

Location Inference for Hidden Population with Online Text Analysis

1 Liu Chuchu^{1,a)*}, Cao Ziqiang^{1,a)}, Lu Xin^{1,2,*}

2 ¹ College of Systems Engineering, National University of Defense Technology, Changsha 410073,
3 China

4 ² School of Software Engineering, Shenzhen Institute of Information Technology, Shenzhen 518172,
5 China

6 ^{a)} Liu C. and Cao Z. contributed equally to this work

7 * **Correspondence:**

8 Liu Chuchu; Lu Xin

9 liuchuchu15@nudt.edu.cn; xin.lu@flowminder.org

10 **Keywords: location inference, hidden population, MSM, text analysis, geographic distribution.**

11 Abstract

12 Understanding the demographics of hidden population, such as men who have sex with men (MSM),
13 sex workers, or injecting drug users, are of great importance for the adequately deployment of
14 intervention strategies and public health decision making. However, due to the hard-to-access
15 properties, e.g., lack of a sampling frame, sensitivity issue, reporting error, etc., traditional survey
16 methods are largely limited when studying such populations. With data extracted from the very active
17 online community of MSM in China, in this study we adopt and develop location inferring methods to
18 achieve a high-resolution mapping of users in this community at national level. The performances of
19 popular inference algorithms are compared to elucidate the most suitable approach. In addition, we
20 propose a new hybrid model, which is proven to achieve the highest accuracy for inferring locations
21 of online users only based on text content. This method is conducive to overcoming the sparse location
22 labeling problem in user profiles, and can be extended to the inference of geo-statistics for other hidden
23 populations.

24 1 Introduction

25 A population is “hidden” when no sampling frame exists and public acknowledgment of membership
26 in the population is potentially threatening (1-3). Examples of hidden populations include men who
27 have sex with men (MSM) (4-6), sex workers (SW) and injecting drug users (IDU). To date, the study
28 of hidden populations has mainly focused on interviews and questionnaire surveys based on offline or
29 online population sampling. While in most cases, these traditional methods are inefficient, limited in
30 sample size and representativeness, and challenged by privacy concerns and reporting errors (7-10).
31 Besides, with concerns about sensitivity and privacy, hidden populations tend to conceal their personal
32 information, including their locations. There are many difficulties in conducting comprehensive studies
33 on demographic characteristics of hidden populations, especially on their geographic distribution.

34 According to recent statistics, approximately 58.8% of the world’s population now use the internet
35 (11), and by 2020, there will be around 30 billion devices connected to each other (12). Nowadays, our
36 daily lives are inseparable from the internet. The numerous data generated on the internet provides
37 opportunities to infer demographic attributes of internet users with computational techniques at an

38 unprecedented scale. Due to social discrimination, hidden populations lack reliable channels for
39 communication and tend not to disclose their information in the real world. Instead, the anonymity of
40 internet provides a good sense of security and has made online social networking prevailing among
41 hidden populations (13, 14), offering alternatives for understanding demographic statistics (gender,
42 age, location, etc.) of traditionally hard-to-access populations.

43 Location-based services (LBS) are incredibly useful across many domains, including personalized
44 services (e.g. local restaurants, hospitals, events), prompting alert, assessment and emergency
45 responses to disease or disasters, as well as detecting security intrusion, etc. (15). However, it is still
46 challenging to obtain user location due to the sparse geo-enabled features in social media. Although
47 users from social networks can fill up profiles with their personal information, the use of these data is
48 however limited as it may be subject to large reporting error and many users may opt to make their
49 personal information hidden to the public. It is reported that, on average only 35% of Facebook users
50 declare locations (16), and a large volume of invalid or low precision locations are often submitted.
51 Regarding other ways of geolocation retrieving, user location given by the IP address is not reliable
52 and needs to be continually updated (17). GPS provides locations with best accuracy and reliability,
53 while many people do not want to disclose their detailed coordinates. For example, in Brazil, under
54 1% of tweets provide GPS data (18). Although there has been a large number of studies on location
55 inference with internet data for general population, e.g., through platforms of Twitter (19, 20) Facebook
56 (21, 22), and Flickr (23, 24), there is however very limited studies with such applications on hidden or
57 hard-to-access populations, as large-scale corpus of hidden populations is generally lacking, and the
58 applicability of algorithms is not known when they are to be generalized to a new population.

59 As a representative of hidden populations, MSM suffers more social pressure and discrimination than
60 the general population in many cultures, and the demographics of them, especially the geographic
61 distribution are of critical importance for public health management and HIV prevention (25, 26).
62 However, it is still difficult to locate the MSM population. To date, relevant study of this population
63 has mainly relied on sampling estimations based on questionnaire surveys. We thereby take MSM as
64 an example in this study, to test and develop appropriate location inferring approaches based on large
65 corpus of online data. We collected a comprehensive dataset in this study, covering 628,360 MSM-
66 related users from the largest sub-community related to MSM topics, gay-bar, in the world's largest
67 Chinese community (i.e., Baidu Tieba (27)). With observations in gay-bar, a user's location is more
68 likely to be mentioned in his/her own posts. Based on users' publicly available posts, we develop a
69 new hybrid method to infer the geographic locations for online MSM. Compared with other approaches
70 that consider only text content, the hybrid method proposed in this paper can achieve the best
71 performance, improving the inferring accuracy from 50.3% to 71.3%.

72 **2 Material and Methods**

73 **2.1 Data collection**

74 Baidu Tieba is the largest Chinese online community and it consists of a variety of sub-communities
75 on different topics and gathers massive numbers of user groups with different interests. The gay-bar is
76 the largest sub-community related to MSM in Baidu Tieba, which serves as a community for
77 homosexual friends to seek partners and chat. As of June 14, 2018, gay-bar had 4.67 million followers
78 and more than 300 million posts. The massive user data generated in this open community is of great
79 significance for analyzing the demographic attributes of the MSM population in China. In this study,
80 we developed a web crawler to collect from the latest to the oldest possible posts, as well as all

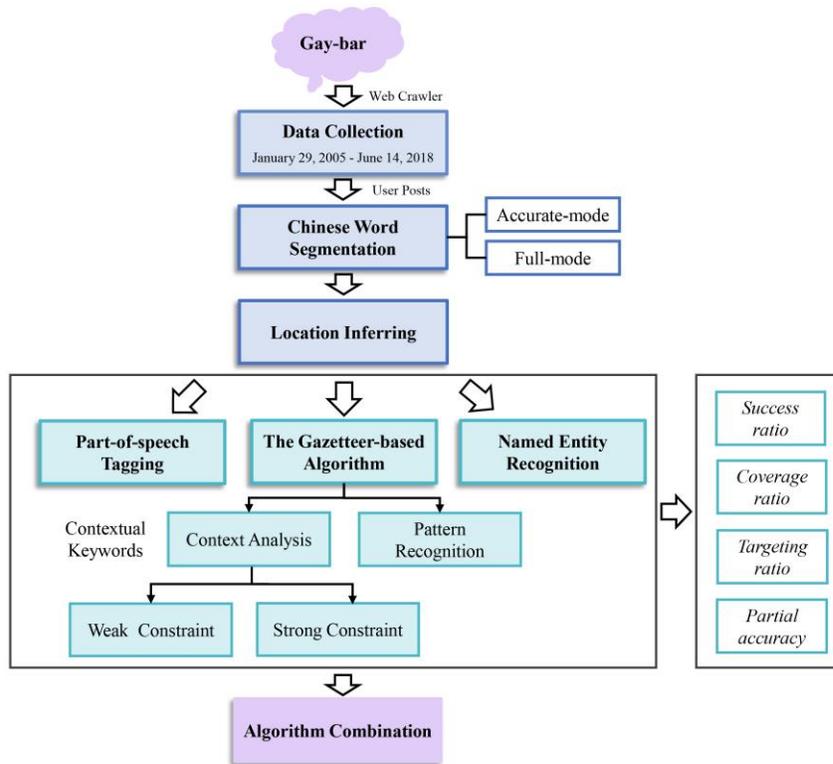
81 following comments and replies under these posts, for a total of 13,023,367 records and 628,360 unique
82 users. All the data was distributed from January 29, 2005, to June 14, 2018.

83 We also obtained public profile data of active users. A total of 359,438 records was collected, including
84 user name, sex, self-reported location, latitude and longitude coordinate located by GPS, etc. However,
85 among all users, only 8.6% (30,948) declare the location field in their profiles, indicating that over
86 91% gay-bar users did not fill in their locations when registering. In addition, the authenticity and
87 accuracy of these reported locations needs to be considered as large reporting error may occur due to
88 the privacy concerns of users. In addition, when the user is using a mobile device, he can opt to reveal
89 his GPS coordinates publicly. Among all 359,438 users, we find that about 10.9% had doing so, adding
90 to a complementary dataset for algorithm validation.

91 **2.2 Location inferring method with online text analysis**

92 It is a conventional and direct approach to obtain user locations through user profiles or GPS
93 coordinates of mobile devices, however very few data are available from these two sources. Due to the
94 privacy concerns, MSM-related users are more reluctant to fill in the location field and refuse to give
95 applications access to their GPS information. Among all records collected this time, only about 10%
96 of user location data (including GPS coordinates) is not empty, which is significantly lower than
97 popular platforms, such as 35% in Facebook (16).

98 Observing the text content posted by gay-bar users, we find that the subject of most posts in this
99 community is related to offline dating. Most users tend to actively reveal their locations when posting
100 online in order to find friends with close distance, which plays a key role in meeting with each other
101 and finally transferring the relationship from the internet to the real world. Therefore locations revealed
102 in these posts can promise a great authenticity, and provide a new perspective for us to evaluate
103 locations of the MSM population as well as their geographic distribution. Based on the text analysis of
104 gay-bar posts, we evaluate and compare the performances of mainstream location inference algorithms
105 to solve the online locating problem of Chinese MSM population. To improve the inference accuracy,
106 we propose a new hybrid method, which integrates different algorithms by voting on the inference
107 results, and is proven to guarantee a higher accuracy (71.3 %) on the location inferring of MSM in
108 online social communities. The workflow and links between algorithms are presented in Figure1.



109

110

Figure 1. The workflow and links between location inferring algorithms.

111 2.2.1 Training and test data

112 In this study, all user posts collected from gay-bar (including posts, comments and replies), as the
 113 corpus input of the inferring algorithm, are used to extract specific location information. There are a
 114 total of 628,360 unique users in the dataset, of which 156,777 have provided geographic information
 115 in their posts at certain degree. However, in many cases these geolocation data are not about where
 116 they stay but discussions about POI (point of interests) such as tourist attractions. In order to effectively
 117 distinguish whether the geographic information appearing in a user's posts represents his actual
 118 location, we experiment with multiple approaches and compare effects of different algorithms with a
 119 testing set. The test data is composed of randomly selected 10,000 users' posts. And each user is
 120 manually labelled with his actual geographic location by reading the post content.

121 2.2.2 Location inferring algorithms

122 Current online-corpus location inferring algorithms mainly considers the text associated with each user,
 123 including traditional gazetteer-based method (identifying geographical names from the gazetteer) (28),
 124 part-of-speech tagging (recognizing geographical terms in a corpus based on the part of speech of
 125 words, according to both their definitions and contexts) (29) (30) and named entity recognition (NER)
 126 (identifying and classifying words in a text as pre-defined entity classes, i.e., persons, locations,
 127 organizations, quantities, etc.) (31, 32). Due to the posting characteristics of gay-bar users, we mainly
 128 infer locations at the province-level and city-level. In the gazetteer-based location inferring method,
 129 we use the popular Chinese word segmentation module, Jieba (33), which is more suitable for Chinese
 130 text analysis, to cut the posts into the most accurate segmentations. We also consider unigrams of the
 131 post text, and remove all punctuation, stop words and URLs. In this method, the gazetteer words used
 132 to match the user location belong to the geographical gazetteer of China (34), where province-level
 133 gazetteer words contain provinces, municipalities, autonomous regions, Hong Kong and Macao; city-

134 level gazetteer words contain prefectural-level city, municipalities, Hong Kong and Macao. Besides,
 135 in the part-of-speech tagging method, the part-of-speech recognition function of Jieba is employed,
 136 where the geographical terms are identified as a particular label ("ns"). And the named entity
 137 recognition method is realized by using the Chinese geographic name recognition based on the HMM-
 138 Viterbi algorithm (35).

139 2.2.3 Improvement on existing algorithms

140 When a user mentions a geographic name in the post, we find that it does not necessarily refer to the
 141 location where the user is or has been or will go, that is, his track. To improve the accuracy of the
 142 inference algorithm, we try to combine other approaches in natural language processing with the
 143 location extraction, such as context analysis and pattern recognition.

144 2.2.3.1 Context keyword analysis

145 The context keyword analysis is added to determine whether the geographic information in the post
 146 refers to a user's track. In this paper, eight Chinese keywords, i.e., ‘坐标’ ('coordinate' in English),
 147 ‘定位’ ('location'), ‘在’ ('in'), ‘是’ ('am'), ‘同’ ('same'), ‘求’ ('seek'), ‘人’ ('from'),
 148 and ‘交友’ ('dating') are selected to filter the geolocation words. The appearance of a context
 149 keyword in the post is then considered as an improved likelihood of referring actual location for the
 150 user.

151 2.2.3.2 Pattern recognition on post sentence

152 In order to further restrict the syntax patterns in user corpus and strengthen the filtering rules, the idea
 153 of pattern recognition is introduced. According to the expression characteristics of gay-bar users, 7
 154 typical modes are defined, i.e., keywords before location (former keywords), keywords after location
 155 (later keywords), global keywords, individual location, individual location with punctuations or
 156 symbols, location with modal particle and province with city. These modes can cover most of syntax
 157 that MSM use when referring to their locations. Keywords employed in pattern recognition are shown
 158 in Table 1.

159 Table 1. Chinese keywords employed in pattern recognition.

	Former keywords	Later keywords	Global keywords
In Chinese	坐标, 定位, 同, 在, 从, 去, 是, 也是, 求, 就是, 大...	人, 上学, 上班, 有, 有吗, 附近, 的, 滴, 是, 加...	交友, 大学, 学院, 公司, 同城, 私聊...
In English	Coordinate, location, same, in, from, come, am, also be, seek, big (usually describe someone's own place), etc.	Person, go to school, work, have, any, close, is, add, follow, etc.	Making friends, university, college, company, same city, private chat, etc.

160

161 2.2.4 Algorithm measurements

162 Since a user's location may consist of several geographic places, such as his hometown, the province
 163 or city where he works or studies. Meanwhile the same user is likely to migrate to different places at

164 different time. A set is used to record locations of each user. Therefore we cannot simply measure the
165 algorithm results with the precision and recall used in binary classification problems. In order to
166 comprehensively judge the performances of different approaches, four new indexes are defined in this
167 paper, namely absolute accuracy (success rate, S for short), coverage ratio (C), targeting ratio (T), and
168 partial accuracy (P).

169 The absolute accuracy measures whether the algorithm result is exactly the same as the manually
170 labeled result, $S = N_s/N_{all}$, where N_{all} is the number of all test users, and N_s is the number of users
171 whose locations are accurately inferred.

172 The coverage ratio measures the comprehensiveness of an algorithm. When the results from manual
173 annotation are completely included in the inference results, the algorithm are believed to achieve a
174 complete coverage. It is defined as $C = N_c/N_{all}$, where N_c is the number of users whose locations are
175 totally excavated by the algorithm.

176 The targeting ratio is used to determine whether all geographic words recognized by the algorithm are
177 correct, $T = N_T/N_{all}$, where N_T is the number of users whose inferred results from the algorithm are
178 all correct even if some manual labels may not be covered.

179 The partial accuracy is used to judge whether inference locations have any intersections with manual
180 labels, $P = N_p/N_{all}$, where N_p is the number of users whose locations are partly inferred.

181 By comparing the inference set evaluated by the algorithm with the manual label set, these four
182 measurements are employed to determine the algorithm with the best performance. S is the most
183 important indicator used to measure the algorithm effect.

184 **3 Results**

185 **3.1 Performances of different text-based location inferring algorithms**

186 In this study, all online posts from 156,777 users who mentioned geographic information are used as
187 the corpus input in location inferring algorithms. As the first step, three mainstream algorithms are
188 applied to determine the most suitable approach regarding the location inference for online hidden
189 population. The performances of three different approaches, i.e., the gazetteer-based method, part-of-
190 speech tagging and Chinese NER, are compared. Accuracy of the inference results from different
191 algorithms is shown in Table 2.

192 As we can see, the gazetteer-based method achieves the highest accuracy on all measurements,
193 suggesting that it is more suitable for extracting the location information from short texts, e.g., user
194 posts on gay-bar. Other algorithms, such as part-of-speech tagging and NER, which are more widely
195 used in location inferring from Chinese texts nowadays, are not so effective than the traditional
196 gazetteer-based method instead. The latter proves to be the simplest, fastest, and most effective.

197 However, the success rate (S) obtained by the gazetteer-based algorithm is only 0.503, which means
198 that only 50.3% users' locations are fully identified without any errors. In order to further improve the
199 performance of this algorithm, context analysis and pattern recognition is introduced to the gazetteer-
200 based method. From Table 3, we can see that the addition of contextual keyword analysis can improve
201 the success rate and targeting ratio of the algorithm, whereas the stricter rules of pattern recognition do
202 not achieve a good performance. This is because that the latter method defines more specific grammar,

203 syntax, keywords as well as keyword positions to filter text, with more restrictions on user posting.
 204 Due to the variety of Chinese expressions, especially in internet language, the syntax patterns used by
 205 online MSM are so flexible, leading to a lower accuracy of the algorithm with pattern recognition.

206 Overall, for the social network whose users mainly focus on making friends or dating, such as gay-bar,
 207 there are numerous obvious geographic information exposed in the short text of user posts. And the
 208 gazetteer-based method combined with the contextual keyword analysis is more effective in user
 209 location inference, by which over 51.2% users' tracks can be absolutely correctly recognized.

210 Table 2. The performances of three location inferring algorithms.

	Gazetteer	Part-of-speech	Chinese NER
<i>S</i>	0.503	0.352	0.487
<i>C</i>	0.932	0.892	0.927
<i>T</i>	0.518	0.370	0.502
<i>P</i>	0.966	0.945	0.964

211 Table 3. The performance of the gazetteer-based method after strengthen filtering rules.

	Gazetteer	Gazetteer with context analysis	Gazetteer with pattern recognition
<i>S</i>	0.503	0.512	0.493
<i>C</i>	0.932	0.929	0.733
<i>T</i>	0.518	0.528	0.667
<i>P</i>	0.966	0.965	0.800

212

213 3.2 Improvement of the gazetteer-based algorithm

214 The comparison above illustrates that the traditional gazetteer-based location inferring method with
 215 context analysis can achieve better performance in the gay-bar dataset. However, the highest accuracy
 216 is still at a relatively low level, therefore, we aim to further improve the algorithm by considering more
 217 conditions.

218 3.2.1 Different constraints on contextual keywords

219 In this section, we try to change the way which keywords are constrained to improve the algorithm
 220 performance. Two constraints (weak or strong) are mainly considered. The weak constraint is that as
 221 long as any keyword appears in any post of a user, all geographic words in his posts are considered to
 222 be his possible locations. For analysis above, the weak constraint on the context keywords is used by
 223 default. We attempt to replace the weak constraint with the strong constraint to explore whether the
 224 algorithm accuracy would be improved as a result. In detail, the geographic word must appear with any
 225 context keyword in a same post, then this geographical term can be considered to be possible user
 226 location. The performance of the algorithm with different constraints is shown in Table 4.

227 It can be seen that compared with weak constraint, the success rate (0.542) and targeting ratio (0.614)
 228 of the algorithm both increase after introducing the strong constraint, while the coverage ratio would
 229 decrease, i.e., the comprehensiveness of the algorithm results has reduced. The strong constraint of
 230 keywords is helpful to improve the correctness of location inference, which is more suitable for
 231 situations aiming at locating accuracy. Although the accuracy of weak constraint is relatively low, the
 232 algorithm can achieve more comprehensive results ($C = 0.929$). And user actual locations are mostly
 233 covered by the results of inference algorithm, which is more conducive to identifying all geographic
 234 tracks of MSM.

235 Table 4. The performance of the gazetteer-based method with strong constraint.

	Gazetteer	Gazetteer with context analysis (weak)	Gazetteer with context analysis (strong)
<i>S</i>	0.503	0.512	0.542
<i>C</i>	0.932	0.929	0.762
<i>T</i>	0.518	0.528	0.614
<i>P</i>	0.966	0.965	0.866

236

237 3.2.2 Expanding contextual keyword set

238 As the context analysis applied in the location inferring algorithm above only utilizes eight keywords,
 239 i.e., 'coordinate', 'location', 'in', 'am', 'same', 'seek', 'from' and 'date'. By further analyzing the posting
 240 characteristics of gay-bar users, other ten common keywords are added, including 'school', 'work',
 241 'friends', etc., for the purpose of covering more posts referring to user locations. The algorithm
 242 performance after expanding contextual keywords is shown in Table 5.

243 It shows that adding keywords has no significant effect on improving the accuracy of location inferring.
 244 Instead, since more contextual keywords tend to broaden the permission on syntax, constraints on
 245 geographic information are insufficient, which reduces the performance of the algorithm.

246 Table 5. The algorithm performance with keyword augmentation.

	Strong constraint	Strong constraint with more keywords	Strong constraint (Full-mode segmentation)
<i>S</i>	0.542	0.540	0.513
<i>C</i>	0.762	0.780	0.756
<i>T</i>	0.614	0.604	0.582
<i>P</i>	0.866	0.878	0.854

247

248 **3.2.3 Different word segmentation mode**

249 In this study, the word segmentation method used in the gazetteer-based algorithm is from Jieba, a
250 popular Chinese word segmentation package. All previous results are based on the accurate-mode of
251 Jieba, which aims at cutting the Chinese sentence most accurately. We attempt to change the word
252 segmentation method and examine whether the accuracy of the algorithm can be improved. In this
253 section, the full-mode word segmentation method (36) is used, which scans all possible words that can
254 be formed in a sentence. With observing the change of measurements, we compare the effect of the
255 two different word segmentation methods on the location inference, shown in Table 5.

256 The results prove that although the full-mode word segmentation can increase the number of words
257 from the text, it may mislead the location inferring and reduce the accuracy of the algorithm as the
258 result of the word ambiguity. Compared with the full-mode method, the accurate-mode word
259 segmentation is more helpful for the gazetteer-based location inference.

260 **3.3 Hybrid voting algorithm on location inference (HVA-LI)**

261 To further improve the accuracy of location inference for hidden population, while maintaining the
262 superiority of existing algorithms, we adopt the ensemble learning approach and develop a hybrid
263 voting algorithm, called HVA-LI, by allowing different approaches to vote to determine the best
264 inference results. The main goal of this hybrid method is by setting multiple filters to improve the
265 inference accuracy.

266 **3.3.1 Gazetteer-based algorithm and part-of-speech tagging (Gazetteer & PT)**

267 In addition to the gazetteer-based algorithm with context analysis (strong constraint), which achieves
268 the highest success rate, different inference algorithms are considered to be introduced to work
269 together. In this section, the most superior algorithm (S-Gazetteer, for short) tend to be combined with
270 the part-of-speech tagging algorithm (PT, for short). Both algorithms calculate each user's corpus to
271 extract possible locations and determine the final output by voting together. Locations that both
272 algorithms agree on are considered the final inference results.

273 **3.3.2 Gazetteer-based algorithm and NER (Gazetteer & NER)**

274 In order to further verify the superiority of the ensemble learning approach, the Gazetteer algorithm is
275 designed to effectively combine with the NER algorithm. Similarly, all final outputs are determined by
276 two algorithms' voting on their results. However, the difference is that in this section two different
277 combinations are carried out. Firstly, without the introduction of context analysis, a simple ensemble
278 of the basic gazetteer-based method (Gazetteer) and NER is employed as the basic benchmark. We
279 then add the context analysis with strong constraint, which is proven to be the best strategy to improve
280 the algorithm accuracy in the gazetteer-based methods, into the simple ensemble algorithm to examine
281 the changes in algorithm performance.

282 **3.3.3 Performances of the algorithm combinations**

283 The performances of different algorithm combinations is shown in Table 6. It can be seen that the
284 success rates of all inference algorithms have improved, compared with the best algorithm before
285 combination (S-Gazetteer, $S = 0.542$). The hybrid voting algorithm (HVA-LI) can be considered more
286 effective in location inferring for online MSM. Significantly, it can be seen that the ensemble of the
287 basic gazetteer-based method and NER can achieve a much higher accuracy, and the success rate for
288 location inference is up to 0.713, approximately 32% higher than the best algorithm without
289 combination. Surprisingly, however, the introduction of context analysis does not improve the

290 effectiveness of the algorithm. We find that the success rate ($S = 0.586$) decreases dramatically in
 291 comparison to the former simple combination. This may be because the contextual keywords destroy
 292 the integrity of post texts, hence degrades the performance of NER as well as the accuracy of the
 293 algorithm.

294 Table 6. The performance of HVA-LI.

	S-Gazetteer & PT	Gazetteer & NER	S-Gazetteer & NER
<i>S</i>	0.578	0.713	0.586
<i>C</i>	0.756	0.713	0.767
<i>T</i>	0.666	0.914	0.668
<i>P</i>	0.883	0.914	0.881

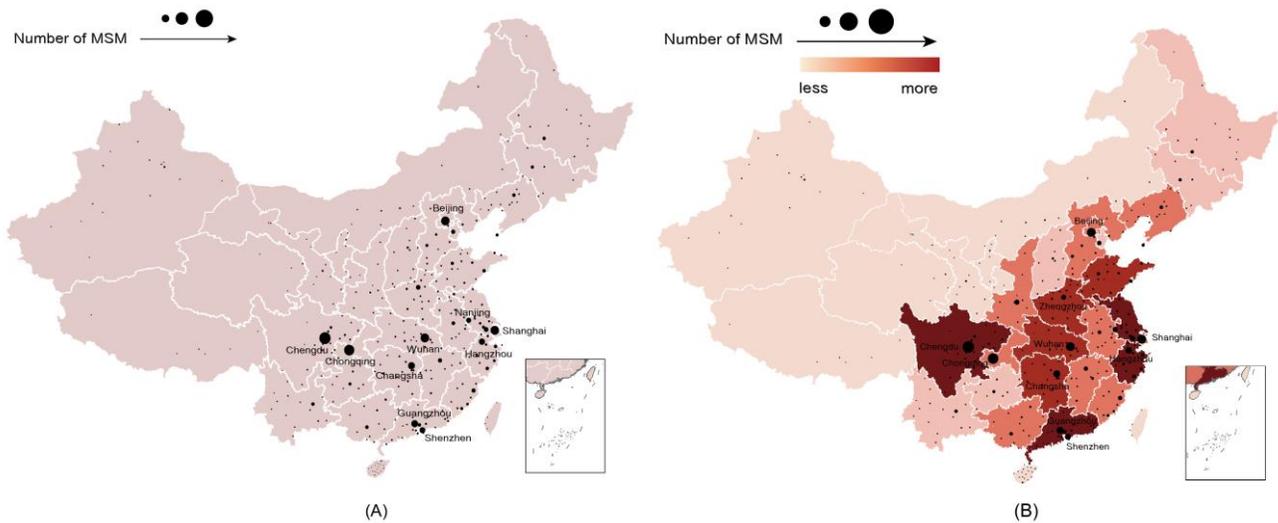
295

296 3.4 Distribution of online MSM in China

297 3.4.1 MSM distribution extracted from user profiles

298 The public profile data provides us with a direct way to obtain users' locations. We collected all
 299 possible profiles of active users in gay-bar, for a total of 359,438 records, in which about 10% of user
 300 location data (including GPS coordinates) is not empty. Figure 2A shows the city distribution of user
 301 geolocations reported in profiles of gay-bar users. It can be seen that most of the users are from cities
 302 of Chengdu, Chongqing, Wuhan, Shanghai, Beijing, Changsha, and Guangzhou. And we find that the
 303 number of users is not entirely related to the population of the city.

304 The city distribution and province distribution resolved from the latest recorded latitude and longitude
 305 coordinates of gay-bar users are shown in Figure 2B. It can be seen that the top 5 provinces are
 306 Guangdong, Sichuan, Jiangsu, Zhejiang and Hunan. The corresponding top 5 cities are Chengdu,
 307 Chongqing, Shanghai, Wuhan and Beijing with most gay-bar users, which is basically the same as the
 308 location distribution filled in user profiles above. Through further analysis, it is found that those users
 309 who were willing to provide their GPS coordinates were more likely to fill in the location fields
 310 correctly in the registration profiles.

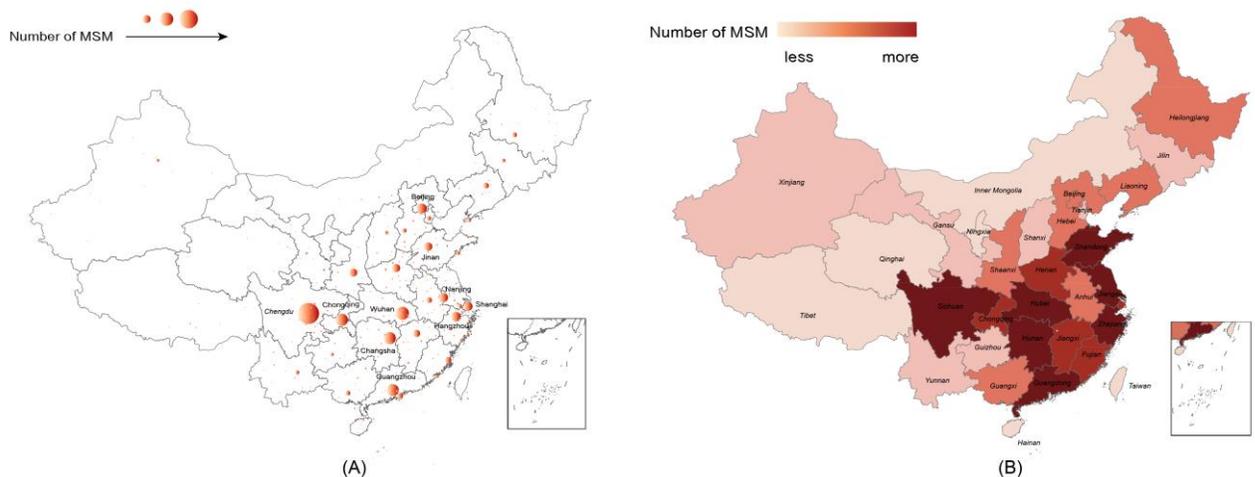


311

312 Figure 2. (A) The city distribution of gay-bar users from location fields in profiles; (B) the location distribution
 313 extracted from the GPS coordinates of gay-bar users on both the city-level and the province-level.

314 **3.4.2 MSM distribution inferred by the Gazetteer & NER algorithm**

315 According to recent locations disclosed in text content posted by gay-bar users, we use the Gazetteer
 316 & NER algorithm to estimate the geographic distribution of online MSM in China, which covers over
 317 156,000 MSM-related users. It makes up for the missing data, as well as the incomplete and inaccurate
 318 location information on user profiles, providing a good solution to the hard-to-access properties of
 319 hidden populations. As shown in Figure 3, it can be seen that most MSM-related users are from cities
 320 of Chengdu, Wuhan, Chongqing, Changsha, Guangzhou and Beijing. And the top 5 provinces with
 321 most relevant users are Sichuan, Guangdong, Zhejiang, Jiangsu and Hunan. These results are consistent
 322 with statistics of Chinese MSM population in recent studies (37, 38), indicating that MSM are mainly
 323 distributed in large cities in the eastern and southwestern China which is more economically developed
 324 and culturally open. Compared with the location distribution extracted from user profiles, MSM
 325 distribution inferred by the algorithm is more accurate. In addition to covering more users, our
 326 algorithm can infer users' more recent whereabouts revealed in their posts, while location fields filled
 327 in profiles when users register can be out of date and less reliable.



328

329 Figure 3. The geographic distribution of gay-bar users inferred from the published posts. (A) City-level distribution;
330 (B) Province-level distribution.

331 4 Conclusion & Discussion

332 For social networking platforms orient the online-to-offline dating, users tend to expose their
333 geographic information, as well as other basic demographic characteristics in the text content of their
334 posts, offering an unprecedented opportunity for the statistical inference on demographic attributes of
335 relevant populations. In this study, we try to compare mainstream location inference algorithms and
336 develop more efficient approaches to infer the geolocation distribution of the hidden population.
337 Among the popular location inferring methods, the classic gazetteer-based algorithm has achieved
338 better results. Meanwhile this algorithm has other advantages, such as fast calculation speed and easy
339 implementation. We have proposed a few amendments to the gazetteer-based algorithm by introducing
340 the context analysis as well as the strong constraint on contextual keywords. In addition, we develop a
341 hybrid voting algorithm (HVA-LI) by allowing different approaches to vote to determine the best
342 inference results, which guarantees a more effective way on location inference for hidden population.
343 Significantly, the hybrid algorithm consisting of the simple gazetteer-based method and named entity
344 recognition (NER) is proven to be the best choice to deal with inferring users' locations disclosed in
345 short texts on online communities, which achieves the best accuracy ($S = 0.713$) on the MSM-related
346 dataset, much higher than those from existing popular algorithms, whose best accuracy is 0.503. The
347 absolute accuracy of all algorithms utilized in this study is shown in Figure 4.

348 In summary, in order to expand the availability of the geolocation information of users in social
349 networks, especially for online hidden population, we have explored the possibility of location
350 inferring by analyzing textual content posted by online users. And we have proposed a more effective
351 hybrid algorithm, i.e., the Gazetteer & NER algorithm, to largely increase the accuracy of location
352 inference for hidden population and to overcome the sparseness problem of dealing with user profile
353 data. The methodology used in this work can also be extended to the location inferring of other
354 populations. In addition, associating the geographic region in user posts with their temporal series can
355 also provide important clues regarding travel, migration and displacement.

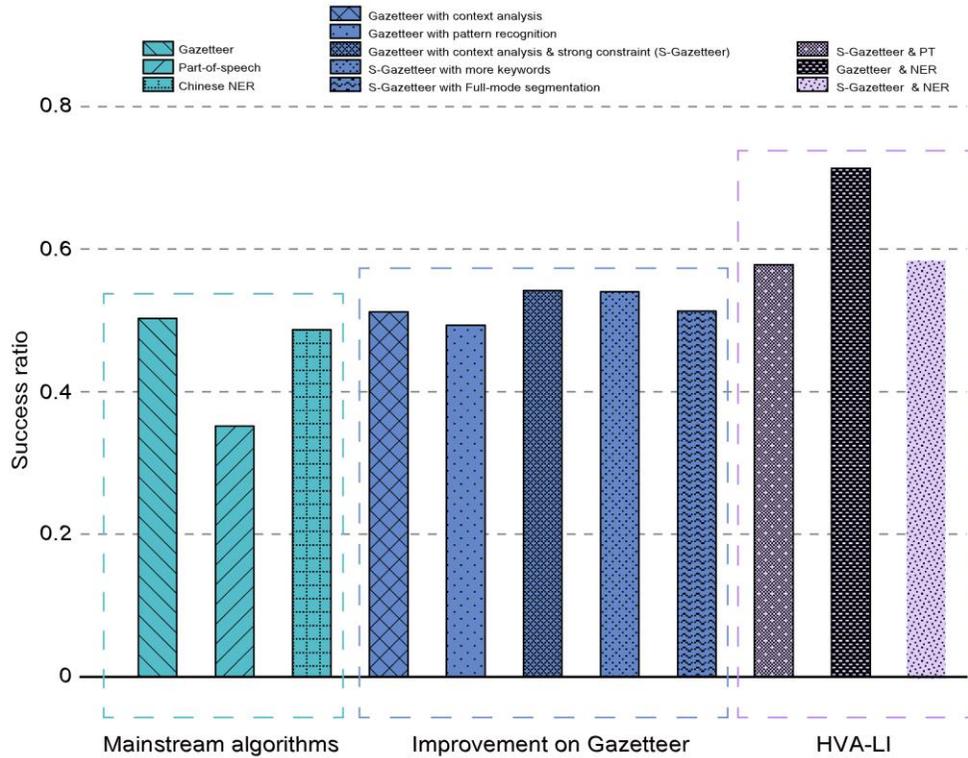


Figure 4. The absolute accuracy of all algorithms.

356
357

358 5 Declarations

359 Ethics approval and consent to participate

360 The study and Liu (2018, 2019) [1, 2] were supported by the Natural Science Foundation of China (91546203,
361 71771213) and approved by the Medical Ethical Committee of the Institutional Review Board (IRB) at Peking
362 University (IRB00001052–16016). This study uses only publicly available data and does not involve any physical,
363 social or legal risks to the users, whose personal identifying information are either unknown to the public or
364 removed if they choose to post online.

365 Consent for publication

366 Not applicable.

367 Availability of data and material

368 All data analyzed in this study are publicly available, all posts in the datasets can be collected on the website of
369 Baidu Tieba (<https://tieba.baidu.com/f?kw=gay>). Other data that support the findings in this study are available
370 from the corresponding author on reasonable request.

371 **Competing interests**

372 The authors declare that they have no competing interests.

373 **Funding**

374 Chuchu Liu and Ziqiang Cao are partially supported by the National Natural Science Foundation of China
375 (82041020, 71771213, 71690233). Xin Lu acknowledges the National Natural Science Foundation of China
376 (71901067, 71790615, and 91846301) and the Hunan Science and Technology Plan Project (2017RS3040, 2018JJ1034).

377 **Authors' contributions**

378 Methodology and experiment, C.L.; Data annotation and visualization, Z.C.; Funding acquisition, X.L.; Writing—
379 original draft, C.L.; Writing—review and editing, X.L.

380 **Acknowledgements**

381 Not applicable.

382 **6 References**

383 1. Liu C, Lu X. Network Evolution of a Large Online MSM Dating Community: 2005 – 2018.

384 (2019) 16(22):4322.

385 2. Liu C, Lu X. Analyzing hidden populations online: topic, emotion, and social network of
386 HIV-related users in the largest Chinese online community. (2018) 18(1):2.

387 3. Lu X, Bengtsson L, Britton T, Camitz M, Kim BJ, Thorson A, et al. The sensitivity of
388 respondent - driven sampling. (2012) 175(1):191-216.

389 4. Jie W, Ciyong L, Xueqing D, Hui W, Lingyao H. A syndemic of psychosocial problems
390 places the MSM (men who have sex with men) population at greater risk of HIV infection. (2012)
391 7(3).

392 5. Berghe WV, Nöstlinger C, Hospers H, Laga M. International mobility, sexual behaviour and
393 HIV-related characteristics of men who have sex with men residing in Belgium. (2013) 13(1):968.

394 6. Huang G, Cai M, Lu X. Inferring Opinions and Behavioral Characteristics of Gay Men with
395 Large Scale Multilingual Text from Blued. (2019) 16(19):3597.

396 7. Lu X. Linked ego networks: improving estimate reliability and validity with respondent-
397 driven sampling. (2013) 35(4):669-85.

398 8. Jia Z, Mao Y, Zhang F, Ruan Y, Ma Y, Li J, et al. Antiretroviral therapy to prevent HIV
399 transmission in serodiscordant couples in China (2003 – 11): a national observational cohort study.
400 (2013) 382(9899):1195-203.

401 9. Magnani R, Sabin K, Saidel T, Heckathorn D. Review of sampling hard-to-reach and hidden
402 populations for HIV surveillance. (2005) 19:S67-S72.

403 10. Lu X, Malmros J, Liljeros F, Britton T. Respondent-driven sampling on directed networks.
404 (2013) 7:292-322.

405 11. World Internet Users Statistics and 2020 World Population Stats [Internet World Stats].
406 (2020) [updated March 3,2020]. Available from: <https://www.internetworldstats.com/stats.htm>.

407 12. Nordrum A. Popular internet of things forecast of 50 billion devices by 2020 is outdated
408 (2016). (2017).

409 13. Bien CH, Best JM, Muessig KE, Wei C, Han L, Tucker JD, et al. Gay apps for seeking sex
410 partners in China: implications for MSM sexual health. (2015) 19(6):941-6.

- 411 14. Young LE, Michaels S, Jonas A, Khanna AS, Skaathun B, Morgan E, et al. Sex behaviors as
412 social cues motivating social venue patronage among young black men who have sex with men.
413 (2017) 21(10):2924-34.
- 414 15. Hinds J, Joinson AN. What demographic attributes do our digital footprints reveal? A
415 systematic review. (2018) 13(11).
- 416 16. Gundecha P, Barbier G, Liu H, editors. Exploiting vulnerability to secure user privacy on a
417 social networking site. *Proceedings of the 17th ACM SIGKDD international conference on*
418 *Knowledge discovery and data mining*; 2011.
- 419 17. Rodrigues E, Assunção R, Pappa GL, Renno D, Meira Jr W. Exploring multiple evidence to
420 infer users' location in Twitter. (2016) 171:30-8.
- 421 18. Davis Jr CA, Pappa GL, de Oliveira DRR, de L. Arcanjo F. Inferring the location of twitter
422 messages based on user relationships. (2011) 15(6):735-51.
- 423 19. Ajao O, Hong J, Liu W. A survey of location inference techniques on Twitter. (2015)
424 41(6):855-64.
- 425 20. Pontes T, Magno G, Vasconcelos M, Gupta A, Almeida J, Kumaraguru P, et al., editors.
426 Beware of what you share: Inferring home location in social networks. *2012 IEEE 12th International*
427 *Conference on Data Mining Workshops*; 2012: IEEE.
- 428 21. Chaabane A, Acs G, Kaafar MA, editors. You are what you like! information leakage through
429 users' interests. *Proceedings of the 19th Annual Network & Distributed System Security Symposium*
430 *(NDSS)*; 2012: Citeseer.
- 431 22. Backstrom L, Sun E, Marlow C, editors. Find me if you can: improving geographical
432 prediction with social and spatial proximity. *Proceedings of the 19th international conference on*
433 *World wide web*; 2010.
- 434 23. Popescu A, Grefenstette G, editors. Mining user home location and gender from flickr tags.
435 *Fourth International AAAI Conference on Weblogs and Social Media*; 2010.
- 436 24. Zheng D, Hu T, You Q, Kautz H, Luo J, editors. Towards lifestyle understanding: Predicting
437 home and vacation locations from user's online photo collections. *Ninth International AAAI*
438 *Conference on Web and Social Media*; 2015.
- 439 25. Beyrer C, Baral SD, Van Griensven F, Goodreau SM, Chariyalertsak S, Wirtz AL, et al.
440 Global epidemiology of HIV infection in men who have sex with men. (2012) 380(9839):367-77.
- 441 26. Qi J, Zhang D, Fu X, Li C, Meng S, Dai M, et al. High risks of HIV transmission for men
442 who have sex with men—a comparison of risk factors of HIV infection among MSM associated with
443 recruitment channels in 15 cities of China. (2015) 10(4).
- 444 27. The introduction of Baidu Tieba [wikipedia]. [cited 2019 22 November]. Available from:
445 https://en.wikipedia.org/wiki/Baidu_Tieba.
- 446 28. Wang X, Xu M, Ren Y, Xu J, Zhang H, Zheng N. A Location Inferring Model Based on
447 Tweets and Bilateral Follow Friends. (2014) 9(2):315-21.
- 448 29. Derczynski L, Ritter A, Clark S, Bontcheva K, editors. Twitter part-of-speech tagging for all:
449 Overcoming sparse and noisy data. *Proceedings of the International Conference Recent Advances in*
450 *Natural Language Processing RANLP 2013*; 2013.
- 451 30. Owoputi O, O' Connor B, Dyer C, Gimpel K, Schneider N, Smith NA, editors. Improved
452 part-of-speech tagging for online conversational text with word clusters. *Proceedings of the 2013*
453 *conference of the North American chapter of the association for computational linguistics: human*
454 *language technologies*; 2013.
- 455 31. Lozano MG, Schreiber J, Brynielsson J. Tracking geographical locations using a geo-aware
456 topic model for analyzing social media data. (2017) 99:18-29.
- 457 32. Gelernter J, Mushegian N. Geo - parsing messages from microtext. (2011) 15(6):753-73.

- 458 33. Project description [pypi]. [cited 2019 29 November]. Available from:
459 <https://pypi.org/project/jieba/>.
- 460 34. The geographical gazetteer of China [cited 2019 20 November]. Available from:
461 <http://www.china.com.cn/ch-quhua/>.
- 462 35. Yu H-k, Zhang H-p, Liu Q, Lv X-q, Shi S-c. Chinese named entity identification using
463 cascaded hidden Markov model. (2006) 27(2):87.
- 464 36. Module description on Chinese word segmentation of Jieba [cited 2019 29 November].
465 Available from: <https://github.com/fxsjy/jieba>.
- 466 37. Hu M, Xu C, Wang J. Spatiotemporal Analysis of Men Who Have Sex With Men in
467 Mainland China: Social App Capture-Recapture Method. (2020) 8(1):e14800.
- 468 38. Huang D, Wang J, Yang T. Mapping the Spatial - Temporal Distribution and Migration
469 Patterns of Men Who Have Sex with Men in Mainland China: A Web-Based Study. (2020)
470 17(5):1469.
- 471