

Comparing the Effectiveness of the Choropleth Map with a Hexagon Tile Map for Communicating Cancer Statistics in Australia

Stephanie Kobakian¹ and Dianne Cook²

Queensland University of Technology and Monash University

1. Introduction

This study compares the effectiveness of the spatial display, a hexagon tile map, against the standard, a choropleth map, for communicating information about disease statistics. The choropleth map is the traditional method for visualizing aggregated statistics across administrative boundaries. The hexagon tile map builds on existing displays, such as the cartogram, and tessellated hexagon displays. A hexagon tile map forgoes the familiar boundaries, in favor of representing each geographic unit as an equally sized hexagon, placed approximately in the correct spatial location. It differs in the relaxed requirement to have connected hexagons and allows sparsely located hexagons. This type of display may be useful for other countries, and other purposes. The algorithm to construct a hexagon tile map is available in the R package `sugarbaq` (Kobakian, Cook & Duncan 2023).

16 The hexagon tile map was designed for Australia, motivated by a need to display
17 spatial statistics for the Australian Cancer Atlas. None of the existing approaches for creating
18 cartograms or hexagon tiling perform well for the Australian landscape, which has vast open
19 spaces and concentrations of population in small regions clustered on the coastlines.

The Australian Cancer Atlas (Cancer Council Queensland 2018) is an online interactive web tool created to explore the burden of cancer on Australian communities. There are many cancer types to be explored individually or aggregated. The Australian Cancer Atlas allows users to explore the patterns in the distributions of cancer statistics over the geographic space of Australia. It uses a choropleth map display and diverging color scheme to draw attention to relationships between neighboring areas. The hexagon tile map may be a useful alternative display to enhance the atlas.

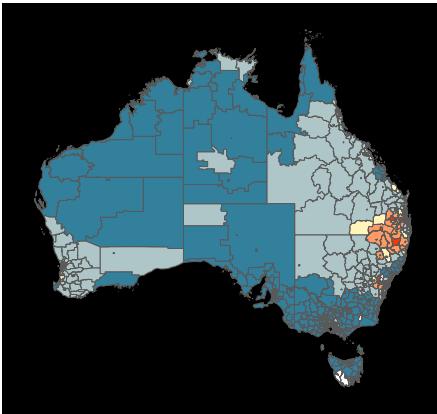
The experiment was conducted using the lineup protocol, a visual inference procedure (Wickham et al. 2010), to objectively test the effectiveness of the two displays.

¹ Science and Engineering Faculty, Queensland University of Technology, Brisbane, Queensland

² Department of Econometrics and Business Statistics, Monash University, Clayton, Victoria

Email: stephanie.kobakian@gmail.com

a. choropleth map



b. hexagon tile map

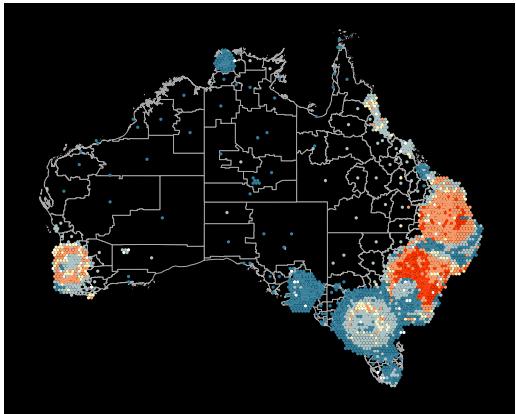


Figure 1. Thyroid incidence among females across the Statistical Areas of Australia at Level 2, displayed using a choropleth (left) and a hexagon tile map (right). Blue indicates lower than average, and red indicates higher than average incidence. The choropleth suggests high incidence is clustered on the east coast but misses the high incidence in Perth and a few locations in inner Melbourne visible in the hexagon tile map.

29 The paper is organised as follows. The next section discusses the background of
 30 geographic data display and visual inference procedures. Section 3 describes the methods
 31 for conducting the experiment and analysing the results. The results are summarized in the
 32 Section 4.

33

2. Background

34 2.1. Spatial data displays

35 Spatial visualisations communicate the distribution of statistics over geographic
 36 landscapes. The choropleth map (Tufte 1990), (Skowronnek 2016) is a traditional display. It is
 37 used to present statistics that have been aggregated on geographic units. Creating a choropleth
 38 map involves drawing polygons representing the administrative boundaries, and filling with
 39 colour mapped to the value of the statistic. The choropleth map places the statistic in the
 40 context of the spatial domain, so that the reader can see whether there are spatial trends,
 41 clusters or anomalies. This is important for digesting disease patterns. If there is a trend it
 42 may imply that the disease is spreading from one location to another. If there is a cluster, or
 43 an anomaly, there may be a localized outbreak of the disease. Aggregating the statistic on
 44 administrative units, provides a level of privacy to individuals, while allowing the impact of
 45 the disease on the community to be analyzed.

46 The choropleth map is an effective spatial display if the size of the geographic
47 units is relatively uniform. This is not the case for most countries. Size heterogeneity in
48 administrative units is particularly extreme in Australia: most of the landscape of Australia is
49 sparsely settled, with the population densely clustered into the narrow coastal strips. Figure
50 1 shows the choropleth map of thyroid cancer rates in Australia. The choropleth map focuses
51 attention on the geography, and for heterogeneously sized areas it presents a biased view of
52 the population-related distribution of the statistic (Kochmoud & House 1998). *Land doesn't get*
53 *cancer, people do* – a more effective way to communicate the spatial distributions of cancer
54 statistics is needed.

55 A cartogram is a general solution for better displaying a population-based statistic.
56 It transforms the geographic map base to reflect the population in the geographic region
57 while preserving some aspects of the geographic location. There are several cartogram
58 algorithms (Dorling 2011; Kochmoud & House 1998); each involves shifting the boundaries
59 of geographic units, using the value of the statistic to increase or decrease the area taken
60 by the geographic unit on the map. The changes to the boundaries result in cartograms that
61 accurately communicate population by map area for each of the geographic units but can
62 result in losing the familiar geographic information. For Australia, the transformations warp
63 the country so that it is no longer recognizable.

64 Alternative algorithms make various trade-offs between familiar shapes and
65 representation of geographic units. The non-contiguous cartogram method (Olson 1976)
66 keeps the shapes of geographic units intact and changes the size of the shape. This method
67 disconnects areas creating empty space on the display losing the continuity of the spatial
68 display of the statistic. The Dorling cartogram (Dorling 2011) represents each unit as a
69 circle, sized according to the value of the statistic. The neighbour relationships are mostly
70 maintained by how the circles touch. A similar approach was pioneered by Raisz (1963),
71 using rectangles that tile to align borders of neighbours (Monmonier 2005). There have been
72 thorough reviews of the array of methods, as suitable for cancer atlas displays (Kobakian,
73 Cook & Roberts 2020), (Skowronnek 2016).

74 The hexagon tile map algorithm, automatically matches spatial regions to their nearest
75 hexagon tile, from a grid of tiles. It has the effect of spreading out the inner city areas while
76 maintaining the spatial locations or regions in remote areas. The algorithm is available in
77 the R package, sugarbag (Kobakian, Cook & Duncan 2023). Figure 1 shows the hexagon
78 tile map. Colour maps from substantially below average (blue) to substantially above average
79 (red) rates. The inner city areas have expanded, making it possible to see the cancer incidence
80 in the small, densely populated areas. Remote regions are represented by isolated hexagons,
81 which is not ideal, but maintains the spatial location of these data values. It is of interest

82 to know how well the spatial distribution is perceived for this display, in comparison to the
83 choropleth.

84 **2.2. Visual Inference**

85 To assess the effectiveness of the hexagon tile map, the lineup protocol (Wickham et al.
86 2010; Buja et al. 2009) from visual inference procedures is employed. The approach mirrors
87 classical statistical inference. The procedures for doing a power comparison of competing
88 plot designed, outlined in Hofmann et al. (2012), are followed.

89 In classical statistical inference hypothesis testing is conducted by comparing the value
90 of a test statistic on a standard reference distribution, computed assuming the null hypothesis
91 is true. If the value is extreme, the null hypothesis is rejected, because the test statistic value
92 is unlikely to have been so extreme if it was true. In the lineup protocol, the plot plays the
93 role of the test statistic, and the data plot is embedded in a field of null plots. Defining the plot
94 using a grammar of graphics (Wickham 2009) makes it a functional mapping of the variables
95 and thus, it can be considered to be a statistic. With the same data, two different plots can be
96 considered to be competing statistics, one possibly a more powerful statistic than the other.

97 To do hypothesis testing with the lineup protocol requires human evaluation. The human
98 judge is required to identify the most different plot among the field of plots. If this corresponds
99 to the data plot – the test statistic – the null hypothesis is rejected. It means that the data plot
100 is extreme relative to the reference distribution of null plots.

101 The null hypothesis is explicitly provided by the grammatical plot description. For
102 example, if a histogram is the plot type being used, the null might be that the underlying
103 distribution of the data is a Gaussian. Null data would be generated by simulating from
104 a normal model, with the same mean and standard deviation as the data. In practice, the
105 null hypothesis used is generic, such as *there is NO structure or a pattern in the plot*, and
106 contrasted to an alternative that there is structure.

107 The chance that an observer picks the data plot out of a lineup of size m plots
108 accidentally, if the null hypothesis is true is $1/m$. With K observers, the probability of k
109 randomly choosing the data plot, roughly follows a binomial distribution with $p = 1/m$.
110 Figure 2 shows a lineup of the hexagon tile map, of size $m = 12$. Plot 3 is the data plot,
111 and the remaining 11 are plots of null data.

112 To determine the effectiveness of a type of display, this probability is less relevant than
113 the overall proportion of observers who pick the data plot, k/K . The power of the test statistic
114 (data plot) is provided by this proportion. Power in a statistical sense is the ability of the
115 statistic to *produce a rejection* of the null hypothesis if it is indeed *not true*. With the same

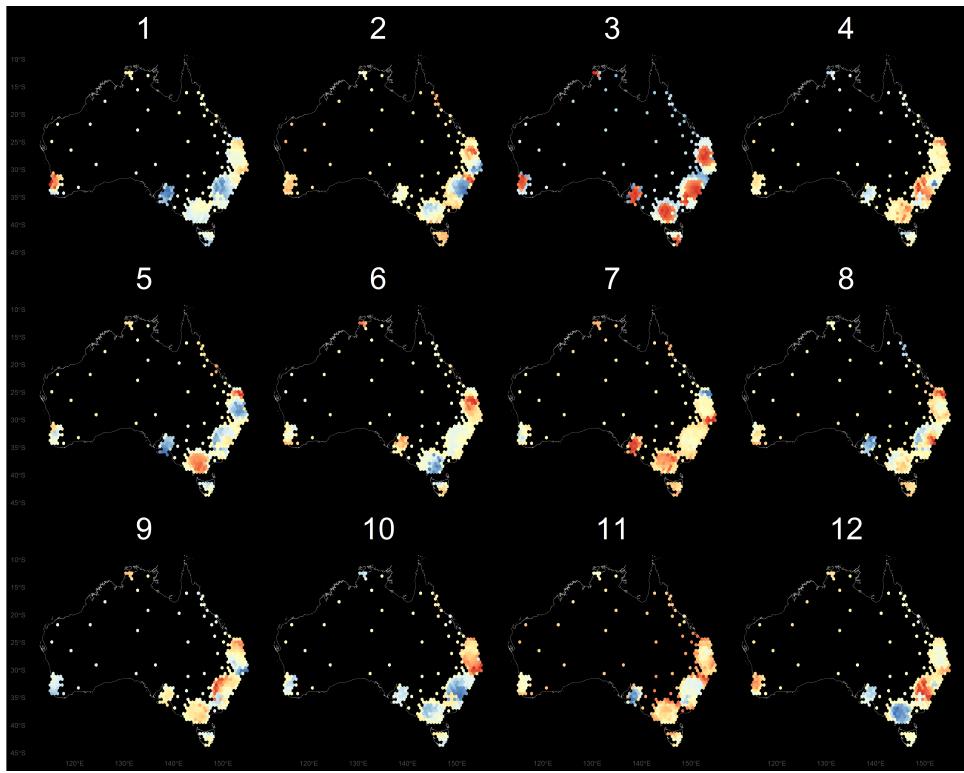


Figure 2. This lineup of twelve hexagon tile map displays contains one map with a real population-related structure. The rest are null plots that contain spatial correlation between neighbours.

116 data plotted using two different displays, the display with the highest proportion of people
117 who choose the data plot would be considered to be the most powerful statistic.

118

3. Methodology

119 This study aims to answer two key questions about the presentation of spatial
120 distributions:

- 121 1. Are spatial disease trends that impact highly populated small areas detected with higher
122 accuracy, when viewed in a hexagon tile map?
123 2. Are people faster in detecting spatial disease trends that impact highly populated small
124 areas when using a hexagon tile map?

125 Additional considerations when completing this experimental task included the
126 difficulty experienced by participants and the certainty they had in their decision.

127 Australia is used for the study, with Statistical Area 3 (SA3) (Australian Bureau of
128 Statistics 2018) as the geographic units. The results should apply broadly to any other
129 geographic area of interest.

130 **3.1. Experimental factors**

131 The primary factor in the experiment is the plot type. The secondary factor is a trend
132 model. Three trend models were developed, one mirroring a large spatial trend for which the
133 choropleth would be expected to do well, and two with differing levels of inner-city hot spots.
134 These latter two reflect the structure seen in thyroid cancer data (Figure 1). This produces six
135 treatment levels:

- 136 • Map type: *Choropleth, Hexagon tile*
137 • Trend: *South-East to North-West; Locations in three population centres; Locations in*
138 *multiple population centres,*

139 Data is generated for each of the trend models, with four replicates, and each displayed
140 both as a choropleth and as a hexagon tile map, which yields 12 data sets, and 24 data plots.
141 This set of displays is divided in half, providing two sets of 12 displays, Group A and Group
142 B. Participants were randomly allocated to Group A or B. Participants saw a data set only
143 once, either as a choropleth or as a hexagon tile map. Figure 3 summarises the design and the
144 allocation of the displays.

145 **3.2. Generating null data**

146 Null data needs to be data with no (interesting) structure. In most scenarios, permutation
147 is the main approach for generating null plots. It is used to break association between
148 variables while maintaining marginal distributions. This is too simple for spatial data. In
149 spatial data, a key feature is the spatial dependence or smoothness over the landscape. To
150 do something simple, like permuting the values relative to the geographic location would
151 produce null plots which are too chaotic, and the data plot will be recognisable for its
152 smoothness rather than any structure of interest.

153 For spatial data, null data is stationary data, where the mean, variance and spatial
154 dependence are constant over the geographic units. Stationary data is specified by a variogram
155 model (Matheron 1963). Simulating from a variogram model, where the spatial dependence
156 is specified, generates the stationary spatial data used for the null plots. The parameters for
157 the Gaussian model were sill=1, range=0.3 with the variance generated by a standard normal
158 distribution.

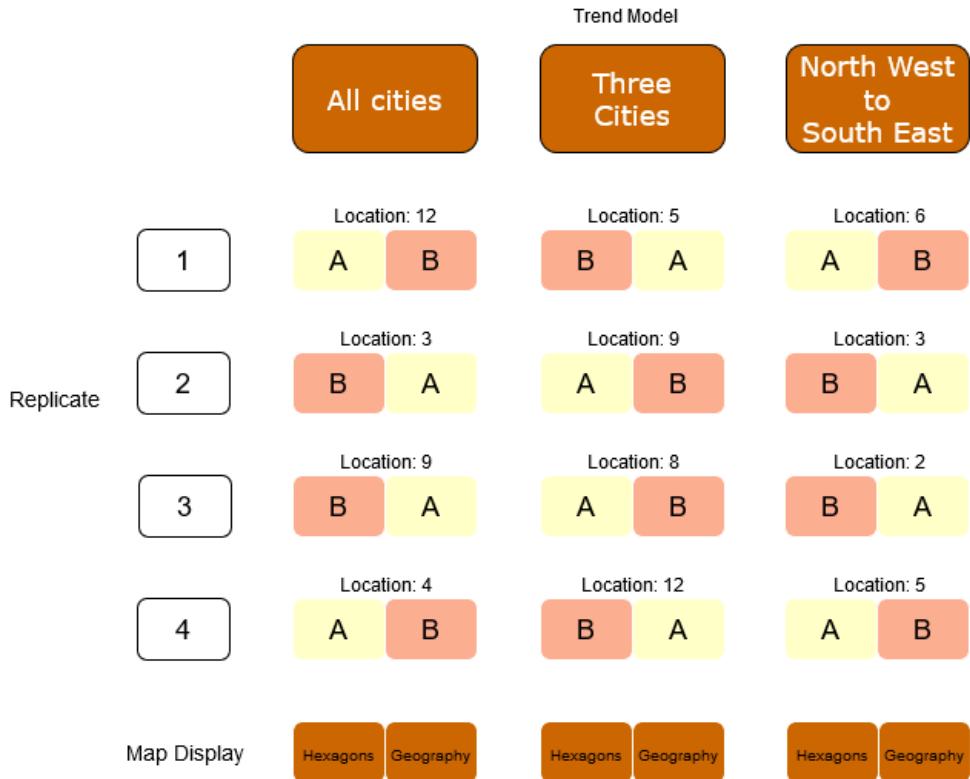


Figure 3. The experimental design used in the visual inference study. Twelve data sets, with four replications of each trend model, are composed into lineups and displayed as either a choropleth map or hexagon tile map. These are organised into two packets of plots, A and B, which are provided to participants to evaluate.

159 The R package `gstat` (Gräler, Pebesma & Heuvelink 2016) was used to simulate 144
 160 null sets, 12 data sets for each plot in a lineup, and 12 sets for 12 lineups.

161 The null model imposed by our hypothesis suggests that neighbours are related. The
 162 randomness induced when generating the null data was smoothed to mirror the practices
 163 employed by the Australian Cancer Atlas statisticians. In these 12 sets of data, each of the
 164 12 maps was smoothed several times to replicate the spatial autocorrelation seen in cancer
 165 data sets presented in the Australian Cancer Atlas, without implementing uncertainty via
 166 transparency.

167 A list of neighbours for each geographic unit was generated to use when smoothing the
 168 distributions. For each geographic unit, the same spatial smoother was applied in each layer
 169 of smoothing. It kept half of the units' previous value and derived the new half as the mean
 170 of the values of its neighbours at the previous layer of smoothing.

171 This smoothing allowed neighbors to be related to each other, but also allowed outliers,
172 and showed distributions similar to the thyroid cancer distribution (Figure 1).

173 **3.3. Generating lineups**

174 For each trend model, four real data displays were created by manipulating the centroid
175 values of each of the SA3 geographic units.

176 The North West to South East (NW-SE) distribution was created using a linear equation
177 of the centroid longitude and latitude values.

178 The All Cities trend model was created using the distance from the centroid of each
179 geographic unit to the closest capital city in Australia, calculated when creating the hexagon
180 tile map using the *sugarbag* (Kobakian, Cook & Duncan 2023) package. 201 of the 336
181 SA3s were considered greater capital city areas, the values of these areas were increased to
182 create red clusters. The amount was chosen to make clusters around the cities visible in the
183 choropleth display even if they were not overtly noticeable.

184 A similar selection process was applied to the Three Cities' trend model. However, for
185 each of the four replicates for the Three Cities trend, a random sample of capital cities was
186 taken from Sydney, Brisbane, Melbourne, Adelaide, Perth, and Hobart. Only the values of
187 the areas nearest to the three cities were increased to create clusters.

188 One of the lineup locations was chosen to embed the real trend model map, in each of
189 the four replicates, for the three trend models. The location was chosen from a sub-sample
190 of the 12 possible locations. The chance of repetition using resampling was introduced to
191 prevent participants from inducing the location by elimination, the locations 1, 7, 10 and 11
192 were not used.

193 As seen in Fig 3, the choropleth and hexagon display used the same location for the
194 real data display of the trend model was added to the spatially correlated null values for each
195 lineup. Each set of lineup data was used to produce a choropleth map lineup and hexagon tile
196 map lineup. These matched pairs were split between Group A and Group B according to the
197 2 x 3 factor experimental design depicted in 3.

198 For each of the 144 individual maps, the values for each geographic area were rescaled
199 to create a similar colour scale from deep blue to dark red within each map. This meant at
200 least one geographic unit was coloured dark blue and at least one was red, in every map
201 display of every lineup.

202 For the geographic NW-SE distribution, this resulted in the smallest values of the trend
203 model (blue) occurring in Western Australia, the North West of Australia, and the largest
204 values of the trend model (red) occurring in the South East. This resulted in Tasmania being
205 coloured completely red.

206 For the population-related displays, the clusters in the cities appeared more red than the
 207 rest of Australia.

208 **3.4. Analysis**

209 **3.4.1. Data Cleaning**

210 The first step in the data cleaning process involved checking that survey responses
 211 collected for each participant were only included once in the data set. The data cleaning
 212 process also involved filtering out participants' who did not provide at least three unique
 213 choices when considering each of the twelve lineups. These participants achieved a detection
 214 rate of 0. If participants had made various plot choices for the 12 displays they saw they were
 215 still included in the dataset.

216 **3.4.2. Descriptive statistics**

217 Basic descriptive statistics were used to contrast the detection rate for the two types
 218 of displays. Comparison was also made across the trend models, contrasting the mean and
 219 standard detection rate for each group, who had seen the different map display types for each
 220 replicate.

221 Side-by-side dot plots were made of accuracy (efficiency) against the map type, faceted
 222 by trend model type.

223 Similar plots were made of the feedback and demographic variables - reason for choice,
 224 reported difficulty, gender, age, education, having lived in Australia - against the design
 225 variables.

226 Plots will be made in R (R Core Team 2019), with the `ggplot2` package (Wickham
 227 2009).

228 **3.4.3. Modelling**

229 The likelihood of detecting the data plot in the lineup can be modelled using a linear
 230 mixed-effects model. The R (R Core Team 2019) `glmer()` function in the `lme4` (Bates
 231 et al. 2015) package implements generalised linear mixed effect models. The model used
 232 includes the two main effects map type and trend model, which gives the fixed effects model
 233 to be:

$$\widehat{y_{ij}} = \mu + \tau_i + \delta_j + (\tau\delta)_{ij} + \epsilon_{i,j}, \quad i = 1, 2; j = 1, 2, 3$$

234 where $y_{ij} = 0, 1$ is the log odds for whether the subject detected the data plot, μ is the
 235 overall mean, $\tau_i, i = 1, 2$ is the map type effect, δ_j is the trend model effect. We are allowing

236 for an interaction between map type and trend model as the response is binary, so a logistic
 237 model was used. As each participant provides results from 12 lineups, this model can account
 238 for each individual participants' abilities as it includes a subject-specific random intercept.

239 The model specifies a logistic link, this means the predicted values from the `glmer`
 240 model should be back-transformed to fit between 0 and 1. The predictions $\hat{p}(\eta)$ are
 241 transformed to be probabilities between 0 and 1 with the link specified below:

$$\hat{p}(\eta) = \frac{e^\eta}{1 + e^\eta}$$

242

$$\eta = f(\tau_i, \delta_j)$$

243 3.5. Web application to collect responses

244 The `taipan` (Kobakian & O'Hara-Wild 2018) package for R was used to create the
 245 survey web application. This structure was altered to collect responses regarding participant's
 246 demographics and their survey responses. The survey app contained three tabs. Participants
 247 were first asked for their demographics their Figure Eight contributor ID, and their consent
 248 to the responses being used for analysis. The demographics collected included participants'
 249 preferred pronouns, the highest level of education achieved, their age range and whether they
 250 had lived in Australia.

251 After submitting these responses, the survey application switched to the tab of lineups
 252 and associated questions. This allowed participants to easily move through the twelve
 253 displays and provide their choice, reason for their choice, and level of certainty.

254 When participants completed the twelve evaluations the survey application triggered
 255 a data analysis script. This created a data set with one row per evaluation, containing the
 256 responses to the three questions. The script also added the title of the image, which indicated
 257 the type of map display, the type of distribution hidden in the lineup, and the location of the
 258 data plot. It also calculated the time taken by participants to view each lineup.

259 Each participant used the internet to access the survey. The data transfer from the web
 260 application to the data set took place using a secure link to the Google sheet used to store
 261 results. The application connected to the Google sheet using the `googlesheets` (Bryan &
 262 Zhao 2018) R package when participants opened the application, and interacted again when
 263 participants chose to submit the survey. At this time it added the participant's responses to
 264 the twelve lineup displays as twelve rows of data in the Google sheet.

265 **3.6. Participants**

266 Participants were recruited from the Figure Eight crowdsourcing platform (Figure Eight
267 Inc 2019) to evaluate lineups. The lineup protocol expects that the participants are uninvolved
268 judges with no prior knowledge of the data, to avoid inadvertently affecting results. Potential
269 participants needed to have achieved level 2 or level 3 from prior work on the platform. All
270 participants were at least 18 years old.

271 Participants were allocated to either group A or group B when they proceeded to the
272 survey web application. There were 92 participants involved in the study. All participants
273 read introductory materials and were trained using three test displays, to orient them to the
274 evaluation task. All participants who completed the task were compensated \$AUD5 for their
275 time, via the Figure Eight payment system.

276 A pilot study was conducted in the working group of the Econometrics and Business
277 Statistics Department of Monash University. This allowed us to estimate the effect size, and
278 thus decide on the number of participants to collect responses from.

279 **3.7. Demographic data collection**

280 Each participant answered demographic questions and provided consent before
281 evaluating the lineups.

282 Demographics were collected regarding the study participants:

- 283 • Gender (female / male / other),
284 • Education level achieved (high school / bachelors / masters / doctorate / other),
285 • Age range (18-24 / 25-34 / 35-44 / 45-54 / 55+ / other)
286 • Lived at least for one year in Australia (Yes / No)

287 Participants then moved to the evaluation phase. The set of images differed for Group A
288 and Group B. After being allocated to a group, each individual was shown the 12 displays in
289 randomised order.

290 Three questions were asked regarding each display:

- 291 • Plot choice
292 • Reason
293 • Difficulty

294 After completing the 12 evaluations, the participants were asked to submit their
295 responses.

296

4. Results

297 Responses from 92 participants were collected. Five participants did not provide more
 298 than three unique choices for the twelve lineups, and their data was removed. Set A was
 299 evaluated by 42 participants, and 53 evaluated set B. This resulted in 1104 evaluations,
 300 corresponding to 92 subjects, each evaluating 12 lineups, that were analysed on accuracy
 301 and speed. The certainty and reasons of subjects in their answers are also examined.

302 **4.1. Participant demographics**

303 Of the 92 participants, 67 were male, and 25 female. Most participants (56) had a
 304 Bachelors degree, 13 had a Masters degree, and the remaining 23 had high school diplomas.

305 **4.2. Accuracy**

306 Figure 4 displays the average detection rates for the two types of plot separately for each
 307 trend model. Each trend model was tested using four repetitions, evaluations on the same
 308 data set were seen as either choropleths or hexagon tile maps by each group as specified in
 309 Figure 3; the detection rates for each display are connected by a line segment. The Three
 310 Cities and All Cities trend models shown in the hexagon tile map allowed viewers to detect
 311 the data plot substantially more often than the choropleth counterparts. One replicate for the
 312 All Cities group had similar detection rates for both plot types, the rate of detection using the
 313 choropleth map was much higher than other replicates. Surprisingly, participants could also
 314 detect the gradual spatial trend in the NW-SE group from the hexagon tile map. We expected
 315 that the choropleth map would be superior for the type of spatial pattern, but the data suggests
 316 the hexagon tile map performs slightly better, or equally as well.

317 Table 1 shows the means and standard deviations of the detection rate for each type
 318 of plot and each trend model. This also gives the standard deviations, the smallest standard
 319 deviation for all sets of replicates was the Three Cities trend model shown in a Choropleth
 320 display. This group of displays had a very small detection rate of 0.04. The mean detection
 321 rate for the Three Cities trend model shown as choropleth map lineups was also the smallest
 322 at 0.40. The North-West to South-East (NW-SE) trend model unexpectedly had a higher mean
 323 detection rate for the hexagon tile map displays, but the difference in the means of detection
 324 rate was only 0.10.

325 Table 2 presents a summary of the generalised linear mixed effects model, testing the
 326 effect of plot type and trend model on the detection rate. The results support the summary
 327 from Figure 4 and all parameters are statistically significant despite the large standard
 328 deviations observed in Table 1. Overall, the hexagon tile map performs marginally better
 329 than the choropleth for all trend models, which is a pleasant surprise. Allowing for the

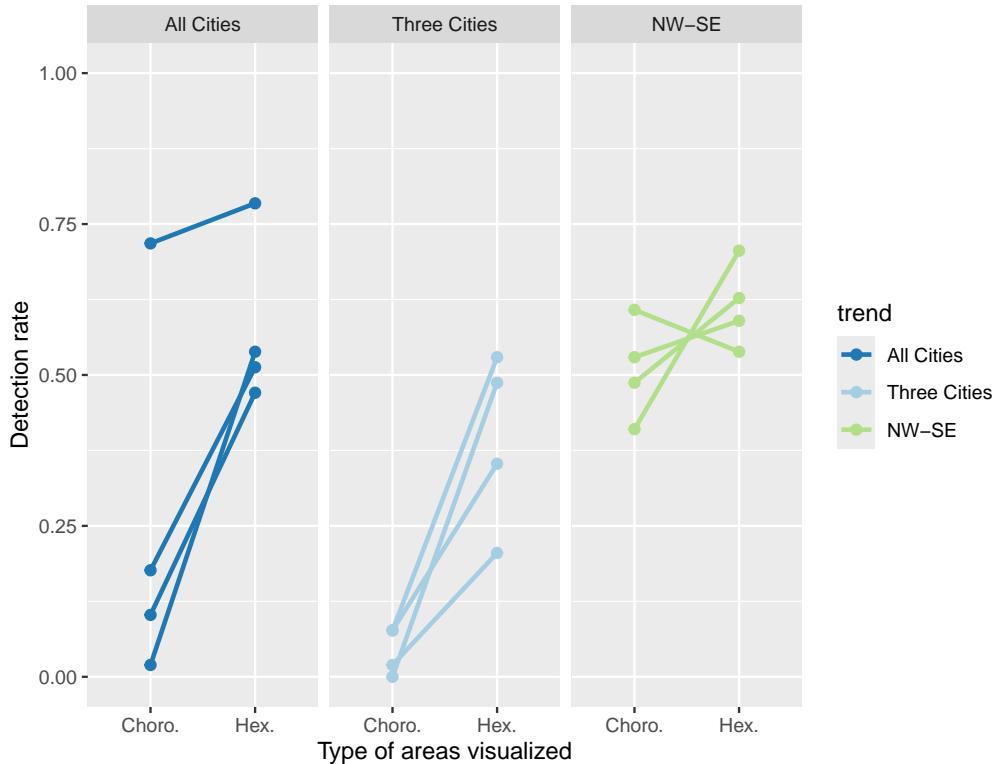


Figure 4. The detection rates achieved by participants are contrasted when viewing the four replicates of the three trend models. Each point shows the probability of detection for the lineup display, the facets separate the trend models hidden in the lineup. The points for the same data set shown in a choropleth or hexagon tile map display are linked to show the difference in the detection rate.

330 interaction effect, the difference in detection rate decreases for population-related displays
 331 for a choropleth map lineup but increases for a hexagon tile map display. The log odds of
 332 detection shown in Table 2 can be back-transformed after taking the sum of all terms for
 333 the trend and type of display that are of interest. For the NW-SE distribution, the predicted
 334 detection rate for the hexagon tile map display increases the predicted probability of detection
 335 to 0.63 from 0.52 for choropleths, this is almost exactly the difference seen in the table of
 336 means and is significant only at the 0.05 level.

337 When a choropleth map display is used, the predicted detection rate for the Three Cities
 338 trend, 0.03; this is extremely low, especially compared to the NW-SE trend of 0.52. When the
 339 All Cities trend is presented in a choropleth display the predicted probability of detection is
 340 0.22. The hexagon tile map has a substantially high detection rate for the display of a Three
 341 Cities trend 0.39 and All Cities trend 0.59.

342 **4.3. Speed**

343 Figure 5 shows horizontally jittered dot plots to contrast the time taken by participants
 344 to evaluate each lineup when viewing each type of display. The time is also separated by
 345 the trend model and whether the data plot was detected or not detected. The time taken to
 346 complete an evaluation ranged from milliseconds to 60 seconds. The average time taken
 347 for each type of display is shown as a large coloured dot on each plot. When considering the
 348 heights of the green and orange dots, there is little difference in the average time taken to read
 349 a choropleth or hexagon tile map. Comparing the same coloured dot across each trend model
 350 row, there is a slight increase in the time taken to correctly detect the data plot in the hexagon
 351 tile map lineup, but little difference in evaluation time for the choropleth display. However,
 352 there were substantially fewer correct detections for choropleth lineups for the Three cities
 353 and All Cities trends.

354 **4.4. Certainty**

355 Participants provided their level of certainty regarding their choice using a five-point
 356 scale. Unlike the accuracy and speed of responses that were derived during the data
 357 processing phase, this was a subjective assessment by the participant prompted by the
 358 question: ‘How certain are you about your choice?’. Figure 6 shows the number of times
 359 participants provided each level of certainty. This was separated for each combination of trend
 360 models and display type, and colored depending on whether a participant correctly detected
 361 the data plot in the lineup. Participants often chose 4 or 5 when viewing the population-
 362 related trends in the choropleth display, even though they were often incorrect when viewing
 363 an All Cities trend and overwhelmingly incorrect for the Three Cities trend. This shows
 364 overconfidence in their detection ability when using a choropleth map display. Participants
 365 were less likely to be certain when their choice was incorrect and they were viewing a
 366 hexagon tile map. For each trend model, participants were more likely to doubt their choice
 367 and choose 1 or 2 in the hexagon tile map displays, even though many had made the correct
 368 choice.

369 **4.5. Reason**

370 Participants were asked why they had made their plot choice and were able to select from
 371 a set of suggested reasons. “Color trend across the areas” was the most common selection for
 372 NW-SE trend displays.

373 The reasons chosen by participants from the list provided to them varied more when
 374 viewing choropleth displays than the hexagon tile map. The hexagon tile map displays
 375 resulted in “Clusters of colour” as the most common choice made by participants.

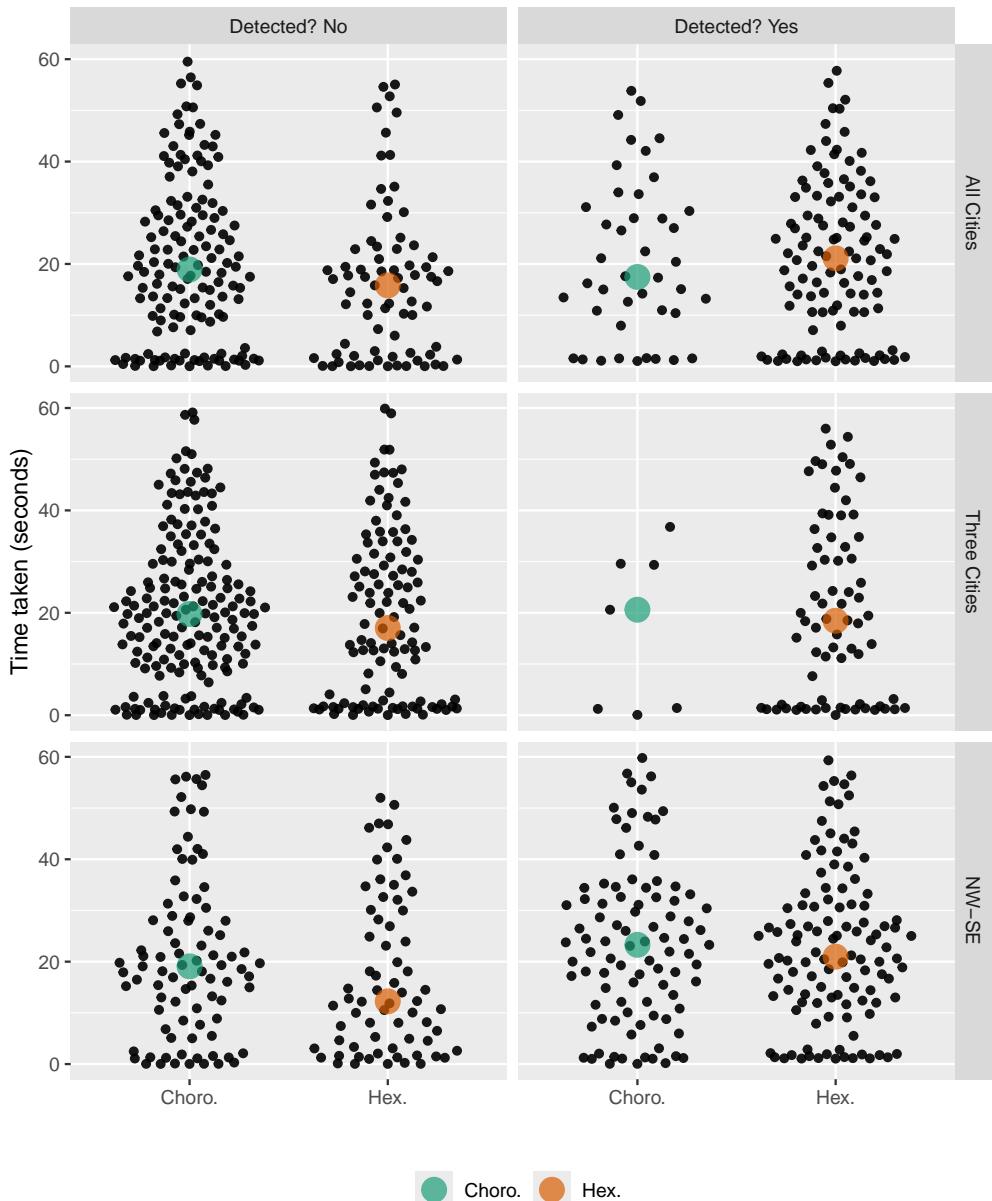


Figure 5. The distribution of the time taken (seconds) to submit a response for each combination of trend, whether the data plot was detected, and type of display, shown using horizontally jittered dot plots. The coloured point indicates the average time taken for each plot type. Although some participants take just a few seconds per evaluation, and some take as much as 60 seconds, there is very little difference in time taken between plot types.

376 The choice “None of these reasons” was used as the default value to minimise noise
 377 from participants who did not select a response.

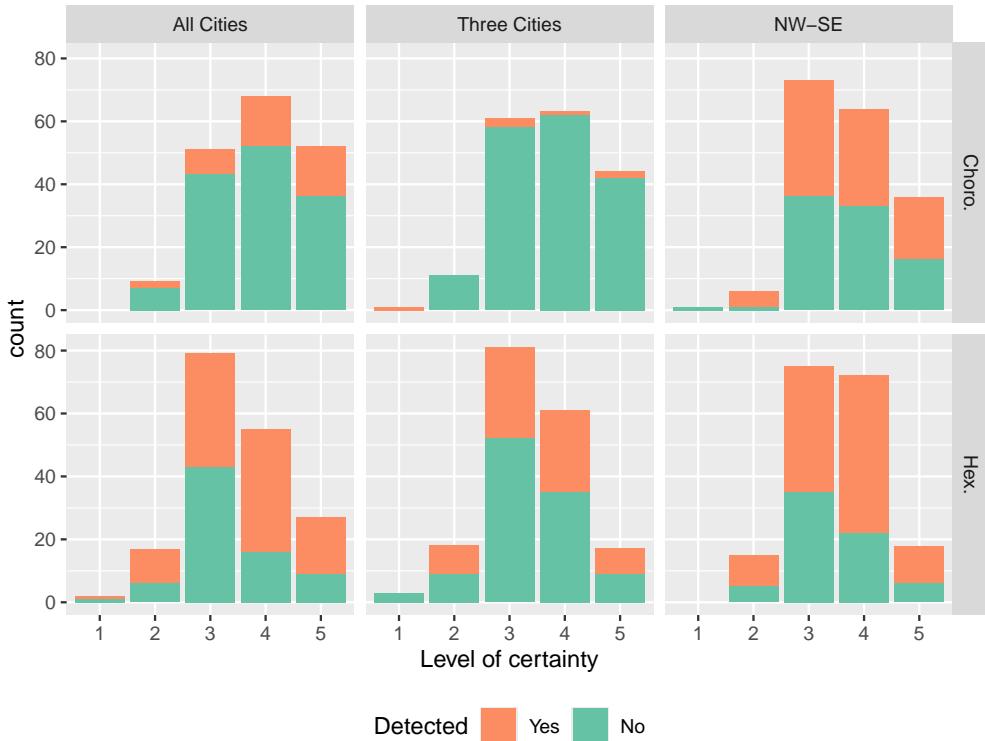


Figure 6. The number of times each level of certainty was chosen by participants when viewing hexagon tile map or choropleth displays. Participants were more likely to choose a high certainty when considering a Choropleth map. The mid value of 3 was the default certainty, it was chosen most for the Hexagon tile map displays.

378

5. Discussion

379 The intention of this study was to contrast the use of the choropleth map and the hexagon
 380 tile map. The visual inference lineup protocol was employed to contrast the effectiveness
 381 of the displays. The results have shown that overall the use of the hexagon tile map
 382 display allows participants to find the data plot in the lineup more often. Using the visual
 383 inference protocol this result can be extended to show that it is a valid alternative display to
 384 communicate spatial distributions of population related data.

385 We expected that the choropleth map would be superior for communicating the spatial
 386 pattern of geographic distributions. The data suggest that the participants perform slightly
 387 better or equally as well for each replicate in each trend model across the two displays. Table
 388 1 shows that the difference in the mean detection rate for the two trend models was 0.10.

389 The differences are seen in Figure 4 and Table 1 are reflected in the model results.
 390 Surprisingly the difference in the geographic distribution was significant at the 0.05 level. It
 391 also showed that the hexagon tile map display performs marginally better than the choropleth

392 for all trend models. Unexpectedly the detection rate suffers when using a choropleth map to
393 display population-related distributions.

394 While the significance of the difference in detection was the key focus of this
395 experiment, the secondary focus was the time taken by participants. It was expected that the
396 participants might take longer to consider the hexagon tile map distribution but would be able
397 to detect the data plot in the lineup. The bimodal distributions seen in Figure 5 showed very
398 little difference in the mean evaluation times. As the maximum time of all of the distributions
399 approached 60 seconds it cannot be said that the participants took longer to evaluate the
400 hexagon tile map displays.

401 The responses to the questions asked of participants included the reason for their choice
402 and the certainty around their choice. Figure 3 shows high levels of certainty of 4 and 5
403 were chosen by participants when looking at the population distributions in a choropleth map
404 display showing that they were overconfident when attempting to find the real data plot in the
405 choropleth map displays. Participants performed better on the NW-SE distribution shown in
406 the choropleth display and were reasonably confident about their decisions. The high levels of
407 the mid-range value of 3 could indicate that the participant did not want to provide a response,
408 as this was the default value. Those who chose level 4 or 5 were equally likely to be correct
409 for the three cities lineups, but more likely to be correct than incorrect for the other two trend
410 models.

411 The colour scaling applied in Three cities and All cities displays resulted in the rural
412 areas of the real data plot appearing more blue or yellow than the other plots in the lineups.
413 Due to the consistent colouring of rural areas in a choropleth display, the choice “All areas
414 have similar colours” was the most common reason for a participant’s choice. The All Cities
415 displays coloured the inner-city areas of all capital cities redder, this was observable to
416 participants and explains the equal choice of the city clusters or rural colour consistency.
417 Choosing “Clusters of colour” was expected when participants viewed the Hexagon tile map
418 display of the All Cities and Three Cities distributions. It was unexpected that it was also the
419 most common reason for the NW-SE hexagon tile map displays. Due to the spatial covariance
420 introduced in the smoothing, groups of similarly coloured hexagons were present in all of the
421 hexagon tile map displays. All Cities and Three Cities distributions of real data trends had
422 distinctly different patterns or red inner-city areas, while some of the plots in each lineup may
423 have shared similar features.

424

6. Conclusion

425 The choropleth map display and the tessellated hexagon tile map have been contrasted
 426 using the lineup protocol. The hexagon tile map was significantly more effective for spotting
 427 a real population-related data trend model hidden in a lineup.

428 The hexagon tile map display should be considered as an alternative visualization
 429 method when communicating distributions that relate to the population across a set of
 430 geographic units. As an additional display to the familiar choropleth map, cancer atlas
 431 products may benefit from the opportunity to allow exploration via an alternative display. The
 432 spatial distributions used to test these displays were inspired by the real spatially smoothed
 433 estimates of the cancer burden on Australian communities. However, this technique may be
 434 extended to other population-related distributions, such as other diseases.

435 The increasing population densities of capital cities despite large land area exacerbates
 436 the difference between the smallest and largest communities. The population density structure
 437 of Australia can be considered similar to that of Canada, New Zealand and many other
 438 countries. Therefore, this display is not only relevant to Australia but all nations or population
 439 distributions that experience densely populated cities separated by vast rural expanses.

440

7. Acknowledgment

441 The authors would like to thank the Australian Cancer Atlas team for discussions
 442 regarding alternative spatial visualizations, and Professor Kerrie Mengersen and Dr Earl
 443 Duncan for regular meetings filled with suggestions and comments. Mitchell O'Hara-Wild
 444 was a co-developer of the *taipan* (Kobakian & O'Hara-Wild 2018) R package for image
 445 tagging, used as the base for the web app constructed to collect participant evaluations
 446 of lineups. We are thankful for the NUMBATs (Non-Uniform Monash Business Analytics
 447 Team) for participating in the pilot study that helped to assess the experimental design and
 448 determine an appropriate sample size for the study.

449 The source code to produce this document can be found on GitHub. Supplementary
 450 materials have been included to discuss the survey procedures and the lineups that were used.
 451 The full set of images can be found here, too.

452 The supplementary material contains:

- 453 • Additional analysis of the experimental results
- 454 • Survey procedure including training materials for the participants
- 455 • 24 lineups as images, that were used in the experiment
- 456 • 12 data sets were used to construct the lineups

457 The analysis of the work was completed in R (R Core Team 2019) with the use of the
458 following packages:

- For document creation: `rmarkdown` (Xie, Allaire & Grolemund 2018), `knitr` (Xie 2015).
 - For lineup creation and data analysis: `tidyverse` (Wickham et al. 2019), `nullabor` (Wickham et al. 2018), `ggthemes` (Arnold 2019), `RColorBrewer` (Neuwirth 2014).
 - For image displays: `cowplot` (Wilke 2019), `png` (Urbanek 2013), `grid` (Murrell 2002).
 - For modelling and presentation of models: `gstat` (Gräler, Pebesma & Heuvelink 2016), `lme4` (Bates et al. 2015), `kableExtra` (Zhu 2019).

468 Ethics approval for the online survey was granted by QUT's Ethics Committee (Ethics
469 Application Number: 1900000991). All applicants provided informed consent in line with
470 QUT regulations prior to participating in this research.

References

- 472 ARNOLD, J.B. (2019). *ggthemes: Extra Themes, Scales and Geoms for 'ggplot2'*. URL <https://CRAN.R-project.org/package=ggthemes>. R package version 4.2.0.

473 AUSTRALIAN BUREAU OF STATISTICS (2018). Australian Statistical Geography Standard (ASGS). URL <https://www.abs.gov.au/statistics/statistical-geography/australian-statistical-geography-standard-asgs>.

474 BATES, D., MÄCHLER, M., BOLKER, B. & WALKER, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* **67**, 1–48. doi:10.18637/jss.v067.i01.

475 BRYAN, J. & ZHAO, J. (2018). *googlesheets: Manage Google Spreadsheets from R*. URL <https://CRAN.R-project.org/package=googlesheets>. R package version 0.3.0.

476 BUJA, A., COOK, D., HOFMANN, H., LAWRENCE, M., LEE, E.K., SWAYNE, D.F. & WICKHAM, H. (2009). Statistical inference for exploratory data analysis and model diagnostics. *Philosophical Transactions of the Royal Society, A (Invited)* **367**, 4361–4383. doi:10.1098/rsta.2009.0120.

477 CANCER COUNCIL QUEENSLAND (2018). Australian Cancer Atlas, publisher = Queensland University of Technology, Cooperative Research Centre for Spatial Information, issue = Version 09, url=<https://atlas.cancer.org.au>, accessed = Jan 12 2020.

478 DORLING, D. (2011). *Area Cartograms: Their Use and Creation*, vol. 59, chap. 3.7. John Wiley & Sons, Ltd, pp. 252–260. doi:10.1002/9780470979587.ch33.

479 FIGURE EIGHT INC (2019). The essential high-quality data annotation platform. URL <https://www.figure-eight.com/>.

480 GRÄLER, B., PEBESMA, E. & HEUVELINK, G. (2016). Spatio-temporal interpolation using gstat. *The R Journal* **8**, 204–218. URL <https://journal.r-project.org/archive/2016/RJ-2016-014/index.html>.

481 HOFMANN, H., FOLLETT, L., MAJUMDER, M. & COOK, D. (2012). Graphical tests for power comparison of competing designs. *IEEE Transactions on Visualization and Computer Graphics* **18**, 2441–2448.

482 KOBAKIAN, S., COOK, D. & DUNCAN, E. (2023). A hexagon tile map algorithm for displaying spatial data. *The R Journal* **15**, 6–16. doi:10.32614/RJ-2023-021. <https://doi.org/10.32614/RJ-2023-021>.

- 498 KOBAKIAN, S., COOK, D. & ROBERTS, J. (2020). Mapping cancer: the potential of cartograms and
 499 alternative map displays. *Annals of Cancer Epidemiology* **4**. URL <https://ace.amegroups.com/article/view/6040>.
- 501 KOBAKIAN, S. & O'HARA-WILD, M. (2018). *taipan: Tool for Annotating Images in Preparation for Analysis*. URL <https://CRAN.R-project.org/package=taipan>. R package version 0.1.2.
- 503 KOCHMOUD, C. & HOUSE, D. (1998). A Constraint-based Approach to Constructing Continuous Cartograms. In *Proc. Symp. Spatial Data Handling*. pp. 236–246.
- 505 MATHERON, G. (1963). Principles of geostatistics. *Economic Geology* **58**, 1246–1266. URL <http://dx.doi.org/10.2113/gsecongeo.58.8.1246>.
- 507 MONMONIER, M. (2005). Cartography: Distortions, World-views and Creative Solutions. *Progress in Human Geography* **29**, 217–224. doi:10.1191/0309132505ph540pr. URL <https://doi.org/10.1191/0309132505ph540pr>. <https://doi.org/10.1191/0309132505ph540pr>.
- 510 MURRELL, P. (2002). The grid graphics package. *R News* **2**, 14–19. <Https://journal.r-project.org/articles/RN-2002-010/>.
- 512 NEUWIRTH, E. (2014). *RColorBrewer: ColorBrewer Palettes*. URL <https://CRAN.R-project.org/package=RColorBrewer>. R package version 1.1-2.
- 514 OLSON, J.M. (1976). Noncontiguous Area Cartograms. *The Professional Geographer* **28**, 371–380. doi:10.1111/j.0033-0124.1976.00371.x. URL <https://doi.org/10.1111/j.0033-0124.1976.00371.x>. <https://doi.org/10.1111/j.0033-0124.1976.00371.x>.
- 517 R CORE TEAM (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <Https://www.R-project.org/>.
- 519 RAISZ, E. (1963). Rectangular Statistical Cartograms of the World. *Journal of Geography* **35**, 8–10. doi:10.1080/00221343608987880.
- 521 SKOWRONNEK, A. (2016). Beyond Choropleth Maps – A Review of Techniques to Visualize Quantitative Areal Geodata. URL Https://alsino.io/static/papers/BeyondChoropleths_AlsinoSkowronnek.pdf.
- 524 TUFTE, E.R. (1990). *Envisioning Information*. Graphics Press.
- 525 URBANEK, S. (2013). *png: Read and write PNG images*. URL <https://CRAN.R-project.org/package=png>. R package version 0.1-7.
- 527 WICKHAM, H. (2009). *ggplot2: elegant graphics for data analysis*. Springer New York. URL <Http://had.co.nz/ggplot2/book>.
- 529 WICKHAM, H., AVERICK, M., BRYAN, J., CHANG, W., McGOWAN, L.D., FRANÇOIS, R., GROLEMUND, G., HAYES, A., HENRY, L., HESTER, J., KUHN, M., PEDERSEN, T.L., MILLER, E., BACHE, S.M., MÜLLER, K., OOMS, J., ROBINSON, D., SEIDEL, D.P., SPINU, V., TAKAHASHI, K., VAUGHAN, D., WILKE, C., WOO, K. & YUTANI, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software* **4**, 1686. doi:10.21105/joss.01686.
- 534 WICKHAM, H., CHOWDHURY, N.R., COOK, D. & HOFMANN, H. (2018). *nullabor: Tools for Graphical Inference*. URL <https://CRAN.R-project.org/package=nullabor>. R package version 0.3.5.
- 537 WICKHAM, H., COOK, D., HOFMANN, H. & BUJA, A. (2010). Graphical inference for infovis. *IEEE Transactions on Visualization and Computer Graphics (Proc. InfoVis '10)* **16**, 973–979.
- 539 WILKE, C.O. (2019). *cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'*. URL <Https://CRAN.R-project.org/package=cowplot>. R package version 1.0.0.
- 541 XIE, Y. (2015). *Dynamic Documents with R and knitr*. Boca Raton, Florida: Chapman and Hall/CRC, 2nd edn. URL <Https://yihui.org/knitr/>. ISBN 978-1498716963.
- 543 XIE, Y., ALLAIRE, J. & GROLEMUND, G. (2018). *R Markdown: The Definitive Guide*. Boca Raton, Florida: Chapman and Hall/CRC. URL <Https://bookdown.org/yihui/rmarkdown>. ISBN 9781138359338.

- 546 ZHU, H. (2019). *kableExtra: Construct Complex Table with 'kable' and Pipe Syntax.* URL <https://CRAN.R-project.org/package=kableExtra>. R package version 1.1.0.
- 547

Table 1. The mean and standard deviation of the rate of detection for each trend model, calculated for the choropleth and hexagon tile map displays.

Type	NW-SE	Three Cities	All Cities
Choro.	0.52 (0.50)	0.04 (0.19)	0.23 (0.42)
Hex.	0.62 (0.49)	0.40 (0.49)	0.58 (0.49)

Table 2. The model output for the generalised linear mixed effect model for detection rate. This model considers the type of display, the trend model hidden in the data plot, and accounts for contributor performance.

Term	Est.	Sig.	Std. Error	P val
Intercept	-1.27	***	0.19	0.00
Hex.	1.63	***	0.24	0.00
Three Cities	-2.07	***	0.43	0.00
All Cities	1.34	***	0.24	0.00
Hex:Three Cities	1.28	**	0.48	0.01
Hex:All Cities	-1.16	***	0.33	0.00

Table 3. The number of participants that selected each reason for their choice of plot when looking at each trend model shown in choropleth and hexagon tile maps. The facets show whether or not the choice was correct.

Trend	Detected	Choro.	Hex.
All Cities	No	trend	clusters
	Yes	clusters, consistent	clusters
Three Cities	No	trend	clusters
	Yes	consistent	clusters
NW-SE	No	trend	clusters
	Yes	trend	clusters