

Tracking Drinking Behavior from Twitter Data

COSC 526 Course Project Report

Zhifei Zhang Yajun Wang Yongli Zhu

Department of Electrical Engineering & Computer Science
The University of Tennessee, Knoxville

Abstract

This report presents the work of a tracking drinking behavior from Twitter data set. In this work, the *Latent Dirichlet Allocation* method is adopted with the assistance of human-created key words dictionary. Two main behaviors are extracted: the user device platform usage distribution and user location distribution. From the results, the IOS platform based device got the No.1 place; while the location behavior indicates a similar distribution pattern in accordance with the whole Twitter uses distribution. Those information can be utilized to instruct more effective advertising and more user-oriented apps developing for smart phone/tablet groups.

1 Introduction

Big data is everywhere people look these days. Businesses are falling all over themselves to hire data scientists, privacy advocates are concerned about personal data and control, and technologists and entrepreneurs scramble to find new ways to collect, control and monetize data. It is universally acknowledged that data is powerful and valuable.

Data Mining is an analytic process designed to explore big data in search of consistent patterns or systematic relationships between variables, and then to validate the findings by applying the detected patterns to new subsets of data. The ultimate goal of data mining is prediction - and predictive data mining is the most common type of data mining and one that has the most direct applications^[1].

Platforms such as Google, Facebook and Twitter are the new factory floor, and online users, who leave digital crumbs as they browse the web and tap into social networks, generate big data that can be bought and sold. Every tweet tweeted, badge unlocked, website searched and "Like" button clicked add to the growing inventory of user information. This rich source of social data is a great point for social data mining because of its inherent openness for public consumption, clean and well-documented API, rich developer tooling, and broad appeal to users from every walk of life. Twitter data is particularly interesting because Tweets happen at the "speed of thought" and are available for consumption as they happen in near real time, represent the broadest cross-section of society at an international level, and are so inherently multifaceted. Tweets and Twitter's "following" mechanism link people in a variety of ways, ranging from short but often meaningful conversational dialogues to interest graphs that connect people and the things that they care about.

One behavior that has the most frequency and periodicity is drinking. Through the social network, people are easy to show their emotions and feelings. Such information as positive or negative opinions about drinking, preference platform and locations of drinkers can be obtained.

Tracking people's drinking behavior from Twitter data allows government not only to identify effect of certain event on the public but also to take actions to prevent unexpected traffic accident. Furthermore, in a world of endless information sharing, consumers have become the product. Data mining then can be applied to sort it, package it, market it — and companies use it to better target wine customers.

52

53 2 Algorithm Description

54

55 2.1 Latent Dirichlet Allocation

56 1) Basic idea

57 In natural language processing, Latent Dirichlet Allocation (LDA) is a generative model that
58 allows sets of observations to be explained by unobserved groups that explain why some parts
59 of the data are similar. For example, if observations are words collected into documents, it
60 posits that each document is a mixture of a small number of topics and that each word's
61 creation is attributable to one of the document's topics. LDA is an example of a topic model and
62 was first presented as a graphical model for topic discovery by David Blei, Andrew Ng, and
63 Michael Jordan [2].

64 In LDA, each document may be viewed as a mixture of various topics. This is similar to
65 probabilistic latent semantic analysis, except that in LDA the topic distribution is assumed to
66 have a *Dirichlet prior*. In practice, this results in more reasonable mixtures of topics in a
67 document.

68 For example, an LDA model might have topics that can be classified as CAT_class and
69 DOG_class. A topic has probabilities of generating various words, such as milk, meow, and
70 kitten, which can be classified and interpreted by the viewer as "CAT_class". Naturally, the
71 word cat itself will have high probability given this topic. The DOG_class topic likewise has
72 probabilities of generating each word: puppy, bark, and bone might have high probability.
73 Words without special relevance, such as the "bird", will have roughly even probability
74 between classes (or can be placed into a separate category).

75

76 2) Mathematic model

77 With plate notation, the dependencies among the many variables can be captured concisely.
78 The boxes are "plates" representing replicates. The outer plate represents documents, while
79 the inner plate represents the repeated choice of topics and words within a document. M
80 denotes the number of documents, N the number of words in a document. Thus:

81 α is the parameter of the Dirichlet prior on the per-document topic distributions,

82 β is the parameter of the Dirichlet prior on the per-topic word distribution,

83 θ_i is the topic distribution for document i,

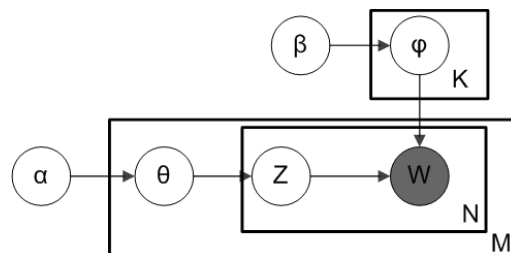
84 ϕ_k is the word distribution for topic k,

85 z_{ij} is the topic for the jth word in document i,

86 w_{ij} is the specific word

87 The w_{ij} are the only observable variables, and the other variables are latent variables.
88 Mostly, the basic LDA model will be extended to a smoothed version to gain better results.
89 The plate notation is shown in the Fig.1, where K denotes the number of topics considered in
90 the model and:

91



92

93

Fig.1 Plate notation representing the LDA model

94

ϕ is a $K \times V$ (V is the dimension of the vocabulary) Markov matrix each row of which

95 denotes the word distribution of a topic.
96 The generative process behind is that documents are represented as random mixtures over
97 latent topics, where each topic is characterized by a distribution over words. LDA assumes
98 the following generative process for a corpus D consisting of M documents each of length N_i :
99 Step1. Choose $\theta_i \sim \text{Dir}(\alpha)$, where $i \in \{1 \dots M\}$ and $\text{Dir}(\alpha)$ is the Dirichlet distribution for
100 parameter α ;
101 Step2. Choose $\phi_k \sim \text{Dir}(\beta)$, where $k \in \{1 \dots K\}$ and $\text{Dir}(\beta)$ is the Dirichlet distribution for
102 parameter β .
103 Step3. For each of the word positions i, j , where $j \in \{1 \dots N\}$, and $i \in \{1 \dots M\}$
104 (a) Choose a topic $z_{ij} \sim \text{Multinomial}(\theta_i)$.
105 (b) Choose a word $w_{ij} \sim \text{Multinomial}(\phi_{z_{ij}})$.
106 Note that the Multinomial distribution here refers to the Multinomial with only one trial. It is
107 formally equivalent to the categorical distribution.
108

109 2.2 Sentiment analysis

110 Sentiment is the attitude, opinion or feeling toward something, such as a person,
111 organization, product or location. Sentiment analysis (also known as opinion mining) refers
112 to the use of natural language processing, text analysis and machine learning techniques to
113 identify and extract subjective information in source materials^[3].

114 Sentiment = feelings / Attitudes / Emotions / Opinions.

115 It is a *subjective impressions*, not *facts*.

116 For example, some typical scenarios that Sentiment analysis can apply are:

- 117 • is this product review positive or negative?
- 118 • based on a sample of Tweets, how are people responding to this ad campaign/product
119 release/news item?
- 120 • how have bloggers' attitudes about the president changed since the election?

121 In our project, the drinking behavior Tweets are classified by two sentiment categories:
122 positive and negative. The “positive” sentiment represents the drinker was in happy/excited
123 mood; otherwise, angry/upset/disappointed for “negative” category.

124 For example, “Hi, guys, I got a job offer today, let’s celebrate it and have a drink!” reflects
125 positive mood of the Twitter user. “Oh my Gosh, my Brazil lost 0-7 to German at our fifth
126 drink round...” indicate he/she was shocked and disappointed by the soccer result, thus
127 belongs to “negative” category.

128 The sentiment analysis is based on the LDA algorithm and human-assisted pre-classification,
129 i.e., we pre-define a list of highly-representative words dictionary for positive and negative
130 sentiments, then apply this dictionary during the LDA algorithm.

131

132 3 Implementation

133 This section implements the tracking algorithm of drinking behavior. Roughly, the method
134 involves three steps: 1) generate corpus for LDA, 2) pick drinking related words from LDA,
135 and 3) extract drinking related Tweets from the raw Twitter data set. The data set we play on
136 is an about 20GB Twitter data, most of which are related to alcohol, but not all of them are
137 tightly concern to drinking. Therefore, the 20GB data set need to be refined into a smaller
138 and more drinking specific data set. Then, a series of analysis, such as sentiment estimation
139 and statistical analysis, can be performed in parallel using the Hadoop.

140

141

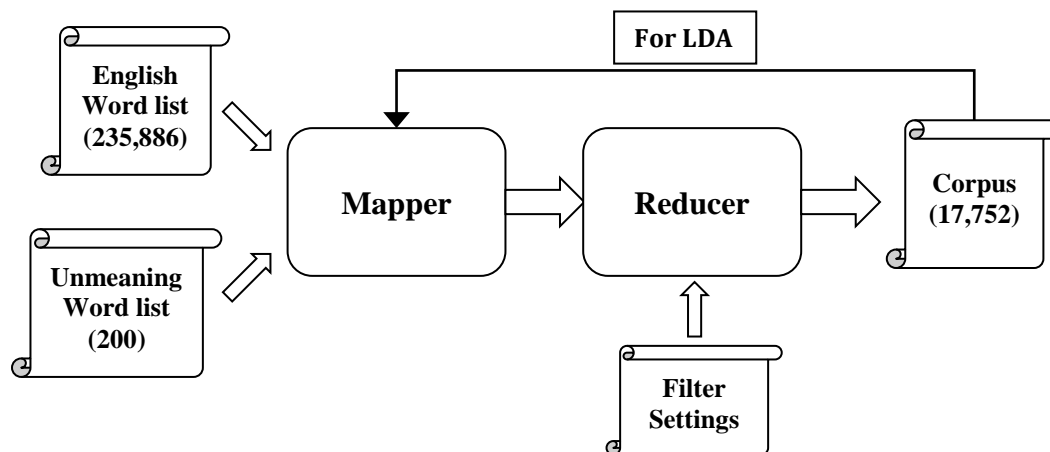
142 **3.1 Data preparation**

143 The data set we use is Twitter data from March 27th to May 1st in 2014. This is a pre-filtered
144 data set, which is mostly related to alcohol. Total account of the Tweets is around 600
145 million. When we check into the data, not all Tweets are so related to alcohol or drinking. So,
146 extracting drinking related Tweets is necessary to squeeze valuable parts from the raw data
147 set. The raw Twitter data is stored in zipped JSON format on the distributed file system
148 (DFS) of Hadoop. Noted that it is unnecessary to unzip those zipped JSON files because
149 Hadoop has unzip them automatically. In addition, Hadoop reads those files line by line, thus
150 we need to check the JSON format for each line first and then parse it using the GSON
151 package provided by Google.

152
153 **3.2 Developing environment**

154 Given the relatively large data set (about 600 million Tweets), we shrink it into a much
155 smaller subset (over 11 million Tweets) by removing those Tweets less related to drinking^[4].
156 This subset focuses more on drinking and yields more reliable analytical results. LDA is
157 applied here to provide drinking related topics, as well as relevant words. However, LDA
158 needs a corpus as initial input, thus we have to go through the raw data set first to collection
159 a corpus that is desired to be as small and representative as possible. In order to obtain more
160 meaningful words, we compare each word from the Tweets with the English word list and
161 unmeaning word list, respectively. The English word list consists of 235,886 English words,
162 and the unmeaning word list stores those words that may not contribute to topic splitting,
163 e.g., “you”, “is”, “on”, “will”, etc. The two word list work like filters, which filter
164 non-English and unmeaning word in the mapper phase and pass those meaningful English
165 words to the reducer phase. Fig.2 demonstrates the process of collecting corpus. The filter
166 settings in reducer phase filters those words with low frequency. Obviously, a word only
167 occurs several times in a huge amount of Tweets contributes little to topic categorization.
168 The corpus eventually obtained involves 17,752 words, all of which are throw into LDA
169 using mapper and reducer again to estimate the hyper parameters [4], namely α , β and γ .

170



171

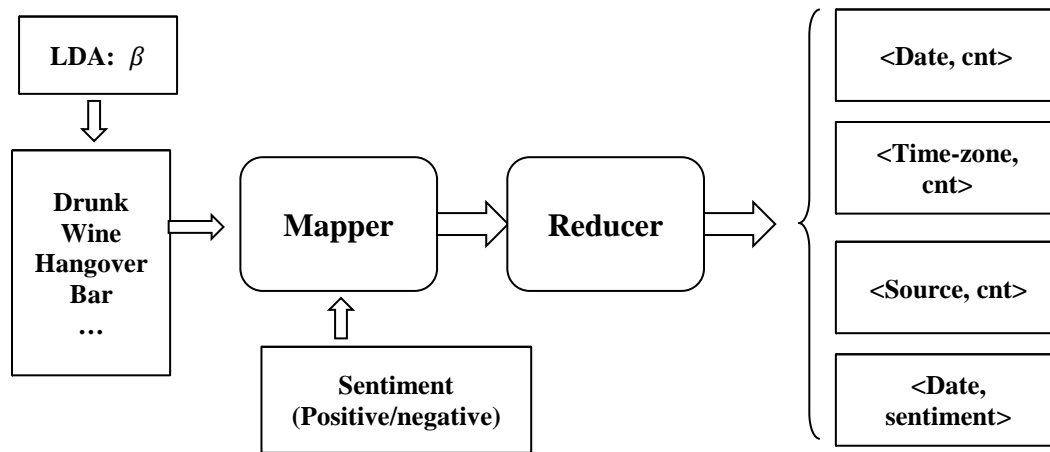
172 Fig.2 Flowchart of collecting corpus for LDA using Hadoop (mapper and reducer)

173

174 Since β describes the relation between topics and words, we pick those word more related
175 to the drinking topic according to β . As list in Fig.3, drunk, wine, hangover and bar are
176 words related to drinking with relatively high probability. Based on the drinking related
177 words, we go through the whole raw data set again to extract the Tweets including one or
178 more drinking related words. At the same time, a sentiment word list is provide to separate
179 positive sentiment from negative ones. All of above is implemented in mapper, and the
180 reducer finally output four sets of data files recording created date of a Tweet, time-zone,
181 source and sentiment. Specifically, <date, cnt> records the count of drinking related Tweets

182 on each corresponding date. By the same token, <time-zone, cnt> and <source, cnt> record
 183 the count of drinking related Tweets for corresponding time-zone and source (platform a
 184 Tweet is sent from), respectively.

185



186

187

188 Fig.3 Flowchart of extracting drinking related Tweets using Hadoop (mapper and reducer).

189

190 The output data file <date, sentiment> stores the count of Tweets with positive or negative
 191 sentiment on each date. Actually, this data file is contracted by two parts---<date, positive
 192 sentiment> and <date, negative sentiment>.

193 In practice, convergence of LDA require tens of iterations which may cost a couple of days.
 194 So we only iterated five times and then pick up drinking related word manually based on the
 195 hyper parameter β . In addition, the sentiment word list is borrowed from some existing
 196 works. In sentiment estimation, both unigram and bigram are employed. Simply speaking, we
 197 try to find isolate and adjacent sentiment words in a Tweet. Isolate sentiment word (unigram)
 198 refers to a word not connected with any other sentiment word; adjacent sentiment words
 199 (bigram) refer to two continuous sentiment words, for example, “don’t like” and “never
 200 hate”. The first example is in “negative + positive” format, so it yields a negative sentiment.
 201 The second example, however, express a positive sentiment since it is in the format of
 202 “negative + negative”. Similarly, the format of “positive + positive” should yield positive
 203 sentiment. For unigram, a signal word represent a sentiment.

204 In a Tweet, it may involves both positive and negative sentiment words. We use a weighted
 205 sentiment to give the final sentiment of the Tweet. First of all, sentiment of a Tweet is
 206 quantified from -1 to 1, where -1 denotes negative and 1 denote positive. Currently, we just
 207 simply set the weight of positive sentiment word as 1 and negative word as -1. If a positive
 208 word(s) occur in the Tweet, plus one to the overall sentiment; if a negative word(s) appears,
 209 subtract one from the overall sentiment. The final overall sentiment is considered as the
 210 sentiment of a Tweet. In a more professional way, the weight of each sentiment word should
 211 vary depending on how positive/negative they are.

212

213 **4 Results and analysis**

214

215 **4.1 Sentiment distribution**

216 The distribution of positive and negative drinking related Tweets are shown in the Fig.4.

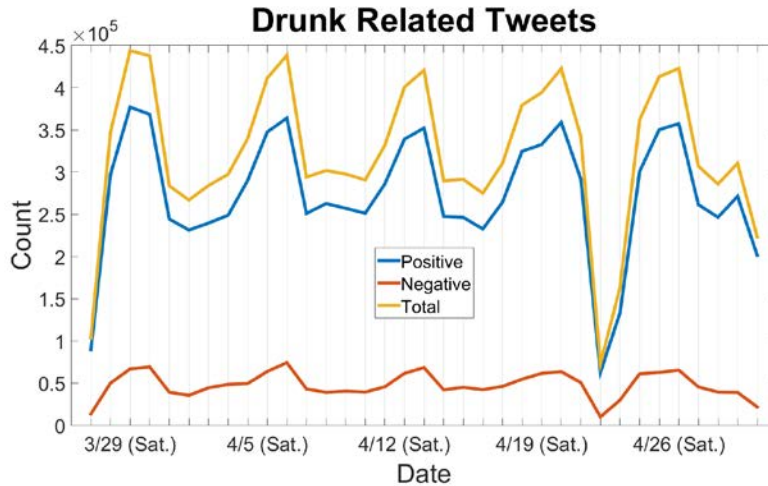


Fig.4 The distribution of positive and negative drinking related Tweets

217

218

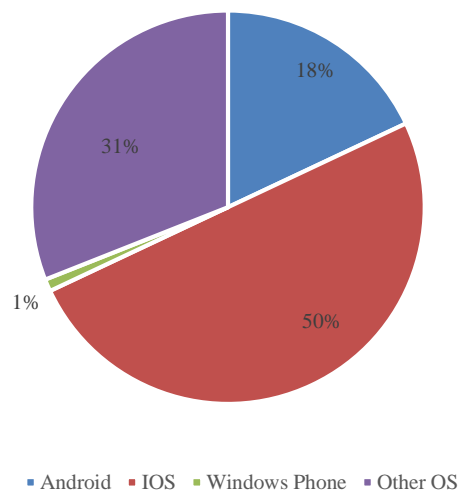
219 We can draw three conclusions from this plots:

- 220 1) The peaks of drinking –related Tweets happens most frequently during weekend,
221 especially Saturday.
- 222 2) Positive Tweets are among the majority of the total drinking related Tweets.
- 223 3) There is one abnormal drops between 04/19 and 04/26. At the beginning, we doubt
224 there may be some “abnormal” events happened on that day. But after searching
225 Google for that “abnormal” period, we cannot find any obvious evidence or valuable
226 information related to this “abnormal drop”. Finally, after going back to the original
227 data set, we found that the reason is very simple, there is one day missing in the
228 data. Thus, this provides us a “by-product” way to discover any missing data.

229

230 4.2 User platform distribution

231 Fig.5 shows the Source of drunk related Tweets. Nearly half of the drunk related Tweets are
232 sent by IOS platform including iPhone, iPad and MacBook. 18% of drunk related Tweets are
233 sent by Android users. Less than 1% of related Tweets came from Windows Phone. Other
234 kinds of software like Instagram and web page generated almost 30% of the related Tweet.



235

236

Fig.5 Source of drunk related Tweets

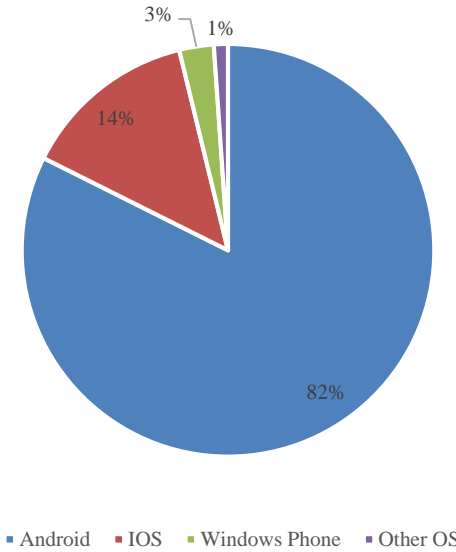


Fig.6 2014 Worldwide Smartphone Market Share (IDC)

237
238

239 In comparison to Fig.6 which shows 2014 Worldwide Smartphone Market Share from
240 International Data Corporation (IDC), the Android platform is the majority component of the
241 markets. The conclusion that people who drink like iPhone more than other brand mobile
242 phones can be draw.

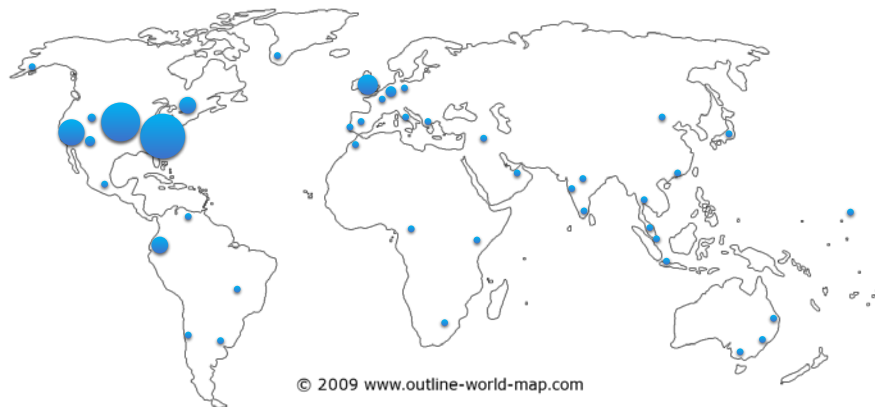
243 At the same time, several studies and some evidence can support this idea. According to
244 research from venture capital firm Battery Ventures, there is some merit to the idea that
245 iPhones are used by a white collar crowd, while Android favors the blue collar set. Extracting
246 information from a set of survey questions, iPhone users are more likely to have flown in an
247 airplane in the past year, drink wine and have investments in the stock market. Android users,
248 on the other hand, take public transportation, prefer beer, consider themselves religious and
249 have eaten at McDonalds.

250 Then, the key word “wine” is searched in both App Store and Google Play. There are 2985
251 Results of “wine” in App Store but only 248 Results of “wine” in Google play. It is obvious
252 that wine or alcohol promotion can be more effective when conducted through people who
253 have iPhones. The same goes to designated driving market.

254
255

4.3 Location distribution

256 The distribution of locations from which Tweets were sent are shown in Fig.7 and Table 1.



257
258

Fig.7 Location distribution of drinking related Tweets

Table 1 46 zones over 10,000 drinking related Tweets

Location	Num.	Location	Num.
'Eastern/Time/(US/&/Canada)'	1349324	'Casablanca'	187121
'Central/Time/(US/&/Canada)'	1141305	'Alaska'	129952
'Pacific/Time/(US/&/Canada)'	783803	'Athens'	102358
'London'	583376	'Brasilia'	92126
'Atlantic/Time/(Canada)'	497300	'Beijing'	74283
'Quito'	490727	'Greenland'	52823
'Amsterdam'	325605	'Chennai'	51933
'Arizona'	291043	'Bangkok'	49231
'Mountain/Time/(US/&/Canada)'	234895	'Edinburgh'	44594
'Hawaii'	198108	'Buenos/Aires'	43854
'Dublin'	38960	'Madrid'	31394
'Singapore'	38276	'Kuala/Lumpur'	29112
'Santiago'	37012	'Paris'	27678
'Sydney'	31784	'Pretoria'	26745
'Mexico/City'	26023	'Caracas'	25682

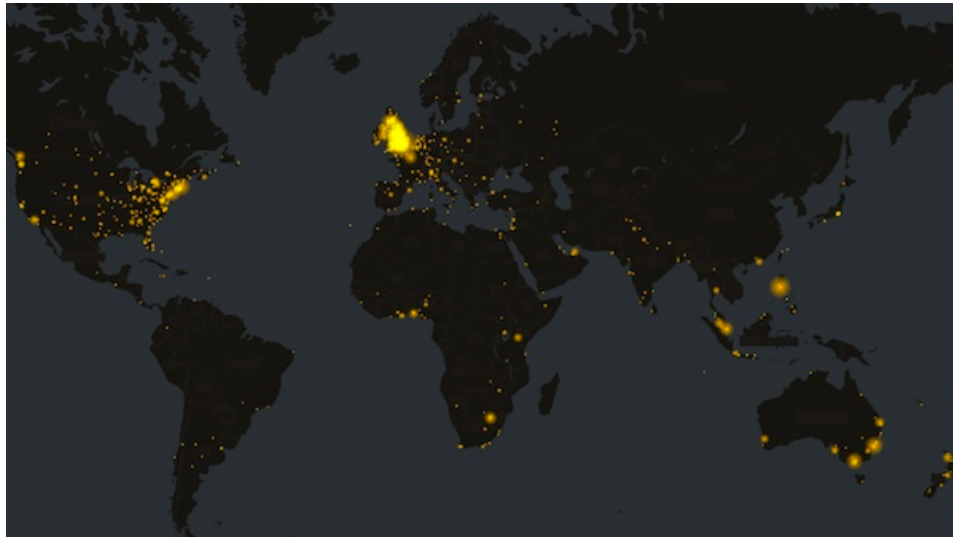


Fig.8 Location distribution of Twitter users

From the Fig.8, one natural conclusion seems to be inferred is, the majority of the alcohol drinkers are from the United States and West Europe. But after serious inspection, we found that this conclusion does not hold rigorously. Because from the Twitter uses distribution (yellow dots), it can be seen that the registered Twitter used located in the U.S. and west Europe. Thus, the blue dots may be not a convincing indicator to support the above conclusion. However, there is still promising application of those location information, for example, the wine manufacturer can produce more wines in advance for some specified location (country/state/province/city) if that place recently has a relatively higher distribution of drinking-related Tweets compared with the remaining area.

273 **4.4 Potential application**

274 Include but not limited to:

- 275 • Public safety: e.g. pre-warning to some tourists.
276 For example, the government/police department can send warning message/notice to
277 tourists who plan to visit an area with high drinking Tweets distribution.
- 278 • Alcohol related business promoting, e.g. giving priority in developing iPhone apps;
279 give priority in advertising investment for those high dinking Tweets areas;
280 increasing stocking amount for those areas, etc.

281

282 **5 Summary**

283 In this project, the drinking behavior of people is analyzed from Twitter by big data mining
284 technic.

285 First of all, 60 million Twitter data is obtained. Then LDA and Human method are applied to
286 get the Mapper and Reducer. Third, 11,255,207 drunk related Tweets are used to do the
287 analysis. In the end, drinking behavior of sentiment, time, source and location is presented.

288 Through big data mining, we may find something hard to be recognized in daily life and
289 identify effect of certain event on the public. At the same time, related marketing can be
290 more targeted. Wine companies may be able to see and predict what their customers like,
291 share, and mention most.

292 But how to interpret the result of data mining is remained to be answered. Because different
293 interpretation may cause different understanding or even misunderstanding. Another issue is
294 that the extracted behavior features can only represent the pubic trend rather than individual
295 characteristics. Both of them can be left as future study topics.

296

297

298 **Acknowledgments**

299 Our team would like to sincerely appreciate the instructor, Dr. Arvind Ramanathan's help and
300 effort in providing us the Tweets data and creating the basic environment of Hadoop in
301 department server for our study convenience.

302 **References**

303 [1] Blei, David M., Ng, Andrew Y., Jordan, Michal I. (2003) Latent Dirichlet Allocation. Journal of
304 Machine Learning Research 3 (2003) 993-1022.

305 [2] Longbing Cao, Philip S. Yu. (2012) Behavior Computing: Modeling, Analysis, Mining and
306 Decision. Springer-Verlag London.

307 [3] Xiaoran An, Auroop R. Ganguly, Yi Fang, et al. (2014) Tracking Climate Change Opinions from
308 Twitter Data. <http://www.cse.scu.edu/~yfang/climate-fang.pdf>

309 [4] Zhai, K., Boyd-Graber, J., Asadi, N., & Alkhouja, M. L. (2012, April). Mr. LDA: A flexible large
310 scale topic modeling package using variational inference in mapreduce. In *Proceedings of the 21st*
311 *international conference on World Wide Web* (pp. 879-888). ACM.