

# Assignment 1 : Sampling from Twitter

**Post Date:** Friday, September 4, 2015

**Due Date:** Friday, September 11, 2015 11:59 PM

**TA Office Hours:** Wednesday, September 9, 2015 4:30 PM - 6:00 PM (LSRC D309)

In this assignment, you will be estimating some properties of twitter users. You find a brief tutorial on how to get information about twitter users using python at (<https://www.cs.duke.edu/courses/fall15/compsci590.4/assignment1/TwitterTut.pdf>). It describes APIs that allow you to query Twitter for various types of information including user profiles, tweets, timelines, etc. However, each user is only allowed to query the website 180 times each 15 minutes. You are advised to form groups of 3. You are also advised to start working on the assignment as soon as possible (reasons will be clear once you understand the assignment).

You are required to use this API to estimate the following properties for the first  $A$  Twitter users, for all  $1 \leq A \leq 5$  million.

- “fraction of users from the US” (*hint: use time zone*)
- “fraction of users with more than 4,000 followers”
- “fraction of users with more than 500 friends”

You are given a file [https://www.cs.duke.edu/courses/fall15/compsci590.4/assignment1/first\\_twitter\\_ids.txt](https://www.cs.duke.edu/courses/fall15/compsci590.4/assignment1/first_twitter_ids.txt) which contains the IDs of the first 5 million users in order of when they joined twitter. A trivial solution to the above problem would be to query twitter 5 million times. However, due to the limit on number of queries to the website this can take a few days. Your goal is to approximate all the fractions within a relative error of  $\epsilon = 0.25$  with high probability ( $\delta = 0.05$ ).

1. How many queries do you need to estimate the fraction of people ( $\mu$ ) satisfying the above properties among the first 5 million users within a relative error of  $\epsilon = 0.25$  with high probability ( $\delta = 0.05$ )? (Hint: you may need to make some reasonable estimate on  $\mu$ )
2. How many queries do you need to estimate the fraction of people ( $\mu$ ) satisfying the above properties among the first  $A$  users, for all  $1 \leq A \leq 5$  million, within a relative error of  $\epsilon = 0.25$  with high probability ( $\delta = 0.05$ )?
3. What is the answer to questions 1 and 2 for  $\epsilon = 0.1$ ?
4. (BONUS) An [article in mashable](#) said that in 2013 that more than half the active users in Twitter came from US, Japan, Indonesia, UK and Brazil, followed by Spain, Saudi Arabia, Turkey Mexico and Russia. How many samples would you need to estimate the fraction of individuals in the first 5 million users who come from these 10 countries, all within a relative error of  $\epsilon = 0.25$  with high probability ( $\delta = 0.05$ )? *No need to answer this bonus question via Twitter API. You are welcome to do so if you are interested.*
5. (BONUS) What epsilon error can we achieve if 10 groups pooled their data for answering questions 1 - 3? *No need to use Twitter API for this bonus question.*

Please submit a PDF report plotting the estimates you got for each of the properties, answering the 4 questions above, and outlining the method used to do the estimation.