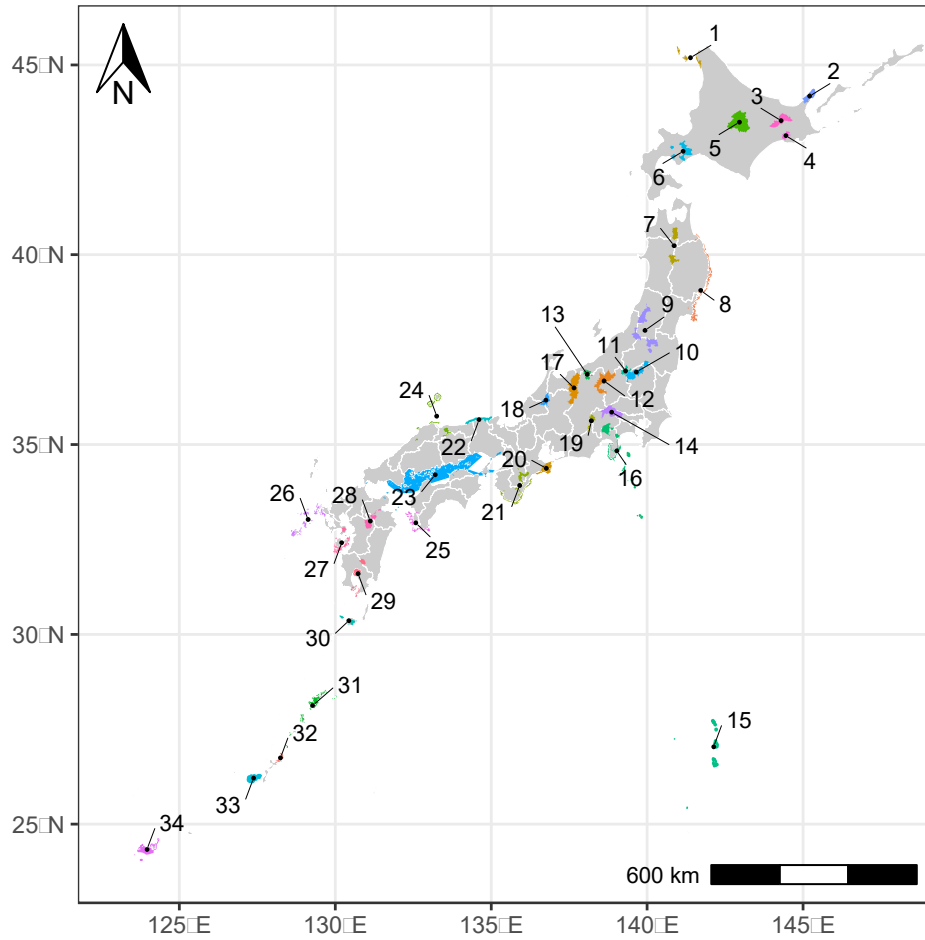


1 **Fig. S1:** Map of the 34 national parks in Japan.



2

- | | | |
|----------------------------|----------------------------|------------------------|
| 1. Rishiri-Rebun-Sarobetsu | 13. Myoko-Togakushi renzan | 25. Ashizuri-Uwakai |
| 2. Shiretoko | 14. Chichibu-Tama-Kai | 26. Saikai |
| 3. Akan-Mashu | 15. Ogasawara | 27. Unzen-Amakusa |
| 4. Kushiroshitsugen | 16. Fuji-Hakone-Izu | 28. Aso-Kuju |
| 5. Daisetsuzan | 17. Chubusangaku | 29. Kirishima-Kinkowan |
| 6. Shikotsu-Toya | 18. Hakusan | 30. Yakushima |
| 7. Towada-Hachimantai | 19. Minami Alps | 31. Amamigunto |
| 8. Sanriku Fukko | 20. Ise-Shima | 32. Yambaru |
| 9. Bandai-Asahi | 21. Yoshino-Kumano | 33. Keramashoto |
| 10. Nikko | 22. San'in kaigan | 34. Iriomote-Ishigaki |
| 11. Oze | 23. Setonaikai | |
| 12. Joshin'etsukogen | 24. Daisen-Okii | |

3

4

5 **Table S1:** Values of the source field of tweets regarded as originating from official Twitter clients
 6 and Foursquare apps.

7

Client type	Source field
Official client	Twitter for iPhone Twitter for Android Twitter Web Client web Twitter for iPad TweetDeck Twitter for Android Tablets Twitter for Windows Twitter for Appli Twitter for BlackBerry® Twitter for Windows Phone Twitter for BlackBerry
Foursquare	Foursquare Foursquare Swarm

8

9

10 **Appendix 1:** Development of custom client applications capturing Twitter data

11 The Twitter data used in this study was obtained by developing small custom client applications. To
12 construct Dataset 1, we developed a small client application which captures target tweets using the
13 ‘tweepy’ module in the Python programming language. The application sampled tweets using the
14 Streaming API of Twitter (statuses/sample endpoint; the archived documentation can be read at
15 [https://web.archive.org/web/20220316014310/https://developer.twitter.com/en/docs/twitter-](https://web.archive.org/web/20220316014310/https://developer.twitter.com/en/docs/twitter-api/tweets/volume-streams/introduction)
16 [api/tweets/volume-streams/introduction](https://web.archive.org/web/20220316014310/https://developer.twitter.com/en/docs/twitter-api/tweets/volume-streams/introduction)) by which we could obtain roughly 1% of tweets in real
17 time. We assumed tweets having Japanese characters were posted by Japanese users and tweets
18 which do not have Japanese characters were discarded using the regular expression (`([一-龠]+[あ-`
19 `ん]+[ア-ヴー]+)`), which matches all three kinds of Japanese characters (i.e. Kanji `[一-龠]`, Hiragana
20 `[あ-ん]`, and Katakana `[ア-ヴー]`). Due to limited storage capacity, this filtering was done in real time
21 and only tweets having Japanese characters were saved. Captured tweet data were stored in text files
22 using JSON (JavaScript Object Notation) format. We adopted JSON text files because it allowed us
23 to manually inspect and fix data corruption if necessary. Also, to reduce the probability of data
24 corruption, the tweets were saved into new files every hour. To maximize continuity of capturing,
25 the application was programmed to automatically restart immediately after occasional crashes of the
26 application or the computer. When a computer restart was needed, we ran the application in a backup
27 system and continued capturing. After obtaining tweet data, we checked and corrected corruption of
28 JSON data caused by occasional network errors and crashes. The data download was conducted
29 from January 1st, 2017, to December 31st, 2022. Note that despite our efforts to maintain continuous
30 capturing, the system experienced downtimes (Table S2). The longest downtime experienced during
31 October to December 2022 was caused by a migration problem for the API update.

32 To obtain tweet data for Dataset 2, another application was developed. Originally, we planned
33 to extract georeferenced tweets from Dataset 1. However, after we started capturing Dataset 1, we

34 noticed that the method could not capture enough georeferenced tweets. Therefore, we developed
35 another application which collects georeferenced tweets with location information in Japan by
36 specifying the bounding box of Japan (122.9336, 20.4253, 153.9864, 45.5572) for the Streaming
37 API (statuses/filter endpoint). This application had similar functionality and operated similarly as the
38 one used for Dataset 1, but the data was stored in new files every 10 minutes. The data collection
39 was conducted from March 19th, 2017 to December 31st, 2022, except during short downtimes
40 (Table S2).
41

42 **Table S2:** Duration of downtime experienced by the applications capturing tweet data. Because we
 43 did not measure actual duration, the durations of downtimes were estimated from the number of files
 44 which we expected to be stored every hour or every 10 minutes for Dataset 1 and Dataset 2,
 45 respectively.
 46

Dataset	Year-month	Downtime (hours for Dataset 1 and mins for Dataset 2)
Dataset 1	2018-04	1
	2021-11	17
	2022-01	20
	2022-11	183
	2022-12	970
Dataset 2	2017-05	10
	2017-08	20
	2017-10	100
	2018-03	60
	2018-04	10
	2018-07	30
	2018-09	90
	2019-07	40
	2019-10	20
	2020-02	30
	2020-05	20
	2020-06	100
	2020-10	80
	2021-11	1100
	2022-01	70
	2022-05	10
	2022-07	50
	2022-08	40
	2022-10	40
2022-11	10	
2022-12	10	

47

48

49 **Appendix 2:** Performance evaluation for the estimation of demographic attributes.

50 To evaluate the reliability of the estimation of demographic attributes, we compared the
51 estimation results with two existing sets of statistics on demographic attributes of Twitter users. The
52 first set of statistics was the results of the Survey on Information and Communication Media Usage
53 Time and Information Behavior conducted by the Institute for Information and Communications
54 Policy, Ministry of Internal Affairs and Communications, Japan
55 (https://www.soumu.go.jp/iicp/research/results/media_usage-time.html, accessed on July 14th, 2023;
56 hereafter, MIC survey). This survey used questionnaires to ask about information and media related
57 activities such as time spent on television, internet, social media, smartphones, etc. It has been
58 conducted every year since 2013 and targeted 1,500 people aged 13-69 in each survey. The
59 respondents of the survey were selected so that the percentages of each sex and age reflect the
60 demographic composition of Japan. We obtained the statistics about demographic attributes (age and
61 sex) and activities (posting and browsing) of Twitter users from 2017 to 2022 and used them for
62 verification (the numbers of respondents for each demographic group are shown in Fig. S2). The
63 second set of statistics used for verification was from the Twitter Usage survey conducted by
64 MyVoice Communications, Inc. (https://myel.myvoice.jp/products/detail.php?product_id=26908,
65 accessed on July 14th, 2023, hereafter MyVoice survey). This survey was conducted every year from
66 2009 to 2020 and contains age and sex as well as resident prefecture, which is not available from the
67 MIC survey, and activities (posting and browsing). We used the survey results from 2017-2020 for
68 this study (the numbers of respondents for each demographic group are shown in Fig. S3).

69 The verification was conducted as follows. First, we calculated the proportion of users who
70 were posting tweets for each demographic group and year using data from the two surveys. The MIC
71 survey separately asked respondents whether they post tweets from smartphones and/or computers.
72 Therefore, we calculated two values of the posting user ratio for each demographic group and year:

73 1) maximum posting user ratio which assumes users who posted tweets using smartphones did not
74 post tweets using computers, that is, the maximum posting user ratio = (number of respondents
75 posting tweets from smartphones + number of respondents posting tweets from computers) / number
76 of respondents; 2) minimum posting user ratio which assumes maximum overlap between those who
77 post on smartphones and computers, that is, the minimum posting user ratio = the larger of the
78 number of respondents posting tweets from smartphones or computers / number of respondents (Fig.
79 S4). For the MyVoice survey, the posting user ratio for each group and year was calculated as the
80 number of respondents who answered that they use Twitter and post tweets (Fig. S5). Next the
81 posting user ratios of each demographic group, year, and survey were multiplied by the population
82 number of the group obtained from the national census conducted in 2020 by the Statistics Bureau,
83 Ministry of Internal Affairs and Communications
84 (<https://www.stat.go.jp/data/kokusei/2020/index.html>, accessed on July 14th, 2023) and the number
85 of tweet-posting users for each demographic group was estimated. Then, using the estimated number
86 of users, we estimated the proportion of each demographic group in the tweet-posting users.

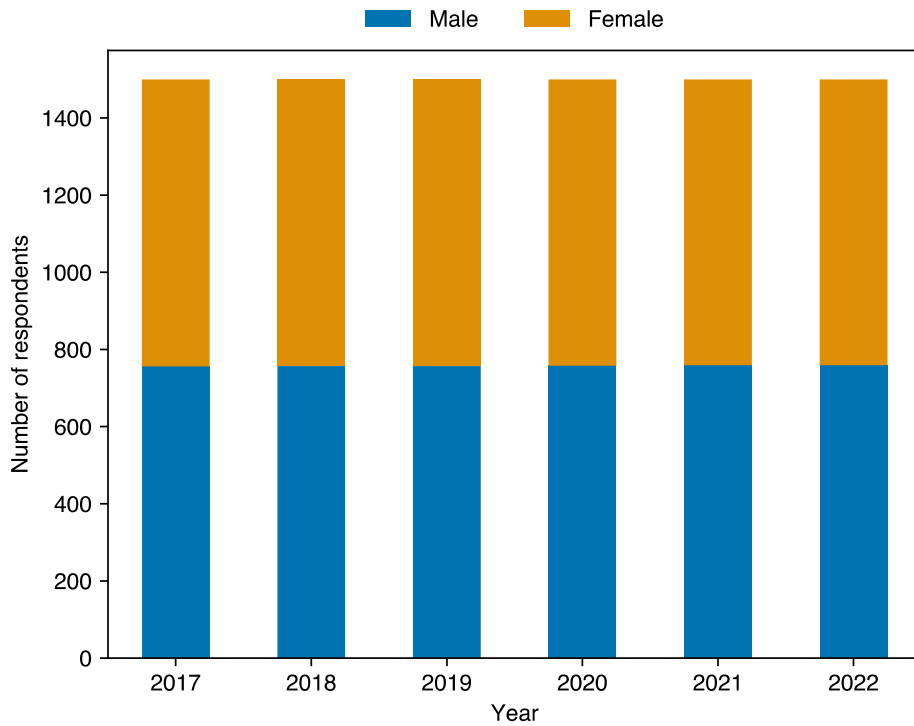
87 The results were compared with the proportion of each demographic group calculated from
88 the estimation results for Dataset 1 because the demographic attributes of the randomly sampled
89 Twitter user would be comparable to the existing statistics. For the verification of age, only the
90 results of the MIC survey were used as the reference because a very small number of respondents in
91 their 10's in the MyVoice survey cause inter-annual fluctuation of the posting ratio of this age group
92 which causes fluctuation of the estimated proportion of Twitter users for each age group. For age
93 and sex, congruency of the estimations based on Twitter to estimations based on other existing
94 statistics was evaluated by calculating Cramér's V (Cramér, 1946) between the estimations based on
95 Twitter and each statistic, for each year. Cramér's V measures association between variables on rows
96 and columns and ranges from 0 to 1 where 0 means no association (i.e., proportions of groups on

97 rows are not different between estimation methods on columns) and 1 means strong association (i.e.,
98 proportions of groups on rows are different between estimation methods on columns). We calculated
99 Cramér's V by creating $2 \times n$ contingency tables containing the estimated number of
100 users/population in each class in a demographic attribute (n = number of classes in the attribute)
101 estimated by Twitter and an existing statistic. Note that when calculating Cramér's V, the estimated
102 Japanese population who post tweets was rounded into nearest integer. For the estimation of
103 prefecture, Pearson's correlation coefficients were calculated between the estimated proportions of
104 users in each prefecture based on Twitter and those based on other statistics.

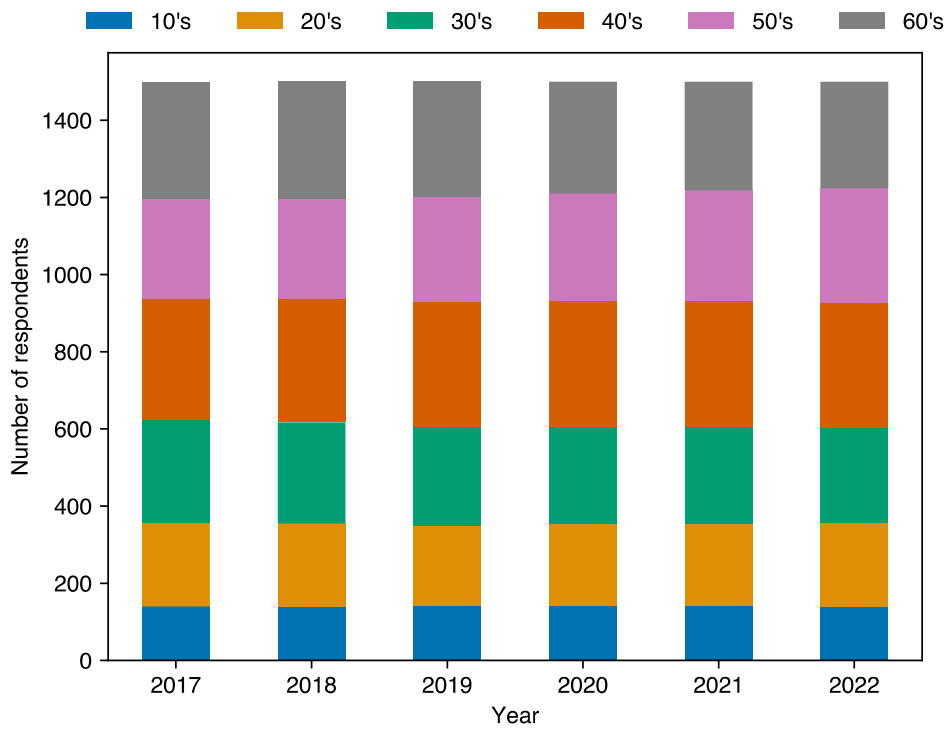
105 The estimated proportions of each sex based on Twitter data were congruent with the
106 estimations based on the existing statistics (Fig. S6a): for all the estimation methods, male users
107 were slightly dominant. Also, all estimated Cramér's V were <0.02 , indicating the differences were
108 very small (Table S3). On the other hand, although the estimated age showed that Twitter users are
109 dominated by younger age groups, the estimation based on the Twitter data showed a stronger
110 dominance compared to existing statistics (Fig. S6b): the estimation based on the Twitter data
111 showed that Twitter is mostly dominated by users less than 50 years old whereas the MIC survey
112 showed 10–20% of Twitter users were over 50. However, all estimated Cramér's V were <0.07 ,
113 indicating the differences were not large taken as a whole (Table S4). The estimated proportions of
114 Twitter users in each prefecture based on Twitter data showed strong correlations with those based
115 on the MyVoice survey for all years (Fig. S6c; Pearson's $R = 0.952, 0.930, 0.939, 0.947, \text{ and } 0.952$
116 for 2017, 2018, 2019, 2020, and all years, respectively and all p values for the coefficients were $<$
117 0.001) though the estimation from Twitter showed a higher proportion of users to be in Tokyo (the
118 cluster of points on the upper right side).

119

120 **Fig. S2:** Number of respondents of the Survey on Information and Communication Media Usage
 121 Time and Information Behavior for each demographic group for each year.



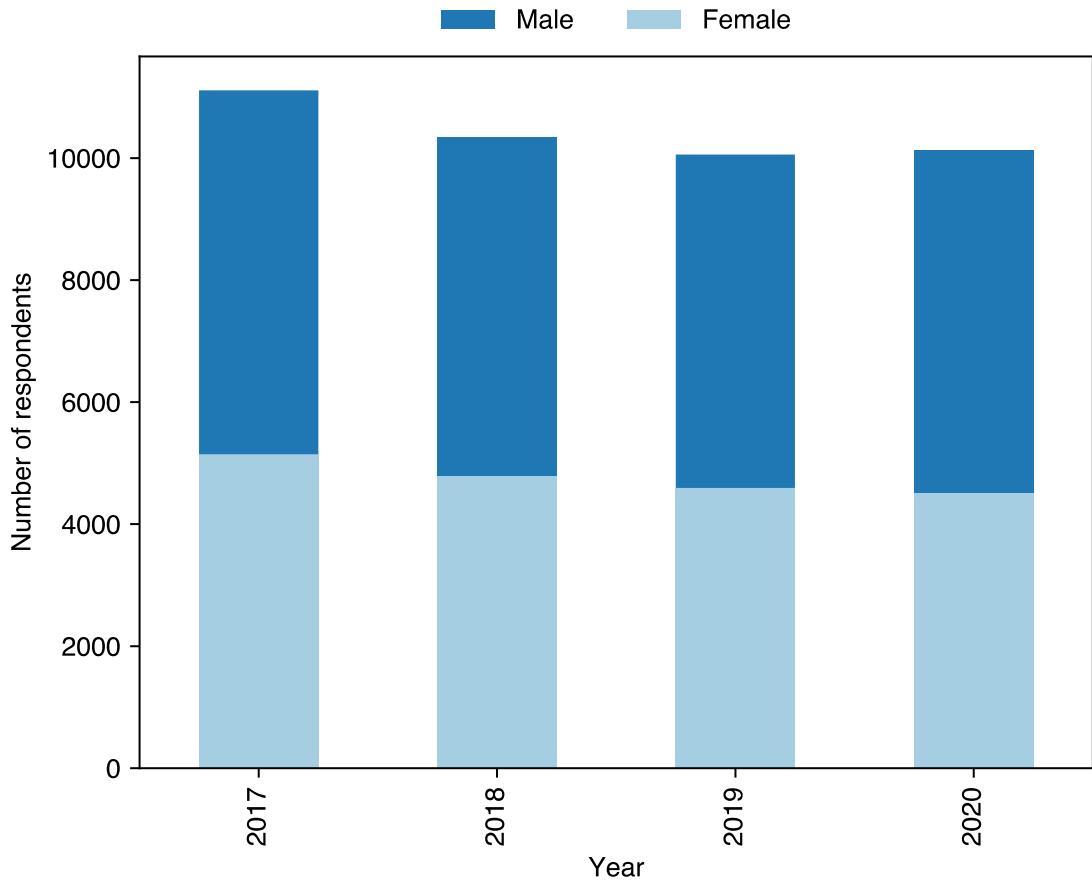
122

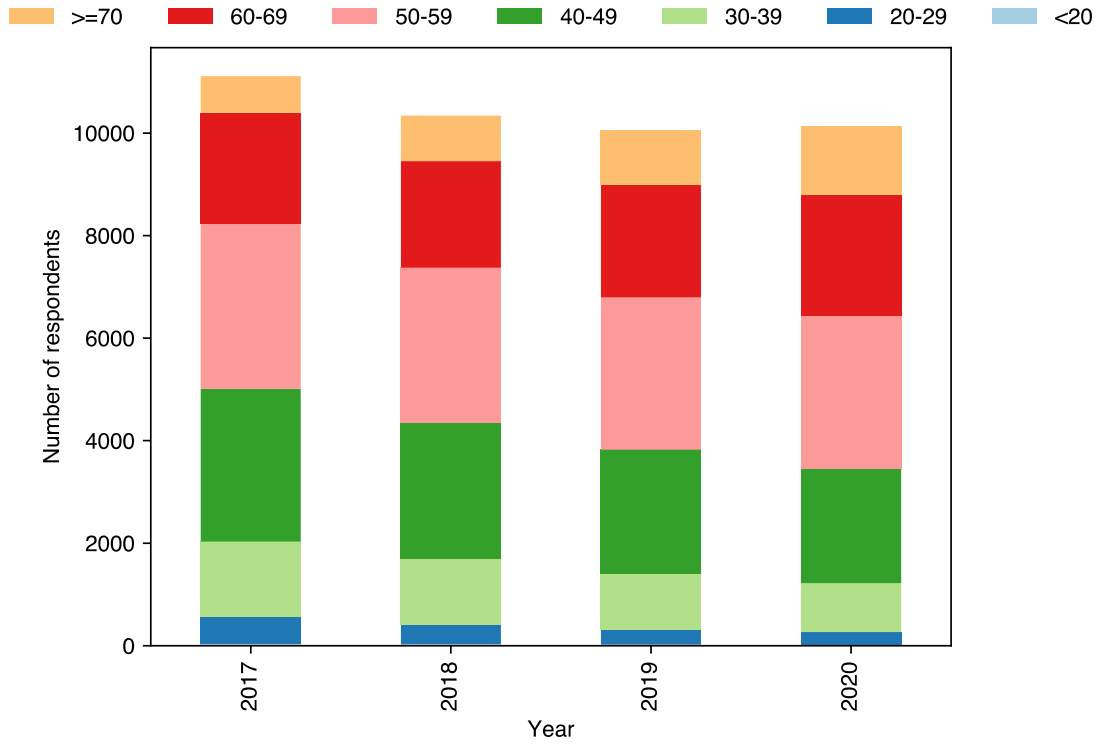


123

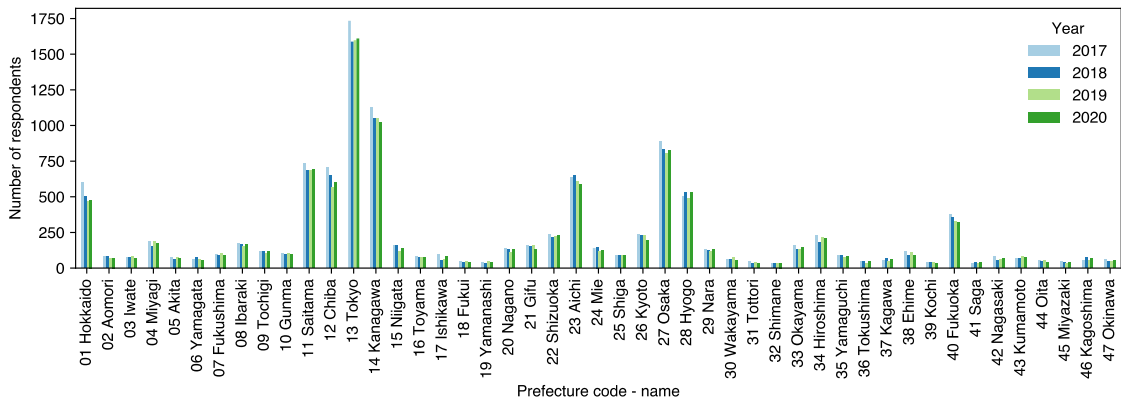
124

125 **Fig. S3:** Number of respondents to the Twitter Usage survey conducted by MyVoice
126 Communications, Inc. for each demographic group for each year.





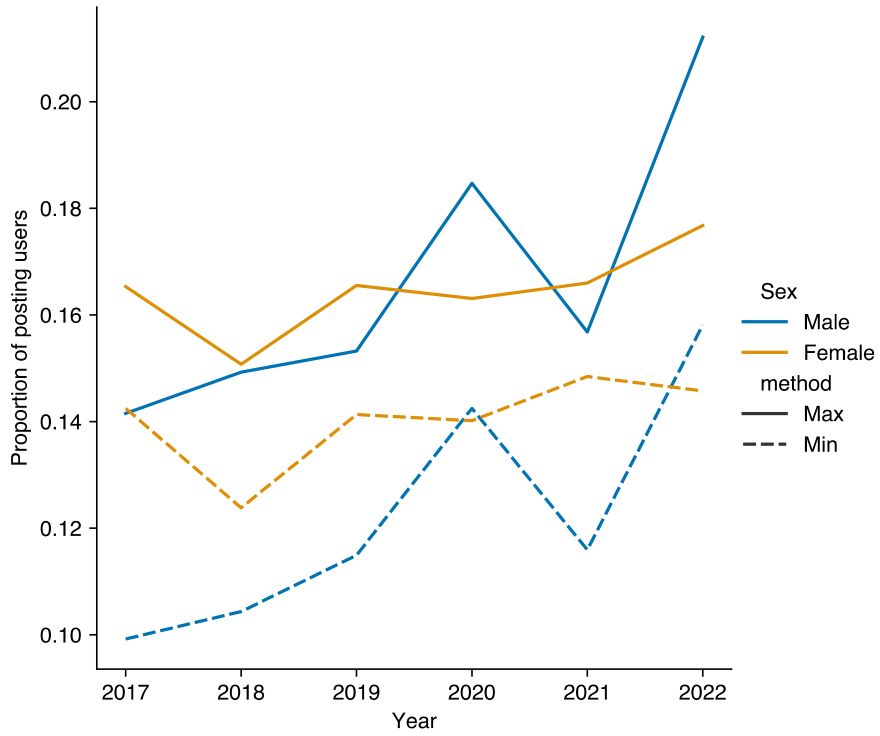
128



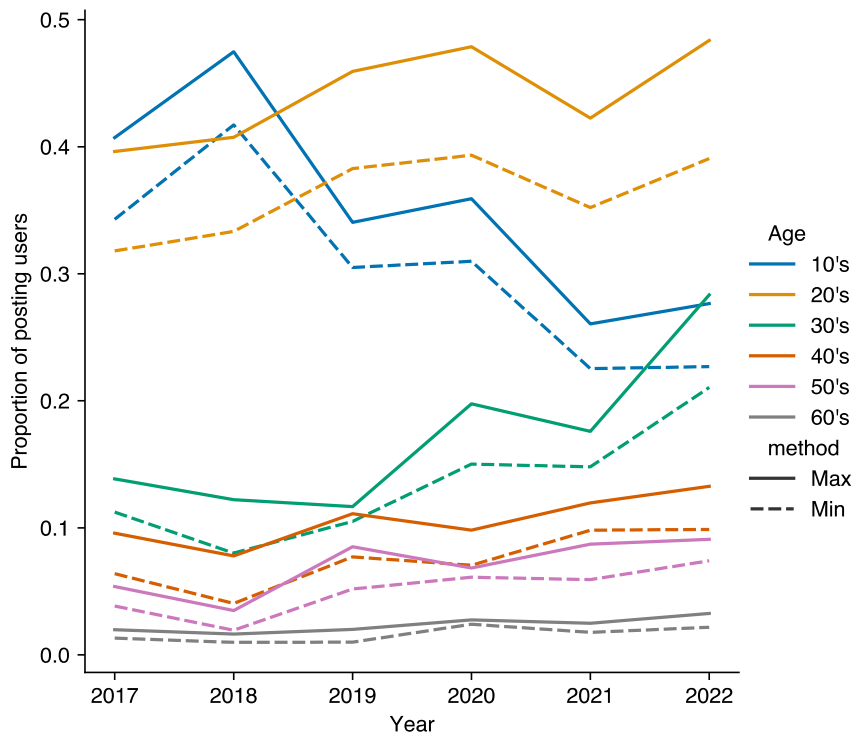
129

130

131 **Fig. S4:** Estimated proportions of users who were posting tweets for each demographic group,
 132 method, and year estimated from the MIC survey.



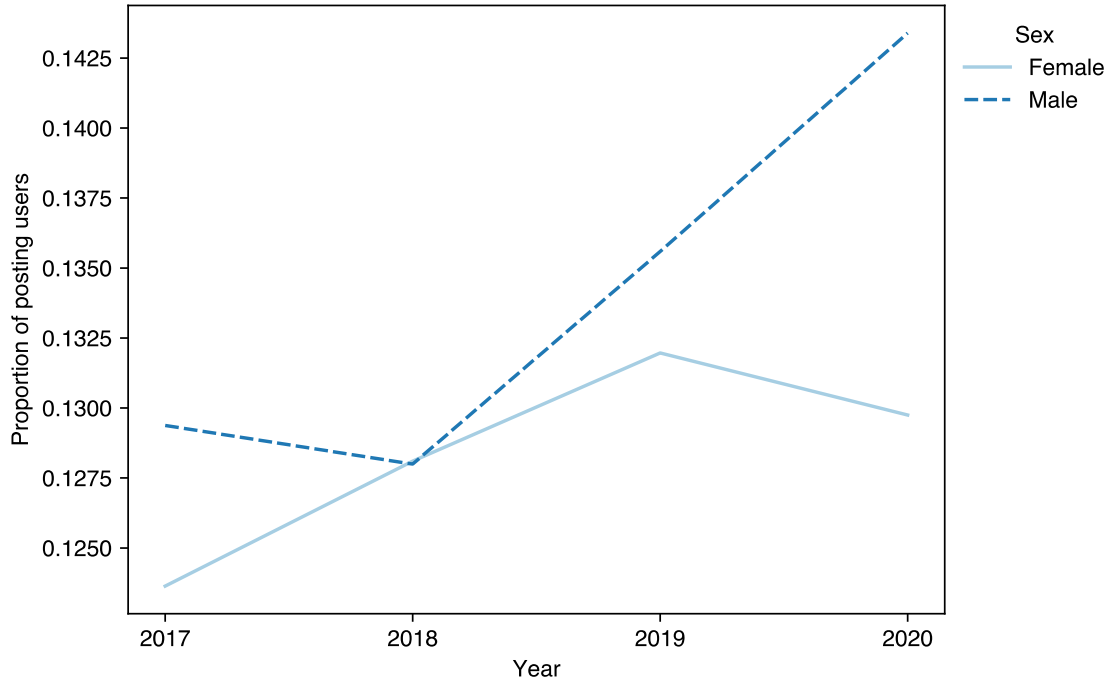
133



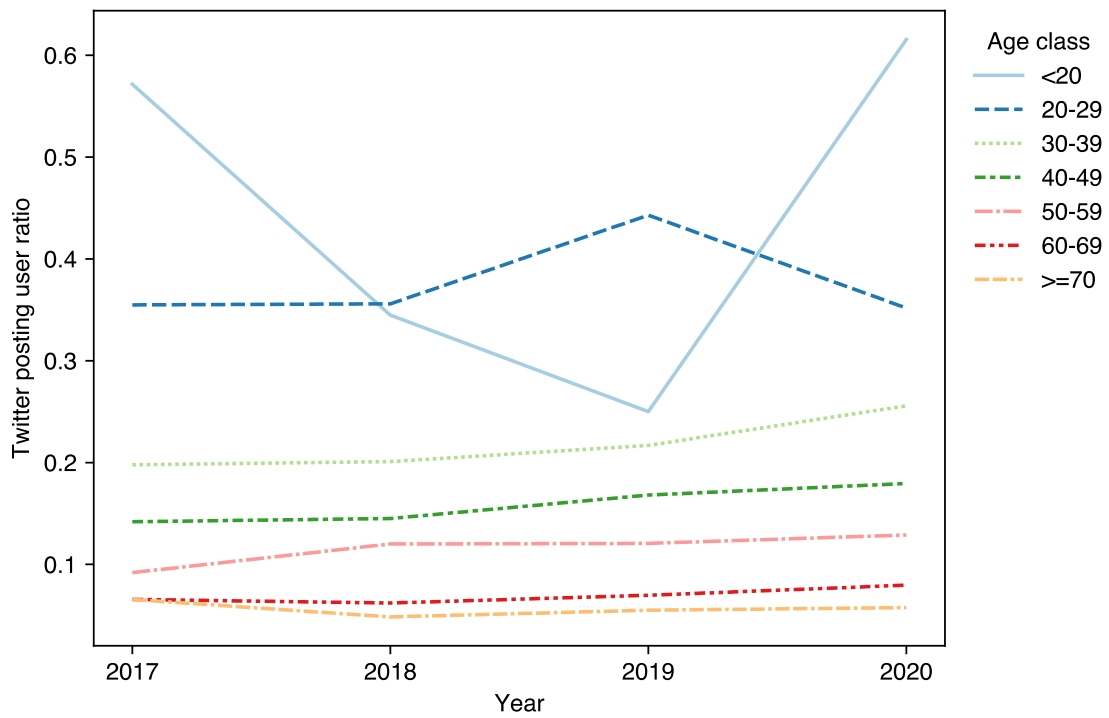
134

135

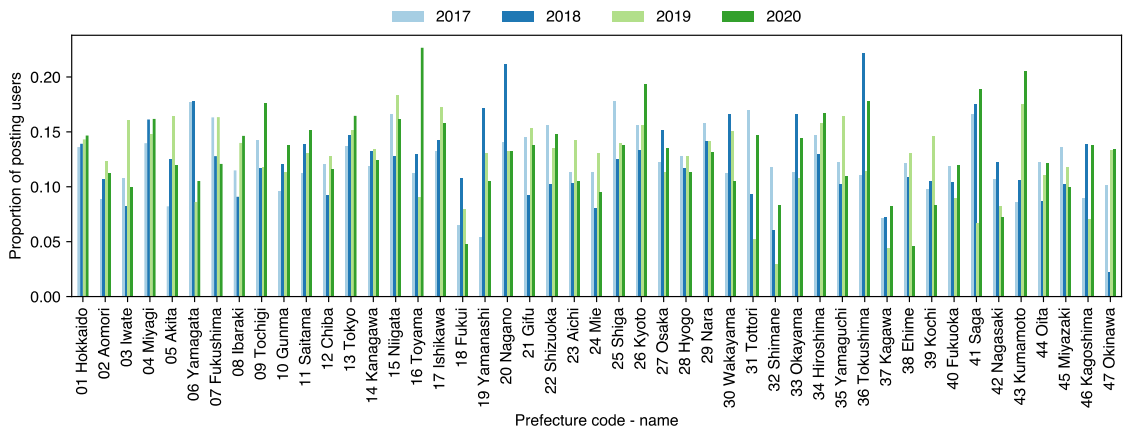
136 **Fig. S5:** Proportions of users who were posting tweets for each demographic group and year
 137 estimated from the MyVoice survey.



138



139

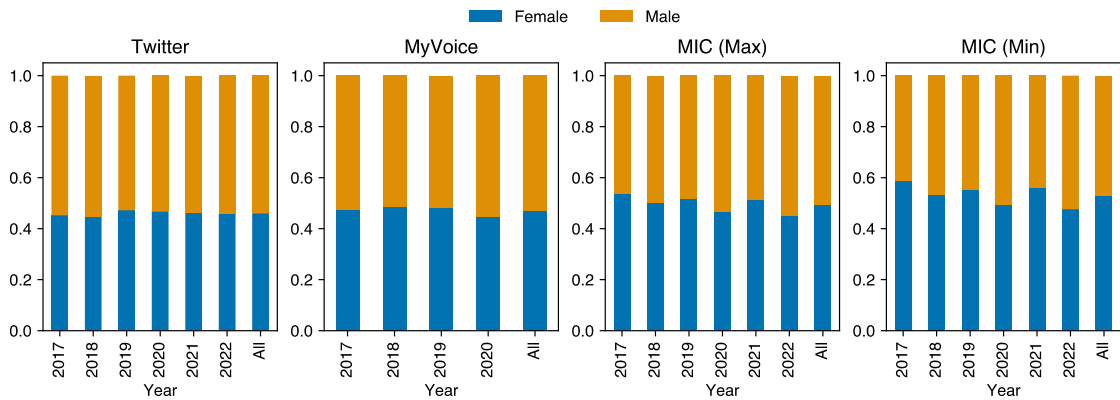


140

141

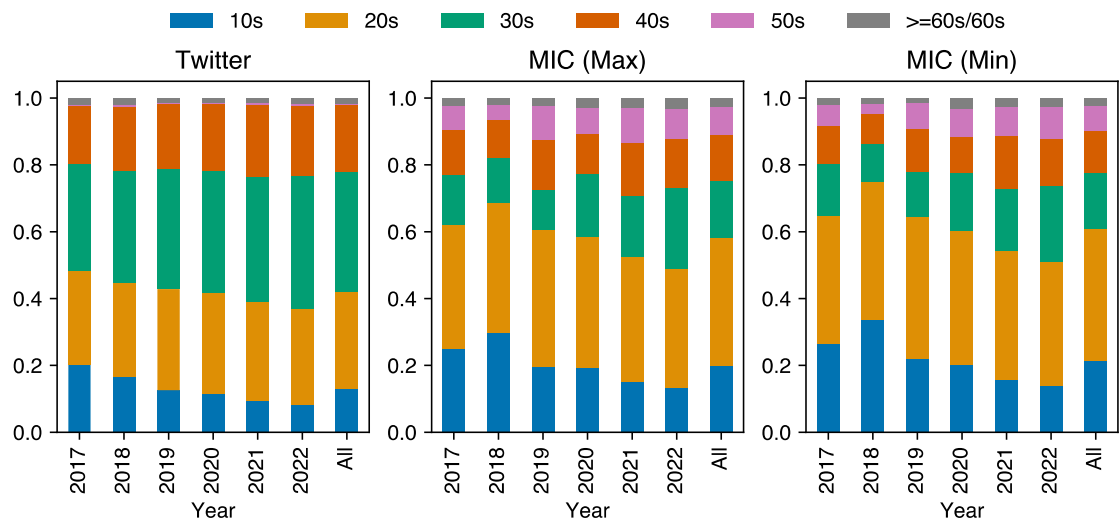
142 **Fig. S6:** Proportions of Twitter users in each demographic group estimated from Twitter and the
 143 existing statistics for each year (a and b) and relationship between proportions of Twitter users in
 144 each prefecture estimated using Twitter and the existing statistics (c). The diagonal line shown on
 145 the panel (c) denotes $y = x$.

146 (a) Sex



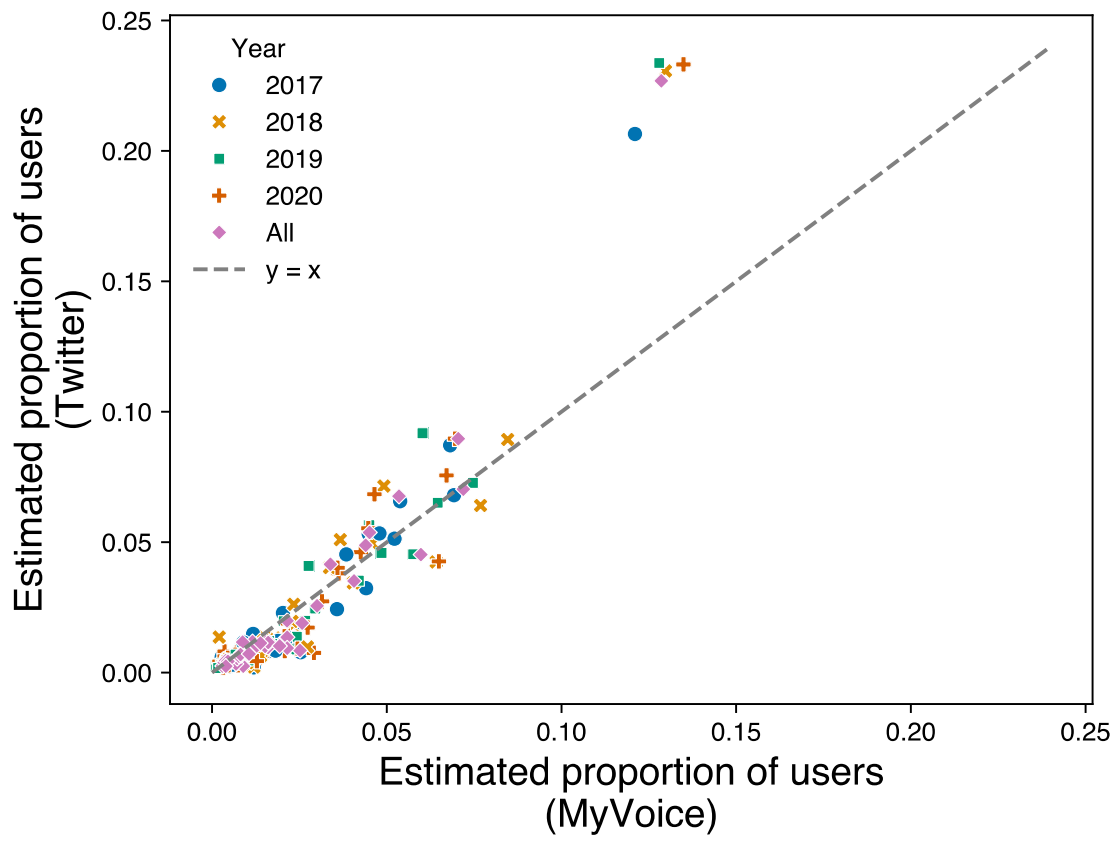
147

148 (b) Age



149

150



154 **Table S3:** Cramér's V statistics for the proportions of each sex estimated based on Twitter and the
155 existing statistics.

Year	MyVoice	MIC (Max)	MIC (Min)
2017	0.0013	0.0062	0.0114
2018	0.0025	0.0042	0.0074
2019	0.0007	0.0034	0.0067
2020	0.0014	0.0001	0.0021
2021	0.0000	0.0036	0.0080
2022	0.0000	0.0005	0.0014
All	0.0015	0.0062	0.0140

156

157

158 **Table S4:** Cramér's V statistics for the proportions of age groups estimated based on Twitter and the
159 existing statistics.

Year	MIC (Max)	MIC (Min)
2017	0.0202	0.0231
2018	0.0267	0.0368
2019	0.0304	0.0315
2020	0.0214	0.0265
2021	0.0229	0.0243
2022	0.0170	0.0207
All	0.0535	0.0627

160

161

162 **Appendix 3:** Performance evaluation of sentiment estimation for Japanese short text using IBM
163 Natural Language Understanding.

164 IBM Watson Natural Language Understanding is a commercially available machine learning service
165 which can estimate sentiment of provided text and returns a sentiment label (positive, neutral, or
166 negative) as well as sentiment score ranging from -1 (negative) to 1 (positive). However, because its
167 estimation performance was not disclosed, we evaluated it by estimating sentiment of an existing
168 corpus with manually attached sentiment labels. Also, we compared its performance with that of
169 other machine learning models and a dictionary-based method.

170 For the machine learning models, we employed *autonlp-japanese-sentiment-59362* (Thakur,
171 2021a) and *autonlp-japanese-sentiment-59363* (Thakur, 2021b) which estimate sentiment of
172 Japanese text and return a sentiment label of the text in binary form (negative/positive) along with its
173 confidence score. As for the dictionary-based method, we used python binding of ML-Ask (Ikegami,
174 2019; Ptaszynski et al., 2009), which assigns a sentiment label of five classes (positive, mostly
175 positive, neutral, mostly negative, and negative) for the analyzed text.

176 For the test data, we employed the *Japanese realistic textual entailment corpus* (Hayashibe,
177 2020) which includes a corpus of 5,553 hotel reviews with manually attached sentiment labels
178 (positive, neutral, and negative). We estimated sentiment of all texts in the corpus and evaluated
179 congruence of estimations and manually attached labels for each method using the Matthews
180 correlation coefficient (MCC; Matthews, 1975), The MCC is a correlation coefficient for discrete
181 predictions ranging from -1 to 1 where 1 indicates complete prediction, 0 indicates random
182 prediction, and -1 indicates complete inverse prediction.

183 Because *autonlp-japanese-sentiment-59362* and *autonlp-japanese-sentiment-59363* do not
184 return neutral labels, we adjusted the threshold value for the confidence score under which the

185 analyzed text was regarded as neutral from 0.5 to 1, and obtained the best threshold which brings the
186 highest value of MCC. As for the ML-Ask, the labels of *mostly positive* and *mostly negative* were
187 regarded as positive and negative, respectively, and MCC was calculated. IBM Natural Language
188 Understanding failed to estimate sentiments of 76 texts which were too short or contained only
189 symbols. Also, ML-Ask could not assign sentiment to 4,019 texts because of their dictionary-based
190 nature. For these cases, texts without sentiment labels were omitted from the calculations of MCC.
191 Python programs and data are available at Zenodo (doi:10.5281/zenodo.11366621) and Mendeley
192 Data (doi:10.17632/2hvh5yzcbd.1).

193 The results showed that IBM Natural Language Understanding had the highest value of MCC
194 among the methods investigated (Table S5). Although ML-Ask also showed a relatively high value
195 of MCC, it could estimate only 27.6% of the text. Therefore, we concluded IBM Natural Language
196 Understanding was the best method for this study.

197 **Table S5:** Matthews correlation coefficient (MCC) for each method/model.

Model/Method	MCC
IBM Watson Natural Language Understanding	0.550
autonlp-japanese-sentiment-59362	0.402
autonlp-japanese-sentiment-59363	0.390
ML-Ask	0.516

198

199

200 **Appendix 4:** Text preprocessing applied before the estimation of sentiment.

201 Before the estimation, we applied a preprocessing algorithm to the text data in each tweet. First,
202 because Foursquare clients add automatically generated text in a fixed format (e.g., “(@
203 LOCATION in CITY, PREFECTURE) <https://t.co/xxxxxxxxxx>”) at the end of the *text* attribute for
204 tweets posted by clients, this text was removed. Next, invalid characters in tweet text which are
205 rarely used in usual sentences and cannot be stored in Microsoft Excel files were removed using the
206 *openpyxl* module in Python. Finally, URLs in the text were removed. If a tweet contained no
207 characters after the preprocessing, it was omitted from the analysis.

208

209 **Appendix 5:** Evaluation and adjustment of spatial auto correlation

210 The spatial autocorrelations of the distributions of the number of posts were evaluated using Moran's
211 I. Moran's I measures spatial autocorrelation using a spatial weight matrix ranging from -1 (negative
212 spatial autocorrelation, i.e., dispersion) to 1 (positive spatial autocorrelation, i.e., clustering) and 0
213 indicates no correlation. The calculation used the *moran.test* function in the *spdep* (Bivand & Wong,
214 2018) package in R. We calculated a spatial weights matrix as an inverse of pairwise geographical
215 distance among the centers of all the grids with the diagonal of the matrix set to 0 using the *spdep*
216 and *geosphere* (Hijmans, 2024) packages in R. To reduce computational requirements, only grids
217 within 5km were considered to have spatial relatedness. The result showed that the spatial
218 distribution of the number of tweets showed weak but significant spatial autocorrelation (Table S6).

219 Because we found significant autocorrelations, we removed them by using a spatial filtering
220 approach using spatial lag variables before calculating the correlation coefficients as follows. 1) We
221 calculated a spatial weights matrix as an inverse of pairwise geographical distance among the centers
222 of all the grids with the diagonal of the matrix set to 0 using the *spdep* (Bivand & Wong, 2018) and
223 *geosphere* (Hijmans, 2024) packages in R. Note that to reduce computational requirements, only
224 grids within 5km were considered to have spatial relatedness. 2) We calculated spatial lag (i.e., the
225 weighted average of values at neighboring locations) for count variables using the weight matrix
226 with the *lag.listw* function in the *spdep* package. 3) Then each count was regressed on its own spatial
227 lag using a linear model and the residuals were extracted. These residuals represent variation that
228 was not explained by local spatial autocorrelation. 4) Using the residuals, correlation coefficients
229 were calculated. Note that the analysis was integrated into a Python script using the *rpy2* module.

230

231 **Table S6:** Moran's I and its two-sided *p*-value for the spatial distributions of the total number of
 232 tweets as well as the numbers of tweets for each sex and age group (a) without or (b) with removing
 233 massively tweeting users.

234

Category	Group	(a)	Moran's I	<i>p</i>	(b)	Moran's I	<i>p</i>
Total	Total		0.104	<0.001		0.103	<0.001
Sex	Male		0.105	<0.001		0.103	<0.001
	Female		0.089	<0.001		0.105	<0.001
Age	<20		0.088	<0.001		0.097	<0.001
	20-29		0.096	<0.001		0.098	<0.001
	30-39		0.108	<0.001		0.109	<0.001
	40-49		0.099	<0.001		0.102	<0.001
	50-59		0.065	<0.001		0.078	<0.001
	>=60		0.038	<0.001		0.071	<0.001

235

236

237 **Table S7:** Number of tweets and users obtained for each dataset and condition.

238

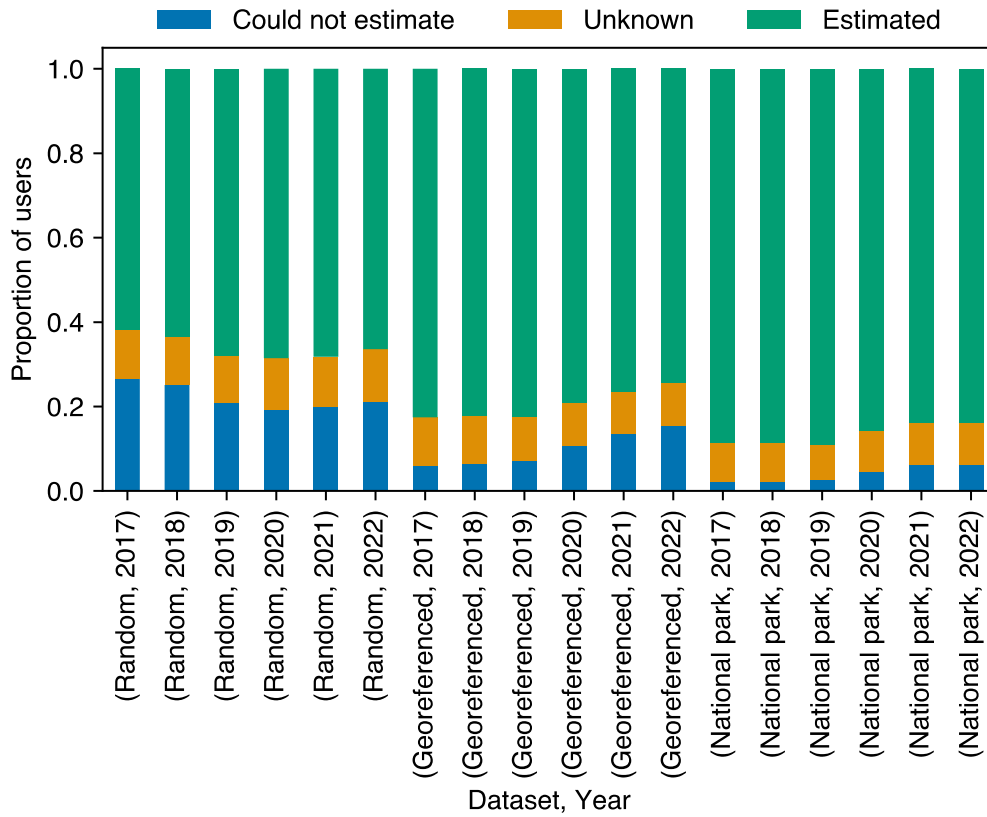
Dataset	Year	Tweets	Users
Randomly sampled Japanese tweets (Dataset 1)	2017	128,182,687	11,698,048
	2018	98,887,809	10,911,155
	2019	83,057,028	10,954,135
	2020	86,980,084	12,561,231
	2021	81,819,279	12,156,577
	2022	74,404,113	11,934,367
	Total	553,331,000	70,215,513
All georeferenced tweets with coordinate or POI	2017	21,464,748	722,011
	2018	23,377,945	729,537
	2019	19,217,940	610,168
	2020	15,384,985	582,433
	2021	15,251,789	538,931
	2022	18,189,046	645,832
	Total	112,886,453	3,828,912
Georeferenced tweets having Japanese characters posted from official Twitter or Foursquare clients having coordinates on the standard 2nd mesh (Dataset 2)	2017	4,331,931	345,129
	2018	5,485,961	386,304
	2019	4,778,818	326,105
	2020	5,533,190	383,847
	2021	6,044,083	371,094
	2022	8,032,706	456,544
	Total	34,206,689	2,269,023
All national park tweets	2017	311,734	72,895
	2018	353,763	78,305
	2019	291,113	64,854
	2020	308,700	64,871
	2021	336,587	61,631
	2022	411,269	76,164
	Total	2,013,166	418,720

239

Dataset	Year	Tweets	Users
National Park tweets with Japanese characters posted from the official the Twitter or Foursquare clients (Dataset 3)	2017	68,518	26,454
	2018	98,269	35,231
	2019	84,432	29,516
	2020	144,446	40,536
	2021	171,512	41,559
	2022	226,544	53,669
	Total	793,721	226,965

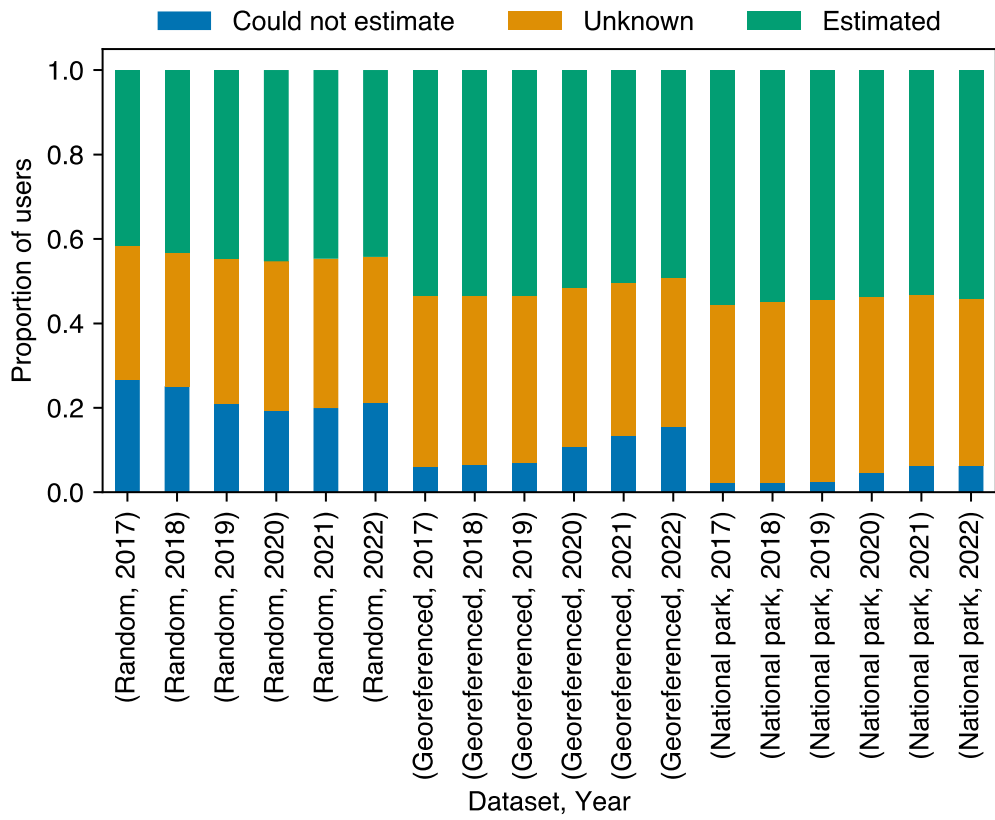
241 **Fig. S7:** Proportion of users whose attributes (a: sex, b: age, and c: prefecture) could not be
 242 estimated (e.g., deleted or private users), users whose attributes could be estimated but the attributes
 243 were unknown (i.e., public users but their attributes could not be estimated), and users whose
 244 attributes were successfully estimated.

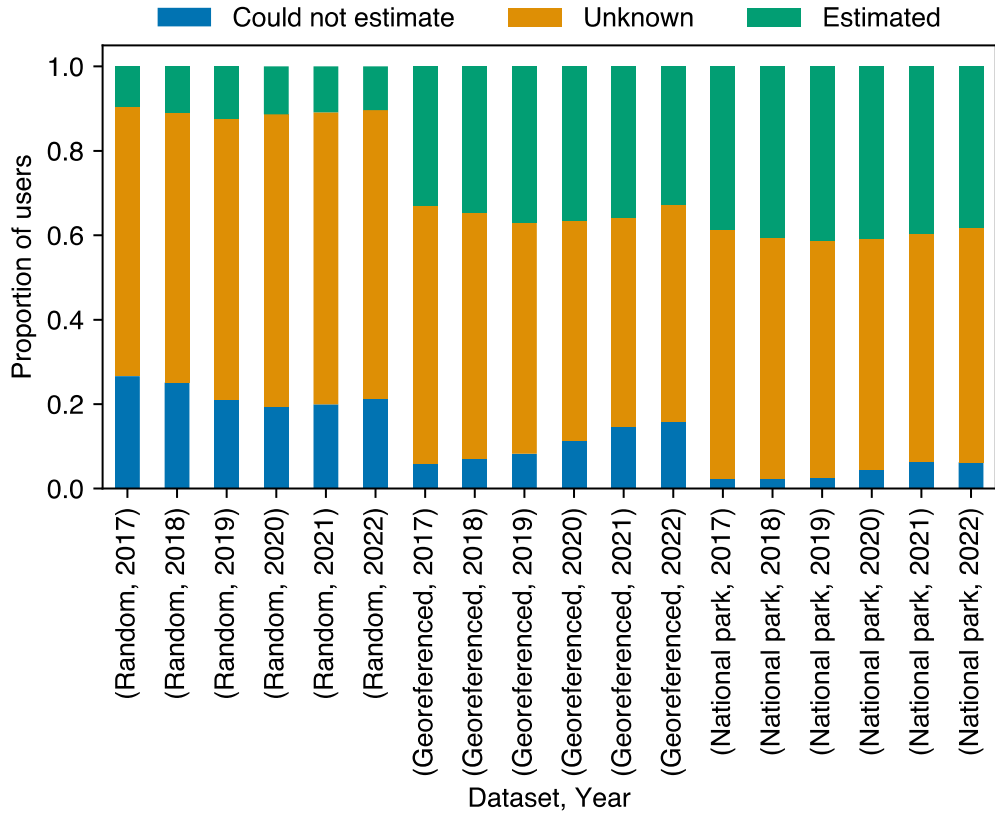
245 (a) Sex



246

247

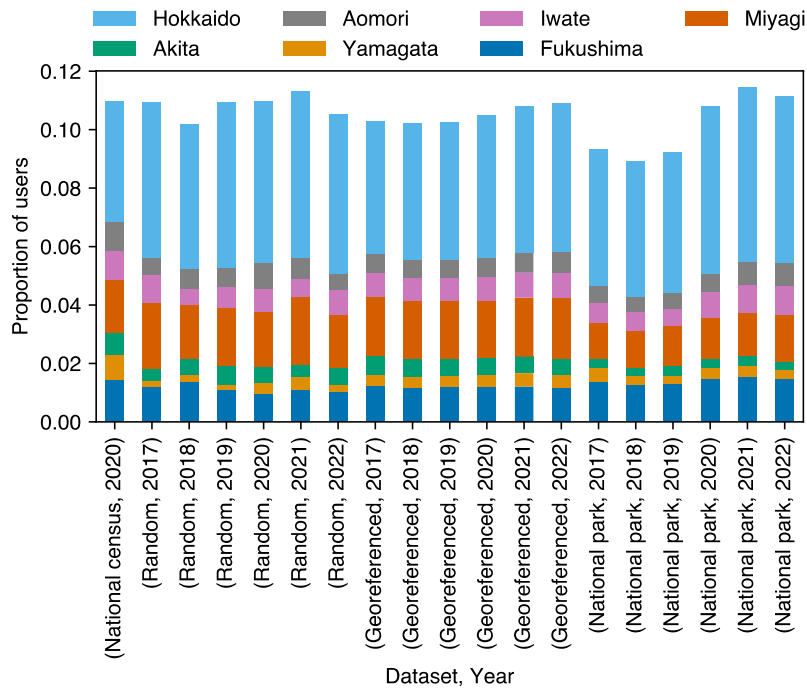




254 **Fig. S8:** Estimated proportions of residential prefectures of the Twitter users for each dataset in each
 255 year for (a) Tohoku and Hokkaido, (b) Kanto, (c) Chubu, (d) Kinki, (e) Chugoku, (f) Shikoku, and
 256 (g) Kyushu and Okinawa regions.

257

258 (a) Tohoku and Hokkaido region

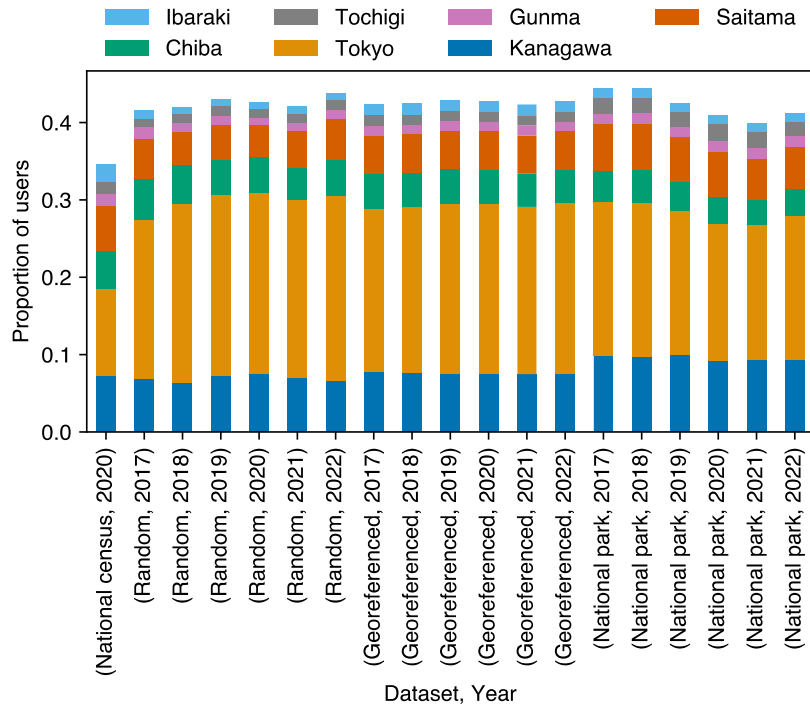


259

260

261

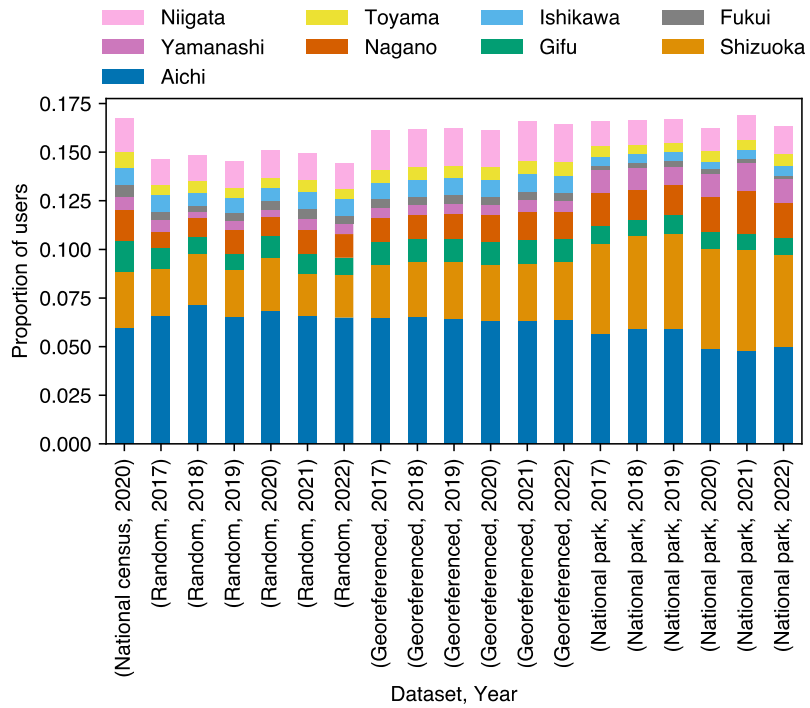
262 (b) Kanto region



263

264

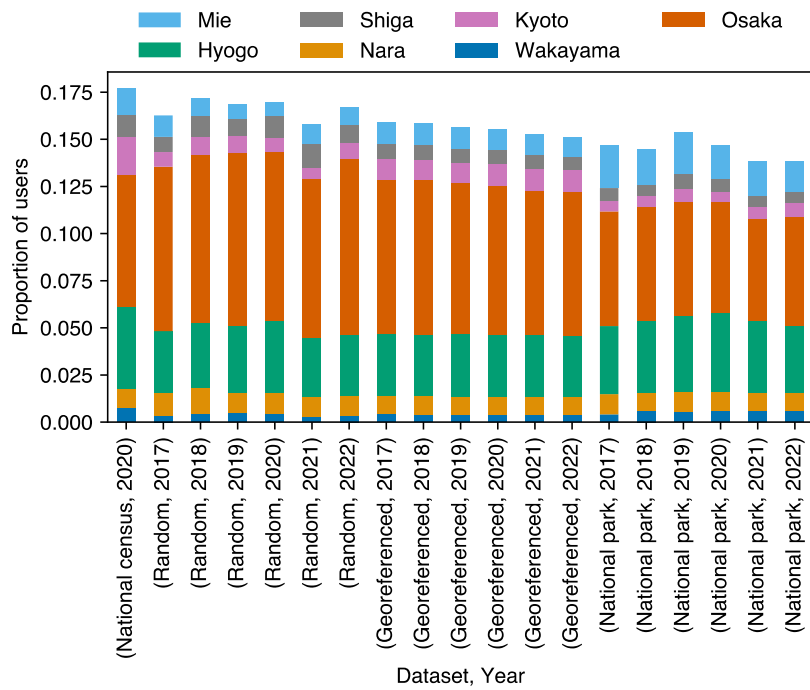
265 (c) Chubu region



266

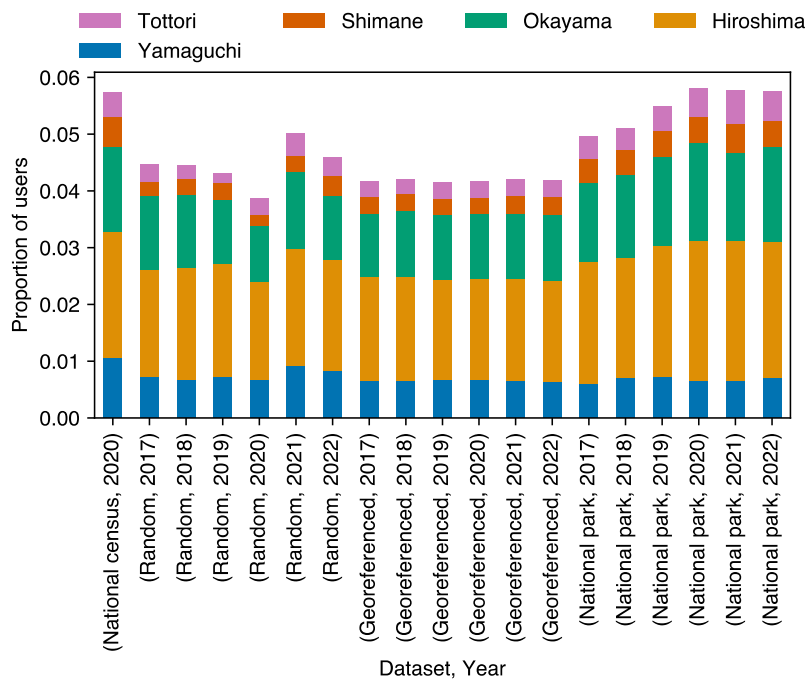
267

268 (d) Kinki region



269

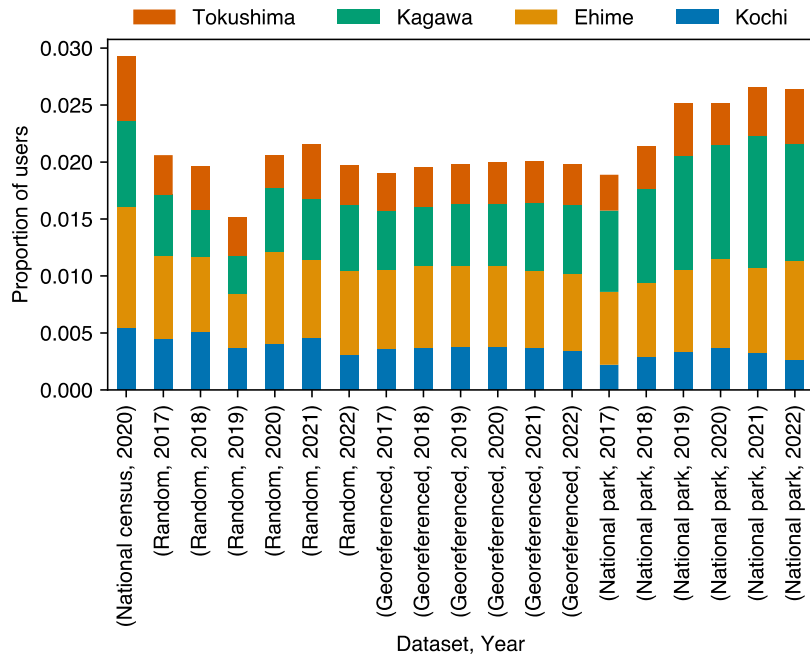
270 (e) Chugoku region



271

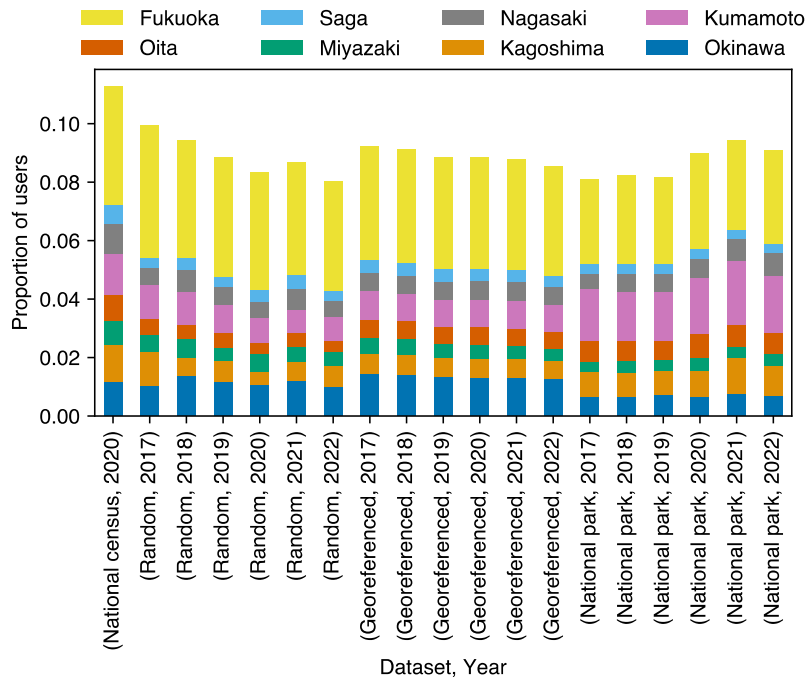
272

273 (f) Shikoku region



274

275 (g) Kyushu and Okinawa region



276

277

278 **Table S8:** Result of the type II analysis of deviance for the GLMM testing differences in travel cost
 279 among client type (Twitter/Foursquare), age, sex, and year. For this model, 155,659 tweets from
 280 34,147 users were used. The marginal R^2 (i.e., performance of explanatory variables) of the model
 281 was 0.001.

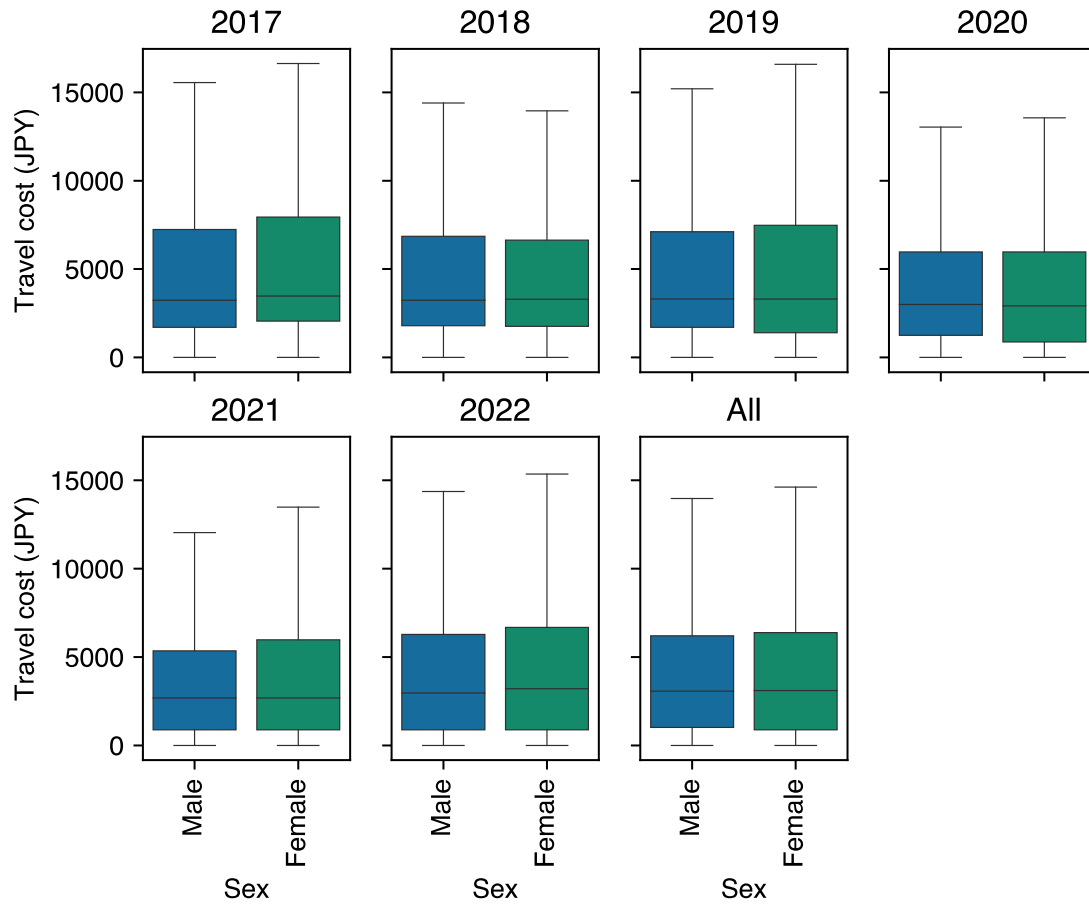
	χ^2	df	p	
Foursquare	0.789	1	0.374404	
Age	19.7569	5	0.001388	**
Sex	0.2882	1	0.591376	
Year	85.6179	5	<0.001	***

282

283

284 **Fig. S9:** Box-whisker plots showing distributions of estimated travel cost for (a) each sex and (b)
285 each age group for each year. Note that to reduce the dominance of massively tweeting users, only
286 one tweet for one day for one grid for each user was randomly selected and used.

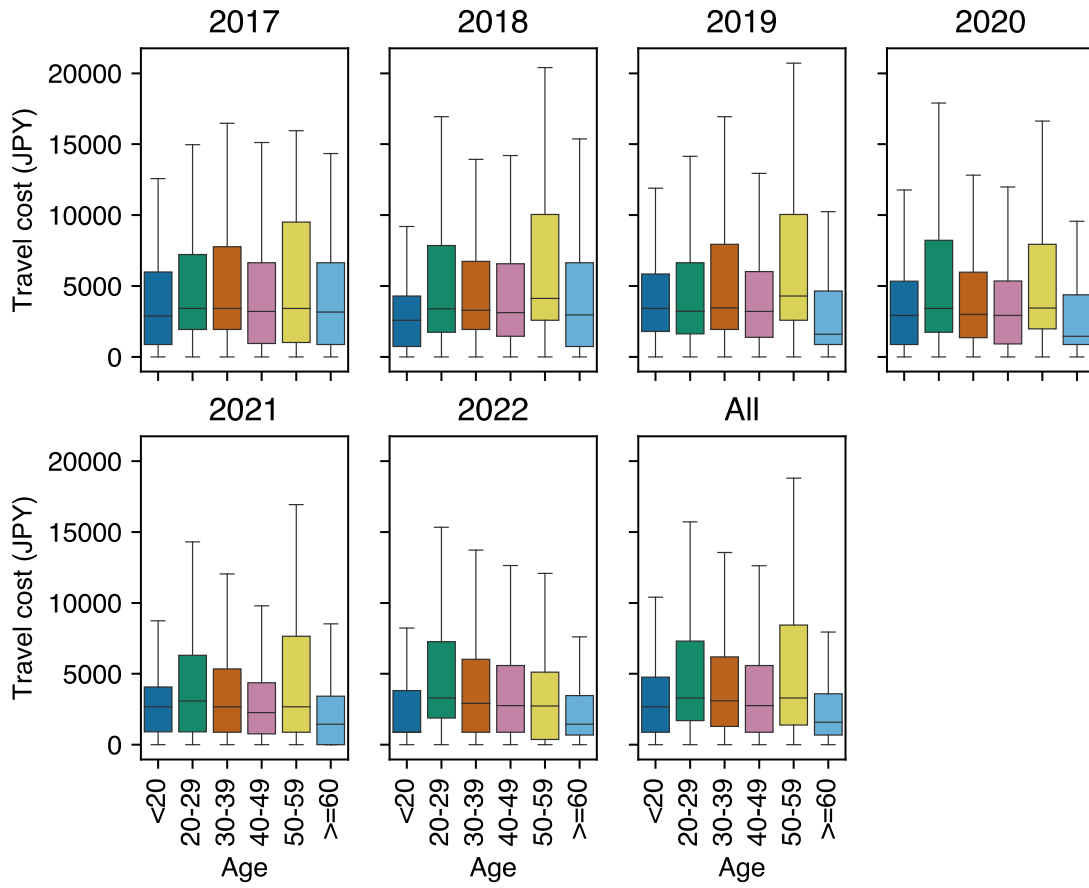
287 (a)



288

289

290 (b)



291

292

293 **Table S9:** Results of the type II analysis of deviance for the GLMM testing differences in sentiment
 294 score among client type (Twitter/Foursquare), age, sex, residential prefecture, and year. For this
 295 model, 150,646 tweets from 33,558 users were used. The marginal R^2 of the model was 0.025.

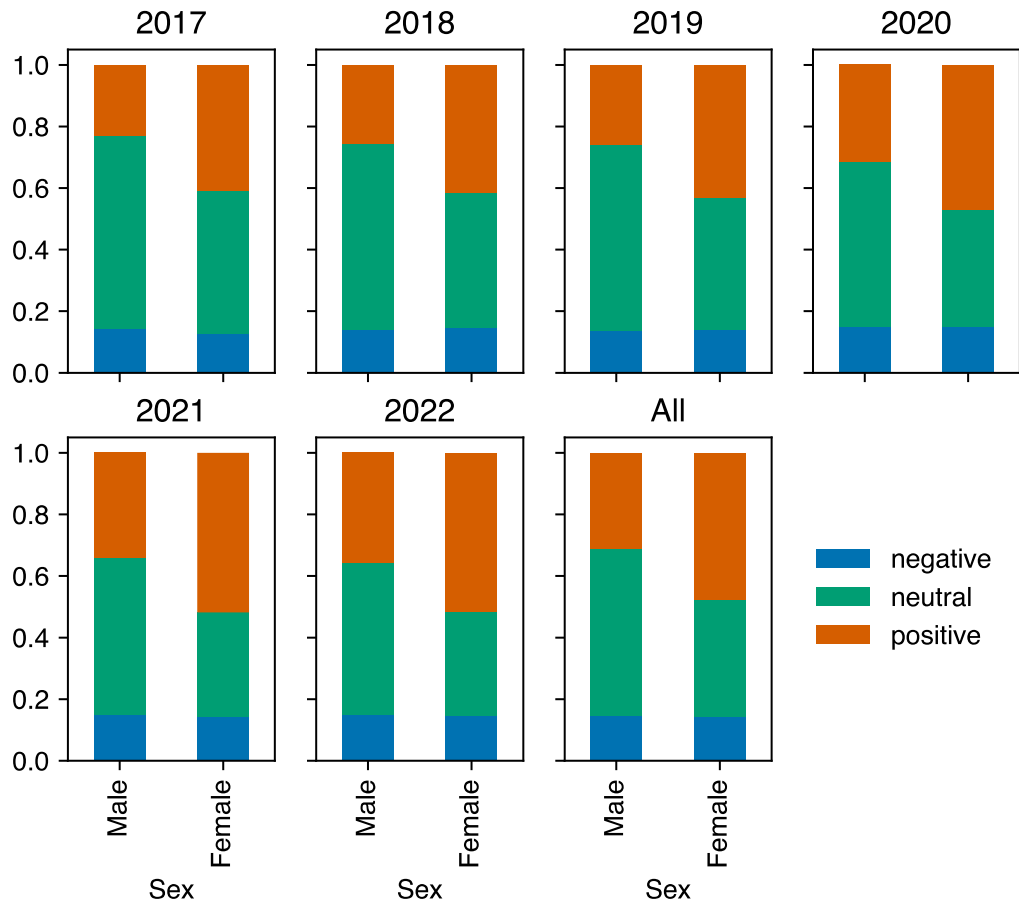
	χ^2	df	p	
Foursquare	797.383	1	<0.001	***
Age	119.57	5	<0.001	***
Sex	491.066	1	<0.001	***
Prefecture	68.644	46	0.01684	**
Year	97.371	5	<0.001	***

296

297

298 **Fig. S10:** Proportion of tweets having positive, neutral, or negative sentiment labels for (a) each sex
299 and (b) each age group for each year. Note that to reduce the dominance of massively tweeting users,
300 only one tweet for one day for one grid for each user was randomly selected and used.

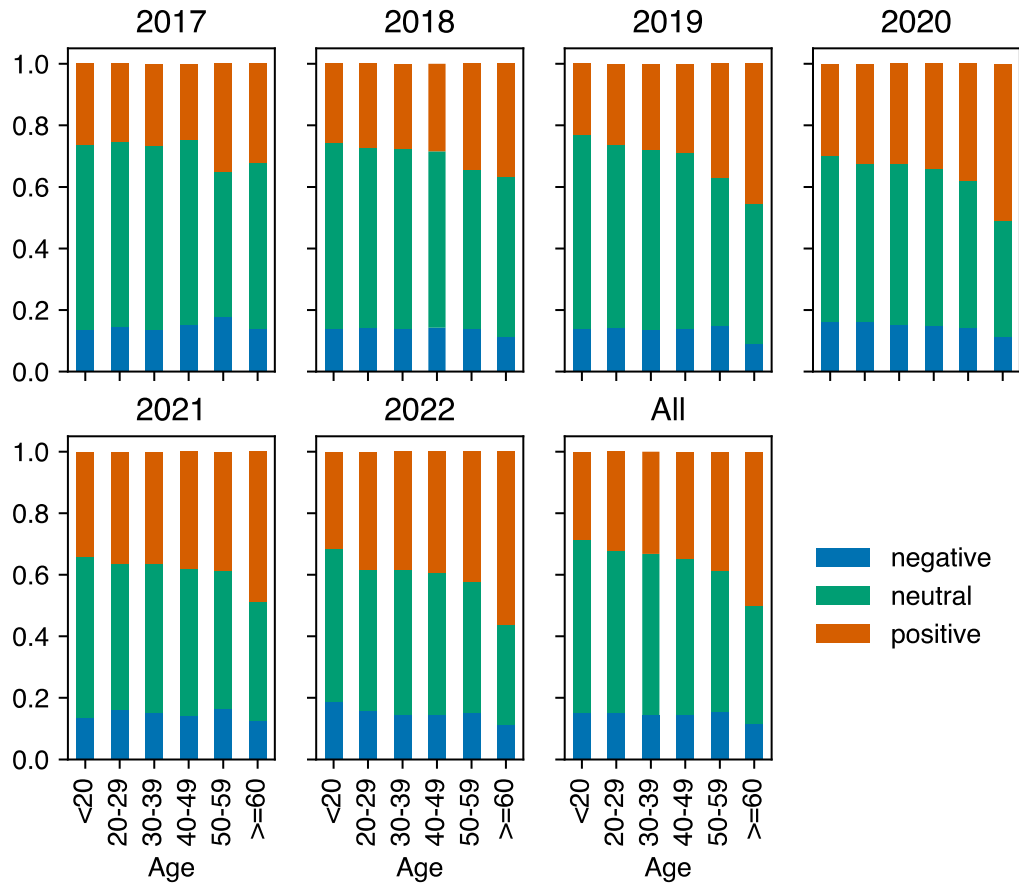
301 (a)



302

303

304 (b)



305

306

307 **Appendix 6:** Prediction of the sentiment labels using a machine learning model.

308 In addition to the model showing effects of the demographic attributes of users on sentiment score,
309 we implemented a machine learning model predicting sentiment labels (negative/neutral/positive)
310 from the demographic attributes. We used the random forest which can automatically incorporate
311 interactions among the explanatory variables implemented in the *randomForest* package (Liaw &
312 Wiener, 2002) with R. For the explanatory variables, we used sex, age, residential prefecture, client
313 type, and year. For the hyperparameters of the model, we used default values.

314 The resultant model could not predict sentiment labels from the explanatory variables. The
315 value of the Cohen's kappa (Cohen, 1960) for the prediction of the model was 0.152, indicating the
316 sentiment labels could not be predicted from the explanatory variables well (Landis & Koch, 1977).
317

318 **Table S10:** Result of the type II analysis of deviance for the GLMM testing effects of age, sex, and
 319 their interaction on the number of tweets per user. For this model 565,641 tweets from 126,494 users
 320 were used. The marginal R^2 of the model was 0.03.

	χ^2	df	p
Sex	616.791	1	<0.001 ***
Age	902.238	5	<0.001 ***
Sex \times Age	8.555	5	0.128

321
 322 **Table S11:** Result of the type II analysis of deviance for the GLMM testing effects of age, sex, and
 323 their interaction on the number of tweets per user after removing the top 1% of users with the highest
 324 number of tweets. For this model 421,789 tweets from 125,036 users were used. The marginal R^2 of
 325 the model was 0.02.

	χ^2	df	p
Sex	569.9959	1	<0.001 ***
Age	653.3263	5	<0.001 ***
Sex \times Age	5.0596	5	0.4086

326

327

328

329 **References**

- 330 Bivand, R. S., & Wong, D. W. S. (2018). Comparing implementations of global and local indicators
 331 of spatial association. *TEST*, 27(3), 716–748. <https://doi.org/10.1007/s11749-018-0599-x>
 332 Cohen, J. (1960). A Coefficient of Agreement for Nominal Scales. *Educational and Psychological*
 333 *Measurement*, 20(1), 37–46. <https://doi.org/10.1177/001316446002000104>
 334 Cramér, H. (1946). *Mathematical methods of statistics*. Princeton university press.
 335 Hayashibe, Y. (2020). Japanese realistic textual entailment corpus. *Proceedings of the Twelfth*
 336 *Language Resources and Evaluation Conference*, 6827–6834.
 337 <https://aclanthology.org/2020.lrec-1.843>
 338 Hijmans, R. J. (2024). *geosphere: Spherical Trigonometry* (Version R package version 1.5-20)
 339 [Computer software]. <https://CRAN.R-project.org/package=geosphere>
 340 Ikegami, Y. (2019). *Pymlask* (Version 0.3.2) [Computer software]. [https://github.com/ikegami-](https://github.com/ikegami-yukino/pymlask/)
 341 [yukino/pymlask/](https://github.com/ikegami-yukino/pymlask/)

- 342 Landis, J. R., & Koch, G. G. (1977). The Measurement of Observer Agreement for Categorical Data.
343 *Biometrics*, 33(1), 159. <https://doi.org/10.2307/2529310>
- 344 Liaw, A., & Wiener, M. (2002). Classification and Regression by randomForest. *R News*, 2(3), 18–
345 22.
- 346 Matthews, B. W. (1975). Comparison of the predicted and observed secondary structure of T4 phage
347 lysozyme. *Biochimica et Biophysica Acta (BBA) - Protein Structure*, 405(2), 442–451.
348 [https://doi.org/10.1016/0005-2795\(75\)90109-9](https://doi.org/10.1016/0005-2795(75)90109-9)
- 349 Ptaszynski, M., Dybala, P., Shi, W., Rzepka, R., & Araki, K. (2009). A System for Affect Analysis
350 of Utterances in Japanese Supported with Web Mining. *知能と情報*, 21(2), 194–213.
351 <https://doi.org/10.3156/jsoft.21.194>
- 352 Thakur, A. (2021a). *Autonlp-japanese-sentiment-59362* [Dataset].
353 <https://huggingface.co/abhishek/autonlp-japanese-sentiment-59362>
- 354 Thakur, A. (2021b). *Autonlp-japanese-sentiment-59363* [Dataset].
355 <https://huggingface.co/abhishek/autonlp-japanese-sentiment-59363>
356