# Problem 1: Data acquisition

Problem Statement: Data collection is an important step that often is quite time consuming and constrained by regulations and privacy and security issues. For instance in order to collect data from patients in a clinical trial, a researcher in an organization has to get prior approval and certification from Internal Review Board (IRB) of the organization. There are many regulatory laws governing what data you collect about common citizens. Students ought to be sensitive to all these when collecting data. The other challenge in data collection is that data comes in different formats (text, html, txt, csv, json etc.) and feature widely varying access methods (web URL, api, hdfs, sqldb etc.). We will learn about this by working on few representative methods for data acquisition given in the handout: http://www.cse.buffalo.edu/~bin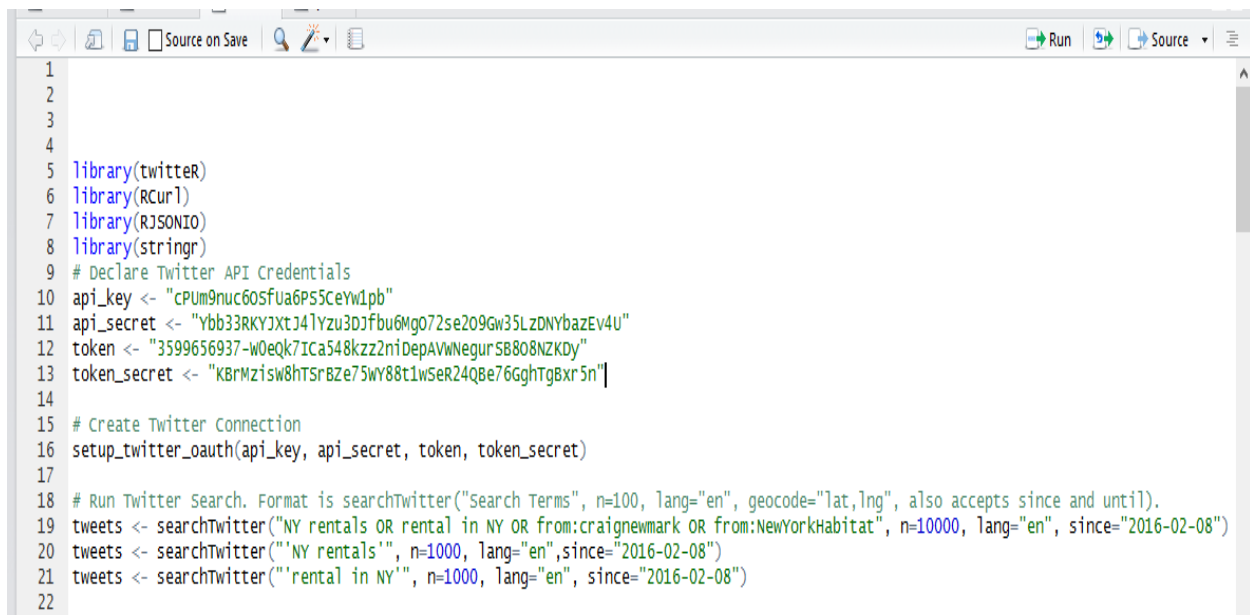a/cse487/spring2016/Lectures/RHandout1.pdf. To that list of methods add an approach for reading json data. For example twitter data is published as JSON objects.

In this project, I have collected data from "**twitter**" on the topic "**Rental Apartments**". I choose "**Twitter**" because of the abundant amount of data that could be acquired on any topic. The topic "**Rental Apartments**" was chosen because it can be used as the data for Problem 4 and 5. The data was collected for a period of approximately **1 month**. This gives us enough data (tweets) using which appropriate analysis can be made.



Fig: Sample Data

Instead of using a python program to retrieve data and then an R script to convert the data to json format, I have used **"twitteR"** package of R. This package is a convenient way of retrieving tweets and converting it to json format. The sample code used to retrieve tweets is shown in the figure below.

Fig: Shows the code for collecting data using "twitter" Package.

The command "library(twitteR)" is used to load the library and all the authorization keys acquired from the Twitter Api is entered. To run the twitter search, the following command is used.

**tweets <- searchTwitter("'NY rentals'", n=1000, lang="en",since="2016-02-08")**

In the above command, the retrieved tweets are stored in "tweets" variable, **"NY rentals"** is the search term, **"n"** is equated to 1000, to retrieve **1000 tweets** from the date **2016-02-08** in **English.**

**twt<-toJSON(tweets)**
**write(twt,file="tweet.json")**

The above command is used to write the tweets in to a json file called "tweet.json". The "toJSON" function converts the data in a variable to Json and then it is written to a file.

The file RP1pgolluin.R is an R script which has retrieves the tweets as well as converts it into a JSON file.