram() {
  wordCram = new WordCram(this)
    //.fromTextFile("../../usconst.txt")
    .fromWebPage("http://curric.dlib.vt.edu/DLcurric/vitae/")
    .withStopWords(StopWords.ENGLISH + " shall")
    .withFont(createFont("../../LiberationSerif-Regular.ttf", 1))
    .sizedByWeight(10, 90)
    .withColors(color(0, 250, 200), color(30), color(170, 230, 200));
}

void draw() {
  if (wordCram.hasMore()) {
    wordCram.drawNext();
  }
}
```

Step 3. Changing to a horizontal layout.

Words that are placed vertically may not be comfortable to read.  The default layout setting in WordCram is a mixed mode that places some words vertically and other words horizontally.   To make all words horizontal:

- Add 'Anglers ang;' (including the semicolon after 'ang') following the 'WordCram wordCram' line.
- Add '.withAngler(ang.horiz())' line following the '.fromWebPage(…)' line that we added in Step 2.
- Save the sketch and run it to see the changes

```
import wordcram.*;
import wordcram.text.*;


WordCram wordCram;
Anglers ang;

void setup() {
  size(700, 500);
  background(255);
  colorMode(RGB);

  initWordCram();
}

void initWordCram() {

  //int everytenmin = (minute() / 10) * 10;
  wordCram = new WordCram(this)
/*      .fromWebPage("http://mule.dlib.vt.edu/~seungwon/"+"libya-"+hour()+".
      .fromWebPage("http://mule.dlib.vt.edu/~seungwon/tweettext/libya-currer
/*      .withStopWords(StopWords.ENGLISH + " RT Libya")    */
      .withAngler(ang.horiz())
      .withFont(createFont("LiberationSerif-Regular.ttf", 1))
      .sizedByWeight(40, 110)
      .withColors(color(50), color(100), color(150), color(200), color(250)
}
```

Step 4. Changing the color setting.
You can make the words either color or gray scale.  If you add a single number inside parentheses of 'color(…)', it is interpreted as a gray scale.  For RGB color mode, please add three numbers separated by commas (e.g., color(150, 255, 150) makes a bright green color). Numbers from 0-255 are allowed.

```
void initWordCram() {

  //int everytenmin = (minute() / 10) * 10;
  wordCram = new WordCram(this)
/*      .fromWebPage("http://mule.dlib.vt.edu/~seungwon/"+"libya-"+hour()+".
      .fromWebPage("http://mule.dlib.vt.edu/~seungwon/tweettext/libya-currer
/*      .withStopWords(StopWords.ENGLISH + " RT Libya")    */
      .withAngler(ang.horiz())
      .withFont(createFont("LiberationSerif-Regular.ttf", 1))
      .sizedByWeight(40, 110)
      .withColors(color(50), color(100), color(150), color(200), color(250)
}
```

Step 5. Adjusting the canvas size, minimum and maximum font sizes.

Canvas size can be adjusted with 'size(width, height)'.

```
import wordcram.*;
import wordcram.text.*;


WordCram wordCram;
Anglers ang;

void setup() {
  size(700, 500);
  background(255);
  colorMode(RGB);

  initWordCram();
}
```

Try different numbers in '.sizedByWeight(minimum font size, maximum font size)'.

```
void initWordCram() {

  //int everytenmin = (minute() / 10) * 10;
  wordCram = new WordCram(this)
/*      .fromWebPage("http://mule.dlib.vt.edu/~seungwon/"+"libya-"+hour()+".
      .fromWebPage("http://mule.dlib.vt.edu/~seungwon/tweettext/libya-curren
/*      .withStopWords(StopWords.ENGLISH + " RT Libya")   */
      .withAngler(ang.horiz())
      .withFont(createFont("Liberati serif-Regular.ttf", 1))
      .sizedByWeight(40, 110)
      .withColors(color(50), color(100), color(150), color(200), color(250)
}
```

Step 6. Exporting to an applet.
Applet is a program written in the Java programming language. It can be included in an HTML page and run on the browser of a client's machine. In Processing, go to 'File' and select 'Export Applet'. Then a pop up window appears. Upload the files in the 'applet' folder to your server. Then, the 'index.html' file will display the word cloud that you created.

*Exercise*
1. Create a word cloud using the New York Times web site (http://nytimes.com/) or your personal homepage URL.
2. Adjust settings such as layout and color based on the given instructions.
3. Another URL is given to each participant, which links to a pre-processed tweet text fetched every five minutes from your archive.
4. Replace the URL in your New York Times (or your homepage) word cloud with the provided one.
5. Compare your dynamic word cloud with the actual tweets (by searching with your key phrase at Twitter.com).
6. Does your word cloud give a snapshot of what is happening in your issue/event?


## 2. Analyzing Tweets with Python Scripts and Term Extraction API

The basic algorithm for text analysis in word cloud creation is to count word frequencies from input texts. For more meaningful content analysis, terms and phrases can be extracted using the term extraction application programming interface (API) and Natural Language Toolkit (NLTK).


## 2.1. Useful Unix/Linux Commands
There are some useful commands to use when working with scripts and APIs in either Linux or Mac OS X systems.

- ssh: Linux Secure Shell
  Provides access to a Linux system using password encryption in the login process.
  **Usage**
  shell> ssh seungwon@newyork.cc.vt.edu  [hit enter key]
  Then, you will be prompted to enter a password

- scp: Secure copy.
  Scp allows to copy files across an ssh connection.
  **Usage**
  shell> scp [[user@]from-host:]source-file [[user@]to-host:][destination-file]
  shell> scp dlrl@georgia.dlib.vt.edu:/home/dlrl/testfile.txt  /home/tmp
  You will be prompted to enter a password.  Once verified, the file 'testfile.txt' is copied to the folder '/home/tmp' in the user's computer.

- ls: List directory contents
  ls lists information about the content of the current directory.  With '-al' option, detailed information is displayed, too.
  **Usage**
  shell> ls  → list content of the current directory

shell> ls –al  /home/seungwon  → list content details of /home/seungwon
directory

- cd: Change directory
  **Usage**
  shell> cd /home/seungwon/tweettext → go to /home/seungwon/tweettext
  directory

- pwd: Print the current working directory
  **Usage**
  shell> pwd  → this will show you the path to the directory that you are in.

- cat: Concatenate files and print on the standard output. It is often used to view the
  content of a file quickly.
  **Usage**
  shell> cat /home/seungwon/tweet.txt  → it shows the content of 'tweet.txt' file on
  the standard output

- vi and vim editors: Command line-based Linux file editor
  **Usage**
  shell> vi <file name>  → this will open the file specified in a vi interface.

  Once the file is open, you can edit/write when it is in input mode.
  **Entering Input Mode**
  a                      Add text after the cursor.
  i                      Insert text before the cursor.

  After finishing edits, push 'esc' key to enter command mode.
  **Command Mode**
  :w             Write the file to disk.
  :q             Quit vi.
  dd             Delete an entire line where the cursor is located
  yy             Copy an entire line where the cursor is located

- crontab: This utility creates, replace, or edit a user's crontab entry.  The entry lists
  when and how often each task should be executed.  It is used to schedule periodic
  background work.
  **Usage**
  shell> crontab –e  [-e option is to edit the invoking user's crontab entry]

To run Python or PHP scripts:
- shell> python <script name>
- shell> php <script name>

## 2.2. Terminology Extraction

*Yahoo! TermExtraction API[11]*
This API allows extracting terms and key phrases from a textual input. This service is limited to 5,000 API calls per IP address per day for non-commercial use. The result is returned in XML format.

*AlchemyAPI*
It provides various services such as named entity extraction, concept tagging, keword/term extraction, sentiment analysis, etc. through their application programming interface (API). The rate limits for a registered account is 1,000 calls per day. It can increase to 30,000 if you get permission from the AlchemyAPI company.

*Exercise*
1. Login to server, virginia.cc.vt.edu, using the account provided.
2. Under your home directory, /home/<account name>/, you will create 4 data files named, input1.txt, input2.txt, input3.txt, and input4.txt. Each of them will include a group of tweets you used for previous exercises.
3. Run a Python script located in your home directory to extract terms from those four input files. For example, you may run the script:

   Shell> python term_extractor.py input1.txt

## 2.3. Natural Language Toolkit (NLTK)
NLTK[12] is a powerful language analysis toolkit for Python. Exercises include identifying bigrams and trigrams using the provided Python script developed using NLTK.

*Exercise*
1. Login to server, virginia.cc.vt.edu, using the account provided.
2. Under your home directory, /home/<account name>/, a Python script "nltk_collocations.py" is located.
3. Run this script to find bigrams and trigrams from the input files, which were used in the previous exercise.

---

[11] http://developer.yahoo.com/search/content/V1/termExtraction.html
[12] http://www.nltk.org/