

New Algorithms For Effectively Visualising Australian Spatio-Temporal Disease Data

A thesis submitted in fulfilment of the
requirement for the degree of
Master of Philosophy (Statistics)

by

Stephanie Rose Kobakian

B.Comm. and B.Eco., Monash University



School of Mathematical Sciences
Science and Engineering Faculty
Queensland University of Technology
Australia
2020

Contents

Statement of Original Authorship	v
Keywords	vii
Abstract	ix
List of Figures	xi
List of Tables	xiii
List of Publications	xv
Declaration	1
Acknowledgements	3
Cancer Applications of Choropleth Maps, and the Potential of Cartograms and Alternative Map Displays	4
Comparing the Effectiveness of the Choropleth Map with a Hexagon Tile Map for Communicating Cancer	4
R packages	4
1 Introduction	7
1.1 The Australian Cancer Atlas	8
1.2 Testing graphical displays	8
1.3 Aims and Objectives	8
1.4 Research Contributions	9
1.5 Thesis Structure	9
2 Cancer Applications of Choropleth Maps, and the Potential of Cartograms and Alternative Map Displays	11
Abstract	12
2.1 Introduction	12
2.2 Traditional approaches for cancer map displays	13
2.3 Contemporary alternatives to choropleth maps	20
2.4 Comparison and critique of alternative displays	28
2.5 User interaction	30
2.6 Conclusion	34

3 An Algorithm For Spatial Mapping Using a Hexagon Tile Map, With Application to Australian Maps	35
Abstract	35
3.1 Introduction	36
3.2 Algorithm	37
3.3 Using sugarbag	45
3.4 Applications	47
3.5 Animation	49
3.6 Conclusion	49
4 Comparing the Effectiveness of the Choropleth Map with a Hexagon Tile Map for Communicating Cancer	51
Abstract	52
4.1 Introduction	52
4.2 Background	53
4.3 Results	66
4.4 Discussion	72
4.5 Conclusion	74
4.6 Supplementary meaterial	74
5 Discussion and Conclusion	77
Bibliography	81
A Appendix	89
A.1 Overall Performance	89
A.2 Lineups	91
A.3 Experiment survey procedure	106
A.4 Ethics Approval	111

Statement of Original Authorship

The work contained in this thesis has not been previously submitted to meet requirements for an award at this or any other higher education institution. To the best of my knowledge and belief, the thesis contains no material previously published or written by another person except where due reference is made.

Signature:

A handwritten signature in cursive script, appearing to read "SR Kobakian".

Date:

© by Stephanie Rose Kobakian (2020), all rights reserved.

Keywords

Keywords: cartogram; choropleth; data science; geospatial statistics; information visualisaton; statistical graphics; tile maps; visual inference

Abstract

Relationships between geographic areas can be communicated using maps. The speed of visual processes allow quick comparisons to be made between neighbouring areas. However, when presenting population-related statistics on the geographic map base, large areas that draw readers' attention can allow too much emphasis to be placed on sparsely populated rural areas. This problem can dramatically impact the interpretation of statistics across Australia, due to the large difference in the land area when the population is concentrated in a few small areas.

In this thesis a new algorithm for laying out spatial areas into a tessellated hexagon tile is developed. This layout is used to create a new visualisation method for displaying disease data for Australia. An experiment was conducted that indicates that this visualisation outperforms the traditional choropleth map for communicating disease data. A software implementation was generated and made publicly available for others to use these techniques. In addition, animation can be used to translate from the choropleth map to the new tessellated hexagon display to help viewers interpret the new display.

List of Figures

Figure	Caption	Page
2.1	A selection of choropleth cancer maps from online atlases that are publicly available.	15
2.2	Common alternatives to maps, showing the same information for the United States of America.	22
2.3	Two alternative displays, tile map (left) and geofaceted map (right), showing state age-adjusted rate of incidence for lung and bronchus in the USA.	27
2.4	Cartograms showing melanoma incidence in Australia.	28
2.5	Interactive controls of displays in publicly available choropleth cancer maps	31
2.6	Two examples of advanced interactivity (and animation) in publicly available choropleth cancer maps.	32
3.1	The geographic shapes of the Statistical Areas of Tasmania at Level 2.	39
3.2	Grid points to create a tilegram.	40
3.3	All possible hexagon locations from the initial grid are shown with blue outlines.	41
3.4	Filter for grid points within a square, then circular, distance for those closest to the centroid.	43
3.5	Filter for grid points within the angle from the focal point to the centroid.	43
3.6	A complete hexagon tile map of Tasmania	44

Figure	Caption	Page
3.7	The Australian Cancer Atlas data has determined the colour of each Statistical Area of Australian at Level 2	45
3.8	A choropleth map of the Statistical Areas of Australia at Level 2.	48
3.9	A hexagon tile map of the Statistical Areas of Australia at Level 2.	49
4.1	A choropleth map of the smoothed average of liver cancer diagnoses for Australian males.	54
4.2	A hexagon tile map of the smoothed average of liver cancer diagnoses for Australian males.	55
4.3	This lineup of twelve hexagon tile map displays contains one map with a real population related structure.	58
4.4	The experimental design used in the visual inference study.	60
4.5	The detection rates achieved by participants are contrasted when viewing the four replicates of the three trend models	68
4.6	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.	70
4.7	The amount of times each level of certainty was chosen by participants when viewing hexagon tile map or choropleth displays.	72

List of Tables

Table	Caption	Page
2.1	A selection of choropleth cancer maps from online atlases.	16
2.2	Common measures for reporting cancer information.	17
2.3	Maps used to present statistics for the United States of America.	23
2.4	Summary of features and constraints of common mapping methods used to display cancer statistics.	29
4.1	The mean and standard deviation of the rate of detection for each trend model, calculated for the choropleth and hexagon tile map displays.	67
4.2	The model output for the generalised linear mixed effect model for detection rate.	69
4.3	The amount of participants that selected each reason for their choice of plot when looking at each trend model shown in Choropleth and Hexagon Tile maps.	72

List of Publications

1. The exploration of the literature regarding current practices for visualising spatial data is presented in Chapter 2 has been submitted to the journal *Annals of Cancer Epidemiology* for publication.

Kobakian S. and Cook, D. and Roberts, J. (2019). Cancer Applications of Choropleth Maps, and the Potential of Cartograms and Alternative Map Displays. Manuscript submitted for publication.

2. The algorithm is documented in Chapter 3 and will be submitted to the *Journal of Statistical Software*.

Kobakian S. and Cook, D. (2020). An Algorithm For Spatial Mapping Using a Hexagon Tile Map, With Application to Australian Maps. Manuscript in preparation.

3. The details of the experiment to test the alternative hexagon tile map is documented in Chapter 4 will be submitted to the *IEEE Transactions of Visualisation and Computer Graphics* under the title “Comparing the Effectiveness of the Choropleth Map with a Hexagon Tile Map for Communicating Cancer”

Kobakian S. and Cook, D. (2020). Comparing the Effectiveness of the Choropleth Map with a Hexagon Tile Map for Communicating Cancer Manuscript in preparation.

4. The code for the algorithm documented in Chapter 3 is currently hosted on CRAN as the package sugarbag.

Kobakian, S. and Cook, D. (2019). sugarbag: Create Tessellated Hexagon Maps. <https://srkobakian.github.io/sugarbag/>, <https://github.com/srkobakian/sugarbag>.

Declaration

I hereby declare that this thesis contains no material which has been accepted for the award of any other degree or diploma at any university or equivalent institution and that, to the best of my knowledge and belief, this thesis contains no material previously published or written by another person, except where due reference is made in the text of the thesis.

This thesis includes three original papers submitted to peer reviewed journals. The core theme of the thesis is spatial visualisations. The ideas, development and writing up of all the papers in the thesis were the principal responsibility of myself, the student, working within the Faculty of Science and Engineering under the supervision of Distinguished Professor Kerrie Mengersen and Dr. Earl Duncan. It was also created under the supervision of the external supervisor Professor Dianne Cook from Monash University.

The papers in Chapters 2, 3, and 4 have been individually formatted for journal submission however, I have renumbered the pages of the submitted papers for cohesion across this thesis.

Acknowledgements

This thesis was made possible due to the opportunity provided by Queensland University of Technology, and the Australian Research Council Centre of Excellence for Mathematical & Statistical Frontiers. This research was also supported by an Australian Government Research Training Program (RTP) Scholarship, and the Cancer Council Queensland.

I would like to acknowledge my supervisors for their tireless work in directing, encouraging and supporting this work.

Professor Dianne Cook (Professor, Department of Econometrics and Business Statistics, Monash University) provided constant support, encouragement and recommendations throughout this degree.

Dr Earl Duncan (Research Associate at ARC Centre of Excellence for Mathematical & Statistical Frontiers, QUT) provided regular time for discussion, editing and commenting throughout this degree.

Professor Kerrie Mengersen (Professor of Statistics, Science and Engineering Faculty, QUT) provided the opportunity to study at Queensland University of Technology, organising and mentoring me through this project.

Working closely with Prof. Dianne Cook enabled the rapid development of my research and writing skills.

Special thanks to Professor Peter Baade and Dr Susanna Cramb.

Cancer Applications of Choropleth Maps, and the Potential of Cartograms and Alternative Map Displays

With my co-authors Prof. Dianne Cook and Jessie Roberts, I would like to acknowledge Dr Susanna Cramb (Strategic Research Fellow, Institute of Health Biomedical Innovation, Queensland University of Technology) and Dr Peter Baade (Senior Research Fellow, Cancer Council Queensland) for providing the opportunity and time to discuss alternative map solutions for presentation in the Australian Cancer Atlas. It was in the development of this online cancer atlas that methods for disease map displays, and visual communication strategies were explored.

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

I would like to thank Mitchell O'Hara-Wild was a co-developer of the *taipan* (Kobakian and O'Hara-Wild, 2018) R package for the web app constructed to collect participant evaluations of lineups.

We are thankful for the NUMBATs (Non-Uniform Monash Business Analytics Team) for participating in the pilot study that helped to assess the experimental design and determine an appropriate sample size for the study.

R packages

Numerous R (R Core Team, 2019) packages were used to produce this thesis:

- **absmapsdata** (Mackey, W. F., 2019)
 - **broom** (Robinson and Hayes, 2019)
 - **cartogram** (Jeworutzki, 2018)
 - **cowplot** (Wilke, 2019)
 - **dplyr** (Wickham et al., 2019)
 - **eechidna** (Cook et al., 2019)
 - **ggplot2** (Wickham, 2016)
 - **ggthemes** (Arnold, 2019)
-

- `gstat` (Pebesma and Graeler, 2019)
 - `kableExtra` (Zhu, 2019)
 - `knitr` (Xie, 2014)
 - `lme4` (Bates et al., 2019)
 - `lubridate` (Grolemund and Wickham, 2011)
 - `nullabor` (Wickham et al., 2018)
 - `magrittr` (Bache and Wickham, 2014)
 - `png` (Urbanek, 2013)
 - `RColorBrewer` (Neuwirth, 2014)
 - `rmarkdown` (Allaire et al., 2019a)
 - `rticles` (Allaire et al., 2019b)
 - `sf` (Pebesma, 2018)
 - `spData` (Bivand, Nowosad, and Lovelace, 2019)
 - `sugarbag` (Kobakian and Cook, 2019)
 - `tibble` (Müller and Wickham, 2019)
 - `tidyR` (Müller and Wickham, 2019)
-

Chapter 1

Introduction

Maps can contribute to the interpretation of spatial distributions of disease occurrence, and help to locate disease clusters. Disease data is commonly aggregated to political areas; privacy is one reason for this and another is that it is the responsibility of the political entity to respond. The typical visualisation for aggregated spatial data is a choropleth map, where areas are coloured according to a numerical value. A choropleth map is the most common display for the presentation of disease data.

Choropleth maps can mislead the map reader, as the attention of the map user is distributed according to the size of the area. Australia provides an extreme example of potential bias and loss of valuable information. In Australia, the geography mismatchs the population, because the communities are densely populated in the inner city areas, especially around the capital cities, and the coastline. There are alternative visualisation methods, like cartograms, that have been developed to correctly focus on the population dense areas. These alternatives should be considered when planning the communication of geospatial statistics, as visualisations should be chosen to best represent the spatial distribution.

This thesis research is motivated by the Australian Cancer Atlas, which presents the spatial patterns of cancer in Australia. The aim of this thesis is to contribute an algorithm that creates effective visualisations for the communication of geospatial population statistics.

1.1 The Australian Cancer Atlas

The Australian Cancer Atlas (ACA) is an online, interactive web tool for exploring the impact of cancer on Australian communities. The prominent display used by the ACA to present incidence rates or excess death rates is the choropleth map. The set of geographic units shown are the Australian Statistical Areas, at Level 2 (SA2s). There are almost 2,200 individual SA2s.

The choropleth map used in the ACA is familiar to the general public of Australia. It is appropriate to use this map display as users can orient themselves on the map and find the geographic areas relevant to them. However, when the intention of the map user is to understand the whole spatial distribution, the information derived from the colours displayed on the map can be misleading. The rural areas in outback Australia are over emphasised, and the densely populated inner city areas are not given enough attention, as they cannot be seen using a choropleth map.

1.2 Testing graphical displays

Visual inference testing will be used to determine if the communication of population geospatial statistics is more effective when using the new alternative hexagon tile map display. Buja et al (Wickham et al., 2010) provide the ‘lineup’ protocol as a formal framework for testing visual statistical methods. The new alternative visualisation method can be tested by implementing this framework. It takes inspiration from a police lineup. The lineup protocol can be used to test if the hexagon tile map is effective, a map displaying a real population based distribution can be hidden in a collection of maps that display null distributions (Roy Chowdhury, 2014).

1.3 Aims and Objectives

This thesis aims to provide a solution to presenting geospatial data regarding populations. It considers the visualisation methods developed over the past two centuries that shift the focus from the geographic map base.

1. *Devising an algorithm for creating hexagon tile maps of Australia:* The algorithm will take geospatial areas and create an alternative visualisation of the spatial distribution.
2. *Test the effectiveness of the hexagon tile map relative to the choropleth map:* The hexagon tile map produced by the algorithm will be contrasted with the traditional choropleth map, applying the same colour methods to represent the data.
3. *Communicating the relationship between the hexmap and choropleth map through animation:* Animations can maximise the benefits of both visualisation methods when communicating to the public. The use of animations may control how people follow a recognisable map of Australia into an alternative visualisation for inference.

1.4 Research Contributions

This research contributes a new algorithm for creating hexagon tile map displays. It contributes an R (R Core Team, 2019) package which implements the algorithm and allows R users to create their own visualisations. It presents a case study that contributes to a growing field of visual inference studies. Applying the lineup protocol to spatial data by comparing a choropleth map to a hexagon tile map display. It also shows how it can be used in practice to effectively communicate population related cancer distributions.

1.5 Thesis Structure

The thesis is structured as follows: Chapter two contains a literature review. The literature reviews considers the current peer reviewed literature and published books that explore spatial distributions of cancer across the globe. It also considers how to evaluate the visualisation methods used for spatial data.

Chapter 3 explores the algorithm to create hexagon tile maps and the code used to create a small example of Tasmania in Australia. Chapter 4 discusses an experiment that uses visual inference. It contains the methods and results that compare the use of a choropleth map and a hexagon tile map. Chapter 5 summarises the contributions of this thesis and possible future work.

Chapter 2

Cancer Applications of Choropleth Maps, and the Potential of Cartograms and Alternative Map Displays

This literature review chapter provides an overview of the key areas of interest in the literature that are relevant to this thesis. This chapter is organised as follows. Section 1 outlines the traditional spatial mapping technique, the choropleth map, and provides design inspiration using examples of online cancer atlases. Section 2 outlines contemporary mapping approaches, suggested as alternatives to the choropleth. Section 3 compares and critiques these alternative displays in light of the strengths and weaknesses of the choropleth method. section 4 considers how users interact with mapping displays in online cancer atlases, and how map creators can direct the attention of uses through animation. Additionally, the following chapters in this thesis each begin with a section to introduce the specific relevant literature.

This chapter was submitted for publication to the journal *Annals of Cancer Epidemiology* for publication. This was intended for an audience of cancer atlas creators to be inspired by current atlases, and to encourage the pursuit of alternative displays.

Abstract

Cancer atlases communicate cancer statistics over geographic domains, typically with a choropleth map. They subdivide these domains into administrative regions such as countries, states, or suburbs. When communicating human-related statistics, the choropleth has a disadvantage in that it draws attention to sparsely populated rural areas to the neglect of small inner city areas. The smaller geographic areas are important to consider if they are densely populated. Alternative map displays, such as a cartogram or a hexagon tile map, can shift the attention of map users from the large rural areas by decreasing their size on the map display. This means alternative displays can be more effective at accurately communicating spatial patterns across spatial areas. It is recommended that alternative displays are included in cancer atlases. In addition, with the ease of today's technology, user interaction with the displays is encouraged. Users should also be able to interactively display different statistics, such as incidence rate or relative incidence, or filtered by demographic variables.

2.1 Introduction

Researchers, health authorities, governments, not-for-profits and the media are common communicators of cancer statistics. They often present statistics to the public as aggregated values for geopolitical areas. Presenting these statistics requires aggregating individual observations for the geographical units, especially for privacy protection, but also for political and policy purposes. Examples of typical geographical units include states, provinces, local government areas, and post/zip codes. It is easy to provide counts or incidence rates of the diagnoses of these areas. This type of data is routinely collected for public health reasons and may be made available to the general public as a service to the community.

To visualize and communicate geospatial cancer statistics over geographic domain, a choropleth map is the common display. Choropleth maps show polygons representing the geographic units, where each polygon is shaded with a color according to the area-specific values of the statistic being conveyed. Visualizing this data is helpful as geographic patterns of disease may be obscured when reported in a table (Moore and Carpenter, 1999).

Providing a visual representation of cancer outcomes allows identification of geographic patterns of the disease that can then be addressed with public health policy and actions. The spatial distribution of the disease incidence can be examined using a choropleth and may reveal a trend in longitude or latitude, or rural vs urban, or coastal vs inland, or even specific hot spots of the disease. One of the key challenges with mapping spatial patterns of disease is the design of visualizations (Exeter, 2016). It is important to consider the strengths and weaknesses of designs, as visualizing diseases on maps is often the first step in exploratory spatial data analysis and helps in the formulation of hypotheses. This paper considers the current visualization techniques to communicate statistics to the public and their applications to cancer statistics. Alternative approaches are posed because they may be more effective than contemporary techniques. The limitations of the visualization methods, highlighting the differences and historic use of these displays is discussed.

The paper is structured as follows. The next section describes the choropleth map, which is the common approach to disease maps and presents examples of atlases in use today and discusses the limitations of the choropleth map. Section 3 describes alternative displays, including the cartogram, which is useful when the map has heterogeneously sized geographic units. Section 4 presents the limitations of the production and use of alternative displays. Disease maps are more useful when made interactive, and common options are described in Section 5, along with a discussion of benefits and disadvantages.

2.2 Traditional approaches for cancer map displays

A choropleth map displays the geographical distribution of data over a set of spatial units by shading areas of a map (Tufte, 1990), (Skowronnek, 2016). Faithful rendering of the geography, when combined with an appropriate color scheme, can reveal spatial patterns among data values. Identifying and explaining spatial structures, patterns, and processes involve considering the individuals and organizing them into representable units of communities (Moore and Carpenter, 1999). Early versions of choropleth maps used symbols or patterns instead of color. Choropleth maps can be used for displaying disease data (Walter, 2001), including cancer data (Bell et al., 2006). In epidemiology, choropleth

maps are often used as a tool to study the spatial distribution of cancer incidence and mortality.

Displaying familiar state boundaries can make a map easier to read (Brewster and Subramanian, 2010) and allow viewers to infer the spatial relationships visually in the data using their mental model of the geography. The map users of disease displays may include researchers, the public, policymakers, and the media (Bell et al., 2006). For these users, the familiarity of the geography is a worthy consideration when presenting results of spatial analysis.

2.2.1 Cancer atlases

A cancer atlas is a map, or collection of maps, representing cancer incidence and mortality for a country, or group of countries. Atlases are key to developing hypotheses regarding areas with unusually high rates, and geographic correlations (d'Onofrio et al., 2016). The data collection methods across regions and the administrative control within regions lends itself to choropleth visualization. Cancer maps and atlases date back to Haviland's maps in 1875, and early work in US cancer atlases appearing in 1971 (Burbank, 1971). The presentation of cancer statistics has increased with greater access to computational power and the availability of geographic information systems software (Exeter, 2016).

Cancer maps are effective tools for communicating incidence, survival, and mortality to a wide range of audiences, including the public and others not trained in statistical analysis. These visualizations enable non-expert audiences to interpret the outputs of sophisticated statistical analysis. Cruickshank (1947) as cited by S. D. Walter (Walter, 2001), discusses using visuals as a 'formal statistical assessment of the spatial pattern'. Overwhelmingly, choropleth maps are visualisations chosen to communicate cancer statistics to members of the public and other non-expert audiences.

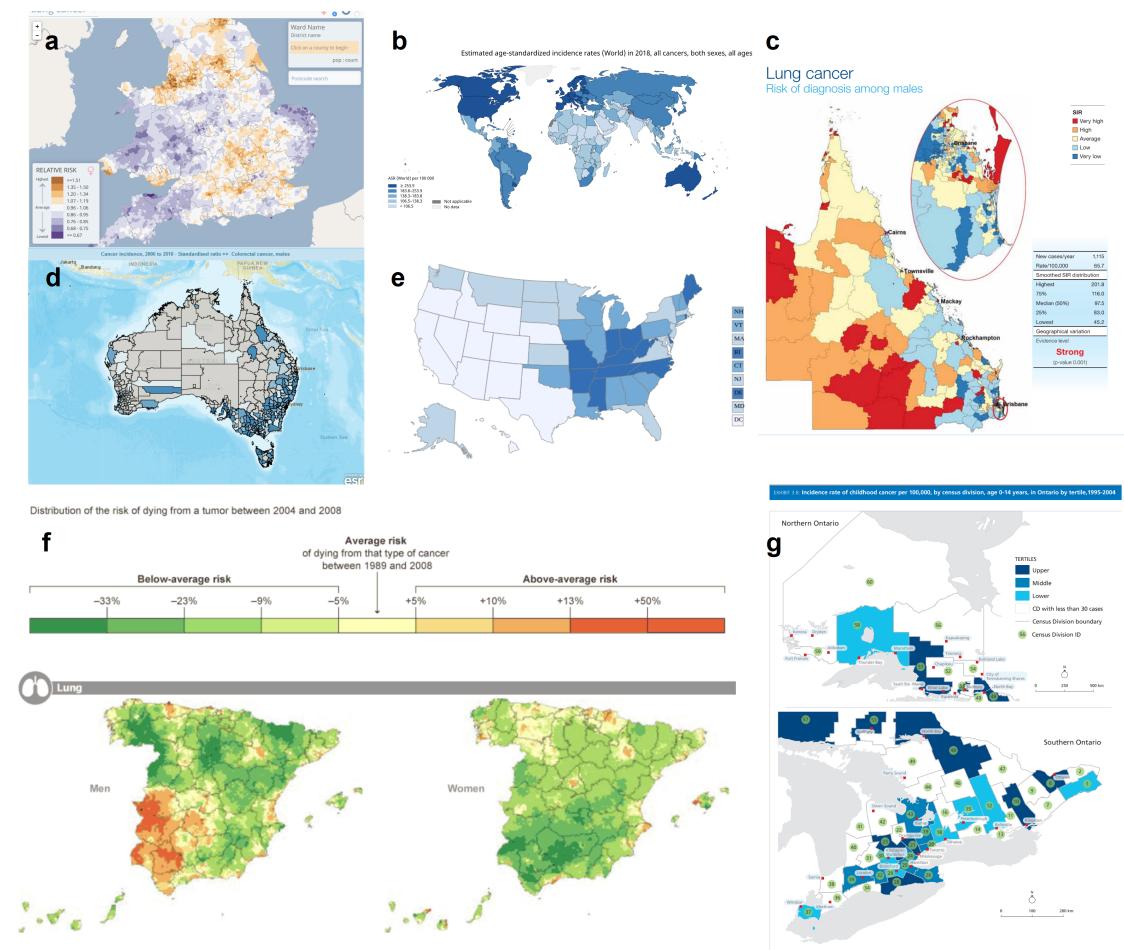


Figure 2.1: A selection of choropleth cancer maps from online atlases that are publicly available. Maps of various countries were chosen: United Kingdom, Australia, Spain, USA, Canada, and display several different colour palettes and legends. These atlases are described in Table 2.1.

Epidemiologists and statisticians have developed the statistics used to communicate the burden of cancer over several decades. Table 2.2 summarizes the measures commonly presented in published cancer atlases. Mortality rates are commonly presented as relative rates of risk across the population and age-adjusted to correct for the higher prevalence of cancers in older populations. As described in Howe (Howe, 1989), Englishman P. Stocks advanced the field of mortality statistics by introducing the standardized mortality ratios in the 1930s, which is an improvement on crude death rates.

Roberts (Roberts, 2019) identified 33 cancer atlases published between 2010 and 2018. Each of these online atlases uses choropleth maps. All except one of these were published by

Table 2.1: A selection of choropleth cancer maps from online atlases.

Fig.	Atlas	Statistic	Data.source
1a	The Environment and Health Atlas of England and Wales	relative risk for women developing lung cancer in England and Wales in 2010	Office for National Statistics (ONS) (England) and from the Welsh Cancer Intelligence and Surveillance Unit (WCISU).
1b	Globocan 2018: Estimated Cancer Incidence, Mortality and Prevalence Worldwide	age standardized incidence rates (per 100,000) for all invasive cancers for both men and women, aggregated at a national level for 2018	World Health Organization's International Agency for Research on Cancer
1c	Atlas of Cancer in Queensland	the relative incidence ratio of lung cancer in males in the state of QLD within Australia based on data from 1998 to 2007	Queensland Cancer Council, Queensland Cancer Registry.
1d	Bowel Cancer Australia Atlas	the percentage of Australian males between 50 - 54 years of age diagnosed with bowel cancer in 2016.	Bowel Cancer Australia.
1e	United States Cancer Statistics: An Interactive Cancer Statistics Website	the incidence rate per 100,000, of all cancer types for men and women in the United States in 2016, aggregated at the state level.	Centers for Disease Control and Prevention, with data from state cancer registries.
1f	Map of Cancer Mortality Rates in Spain	side by side maps of relative risk of lung cancer for men vs women for 2004 to 2008.	Map of cancer mortality rates in Spain.
1g	Atlas of Childhood Cancer in Ontario	the incidence rate of childhood cancers per 100,000 (by census division) for children aged 0-14, in Ontario from 1995 to 2004.	The Paediatric Oncology Group of Ontario Networked Information System.

Table 2.2: Common measures for reporting cancer information.

Measure	Details
1. Count	Crude cancer counts
2. Rate per 100,000	Cancer incidence per 100,000 population
3. IR (Incidence Ratio) NA	$(IR)_i = \frac{(Incidence\ Rate)_i}{Average\ Incidence\ Rate}$, The cancer incidence rate in region i over the average cancer incidence rate for all of the regions
4. Age-Adjusted Rate per 100,000	Standardized by age structure or region
5. Age-Adjusted Relative Risk	Standardized by age structure in each region i
6. SIR (Standardized Incidence Ratio)	Incidence standardized by population at risk in each region i
7. Below or above Expected	An alternative expression of the SIR
8. RER (Relative Excess Risk)	$RER = \frac{(Cancer\ related\ mortality)_i}{Average\ cancer\ related\ mortality}$ Represents the estimate of cancer-related mortality within five years of diagnosis. Also referred to as 'excess hazard ratio'

non-commercial organizations, including not-for-profits, government, research organizations, advocacy groups or government-funded partnerships. Figure 2.1 displays a subset of maps from these atlases, the selection varies in the geographies explored. Figure 2.1b shows Globocan 2018 (World Health Organization's International Agency for Research on Cancer, 2018) which explores Estimated Cancer Incidence, Mortality and Prevalence Worldwide using data sourced from cancer registries of each country. The Bowel Cancer Australia Atlas in Figure 2.1d presents an example of a cancer specific atlas – it shows the average Standardized Incidence Ratio of colorectal cancer for Australian males from 2006 to 2010 (Bowel Cancer Australia, 2016). Like many of the atlases examined, there is a choice of gender displayed in the Bowel Cancer Atlas. Gender is displayed in side-by-side maps in the Map of Cancer Mortality Rates in Spain (Figure 2.1f) (El Pais, 2014).

Resolution of the maps varies greatly. Figure 2.1b shows global information at a national level. The United States Cancer Statistics (U.S. Department of Health and Human Services, Centers for Disease Control and Prevention and National Cancer Institute - Cancer Statistics Working Group, 2019) shows data aggregated at the state level. The Environment

and Health Atlas of England and Wales (Emperial College London - Small Area Health Statistics Unit, 2010) (Figure 2.1a) shows the relative risk for women developing lung cancer at a neighbourhood (small-area) scale. The Atlas of Cancer in Queensland (Figure 2.1c) shows the relative incidence ratio of lung cancer in males for each Statistical Area at Level 2 (Statistics, 2018) in the state of Queensland within Australia (Queensland Cancer Registry, 2011).

Age-specific atlases are less common. Figure 2.1g displays Atlas of Childhood Cancer in Ontario, this communicates the incidence rate of childhood cancers per 100,000 (by census division) for children aged 0-14, in Ontario from 1995 to 2004 (Pediatric Oncology Group of Ontario, 2015).

2.2.2 Additional considerations

Cancer atlases often display supplementary graphs and plots to add more information. Additional materials such as tables, graphs, and text explanations support understanding and inference derived from maps, ensuring the message communicated will be consistent across a range of viewers (Bell et al., 2006). The many displays of statistical summaries, including dot plots, bar plots, box plots, cumulative distribution plots, scatter plots, and normal probability plots, can provide alternative views of the cancer statistics. These can also display supporting statistics such as error, confidence intervals, distributions, sample or population sizes, and standard deviation.

The statistics communicated in atlases are often used to describe differences between areas. This can occur at different levels of aggregation. Aggregation of global health statistics occurs within administrative and arbitrarily defined regions, such as those used by the World Health Organization and the United Nations (Ferlay J, Ervik M, Lam F, Colombet M, Mery L, Piñeros M, Znaor A, Soerjomataram I, Bray F, 2018). World atlases can allow for displays of data aggregated into continents, countries, states, provinces and congressional districts (U.S. Department of Health and Human Services, Centers for Disease Control and Prevention and National Cancer Institute - Cancer Statistics Working Group, 2019). Each population area will probably have a different number of people, which is typically used to calibrate the statistic. Cancer atlases may also communicate the

distribution of the population living in all areas in a table or histogram display (Northern Ireland Cancer Registry, 2011). Atlases can connect the population to the land available to them by communicating population density.

Maps can also be used to focus on demographic strata, such as age and sex. Some of the digital atlases surveyed allow subsets such as males, females, or those aged over 65, to be selected for display. Similarly, socioeconomic indicators, such as unemployment rates, poverty rates, remoteness, and education levels, can be used to filter data, in order to communicate how cancer prevalence varies for different members of society. Few atlases provide this level of detail.

Introducing population and demographic information helps to interpret the rates in areas effectively, but there will still be uncertainty around the rates. To address this, a cancer atlas often communicates uncertainty about the value of a statistic. There are several potential sources of uncertainty: sampling error, errors arising from the disease reporting process (or data collection), and errors arising from the statistical modelling or simulation process. The most common measures used to present uncertainty are credible or confidence intervals (CIs). Displaying the uncertainty associated with reported statistics is a vital feature of a cancer map, but it is difficult to display effectively. The map focuses on displaying the statistic and lacks additional space to represent the uncertainty. Providing an adjacent map or overlaying maps with symbols (Kronenfeld and Wong, 2017) are two common solutions.

2.2.3 Limitations of choropleth displays

Australia presents an extreme case of an urban rural divide. The land mass occupied by urban electoral districts is only 10% of Australia, yet 90% of the population live in these urban areas (Dorling, 2011).

Choropleth maps provide a familiar display, which shows data in a geographically recognisable way. A disadvantage is that the different population and geographical sizes of administrative areas can attract attention to the shades of the underpopulated but large areas (Tufte, 1990). Skowronnek also (Skowronnek, 2016) discusses how choropleth maps

suffer from area-size bias, as they give a ‘stronger visual weight’ to large administrative units. The administrative boundaries used to define regions may limit a choropleth display, as this display unfaithfully represents the disease distribution across the region by obscuring small geographic areas. Sparsely populated rural areas are emphasized, whereas the areas representing inner city communities are very small. This is especially true for Australia.

Choropleth maps colour each geographic unit to allow map users to measure the value of the statistic (Tufte, 1990). Map users contrast the colours in neighbouring areas to understand the spatial distribution. The ColorBrewer system (Harrower and Brewer, 2003) and viridis (van der Walt, S. and Smith, N, 2015) palettes provide effective colour schemes for qualitative, sequential and diverging data. When communicating information using colour, a map creator should use a scheme that has a linear color gradient, with perceptually uniform color spaces that match equal steps in data space with equal steps in the colour space (Madsen, 2019). The use of borders and backgrounds, and their colours, can also change the appearance of the colors representing the value of the statistics (Harrower and Brewer, 2003). These supports can be used to implement a reference point in the colour scheme as well as orient users to the geographic regions.

Inset maps like in Brisbane city in Figure 2.1c of the state of Queensland are commonly used to reduce distorted interpretations, but it is a bandaid remedy. For Australia, many, many inset maps would be needed.

2.3 Contemporary alternatives to choropleth maps

2.3.1 Cartograms

Choropleth maps imply uniformity of data across the geographic space but population densities are unlikely to be uniform (Skowronnek, 2016). Cartographers developed the cartogram to draw the attention to the population by transforming the map (Dougenik, Chrisman, and Niemeyer, 1985). The resulting display can communicate the impact of the disease more accurately across the population, as recorded by the statistic, at the sacrifice of geographic accuracy.

When a map creator desires a uniform population density of the map base, the purposeful distortion of the map space is beneficial. The “population distribution is often extremely uneven”, making a distortion necessary so that population is more faithfully represented as a uniformly distributed background for the statistic to be presented (Dorling, 2011) (Griffin, 1980) (Berry, Morrill, and Tobler, 1964). An area cartogram (Olson, 1976), or population-by-area cartogram (Levison and Haddon Jr, 1965) is produced from the distortion of the geographical shape according to population. Event cartograms (Kronenfeld and Wong, 2017) change the area of regions on a map depending on the amount of disease-related events, rather than population.

Cartograms provide an alternative visualization method for statistical and geographical information. Monmonier (Monmonier, 2018) suggests that map creators can use white lies to create useful spatial displays. It is easy for the reader to disregard the impact of transformations used to create cartograms, for the benefit of reading the statistical distribution more accurately with approximate geographic information. The spatial transformation of map regions relative to the data emphasizes the data distribution instead of land size (Kocmoud and House, 1998). When visualizing population statistics, Dorling considers this design ‘more socially just’ (Dorling, 2011), or honest (Dent, 1972), giving equitable representation and attention to all members of the population and reducing the visual impact of large areas with small populations (Walter, 2001). Howe (Howe, 1989) suggests that ‘cancer occurs in people, not in geographical areas’ and that spatial socio-economic data, like cancer rates, are best presented on a cartogram for urban areas as the population map base avoids allocating ‘undue prominence’ to rural areas (Griffin, 1980).

The creation of cartograms was historically in the hands of professional cartographers (Kraak, 2017). Early approaches by John Hunter and Jonathan Young (1968) and Durham’s wooden tile method, Skoda and Robertson’s (1972) steel ball-bearing approach and Tobler’s (1973) computer programs (Dorling, 2011). Howe (Howe, 1989) discusses the impact of electronic computer-assisted techniques. Geographical information systems allow map creators to produce cartograms and they use these systems depending on ‘the effectiveness, efficiency, and satisfaction of the map products’ (Kraak, 2017).

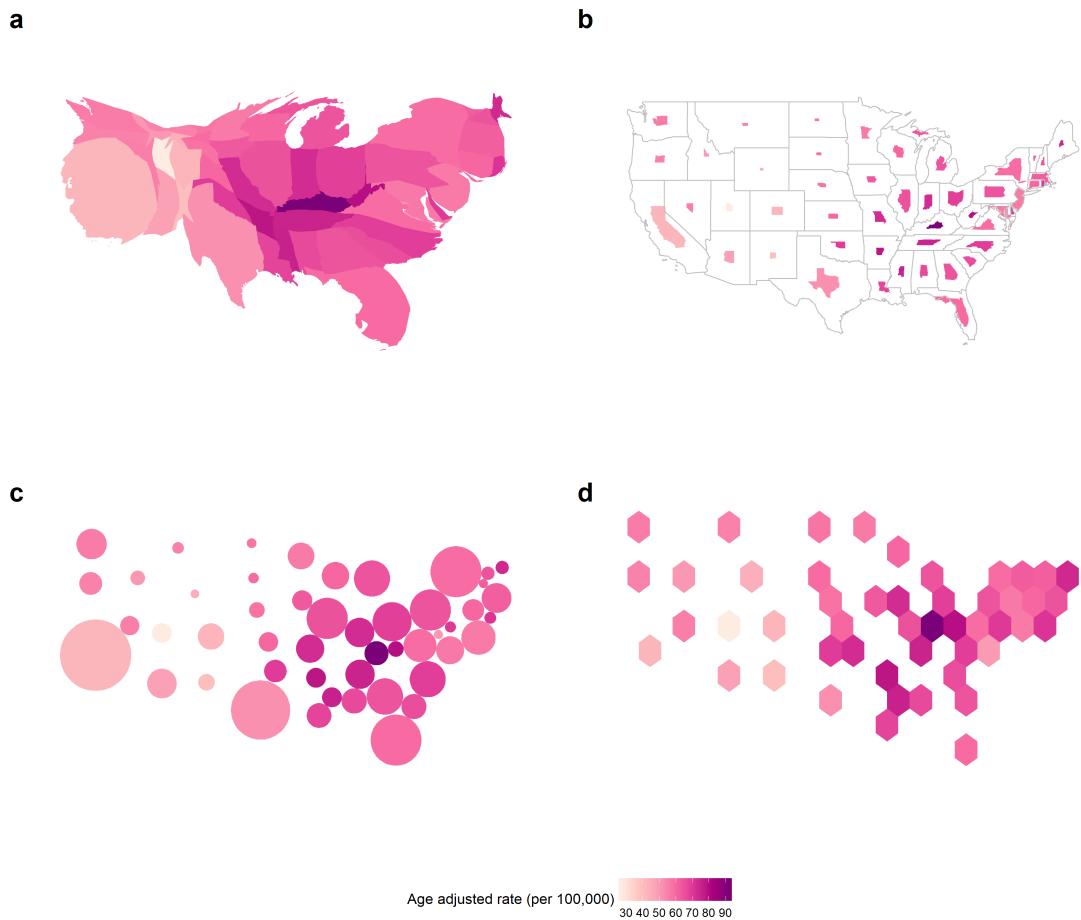


Figure 2.2: Common alternatives to maps, showing the same information for the United States of America: (a) contiguous cartogram, (b) non-contiguous, shape-preserved cartogram, (c) Dorling cartogram (non-contiguous), (d) hexagon tile map (non-contiguous). Maps (a) - (c) are created by resizing and reshaping the states of the USA to match the 2015 population of the state. This provides a better sense of the extent of disease relative to the population in the country and can help ease losing information about physically small but population-dense states. Map creators give each state equal size and thus equal emphasis in (d) the hexagon tile map.

There are two key issues to consider when creating alternative map displays, (1) the intended audience of the map, and (2) its purpose. Nusrat and Kobourov (Nusrat and Kobourov, 2016) provided a framework to investigate implementations of the many algorithms presented, and the “statistical accuracy, geographical accuracy, and topological accuracy”.

Figure 2.2 shows four different cartograms for the same data. The information in Table 2.3 summarizes what can be observed in the four types of cartograms.

Table 2.3: Maps used to present statistics for the United States of America. The colour of each state communicates the average age-adjusted rate of incidence for lung and bronchus for females and males in the United States 2012-2016.

Figure	Map display	Details
2.2a	Contiguous	It has distorted each state's shape according to the population of the state in 2015. The state of California has become much larger because of the large population density. This draws attention to the densely populated North-East region and detracts from the less populated Mid West.
2.2b	Non-contiguous	It maintains the geographic shape of the states, but the size has altered according to the population of the state in 2015. The state of California has remained closer to its original size than its surrounding states. The North-East states have remained closer to their geographical size, for Massachusetts and Connecticut. This draws attention to the densely populated North-East region and the sparse Mid West.
2.2c	Dorling	Circles are used to represent each state, but the population of the state determines the size in 2015. The North-East states remain closer to their neighbors and are slightly displaced from their geographic location. It highlights the sparsity of the population in the Mid West by the distance between the circles at the geographic centroids.
2.2d	Hexagon tile map	A hexagon of equal size represents each state. It is easy to contrast the neighboring states however the North-East regions have been displaced from their geographic location. It highlights the sparsity of the population in the Mid West by the light yellow color, the Age-Adjusted rate in Kentucky is the darkest and its neighbors are similar.

Contiguous

A contiguous cartogram alters the choropleth according to a statistic and maintains connectivity of the map regions. Min Ouyang and Revesz (Min Ouyang and Revesz, 2000) present three algorithms for creating value-by-area cartograms. They implement 'map deformation' to account for the value assigned to each area. Other methods include Tobler's Pseudo-Cartogram Method, Dorling's Cellular Automaton Method (Dorling, 2011), Radial Expansion Method, Rubber Sheet Method, Line Integral Method, Constraint-Based Method (Kocmoud and House, 1998).

Figure 2.2a shows a population contiguous cartogram of the United States. All states are visible and the shape of the United States overall is still recognizable. In contrast, Figure 2.4a shows an Australian contiguous cartogram also based on population. The south east is enlarged, but high population areas are still small, and low population areas are still large on the map. The algorithm doesn't fully reach an optimal configuration where area matches population – Australia is too heterogeneous for the algorithm to handle.

To be able to recognize the significant changes, a reader will usually have to know the initial geography to find the differences in the new cartogram layout (Olson, 1976). The shapes of small areas on a choropleth map and a cartogram are preserved using Tobler's Conformal mapping method. Kocmoud and House (Kocmoud and House, 1998) present this issue as conflicting tasks or aims, to adjust region sizes and retain region shapes.

Non-contiguous

Non-contiguous cartograms prioritize the shapes of the areas instead of connectivity. Each area stays in a similar position to its location on a choropleth map. Displaying the choropleth map base allows map users to make comparisons regarding the change in the area. The addition is the gap between areas, created as each area shrinks or grows according to the associated value of the statistic. Olson (Olson, 1976) discusses the creation of these maps and the significance of the empty areas left between the geographic boundaries and the new shape.

The white space presents the meaningful empty-space property (Keim et al., 2002), (Olson, 1976) but it also distracts the reader from the data, with a low data density (Tufte, 2001).

Dorling

Daniel Dorling presents an alternative display engineered to highlight the spatial distribution and neighbourhood relationships without complex distortions of borders and boundaries (Dorling, 2011):

“If, for instance, it is desirable that areas on a map have boundaries which are as simple as possible, why not draw the areas as simple shapes in the first place?”

He acknowledged the sophistication of contiguous cartograms but critiqued their ‘very complex shapes,’ he answers this with his implementation of maps created using ‘the simplest of all shapes’. Circular cartograms use the same circle shape for every region represented, resized according to the statistic represented or the population. This simple shape may be more effective for understanding the spatial distribution than contiguous cartograms. Contiguous cartograms create ‘nonsense’ shapes that have ‘no meaning’ (Dent, 1972). Both methods applies a gravity model to produce a layout, that avoids overlaps and keep spatial relationships with neighboring areas over many iterations. The circular cartogram is relatively fast to compute.

Raisz (Raisz, 1963) laid the groundwork for this approach in the mid-1930s, drawing rectangular cartograms that provide simple comparisons, effective for correcting misconceptions communicated by geographic maps. Tobler (Tobler, 2004) names and defines these as Value-Area Cartograms. This rectangular display may sacrifice contiguity but allows for tiling where geographic neighbors placed in suitable relative positions also share borders (Monmonier, 2005). Rectangular cartograms communicate bivariate displays of the population by the size of each rectangular, and they use color to communicate a second variable (Kreveld and Speckmann, 2007).

2.3.2 Tile Map

A tile map provides a tessellated display of consistent shapes. A similar method to a rectangular cartogram, represents each geographic area using a square. The squares are tessellated to create a grid. Each area is represented by a square of the same dimensions, each tile is usually one unit of measurement, this could be geographic regions such as states or population-based that use a consistent measure of population for each tile. Regions with over four neighbors require some necessary displacement. The tile map uses color to represent a value of a statistic for each area. A similar method to a rectangular cartogram represents each geographic area using a square of the same dimensions. There are online media sources using this method, these include (Montanaro, 2016), (Kanjana and Mehta, 2016), (Zitner, Yeip, and Wolfe, 2016), (Gamio and D., 2016). Tile maps may be

difficult to create as they are best created manually, they require additional time and care as the number of geographic areas to include increases.

Cano and others (Cano et al., 2015) define the term ‘mosaic cartograms’ for hexagonal tile displays, where the number of tiles for each area or the color of them can communicate the statistic of regions. When using several tiles per region, map makers can adjust the complexity of the boundaries in the resulting display. They can also make a trade-off between boundary complexity and simplicity by the size of the tiles used. A mosaic cartogram employs tessellation to connect the hexagons, triangles or squares used to represent the geographic land mass. Tessellation closely arranges each of the shapes so that the sides of neighbouring shapes align. Tile maps do not have to tessellate completely, this flexibility is helpful if the land mass has islands.

2.3.3 Geofacet

Hafen (Hafen, 2019) introduces the term geofacet to describe a grid display of small plots. The arrangement of tiles mimics the geographic topology. Geofaceting has the functionality that a statistical plot can be constructed in each facet for each geographic area. A tile map can communicate only one value per region in a visualization, while geofaceting is a more flexible visualization for communication as it increases the amount of information displayed. Virtually any type of plot can be shown in the tile, allowing displays of multiple variables or values per geographic entity. Creating the layout of a geofacet is manual, but once created can be used for any data on that geographic base.

2.3.4 Multivariate displays

Pickle and others (W., Carr, and Pearson, 2015) present linked micromap plots to match geographic and statistical data visually, this serves as a solution to multi-dimensionality issues. These maps group areas based on their value for one variable, and additional columns provide displays that contrast the areas in each group by other variables. The display juxtaposes choropleth maps and statistical plots; it shows one map per group of the key separating variable, in a row with each additional statistical plot. Linked micromaps predominantly use the choropleth map for displays of spatial relationships. These maps

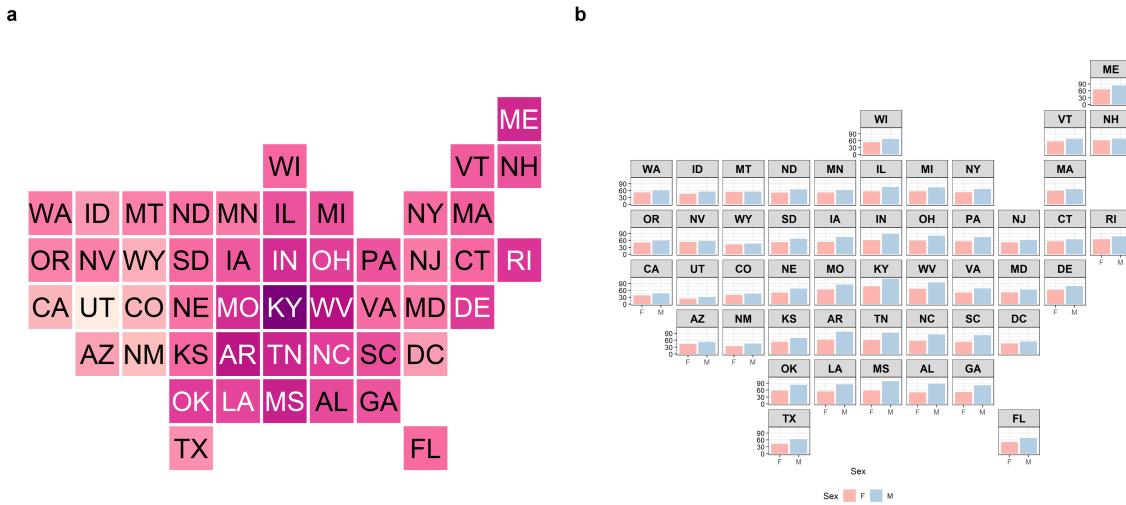


Figure 2.3: Two alternative displays, tile map (left) and geofaceted map (right), showing state age-adjusted rate of incidence for lung and bronchus in the USA. In the tile map, the layout approximates spatial location, with each state being an equal box filled with color representing cancer incidence. The geo-faceted map shows bar charts laid out in a grid approximating the spatial location of the state. The maps show age-adjusted rates for males and females. This display allows the presentation of multiple variables for each geographic area.

show spatial relationships by allotting spatial neighbors to the same group. It is one of several alternative displays that allow maps to become bivariate displays, commonly used to present both an estimate and the associated uncertainty.

Lucchesi and Wikle (Lucchesi and C.K., 2017) present bivariate choropleth maps blend color schemes to convey the intersection of categorized levels of an estimate and the associated uncertainty for each spatial area. They also suggest map pixilation, which breaks each region into small pixels, and allocates values to the individual pixels to create texture. This reflects the uncertainty around the area's estimate by randomly sampling from the confidence interval of the estimate of the area. Animating these displays involves resampling the pixels for each frame. Areas with uncertain values will flicker more dramatically than areas with more certain values.



Figure 2.4: Cartograms showing melanoma incidence in Australia: (a) contiguous, partially population transformed, (b) non-contiguous shape preserved, (c) Dorling, (d) hexagon tile map. The contiguous cartogram has expanded the highly populated areas while preserving the full shapes of rural areas. If it accurately sized areas by population, the country would be unrecognizable. The shape-preserved is unreadable due to the small area sizes. The Dorling cartogram presents all areas but many are difficult to compare. The hexagon tile map provides a reasonable spatial distribution despite having isolated hexagons in the outback areas.

2.4 Comparison and critique of alternative displays

2.4.1 Neither choropleth maps or cartograms perform well for Australia

Figure 2.4 shows four main types of cartograms using melanoma incidence on Australian Statistical Areas at Level 3 (Statistics, 2018). The version of a contiguous cartogram (a) has expanded the highly populated areas while preserving the full shapes of rural areas. It has not fully resolved the population transformation of areas, and if it had accurately sized areas by population, the country would be unrecognizable. The shape-preserved cartogram is unreadable, and it has reduced all areas to tiny spots on the map. Zooming in on a high-resolution output shows it does preserve the shapes. The Dorling cartogram

Table 2.4: Summary of features and constraints of common mapping methods used to display cancer statistics (Y=Yes, N=No, S=Sometimes)

Feature	Choro.	Contig.	Non-contig.	Dorling	Tiles	Geofacets
Spatial distortion	N	Y	Y	Y	Y	Y
Preserves neighbors	Y	Y	Y	S	S	S
Conceals small areas	Y	S	N	N	N	N
Uniform shape	N	N	N	Y	Y	Y
Univariate only	Y	Y	Y	S	S	N
Manual construction	N	N	N	N	Y	Y

and the hexagon tile map provide reasonable displays of the spatial distribution, despite having too much white-space in the outback areas.

2.4.2 Limitations of alternative displays

Cartograms provide the spatial distortion to more accurately convey the statistical distribution, focusing on the human impact of the disease. However, the transformation of contiguous cartograms often occurs at the expense of the shape of areas (Kocmoud and House, 1998), [Olson (1976), (Levison and Haddon Jr, 1965)]. When the population density of the geographic units is highly dissonant with geographic density, the cartogram will lose all spatial context. Dorling (Dorling, 2011) has a cartogram showing the 1966 general election results, which looked very little like the geographical shape of Australia.

Some mix of tiling, faceting or even micromaps, which allow some spatial continuity while also zooming into small areas, are good solutions for difficult geographies. Table 3 summarizes the key criteria for testing maps and alternative displays. Moore and Carpenter (Moore and Carpenter, 1999) and Bell et al. (Bell et al., 2006) provide suggestions and comments to help map creators best communicate their health data and spatial analysis.

2.5 User interaction

One of the concerns of adding too much information to a map is the fear of cognitive overload (McGranaghan, 1993) in which the user reaches an information threshold, beyond which they become confused. It can be a juggling act for a diverse audience, with experts probably preferring more detail (Cliburn et al., 2002) while a simpler display is more broadly readable. Interactivity is a design feature within modern mapping methods that can be used to incorporate additional information and complexity without overloading the user. Effective user-centred interactive actions produce rapid, incremental, and reversible changes to the display (Perin, 2014).

Monmonier (Monmonier, 2018) says that interactivity can be used to allow users to explore the map for more information and provides flexibility for the display. The user can toggle between different variables, map views or even multiple realizations of future scenarios (Goodchild, Buttenfield, and Wood, 1994). This provides additional mechanisms for the users to digest the uncertainty of the available information (MacEachren, 1992), (Van der Wel, Hootsmans, and Ormeling, 1994). When the needs of the audience are changeable and are also the priority, the map creator can allow interactivity for map users to explore a data set through dynamic interactions. This can allow inspection of the data from many views (Dang, North, and Shneiderman, 2001). User interaction with maps helps to understand and interpret the spatial distribution of disease, to validate, explain or explore the presented statistics and their relationships to each other (Carr, Wallin, and Carr, 2000).

Interactivity enables supplementary information to be incorporated into online atlases without cluttering the display. Interactive design features, found in online cancer maps, include tool tips, drop-down menus, data selection, zooming, and panning to allow users to explore the map as they want more information and allow flexibility in the display (Monmonier, 2018). The use of these supports can be found in various online cancer maps and are shown in Figure 2.5 (Roberts, 2019).

Animation, in contrast to interactivity, usually involves pre-computing views and showing these in a sequence. Lin Pedersen (Pedersen, 2018) provides an overview of animation for maps using the R package `ganimate` (Pedersen and Robinson, 2019). Animations are

used to communicate a message by capturing and directing users' attention. It is most often employed to show changes over time. The controls for basic animation are usually placed outside of the plot space (Pedersen, 2018), and the map image is updated/replaced as the animation progresses.

Weather maps are a thoroughly developed examples of animation of spatial displays to communicate information to the general public (Bell et al., 2006). The movement of a weather system will follow a forecasted path. All map users can follow the animated path of the weather system across the geography over a specified period.

The Australian Cancer Atlas (Cancer Council Queensland, Queensland University of Technology, and Cooperative Research Centre for Spatial Information, 2018) provides [tours](#) that change the display to draw users' attention to areas on the map that are relevant to the story. This implementation of animation gives users tools to plan their exploration.

Figure 2.6 shows two examples of more sophisticated interactive maps. The Spanish Cancer map (left) contains a linked display between a choropleth map and time series plots of cancer change. In linked plots, changing values in one display will trigger changes of corresponding elements in another display. Here, the temporal change in the choropleth map can be played out as an animation. Mousing over the time series plots will highlight the line for a particular region. The Canadian Breast Cancer Mortality map (right) has a magnifying glass that allows the user to zoom into small areas. It is easy to control and shows precise details in small areas.

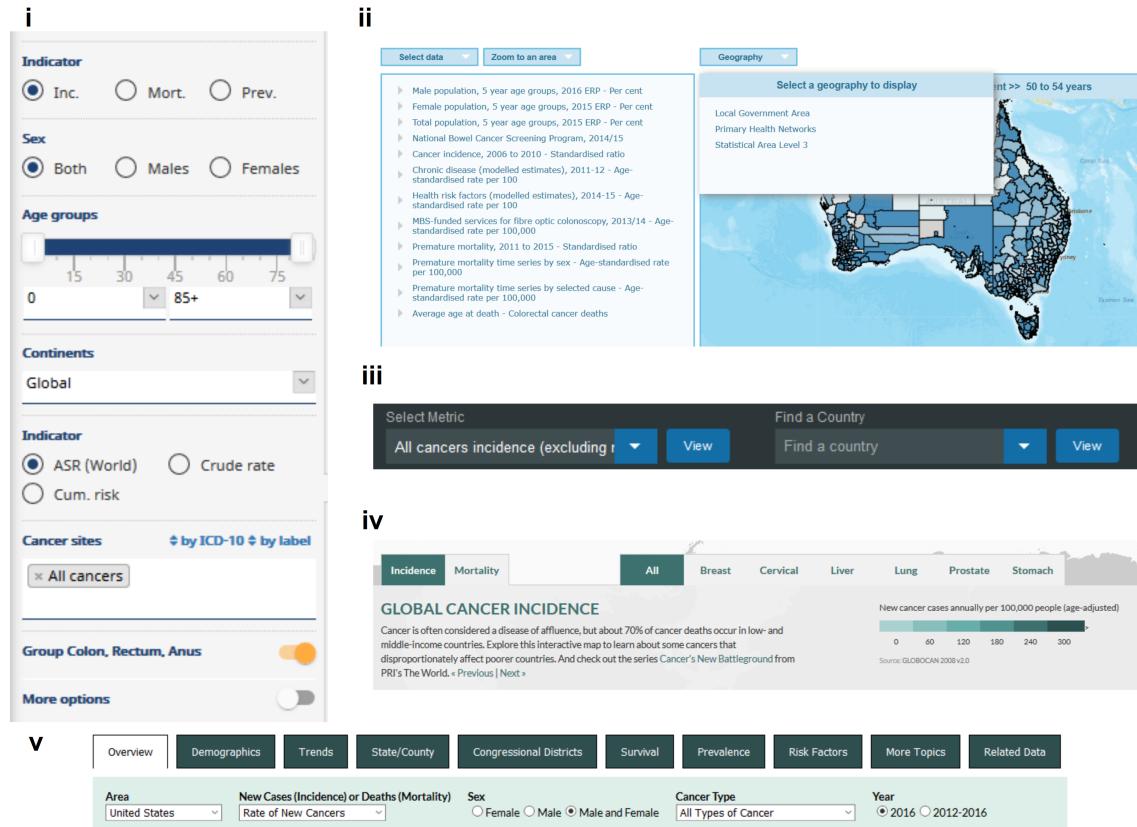


Figure 2.5: Interactive controls of displays in publicly available choropleth cancer maps: (i) GUI controls for statistic, sex, age groups, continents, and cancer types for Globocan 2018 [[@Globocan](#)], (ii) Menus for variable selection and zooming on Bowel Cancer Australia Atlas, (iii) Menus for choosing variables and countries in The Cancer Atlas, (iv) Tabs for different indicators and cancer types in Global Cancer Map, (v) Menus and toggles for variable and subset selection in United States Cancer Statistics: Data Visualizations.

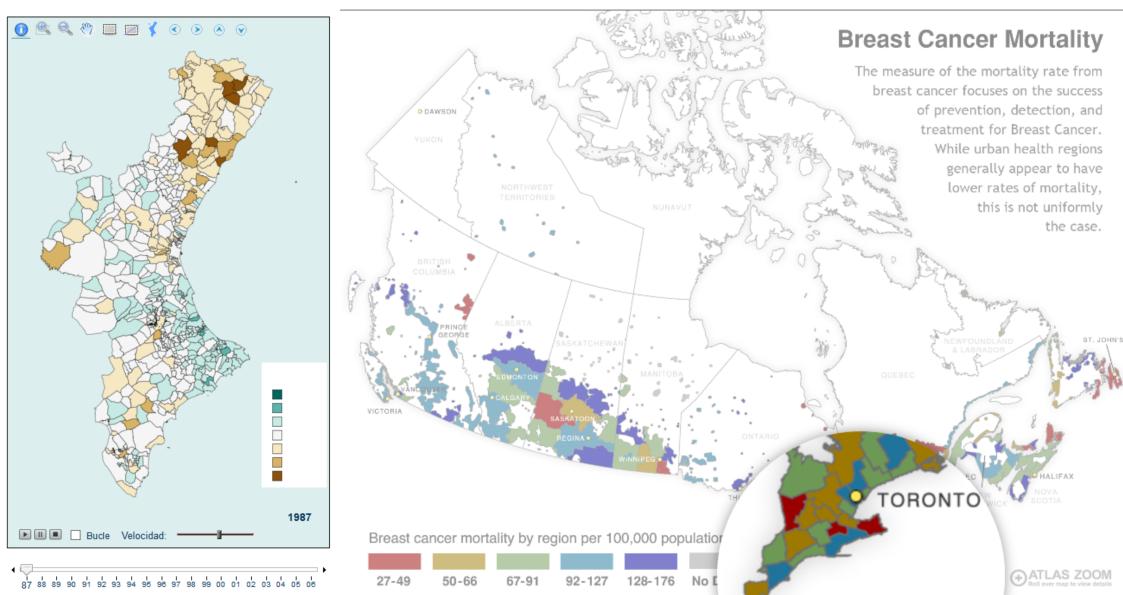


Figure 2.6: Two examples of advanced interactivity (and animation) in publicly available choropleth cancer maps: a. Linked maps and time-series line plots, with temporal animation in Map of Cancer Mortality Rates in Spain, b. A highly responsive magnifying glass on a map of Breast Cancer Mortality in Canada.

2.6 Conclusion

This paper provides an overview of mapping practices as commonly used for cancer atlases and recommends new approaches, such as cartograms and hexagon tile maps that should be adopted going forward. The conventional approach is the choropleth map, and it is widely used. The choropleth map suffers when there are small geographic units, as occurs in Australia where the population is concentrated on the coast, the information about the burden of cancer on those communities can be hidden. Making an inset can clarify congested regions but this breaks the viewers' attention as they shift focus from the map to the inset, and if there are many congested areas, many insets would be needed. The map alternatives implement trade-offs between the familiar shapes, and the importance of the geographic areas in the context of the areas. Given the population or a cancer statistic for each area, the geographic size or shape will change. Alternative displays allow the spatial distribution of cancer data to be digested by map users.

Many statistics are commonly used in cancer displays. The most basic is the incidence rate. It is common to see relative rates which measure how far a region is above or below the average. The purpose of using a relative rate is, perhaps the desire to pinpoint the areas that need attention because they have higher than expected rates. A region might be much higher than average, but it may not be close to a health concern, because all regions have a low incidence. Supplementary materials can allow map users to recognise when this occurs.

Interaction with maps is an important component of public atlases, and is easy to add with today's technology. The purpose is to provide access to more information than is possible to display in a single map, without overwhelming the viewer. Too many choices can similarly overwhelm a viewer, and thus decisions do need to be made about content to provide for accurate and comprehensive communication of information. Similarly, providing ways for users to interact with the display encourages engagement, and creative, efficient, elegant, interactive tools elicit curiosity about the data.

Chapter 3

An Algorithm For Spatial Mapping Using a Hexagon Tile Map, With Application to Australian Maps

This chapter relates to the first research aim as stated in Section 1.3. The chapter introduces the steps of the algorithm. The steps follow the procedures implemented in the `sugarbag` (Kobakian and Cook, 2019) package functions, that allow users to run the algorithm in R (R Core Team, 2019). It uses the Statistical Areas of Australia at Level 2, taking a subset and considering only those located within the island of Tasmania. It also provides an example of how to animate between the choropleth map display and the hexagon tile map.

This chapter will be submitted for publication to the *Journal of Statistical Software* for publication. The steps in this algorithm are implemented in the `sugarbag` (Kobakian and Cook, 2019) package for R (R Core Team, 2019).

Abstract

This algorithm creates a tessellated hexagon display to represent each of the spatial polygons. It allocates these hexagon in a manner that preserves the spatial relationship of the geographic units. It showcases spatial distributions, by emphasising the small geographical regions that are often difficult to locate on geographic maps. Spatial distributions

have been presented on alternative representations of geography for many years. In modern times, interactivity and animation have begun to play a larger role, as alternative representations have been popularised by online news sites, and atlas websites with a focus on public consumption. Applications are increasingly widespread, especially in the areas of disease mapping, and election results.

3.1 Introduction

The current practice for presenting geospatial data is a choropleth map display. These maps highlight the geographic patterns in geospatially related statistics (Moore and Carpenter, 1999). The land on the map space is divided into geographic units, these boundaries are usually administrative, such as states or counties. The units are filled with colour to represent the value of the statistic (Tufte, 1990).

Australian residents are increasingly congregating around major cities, the vast rural areas are often sparsely populated in comparison to the urban centres. In Australia, government bodies such as the Australian Bureau of Statistics (ABS), and the Australian Electoral Commission (AEC) hold the responsibility for the division of the population into geographic units. If it necessary, the AEC may adjust the boundaries of the areas as the population increases. The division of the population into approximately equal population areas results in dramatically different square meterage of the geographic areas. This can give unequal attention to the statistic of each area, this can cause misrepresentation of the spatial distributions of human related statistics in geographic maps.

The solutions to this visualisation problem begin with the geography. Cartograms apply a transformation to the geographic boundaries based on the value of the statistic of interest. These displays result in a distortion of the map space to represent differences in the statistic across the areas (Dougenik, Chrisman, and Niemeyer, 1985). The statistic of interest is used to determine the cartogram layout. When the Australian population is the statistic of interest, the result is a population cartogram. They fail to preserve a recognisable display due to the difference in size of metropolitan and rural areas (Dorling, 2011), (Berry, Morrill, and Tobler, 1964). Contiguous cartograms change the shape of areas, while preserving boundary relationships of neighbours. Non-contiguous cartograms maintain

the geographic shape of each geographic area, but will lose the connection to neighbours as the polygon for each geographic area shrinks or grows.

Alternative maps shift the focus from land area and shape, to the value of the statistics in a group of areas. Alternative mapping methods allow increased understanding of the spatial distribution of a variable across the population, by fairly representing each administrative area. This acknowledges that the amount of residents can be different but recognises that each area, or person within it is equally important.

tile maps, Rectangular cartograms (Kreveld and Speckmann, 2007) and Dorling cartograms (Dorling, 2011), all use one simple shape to represent each area. They place various importance on the preservation of spatial relationships, but all decrease the emphasis on the size of the geographic areas. These alternative map displays focus on the relationship between neighbours, attempting to preserve connections, and disregard the unique shapes of the administrative boundaries.

The `sugarbag` package provides a new algorithm to create tessellated hexagon tile maps. It emphasises the capital cities as population hubs, and emphasises the distances rather than size of large, rural geographic units.

3.2 Algorithm

The algorithm presented in `sugarbag` package operates on a set of simple feature geometry objects, also known as `sf` (Pebesma, 2018) polygons.

There are four steps performed to create a tessellated hexagon tile map. These steps can be executed by the main function, `create_hexmap`, or can be implemented separately for more flexibility. There are parameters used in the process that can be provided by users, if they are not, they will be automatically derived.

1. Create the set of centroids to allocate
 2. Create the grid of hexagons locations to use
 3. Allocate each centroid to an available hexagon
 4. Transform the data for plotting
-

Parameters

The `create_hexmap` function requires several parameters, if they are not provided, the information will be derived from the simple features (`sf`) set of shapes used. Users may choose to only use the `allocate` function when they wish to use a set of centroids, rather than (Pebesma, 2018) polygons.

The following parameters must be provided to ‘`create_hexmap`’:

- `shp`: an `sf` object containing the polygon information
- `sf_id`: name of a column that distinguishes unique areas
- `focal_points`: a data frame of reference locations used to allocate hexagons

Polygon set

The polygon set of Statistical Areas at Level 2 (SA2) (Statistics, 2018) of Tasmania in 2016 is provided with the `sugarbag` package as `tas_sa2`. A single column of the data set is used to identify the unique areas. In this case, the unique SA2 names for each SA2 have been used.

The longitude and latitude centre of the capital cities of Australia are used as focal points to allocate each geographic area around the closest capital city. Hobart will be the common focal point, as this example uses only the areas in the state of Tasmania.

The following parameters will be determined within `create_hexmap` if they are not provided. They are created as they are needed throughout the following example:

- `buffer_dist`: a float value for distance in degrees to extend beyond the geometry provided
- `hex_size`: a float value in degrees for the diameter of the hexagons
- `hex_filter`: amount of hexagons around centroid to consider for allocation
- `width`: the angle used to filter the grid points around a centroid

Create the set of centroid points

A set of centroids may be used directly. The set of polygons should be provided as an `sf` object, this is a data frame containing a geometry column. The `read_shape` function can assist in creating this object for use in R.

The centroids can be derived from the set of polygons using the `create_centroids` function:

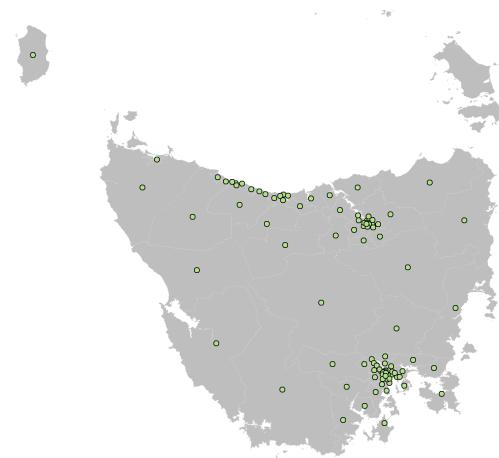


Figure 3.1: The geographic shapes of the Statistical Areas of Tasmania at Level 2. The points show the locations of the centroids of the SA2 areas.

Create the hexagon grid points

A grid is created to ensure tessellation between the hexagons that represent the geographic units on a hexagon tile map.

The grid of possible hexagon locations is made using the `create_grid` function. It uses the centroids, the hexagon size and the buffer distance.

Step 1: Creating a tessellated grid A set of longitude columns, and latitude rows are created to define the locations of the hexagons. The distance between each row and column is the size specified by `hex_size`. Equally spaced columns are created from the minimum longitude minus the buffer distance, up to the maximum longitude plus the buffer distance. Similarly, the rows are created from the latitude values and the buffer distance. A unique hexagon location is created from all intersections of the longitude

columns and latitude rows. Figure 3.2 shows the original grid on the left, to allow for tessellating hexagons, every second latitude row on the grid is shifted right, by half of the hexagon size. The grid for tessellation is shown on the right.

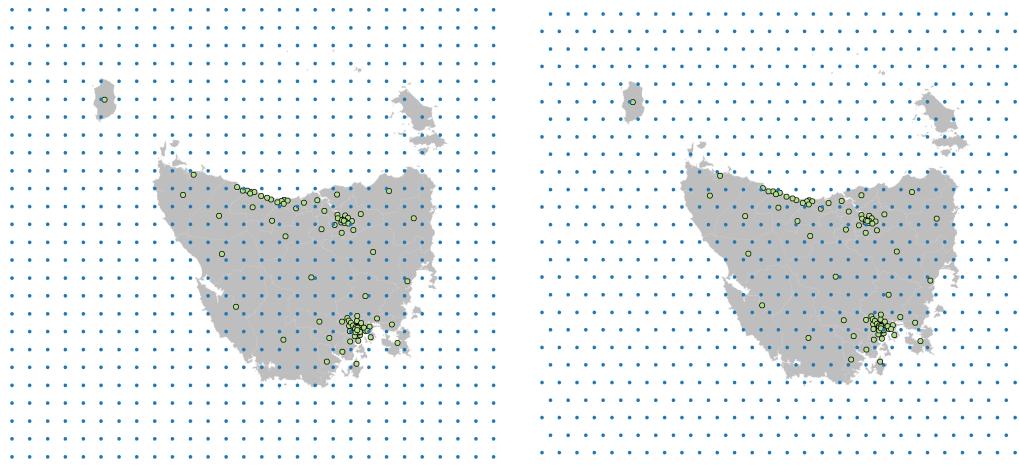


Figure 3.2: Grid points to create a tile map.

Step 2: Rolling windows Not all of the grid points will be used, especially if islands result in a large grid space. To filter the grid for appropriate hexagon locations for allocation, the `create_buffer` function is used by `create_grid`. It finds the grid points needed to best capture the set of centroids on a hexagon tile map.

The closest latitude row and longitude column are found for each centroid location. Then rows and columns of centroids are divided into 20 groups. The amount of rows in each latitude group and the amount of columns in each longitude group are used as the width of rolling windows. The rolling windows can be seen on the bottom and right of the grid shown in Figure 3.3. This will tailor the available grid points to those most likely to be used. It also helps reduce the amount of time taken, as it decreases the amount of points considered for each centroid allocation.

The first rolling window function finds the minimum and maximum centroid values for the sliding window groups of longitude columns and the groups of latitude rows.

The second rolling window function finds the average of the rolling minimum and maximum centroid values, for the longitude columns and latitude rows.

Step 3: Filtering the grid The grid points are kept only if they fall between the rolling average of the minimum and maximum centroid values after accounting for the buffer distance, for each row and column of the grid. The sparsely populated South-West region of National Park has much fewer points available compared to the South-East region containing the city of Hobart.

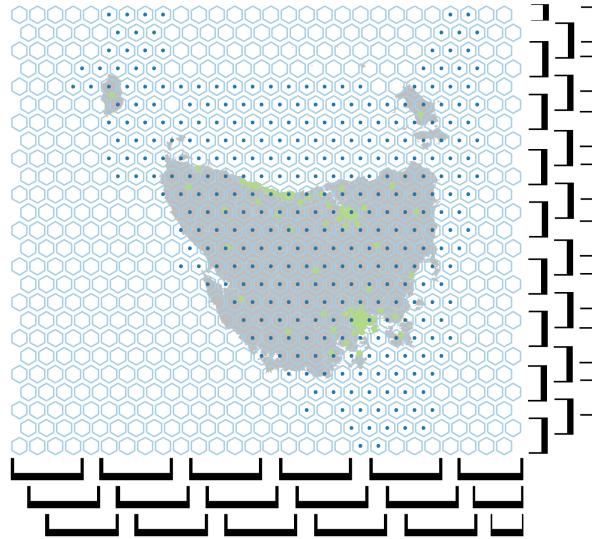


Figure 3.3: All possible hexagon locations from the initial grid are shown with blue outlines. The blue dots show the grid points left to choose from after the buffer step. The rolling windows show the collections of rows and columns used to filter the hexagon locations.

Centroid to focal point distance

The distance between each centroid in the set, and each of the focal points provided is calculated. The name of the closest focal point, and the distance and angle from focal point to polygon centroid is joined to polygon data set. To minimise time taken for this step only one option is provided, Tasmania's capital city Hobart. The order for allocation is determined by the distance between the polygon centroid and it's closest focal point. The points are arranged from the centroid closest to the focal point(s), to the furthest.

Allocate each centroid to a hexagon grid point

Allocation of all centroids takes place using the set of polygon centroids and the hexagon map grid. Centroid allocation begins with the closest centroid to a focal point. This will preserve spatial relationships with the focal point, as the inner city areas are allocated first, they will be placed closest to the capital, and the areas that are further will then be

accommodated. The possible hexagon grid points reduces by one after each allocation, then only those that have not yet been allocated are considered.

The possible hexagon locations to consider for a centroid are determined by the `hex_filter`. This is the maximum amount of hexagons between the centroid and the furthest considered hexagon. It is used to subset possible grid points to only those surrounding the polygon centroid within an appropriate range. A smaller distance will increase speed, but can decrease accuracy when width of the angle increases.

The following example considers the first of the Statistical Areas at Level 2. Within the algorithm, these steps are repeated for each polygon.

Step 1: Filter the grid for unassigned hexagon points Keeping only the available hexagon points prevents multiple geographic units from being allocated to the same hexagon.

Step 2: Filter the grid points for those closest to the centroid A box of possible hexagon locations around the centroid allows only the closest points that are not yet assigned to be considered. The corners of the box may not appear square if the buffer step has already removed unnecessary points from over the ocean.

The algorithm then removes the outer corners of the square, creating a circle of points, by only keeping points within a certain radial distance around the original centroid location.

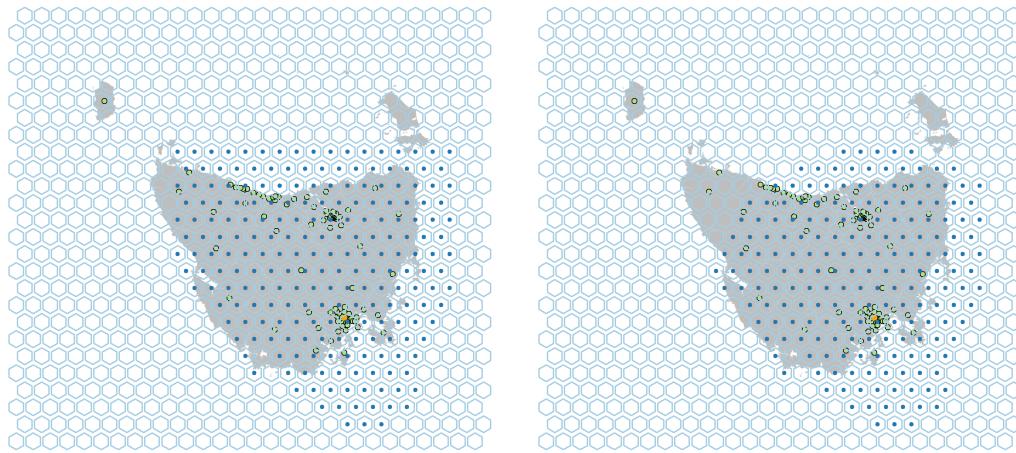


Figure 3.4: Filter for grid points within a square, then circular, distance for those closest to the centroid.

The width parameter is used to take a slice of the remaining points. The width is the amount of degrees used on either side of the angle from the focal point to centroid location. This uses the angle from the closest capital city, to the current centroid as seen in Figure 3.5 . This allows the spatial relationship to be preserved, even when it is allocated to a hexagon that is further from the focal point than the original centroid location.

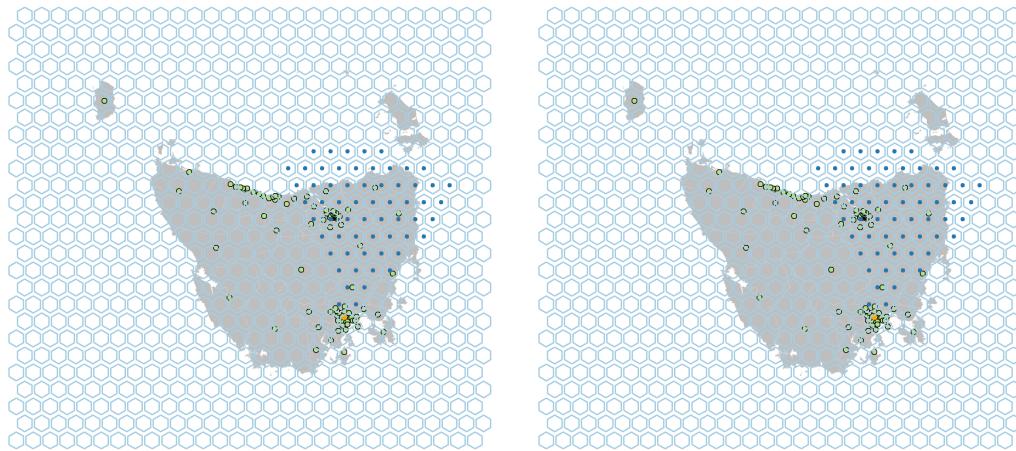


Figure 3.5: Filter for grid points within the angle from the focal point to the centroid.

If no available hexagon grid point is found within the original filter distance and angle, the distance is expanded, only when a maximum distance is reached will the angle expand

to accommodate more possible grid points.

By default the angle filter to hexagon grid points that fall within the bounds of the angle from the focal point to the geographic centroid, plus and minus 30 degrees. This will increase if no points can be found within the `hex_filter` distance. The default angle of 30 was chosen to allow the algorithm to choose hexagons that best maintained the spatial relationship between the focal point and geographic centroid.

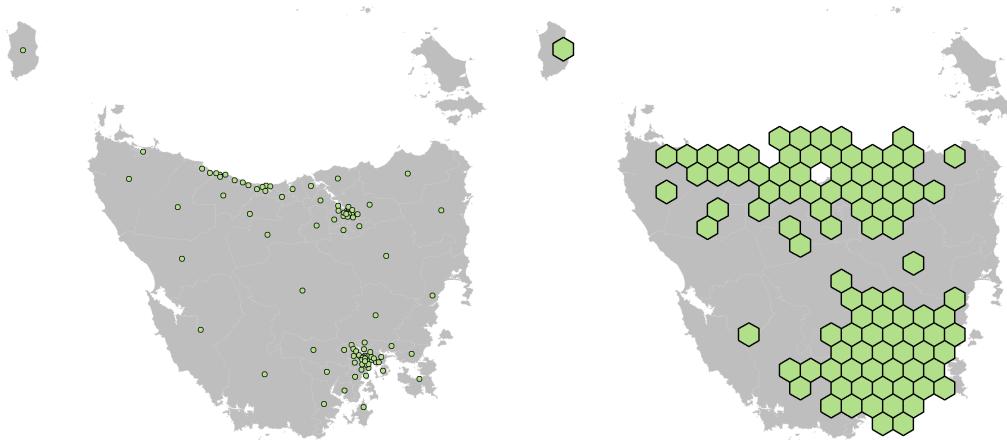


Figure 3.6: A complete hexagon tile map of Tasmania.

A complete hexagon tile map of Tasmania is created by applying the algorithm steps to each centroid. The hexagon tile map visualisation is used below to visualise the Australian Cancer Atlas data. Two views of the same data are produced by filling according to the Lung Cancer Standardised Incidence Rates (SIRs) downloaded from the Australian Cancer Atlas site. This small example in Figure 3.7 shows the group of blue areas in the Hobart CBD more prominently in the hexagon tile map (b). The small red areas visible in the choropleth map (a) along the north coast are much larger in the hexagon tile maps. The hexagon tile map shows less yellow, this no longer overwhelms the map space with the information regarding the rural areas.

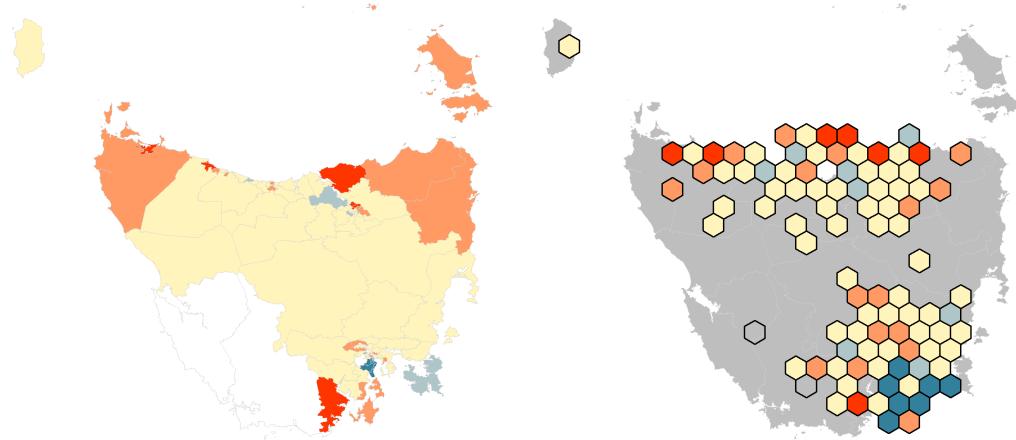


Figure 3.7: The Australian Cancer Atlas data has determined the colour of each Statistical Area of Australian at Level 2. A choropleth map (a) of SIR is paired with a hexagon tile map (b) to contrast the colours that are made obvious when every SA2 is equally represented.

Neighbour relationships

It is possible to consider the neighbouring areas for each SA2, for stronger preservation of the spatial distribution.

An additional step can be included to allow the neighbours that have already been allocated to influence the placement of the current centroid. This requires specifying the `sf` object as the argument for the `use_neighbours` parameter. This calculates neighbours using intersections of their polygons. This occurs for all areas before any allocations begin.

During the allocation of each centroid, the list of neighbours is consulted. If any neighbour was already allocated, the hexagons surrounding the neighbours on the grid are prioritised. For multiple neighbours, the neighbouring hexagon grid points are aggregated and considered in order of distance from the original centroid.

3.3 Using sugarbag

Installation

The package can be installed from CRAN:

and the development version can be install from the GitHub repository:

Load the library into your R session with:

Creating a hexagon tile map

The following code creates the hexagon tile map for all the Statistical Areas at Level 2 in Tasmania.

```
## Load data
data(tas_sa2)

## Create centroids set
centroids <- create_centroids(tas_sa2, "SA2_NAME16")

## Create hexagon grid
grid <- create_grid(centroids = centroids,
                    hex_size = 0.2,
                    buffer_dist = 1.2)

## Allocate polygon centroids to hexagon grid points
hex_allocated <- allocate(
  centroids = centroids,
  hex_grid = grid,
  sf_id = "SA2_NAME16",
  ## same column used in create_centroids
  hex_size = 0.2,
  ## same size used in create_grid
  hex_filter = 10,
  use_neighbours = tas_sa2,
  focal_points = capital_cities %>% filter(points == "Hobart"),
  width = 35,
  verbose = FALSE)
```

```
## Prepare to plot

fort_hex <- fortify_hexagon(data = hex_allocated,
  sf_id = "SA2_NAME16", hex_size = 0.2)

## Make a plot

library(ggplot2)

ggplot(fort_hex) +
  geom_polygon(aes(x=long, y=lat, group=hex_id, fill = lat)) +
  scale_fill_distiller("", palette="PRGn")
```

3.4 Applications

Australian Cancer Atlas

The Australian Cancer Atlas (Cancer Council Queensland, Queensland University of Technology, and Cooperative Research Centre for Spatial Information, 2018) allows estimates derived from the models of SIRs and excess deaths to be downloaded. Figure 3.8 is a choropleth map that uses colour to display the estimated SIRs of melanoma cancer for all persons for each SA2. The Australian choropleth map display draws attention to the expanse of light blue areas across the rural communities in all states. The SA2s around Brisbane stand out as more orange and red. Comparatively, the hexagon tile map display in Figure 3.9 draws attention to contrast of the blue areas in Sydney and Melbourne and the capital city of Brisbane. In both Sydney and Melbourne, the hexagons that represent the SA2 areas in the inner-city areas have lower than average Incidence Rates.

With careful consideration of the choropleth map, the small geographic inner city areas may have been noticed by viewers, but the hexagon tile map display emphasises them. The communities in northern Queensland and the Northern territory do not draw attention because of their size as in the choropleth, but their colour is still noticeably below average when contrasted with the hexagons further south.

To create this choropleth map the SA2 polygons for 2011 from the ABS. The SIRs for each geographic unit are joined to the appropriate polygons.



Figure 3.8: A choropleth map of the Statistical Areas of Australia at Level 2. The colours communicate the value of the estimated SIR of Melanoma for all persons, they range from much lower than average (blue) to much higher than average (red)

To create the hexagon tile map display, the same steps are followed as outlined above:

- Create the set of centroid points
 - Create the hexagon grid points
 - Allocate each centroid to a hexagon grid point
-

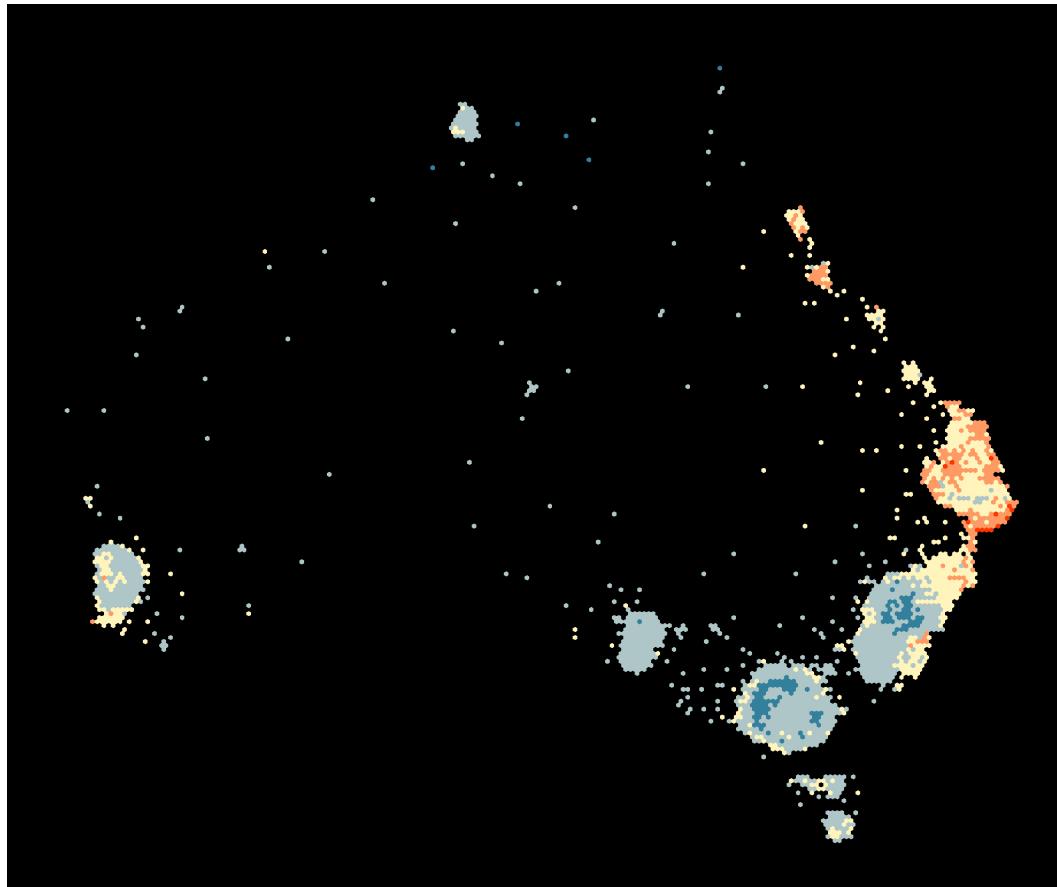


Figure 3.9: A hexagon tile map of the Statistical Areas of Australia at Level 2. The colours communicate the value of the estimated SIR, they range from much lower than average (blue) to much higher than average (red)

3.5 Animation

The `ganimate` (Pedersen and Robinson, 2019) package can be used to make an animation. It requires connecting the polygons for each area in two displays, which can be done using the `sf_id` variable, such as the SA2 name. The animation¹ connecting these two displays will highlight the rapid growth of the inner-city areas, and will decrease the large rural areas. The hexagons that move the furthest will move rapidly in the animation.

3.6 Conclusion

It is possible to use alternative maps to communicate spatial distributions. While a choropleth map display is the current practice spatial visualisation of geographical data. Current methods do not always work for Australia due to the large geographic space between the

¹This animation can be viewed at: <https://sugarbagjss.netlify.com/>

densely populated capital cities. The administrative boundaries may also distract from the statistics communicated using colour.

Alternative maps highlight the value of the statistics across the geographic units. Alternative mapping methods allow increased understanding of the spatial distribution of a variable across the population, by fairly representing each administrative area. This acknowledges that the amount of residents can be different but recognises that each population area is equally important. The solution to this visualisation problem has equally sized areas, with neighbourhood boundary connections. This map algorithm is implemented in the `sugarbag` (Kobakian and Cook, 2019) package written for R (R Core Team, 2019). The `sugarbag` package creates tessellated hexagon tile maps. The Australian application preserves the spatial relationships, emphasising capital cities. The hexagon tile map is a visualisation solution that highlights spatial distributions.

These hexagons equally represent each area. However, the tessellation does not allow the size of the hexagons to represent another variable, similar to the choropleth maps. The algorithm is heavily dependent on the focal points used, as this determines the order of allocation. It works on the assumption that viewers can use directional relationships to identify their neighbourhoods but this can be aided by the animation.

Future work will include refining the algorithm. It would be possible to take a logarithmic function rather than a direct angle to help choose a closer hexagon to the original centroid location, before increasing the width of the angle used to filter the hexagons.

This algorithm has only been tested using single countries, and does not consider definite borders of countries. While the buffer allows extension beyond the furthest centroids, there is no mechanism to protect the borders and ensure centroids are placed within the geographic borders of a country.

This algorithm is an effective start to creating hexagon tile maps for many geographic units.

Chapter 4

Comparing the Effectiveness of the Choropleth Map with a Hexagon Tile Map for Communi- cating Cancer

This chapter tests the performance of the hexagon tile map display created using the algorithm discussed in Section 3. It outlines the lineup protocol method of visual inference that can be used to test the effectiveness of information visualisations. Using a two factor experimental design, the experiment contrasts the performance of participants when they viewed a choropleth map, and their performance when viewing a hexagon tile map. The experiment also considered three types of spatial trends, one geographic trend, and two population related distributions. The results showed that participants did in fact more frequently find the population related distributions when using the hexagon tile map.

This chapter will be submitted to the journal *IEEE Transactions of Visualisation and Computer Graphics*.

Abstract

The choropleth map display is commonly used for communicating spatial distributions across geographic areas. However, when choropleths are used the size of the geographic units will influence the understanding of the distribution derived by map users. The hexagon tile map is presented as an alternative display for visualizing population related distributions effectively. Visual inference is used to measure the power of the hexagon tile map design, and the choropleth is used as a comparison. The hexagon tile map display is tested using a distribution that is directly related to the geography, with values monotonically increasing from the North-West to South-East areas of Australia. This study finds in a hexagon tile map lineup the single map that contains a population related distribution is detected with greater probability than the same data displayed in a choropleth map. These findings should encourage map creators to implement alternative displays and consider a hexagon tile map when presenting spatial distributions of heterogeneous areas.

4.1 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 *sugarbag* (Kobakian and Cook, 2019).

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

The Australian Cancer Atlas (Cancer Council Queensland, Queensland University of Technology, and Cooperative Research Centre for Spatial Information, 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.

The paper is organised as follows. The next section discusses the background of geographic data display and visual inference procedures. The [Methodology](#) section describes the methods for conducting the experiment and analysing the results. The results are summarized in the [Results](#) section.

4.2 Background

4.2.1 Spatial data displays

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

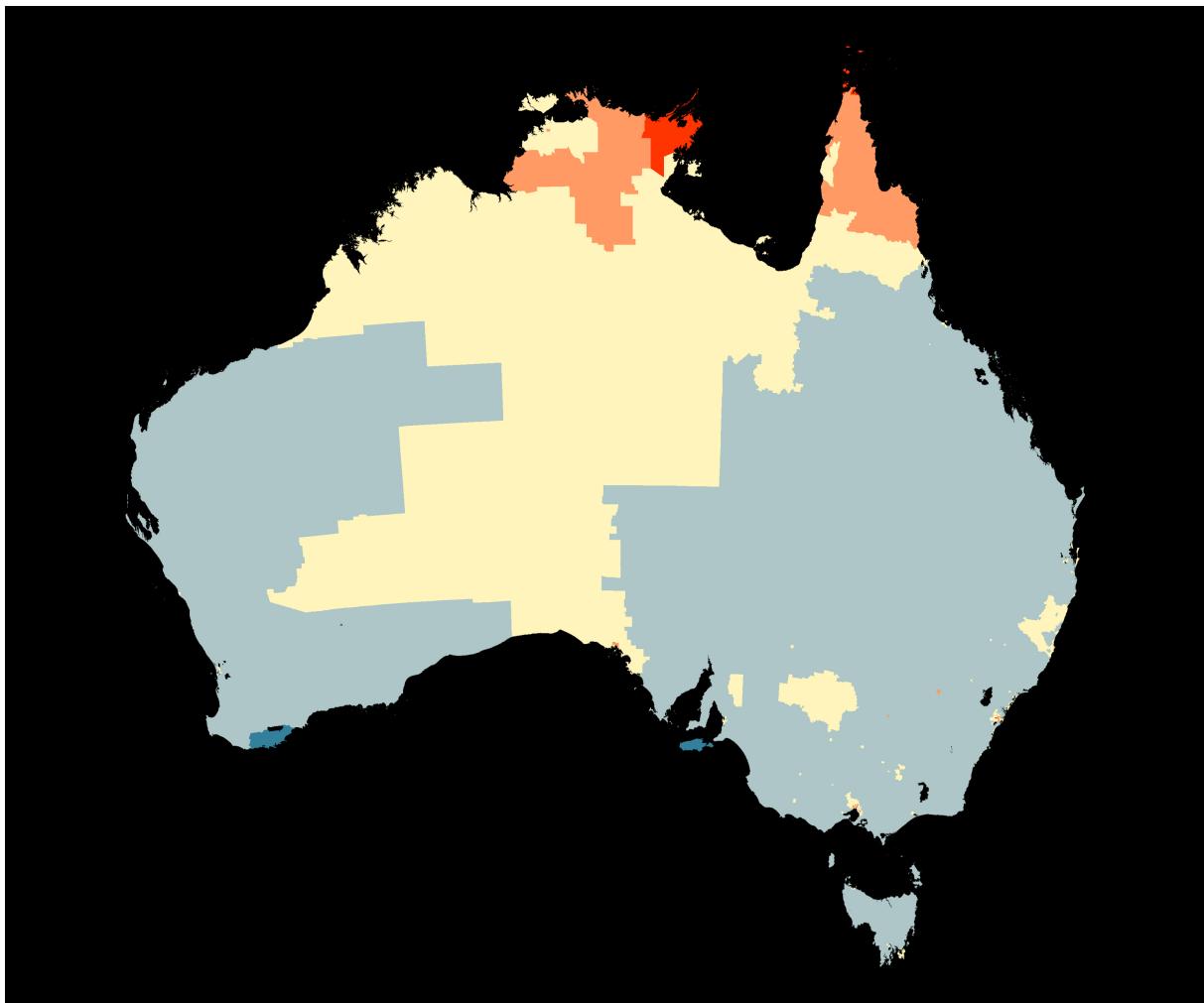


Figure 4.1: A choropleth map of the smoothed average of liver cancer diagnoses for Australian males. The diverging colour scheme uses dark blue areas for much lower than average diagnoses, yellow areas with diagnoses around the Australian average, red shows diagnoses much higher than average. The hexagon tile map shows concentrations of higher than expected liver cancer rates in the cities of Melbourne and Sydney, which is not visible from the choropleth.

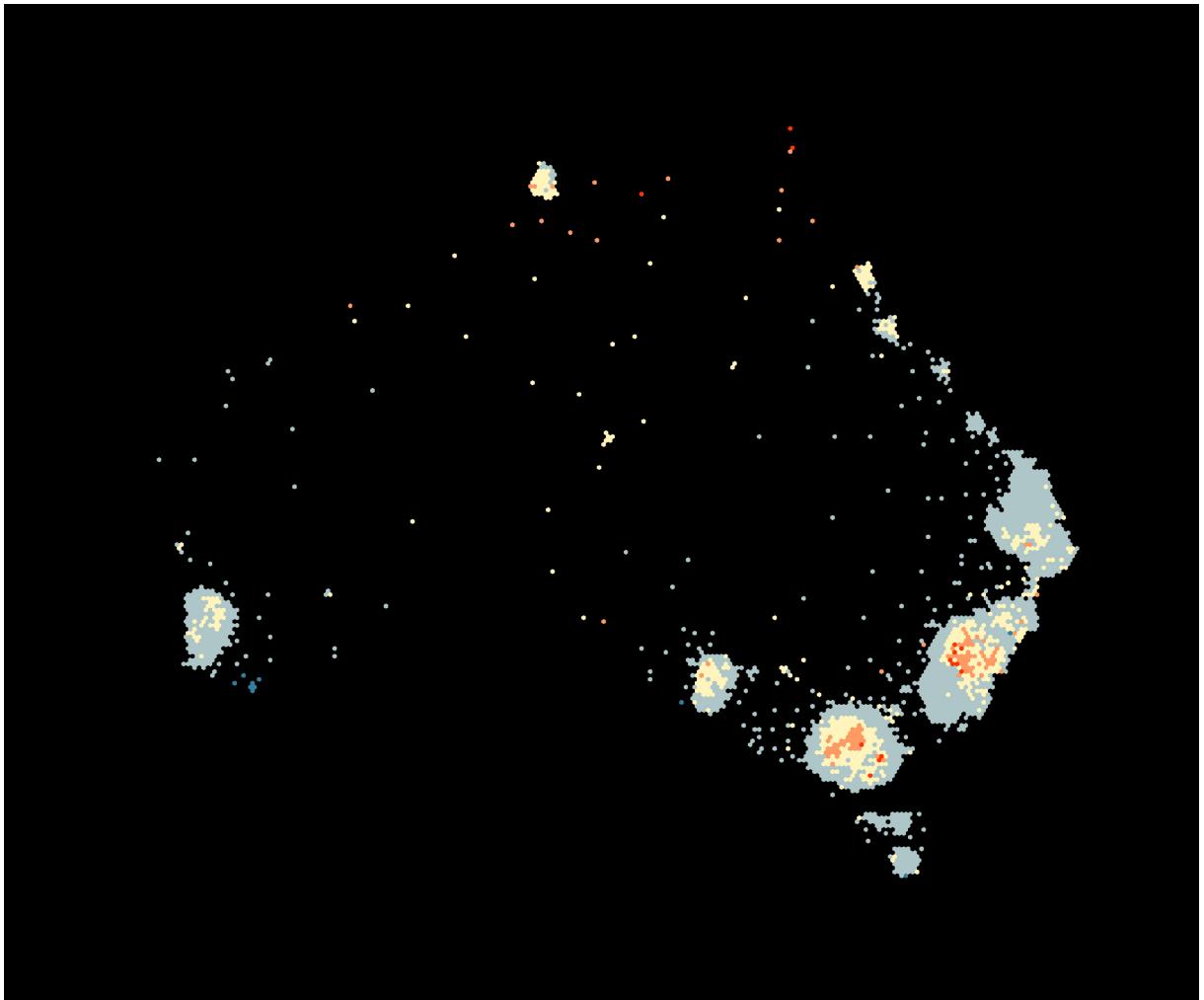


Figure 4.2: A hexagon tile map of the smoothed average of liver cancer diagnoses for Australian males. The diverging colour scheme uses dark blue areas for much lower than average diagnoses, yellow areas with diagnoses around the Australian average, red shows diagnoses much higher than average. The hexagon tile map shows concentrations of higher than expected liver cancer rates in the cities of Melbourne and Sydney, which is not visible from the choropleth.

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

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

Alternative algorithms make various trade offs between familiar shapes and representation of geographic units. The non-contiguous cartogram method (Olson, 1976) keeps the shapes of geographic units intact, and changes the size of the shape. This method disconnects areas creating empty space on the display losing the continuity of the spatial display of the statistic. The Dorling cartogram (Dorling, 2011) represents each unit as a circle, sized according to the value of the statistic. The neighbour relationships are mostly maintained by how the circles touch. A similar approach was pioneered by Raisz (1963), using rectangles that tile to align borders of neighbours (Monmonier, 2005). There have been thorough reviews of the array of methods, as suitable for cancer atlas displays (S., Cook, and Roberts, 2019), (Skowronnek, 2016).

The hexagon tile map algorithm, automatically matches spatial regions to their nearest hexagon tile, from a grid of tiles. It has the effect of spreading out the inner city areas while maintaining the spatial locations or regions in remote areas. The algorithm is available in the R package, sugarbag (Kobakian and Cook, 2019). Figure 4.1 shows the hexagon tile map, along with the choropleth map of liver cancer rates in Australia. Colour maps from substantially below average (blue) to substantially above average (red) rates. The inner city areas have expanded, making it possible to see the cancer incidence in the small, densely populated areas. Remote regions are represented by isolated hexagons, which is not ideal, but maintains the spatial location of these data values. It is of interest to know how well the spatial distribution is perceived for this display, in comparison to the choropleth.

4.2.2 Visual Inference

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

In classical statistical inference hypothesis testing is conducted by comparing the value of a test statistic on a standard reference distribution, computed assuming the null hypothesis is true. If the value is extreme, the null hypothesis is rejected, because the test statistic value is unlikely to have been so extreme if it was true. In the lineup protocol, the plot plays the role of the test statistic, and the data plot is embedded in a field of null plots. Defining the plot using a grammar of graphics (Wickham, 2016) makes it a functional mapping of the variables and thus, it can be considered to be a statistic. With the same data, two different plots can be considered to be competing statistics, one possibly a more powerful statistic than the other.

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

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

The chance that an observer picks the data plot out of a lineup of size m plots accidentally, if the null hypothesis is true is $1/m$. With K observers, the probability of k randomly choosing the data plot, roughly follows a binomial distribution with $p = 1/m$. Figure 4.3 shows a lineup of the hexagon tile map, of size $m = 12$. Plot 3 is the data plot, and the remaining 11 are plots of null data.

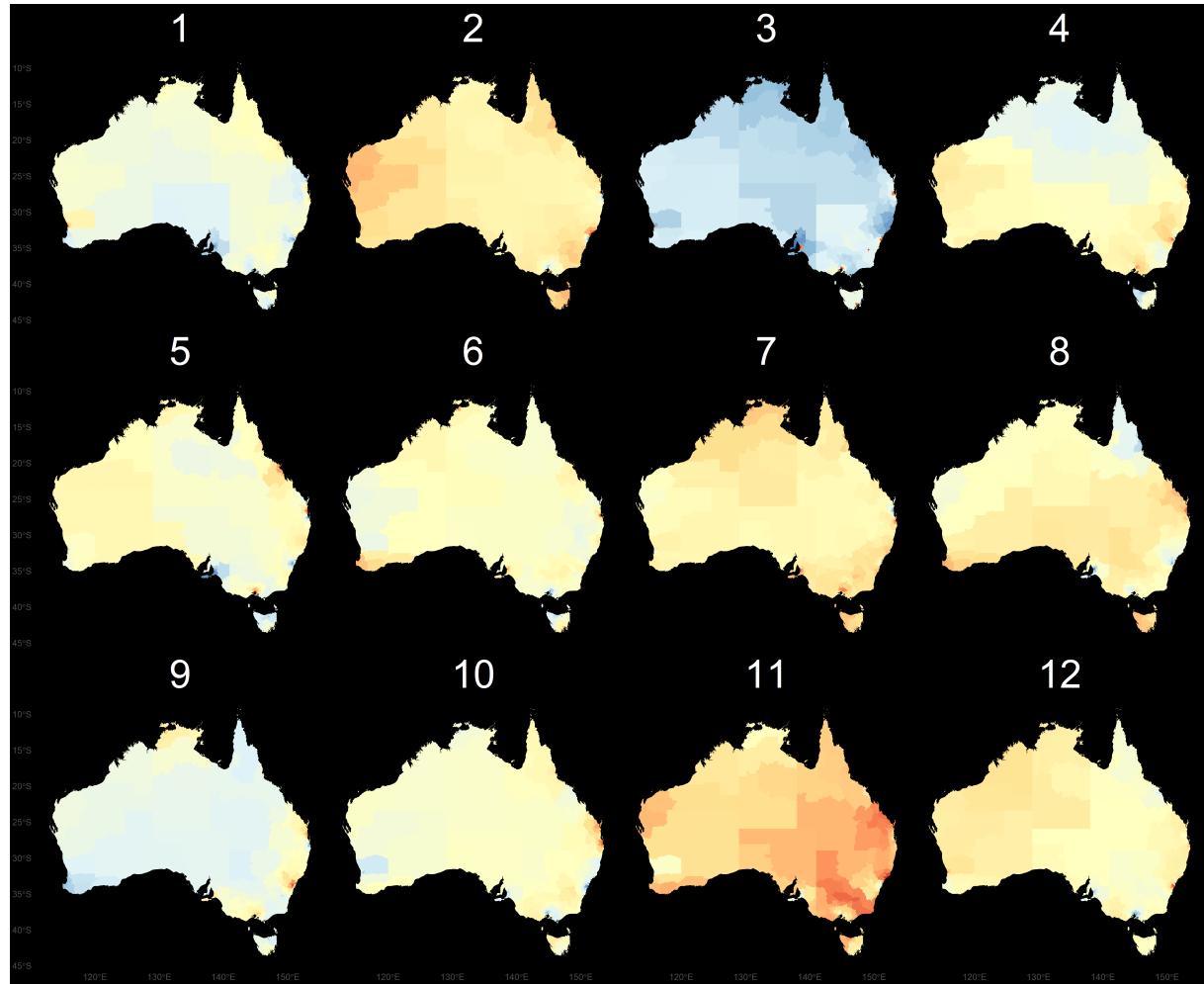


Figure 4.3: This lineup of twelve choropleth displays contains one map with a real population related structure. The rest are null plots that contain spatial correlation between neighbours.

In order to determine the effectiveness of a type of display, this probability is less relevant than the overall proportion of observers who pick the data plot, k/K . The power of the test statistic (data plot) is provided by this proportion. Power in a statistical sense is the ability of the statistic to *produce a rejection* of the null hypothesis, if it is indeed *not true*. With the same data plotted using two different displays, the display with the highest proportion of people who choose the data plot would be considered to be the most powerful statistic.

4.2.3 Methodology

This study aims to answer two key questions around the presentation of spatial distributions:

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

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

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

4.2.4 Experimental factors

The primary factor in the experiment is the plot type. The secondary factor is a trend model. Three trend models were developed, one mirroring a large spatial trend for which the choropleth would be expected to do well, and two with differing level of inner city hot spots. These latter two reflect the structure seen in the liver cancer data (Figure 4.1). This produces six treatment levels:

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

Data is generated for each of the trend models, with four replicates, and each displayed both as a choropleth and as a hexagon tile map, which yields 12 data sets, and 24 data plots. This set of displays is divided in half, providing two sets of 12 displays, Group A and Group B. Participants were randomly allocated to Group A or B. Participants saw a data set only once, either as a choropleth or as a hexagon tile map. Table 4.4 summarises the design and the allocation of the displays.

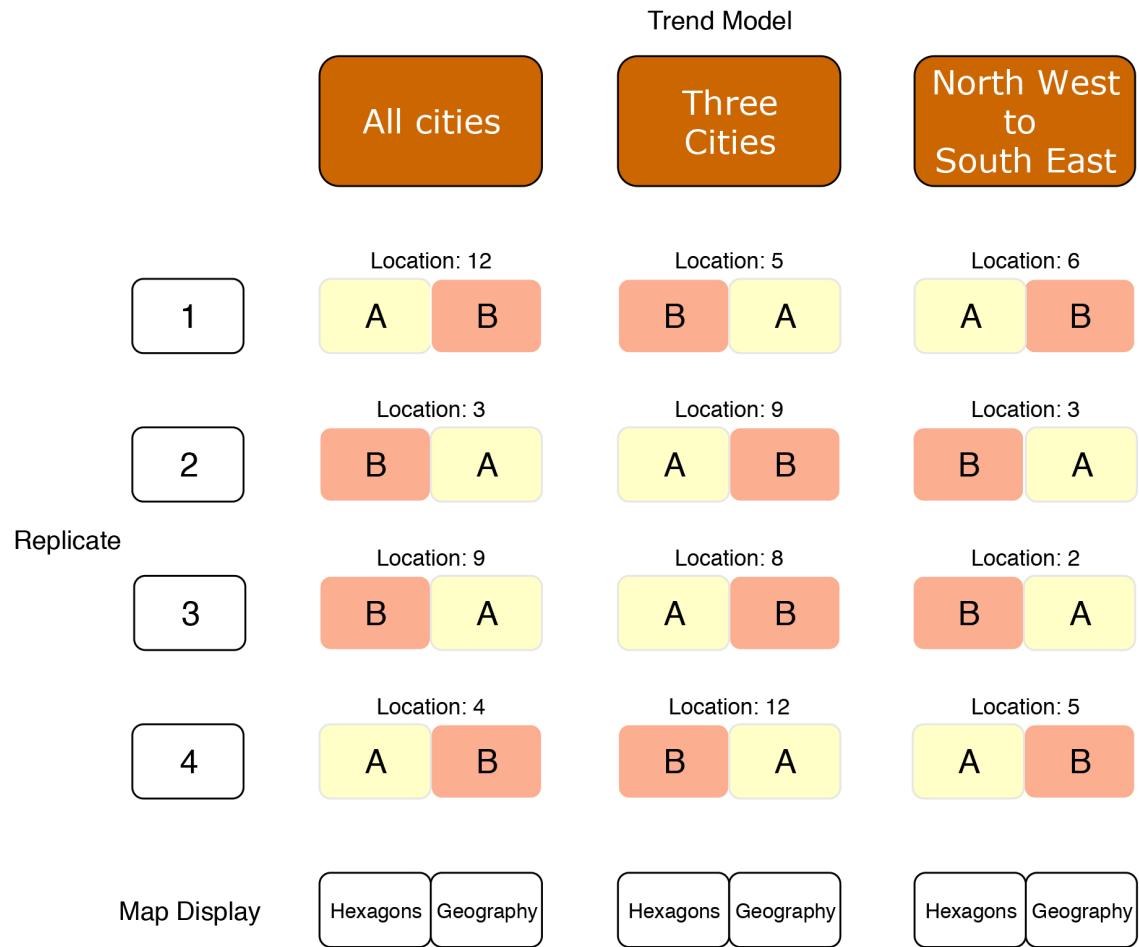


Figure 4.4: The experimental design used in the visual inference study.

4.2.5 Generating null data

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

For spatial data, null data is stationary data, where the mean, variance and spatial dependence are constant over the geographic units. Stationary data is specified by a variogram

model (Matheron, 1963). Simulating from a variogram model, where the spatial dependence is specified, generates the stationary spatial data used for the null plots. The parameters for the Gaussian model were sill=1, range=0.3 with the variance generated by a standard normal distribution.

The R package `gstat` (Pebesma and Graeler, 2019) was used to simulate 144 null sets, 12 data sets for each plot in a lineup, and 12 sets for 12 lineups.

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

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

This smoothing allowed neighbors to be related to each other, but also allowed outliers, and showed distributions similar to the Liver cancer distribution (Figure 4.1).

4.2.6 Generating lineups

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

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

The All Cities trend model was created using the distance from the centroid of each geographic unit to the closest capital city in Australia, calculated when creating the hexagon tile map using the `sugarbag` (Kobakian and Cook, 2019) package. 201 of the 336 SA3s were considered greater capital city areas, the values of these areas were increased

to create red clusters. The amount was chosen to make clusters around the cities visible in the choropleth display even if they were not overtly noticeable.

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

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

As seen in Figure 4.4, the choropleth and hexagon display used the same location for the real data display of the trend model was added to the spatially correlated null values for each lineup. Each set of lineup data was used to produce a choropleth map lineup and hexagon tile map lineup. These matched pairs were split between Group A and Group B according to the 2 x 3 factor experimental design depicted in Figure 4.4.

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

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

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

4.2.7 Analysis

Data Cleaning

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

Descriptive statistics

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

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

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

Plots will be made in R (R Core Team, 2019), with the `ggplot2` package (Wickham, 2016).

Modelling

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

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

where $y_{ij} = 0, 1$ is the log odds for whether the subject detected the data plot, μ is the overall mean, $\tau_i, i = 1, 2$ is the map type effect, δ_j is the trend model effect. We are allowing for an interaction between map type and trend model as the response is binary, so a logistic model was used. As each participant provides results from 12 lineups, this model can account for each individual participants' abilities as it includes a subject-specific random intercept.

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

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

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

4.2.8 Web application to collect responses

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

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

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

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

4.2.9 Participants

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

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

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

4.2.10 Demographic data collection

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

Demographics were collected regarding the study participants:

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

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

Three questions were asked regarding each display:

- Plot choice
- Reason
- Difficulty

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

4.3 Results

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

4.3.1 Participant demographics

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

4.3.2 Accuracy

Figure 4.5 displays the average detection rates for the two types of plot separately for each trend model. Each trend model was tested using four repetitions, evaluations on the same data set were seen as either choropleths or hexagon tile maps by each group as specified in Table 4.4; the detection rates for each display are connected by a line segment. The Three Cities and All Cities trend models shown in the hexagon tile map allowed viewers to detect the data plot substantially more often than the choropleth counterparts. One replicate for the All Cities group had a similar detection rate for both the choropleth and the hexagon tile map. Interestingly, in post-analysis we found that participants chose the data display in the choropleth lineup for reasons unrelated to the All Cities data structure. Participants detected the gradual spatial trend in the NW-SE group equally well from both map types. This was a pleasant surprise; we expected that the choropleth map would be superior for the type of spatial pattern, but the data suggests the hexagon tile map performs equally as well.

Table 4.1 shows the means and standard deviations of the detection rate for each type of plot and each trend model. This also gives the standard deviations, the smallest standard deviation for all sets of replicates was the Three Cities trend model shown in a choropleth display. This group of displays had a very small detection rate of 0.04. The mean detection rate for the Three Cities trend model shown as choropleth map lineups was also the smallest at 0.40. The North-West to South-East (NW-SE) trend model unexpectedly had a higher mean detection rate for the hexagon tile map displays, but the difference in the means of detection rate was only 0.10.

Table 4.1: *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)

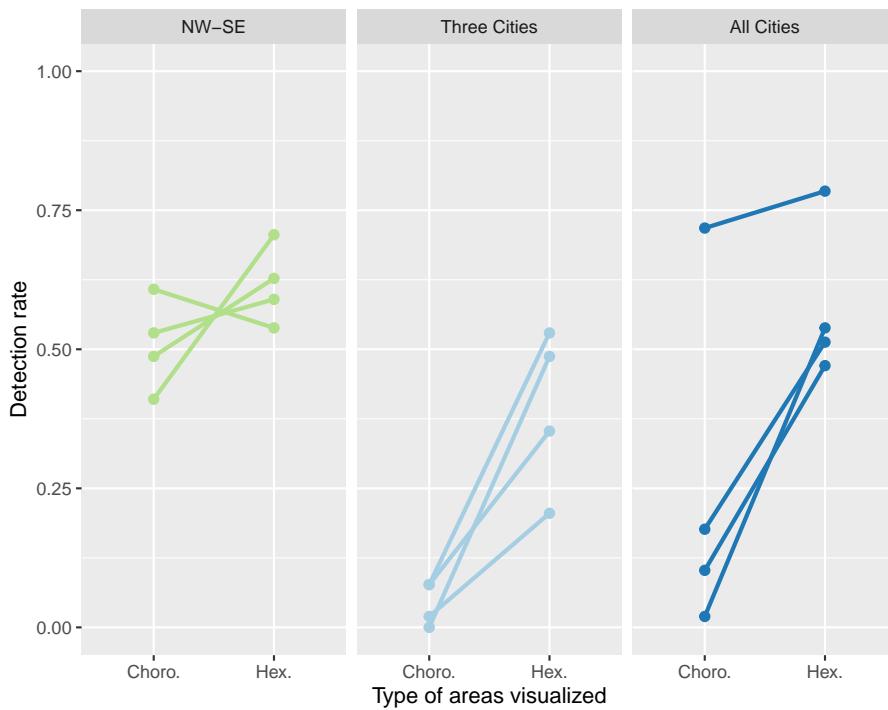


Figure 4.5: The detection rates achieved by participants are contrasted when viewing the four replicates of the three trend models. Each point shows the probability of detection for the lineup display, the facets separate the trend models hidden in the lineup. The points for the same data set shown in a choroleth or hexagon tile map display are linked to show the difference in the detection rate.

Table 4.2 presents a summary of the generalised linear mixed effects model, testing the effect of plot type and trend model on the detection rate. The results support the summary from Figure 4.5 and all parameters are statistically significant despite the large standard deviations observed in Table 4.1. Overall, the hexagon tile map performs marginally better than the choropleth for all trend models, which is a pleasant surprise. Allowing for the interaction effect, the difference in detection rate decreases for population related displays for a choropleth map lineup, but increases for a hexagon tile map display. The log odds of detection show in Table 4.2 can be back transformed after taking the sum of all terms for the trend and type of display that are of interest. For the NW-SE distribution, the predicted detection rate for the hexagon tile map display increases the predicted probability of detection to 0.63 from 0.52 for choropleths, this is almost exactly the difference seen in the table of means and is significant only at the 0.05 level.

When a choropleth map display is used, the predicted detection rate for the Three Cities trend, 0.03; this is extremely low, especially compared to the NW-SE trend of 0.52. When

the All Cities trend is presented in a choropleth display the predicted probability of detection is 0.22. The hexagon tile map has a substantially high detection rate for the display of a Three Cities trend 0.39 and All Cities trend 0.59.

Table 4.2: *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	0.07		0.16	0.67
Hex.	0.46	*	0.22	0.04
Three Cities	-3.41	***	0.42	0.00
All Cities	-1.34	***	0.24	0.00
Hex:Three Cities	2.44	***	0.47	0.00
Hex:All Cities	1.16	***	0.33	0.00

4.3.3 Speed

Figure 4.6 shows horizontally jittered dot plots to contrast the time taken by participants to evaluate each lineup when viewing each type of display. The time are also separated by trend model and whether the data plot was detected or not detected. The time taken to complete an evaluation ranged from milliseconds to 60 seconds. The average time taken for type of display is shown as a large colored dot on each plot. when considering the heights of the green and orange dots, there is little difference in the average time taken to read a choropleth or hexagon tile map. Comparing the same colored dot across each trend model row, there is a slight increase in the time taken to correctly detected the data plot in the hexagon tile map lineup, but little difference in evaluation time for the choropleth display. However, there were substantially less correct detections for choropleth lineups for the Three cities and All Cities trends.

4.3.4 Certainty

Participants provided their level of certainty regarding their choice using a five point scale. Unlike the accuracy and speed of responses that were derived during the data processing phase, this was a subjective assessment by the participant prompted by the question: ‘How certain are you about your choice?’. Figure 4.7 shows the amount of times

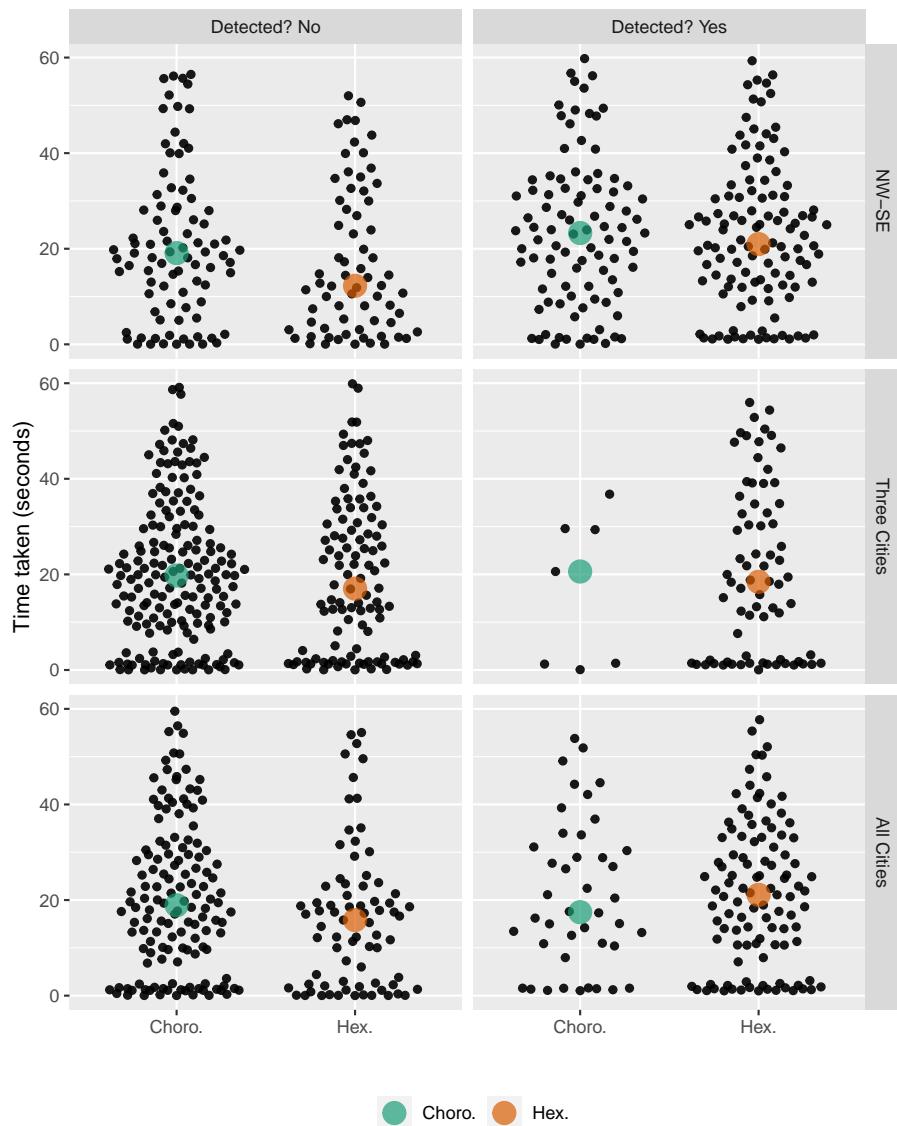


Figure 4.6: The distribution of the time taken (seconds) to submit a response for each combination of trend, whether the data plot was detected, and type of display, shown using horizontally jittered dotplots. The colored point indicates average time taken for each plot type. Although some participants take just a few seconds per evaluation, and some take as much as much as 60 seconds, but there is very little difference in time taken between plot types.

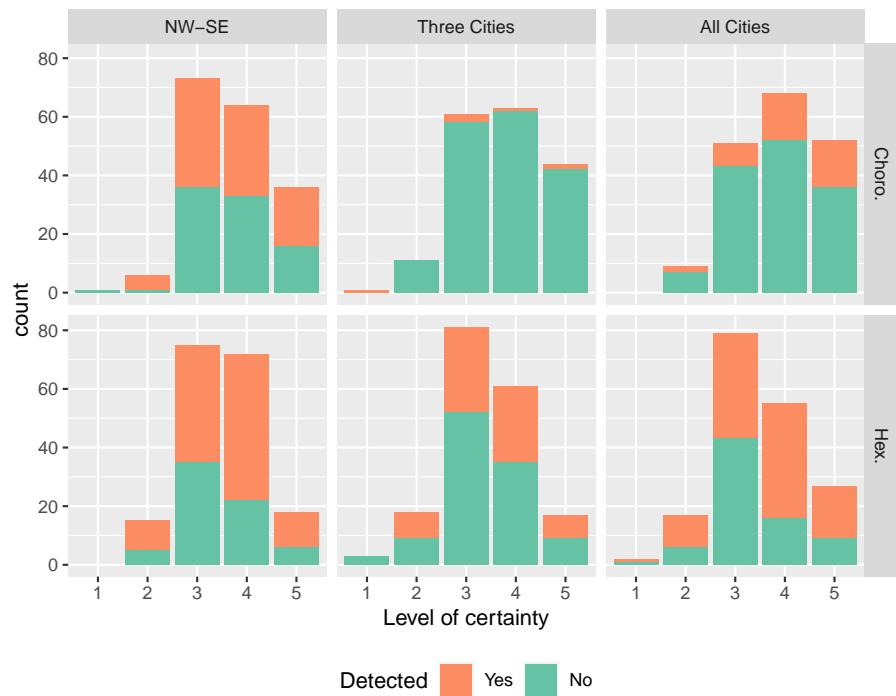


Figure 4.7: The amount of times each level of certainty was chosen by participants when viewing hexagon tile map or choropleth displays. Participants were more likely to choose a high certainty when considering a choropleth map, but more likely to be wrong for the All Cities and Three Cities patterns. Participants were less certain in their responses for the hexagon tile map lineups, perhaps reflecting the lack of familiarity.

participants provided each level of certainty. This was separated for each combination of trend models and display type, and colored depending on whether a participant correctly detected the data plot in the lineup. Participants often chose 4 or 5 when viewing the population related trends in the choropelth display, even though they were often incorrect when viewing an All Cities trend and overwhelmingly incorrect for the Three Cities trend. This shows overconfidence in their detection ability when using a choropleth map display. Participants were less likely to be certain when their choice was incorrect and they were viewing a hexagon tile map. For each trend model, participants were more likely to doubt their choice and choose 1 or 2 in the hexagon tile map displays, even though many had made the correct choice.

4.3.5 Reason

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

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

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

Table 4.3: *The amount 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.
NW-SE	No	trend	clusters
	Yes	trend	clusters
Three Cities	No	trend	clusters
	Yes	consistent	clusters
All Cities	No	trend	clusters
	Yes	clusters, consistent	clusters

4.4 Discussion

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

We expected that the choropleth map would be superior for communicating the spatial pattern of geographic distributions. The data suggest that the participants perform slightly better or equally as well for each replicate in each trend model across the two displays.

Table 4.II shows that the difference in the mean detection rate for the two trend models was 0.10.

The differences seen in Figure 4.5 and Table 4.1 are reflected in the model results. Surprisingly the difference for the geographic distribution was significant at the 0.05 level. It also showed that the hexagon tile map display performs marginally better than the choropleth for all trend models. Unexpectedly the detection rate suffers when using a choropleth map to display population related distributions.

While the significance of the difference in detection was the key focus of this experiment, the secondary focus was the time taken by participants. it was expected that the participants may take longer to consider the hexagon tile map distribution but would be able to detect the data plot in the lineup. The bimodal distributions seen in Figure 4.6 showed very little difference in the mean evaluation times. As the maximum time of all of the distributions approached 60 seconds it cannot be said that the participants' took longer to evaluate the hexagon tile map displays.

The responses to the questions asked of participants included the reason for their choice and the certainty around their choice. Figure 4.3 shows high levels of certainty of 4 and 5 were chosen by participants when looking at the population distributions in a choropleth map display show that they were over confident when attempting to find the real data plot in the choropleth map displays. Participants performed better on the NW-SE distribution shown in the choropleth display and were reasonably confident about their decisions. The high levels of the mid range value of 3 could indicate that the participant did not want to provide a response, as this was the default value. Those who chose level 4 or 5 were equally likely to be correct for the three cities lineups, but more likely to be correct than incorrect for the other two trend models.

The color scaling applied in Three cities and All cities displays resulted in the rural areas of the real data plot appearing more blue or yellow than the other plots in the lineups. Due to the consistent coloring of rural areas in a choropleth display, the choice "All areas have similar colors" was most common reason for a participants choice. The All Cities displays colored the inner-city areas of all capital cities more red, this was observable to

participants and explains the equal choice of the city clusters or rural color consistency. Choosing “Clusters of colour” was expected when participants viewed the hexagon tile map display of the All Cities and Three Cities distributions. It was unexpected that it was also the most common reason for the NW-SE hexagon tile map displays. Due to the spatial covariance introduced in the smoothing, groups of similarly colored hexagons were present in all of the hexagon tile map displays. All Cities and Three Cities distributions of real data trends had distinctly different patterns or red inner-city areas, while some of the plots in each lineup may have shared similar features.

4.5 Conclusion

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

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

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

4.6 Supplementary material

The appendix [A](#) contains:

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

Chapter 5

Discussion and Conclusion

Cancer Atlases are used to develop hypotheses about spatial distributions of cancer statistics (Bell et al., 2006). However, the use of the choropleth map may lead to misinterpretation of the overall distribution. This is because of the overemphasis on the large geographic areas, and the lack of visibility for the small inner-city communities (Dorling, 2011).

The first aim of this thesis was to present an alternative visualisation method for spatial data. This thesis has provided a new algorithm to present spatial distributions of disease data, and includes an R code (R Core Team, 2019) implementation. The spatial data sets with population related distributions will be effectively communicated by this display. The hexagon tile map display will represent each area equally on the map space to effectively convey the spatial distribution. This does not require manual creation of layouts, and the displays are reusable for any data set that uses the same set of geographic units.

The second aim was to test effectiveness of the hexagon tile map relative to the choropleth map. It was expected that the familiar map base would have advantages in communicating geospatial distributions. However, when tested using the lineup protocol the results for the choropleth map and hexagon tile map were extremely close. The hexagon tile map was much more effective for communicating the population related distributions.

To achieve the third aim of this thesis, the hexagon tile map output allows for animation between the choropleth display and the hexagon tile map display. This was achieved through an implementation as part of the `sugarbag` (Kobakian and Cook, 2019) package for R (R Core Team, 2019).

The hexagon tile map visualisation method solves the misrepresentation problem of choropleth display. Especially for geographic data sets that contain disparity between land size and population density. This algorithm is accessible to all R users, it can be applied to any set of areas in an `sf` (Pebesma, 2018) object through the set of simple functions. The tessellation employed in the hexagon tile map algorithm maintains connectedness between neighbouring areas, this draws inspiration from contiguous cartograms (Min Ouyang and Revesz, 2000), rectangular cartograms (Raisz, 1963) and Dorling's circular cartograms (Dorling, 2011). However, the hexagon tile map algorithm does not employ the gravitation pull mathematics that is used to create contiguous cartograms. It also does not iterate on the placement of hexagons. The choice of a consistent shape to be used for all areas draws from rectangular and Dorling cartograms. This encourages map readers to focus on the similarities or difference in the colour between geographic neighbours, and does not distract them with unfamiliar boundaries produced during a contiguous cartogram transformation.

The effectiveness of the hexagon tile map has been tested by the visual inference study. It showed that participants could recognise the data display in the set of null distributions more frequently when viewing a hexagon tile map display. The choropleth map display is still effective for distributions that are directly related to the geography, such as the North-West to South-East distribution used in the study. This has expanded the applications of visual inference studies in a spatial data context by contributing a specific example of testing new graphical displays for geospatial data.

The tile map allocation provided by the algorithm can be used to create animations between a choropleth and hexagon tile map display. Linking the familiar geography to the effective display for understanding the distribution across many heterogeneous geographic regions. Many interactive tools are included in current cancer atlases, these additions allow user driven exploration, but do not guarantee that the spatial distribution

across the geographic space is digested accurately. Animating between a choropleth and a hexagon tile map will allow map users to understand how the small communities of a whole country are affected simultaneously. It also teaches map users how to find areas of interest as their attention is drawn to the capital cities, that may not have caught their attention in the display of the choropleth map. When communicating cancer statistics, there should be a balance between providing people with a familiar landscape and ensuring they interpret the spatial distribution correctly. Animations will communicate a specific message through the capture and direction of users' attention.

Future work will include expanding on the criteria used to evaluate the hexagon tile maps produced by the algorithm. The methods to evaluate the alternative displays have not been thoroughly explored in this thesis, but could be included as functions with the R implementation. This framework will be used to create relevant tests that contrast the use of the map area, and changes in the visual when the parameter of the hexagon tile map algorithm are altered.

The current hexagon tile map creates a template map that can be used to visualise any data set that contains the areas used to create the map. There is the possibility of allowing a bivariate display to incorporate uncertainty by using a colour scheme that operates in two directions, as suggested by Lucchesi and Wikle (Lucchesi and C.K., 2017). The animation methods that allow the colours filling the hexagons to flicker to communicate the uncertainty around an estimate could also be employed. With large hexagons, there is a potential to incorporate geofacets (Hafen, 2019) to create a tessellated display of small visualisations for each geographic unit. These displays become increasingly complex if the visualisation becomes more detailed, or the hexagons become smaller.

The animations created of the Australian Statistical Areas at Level 2 highlight just how many SA2 areas are hidden due to their size in the choropleth display. This animation could be included future iterations of the Australian Cancer Atlas to improve the communication of the spatial distributions of the burden of cancer on Australian communities.

In summary, this work has contributed a new alternative visualisation method to highlight the communities in spatial data sets. This is valuable as the spatial distributions of cancer

burden for different types of cancers largely relates to the population rather than the geography. It also contributes an algorithm to produce these displays for any set of spatial polygons. Through this thesis an open-source R package implementation has been included on the CRAN package repository, with associated examples and documentation for use by any R user. This work has also contributed to the literature of visual inference studies, by using the “lineup” protocol developed by Buja et al. and used by Wickham et al. (2010), and Hofmann et al. (2012). To communicate human related spatial patterns of disease, map creators should consider the use of alternative displays. The hexagon tile map display has proven effective in this thesis for communicating spatial distributions in sets of heterogeneous geographic units. This thesis provides a practical guide for map creators to communicate spatial displays of cancer data in Australia.

Bibliography

- Allaire, J, Y Xie, J McPherson, J Luraschi, K Ushey, A Atkins, H Wickham, J Cheng, W Chang, and R Iannone (2019a). *rmarkdown: Dynamic Documents for R*. R package version 1.15. <https://github.com/rstudio/rmarkdown>.
- Allaire, J, Y Xie, R Foundation, H Wickham, Journal of Statistical Software, R Vaidyanathan, Association for Computing Machinery, C Boettiger, Elsevier, K Broman, K Mueller, B Quast, R Pruij, B Marwick, C Wickham, O Keyes, M Yu, D Emaasit, T Onkelinx, A Gasparini, MA Desautels, D Leutnant, MDPI, Taylor and Francis, O Ögreden, D Hance, D Nüst, P Uvesten, E Campitelli, J Muschelli, ZN Kamvar, N Ross, and R Cannoodt (2019b). *rticles: Article Formats for R Markdown*. R package version 0.13. <https://CRAN.R-project.org/package=rticles>.
- Arnold, JB (2019). *ggthemes: Extra Themes, Scales and Geoms for 'ggplot2'*. R package version 4.2.0. <https://CRAN.R-project.org/package=ggthemes>.
- Bache, SM and H Wickham (2014). *magrittr: A Forward-Pipe Operator for R*. R package version 1.5. <https://CRAN.R-project.org/package=magrittr>.
- Bates, D, M Maechler, B Bolker, and S Walker (2019). *lme4: Linear Mixed-Effects Models using 'Eigen' and S4*. R package version 1.1-21. <https://CRAN.R-project.org/package=lme4>.
- Bell, BS, RE Hoskins, LW Pickle, and D Wartenberg (2006). Current practices in spatial analysis of cancer data: mapping health statistics to inform policymakers and the public. *International Journal of Health Geographics* 5, 49.
- Berry, BJL, RL Morrill, and WR Tobler (1964). Geographic Ordering of Information: New Opportunities. *The Professional Geographer* 16 (4), 39–44.

- Bivand, R, J Nowosad, and R Lovelace (2019). *spData: Datasets for Spatial Analysis*. R package version 0.3.0. <https://CRAN.R-project.org/package=spData>.
- Bowel Cancer Australia (2016). *Bowel Cancer Australia Atlas*. <http://www.bowelcanceratlas.org/> (visited on 09/26/2019).
- Brewster, MB and SV Subramanian (2010). Cartographic Insights into the Burden of Mortality in the United Kingdom: A Review of 'The Grim Reaper's Road Map'. *International Journal of Epidemiology* **39**(4), 1120–1122.
- Bryan, J and J Zhao (2018). *googlesheets: Manage Google Spreadsheets from R*. R package version 0.3.0. <https://CRAN.R-project.org/package=googlesheets>.
- Buja, A, D Cook, H Hofmann, M Lawrence, EK Lee, DF Swayne, and H Wickham (2009). Statistical Inference for Exploratory Data Analysis and Model Diagnostics. *Philosophical Transactions of the Royal Society, A (Invited)* **367**. doi: [10.1098/rsta.2009.0120](https://doi.org/10.1098/rsta.2009.0120), 4361–4383.
- Burbank, F (1971). *Patterns in Cancer Mortality in the United States 1950-67*. NCI, Washington DC: National Cancer Institute Monograph Vol. 33.
- Cancer Council Queensland, Queensland University of Technology, and Cooperative Research Centre for Spatial Information (2018). *Australian Cancer Atlas*. Version 09-2018. <https://atlas.cancer.org.au>.
- Cano, RG, K Buchin, T Castermans, A Pieterse, W Sonke, and B Speckmann (2015). Mosaic Drawings and Cartograms. In: *Computer Graphics Forum*. Vol. 34. 3. Wiley Online Library, pp.361–370.
- Carr, DB, JF Wallin, and DA Carr (2000). Two New Templates for Epidemiology Applications: Linked Micromap Plots and Conditioned Choropleth Maps. *Statistics in Medicine* **19**, 2521–2538.
- Cliburn, DC, JJ Feddema, JR Miller, and TA Slocum (2002). Design and Evaluation of a Decision Support System in a Water Balance Application. *Computers and Graphics* **26**(6), 931–949.
- Cook, D, A Ebert, J Forbes, H Hofmann, R Hyndman, T Lumley, B Marwick, C Sievert, M Sun, D Talagala, N Tierney, N Tomasetti, E Wang, and F Zhou (2019). *echedna: Exploring Election and Census Highly Informative Data Nationally for Australia*. R package version 1.4.0. <https://CRAN.R-project.org/package=echedna>.

- d'Onofrio, A, C Mazzetta, C Robertson, M Smans, P Boyle, and M Boniol (2016). Maps and Atlases of Cancer Mortality: A Review of a Useful Tool to Trigger New Questions. *Ecancermedicalscience* **10**, 670–670.
- Dang, G, C North, and B Schneiderman (2001). Dynamic Queries and Brushing on Choropleth Maps. In: *Proceedings Fifth International Conference on Information Visualisation*, pp.757–764.
- Dent, BD (1972). A Note on the Importance of Shape in Cartogram Communication. *Journal of Geography* **71**(7), 393–401.
- Dorling, D (2011). "Area Cartograms: Their Use and Creation". In: *Concepts and Techniques in Modern Geography (CATMOG)*. Vol. 59, pp. 252–260.
- Dougenik, JA, NR Chrisman, and DR Niemeyer (1985). An Algorithm to Construct Continuous Area Cartograms. *The Professional Geographer* **37**(1), 75–81. eprint: <https://doi.org/10.1111/j.0033-0124.1985.00075.x>.
- El Pais (2014). *Map of Cancer Mortality Rates in Spain*. http://elpais.com/elpais/2014/10/06/media/1412612722_141933.html (visited on 09/26/2019).
- Emperial College London - Small Area Health Statistics Unit (2010). *The Environmental and Health Atlas of England and Wales: National male lung cancer rate*. <http://www.envhealthatlas.co.uk/eha/Breast/> (visited on 09/26/2019).
- Exeter, DJ (2016). Spatial Epidemiology. *International Encyclopedia of Geography: People, the Earth, Environment and Technology: People, the Earth, Environment and Technology*, 1–4.
- Ferlay J, Ervik M, Lam F, Colombet M, Mery L, Piñeros M, Znaor A, Soerjomataram I, Bray F (2018). *Global Cancer Observatory: Cancer Today*. <https://gco.iarc.fr/today>.
- Figure Eight Inc (2019). *The Essential High-Quality Data Annotation Platform*. <https://www.figure-eight.com/>.
- Gamio, L and C D. (2016). *Poll: Redrawing the electoral map*. <https://www.washingtonpost.com/graphics/politics/2016-election/50-state-poll/>.
- Goodchild, M, B Buttenfield, and J Wood (1994). On Introduction to Visualizing Data Validity. *Visualization in geographical information systems*, 141–149.
- Griffin, T (1980). Cartographic Transformation of the Thematic Map Base. *Cartography* **11**(3), 163–174. eprint: <https://doi.org/10.1080/00690805.1980.10438102>.

- Grolemund, G and H Wickham (2011). Dates and Times Made Easy with lubridate. *Journal of Statistical Software* **40**(3), 1–25.
- Hafen, R (2019). *geofacet: 'ggplot2' Faceting Utilities for Geographical Data*. R package version 0.1.10. <https://CRAN.R-project.org/package=geofacet>.
- Harrower, M and CA Brewer (2003). ColorBrewer.org: An Online Tool for Selecting Colour Schemes for Maps. *The Cartographic Journal* **40** (1), 27–37.
- Hofmann, H, L Follett, M Majumder, and D Cook (2012). Graphical Tests for Power Comparison of Competing Designs. *IEEE Transactions on Visualization and Computer Graphics* **18**(12), 2441–2448.
- Howe, G (1989). “Historical Evolution of Disease Mapping in General and Specifically of Cancer Mapping”. In: *Cancer mapping*. Springer, pp.1–21.
- Jeworutzki, S (2018). *cartogram: Create Cartograms with R*. R package version 0.1.1. <https://CRAN.R-project.org/package=cartogram>.
- Kanjana, J and D Mehta (2016). *Who will win the presidency?* <https://fivethirtyeight.com/2016-election-forecast/>.
- Keim, D, S North, C Panse, and J Schneidewind (2002). Efficient Cartogram Generation: A Comparison. eng. In: *IEEE Symposium on Information Visualization, 2002. INFOVIS 2002*. Vol. 2002. IEEE, pp.33–36.
- Kobakian, S and D Cook (2019). *sugarbag: Create Tessellated Hexagon Maps*. R package version 0.1.0. <https://CRAN.R-project.org/package=sugarbag>.
- Kobakian, S and M O’Hara-Wild (2018). *taipan: Tool for Annotating Images in Preparation for Analysis*. R package version 0.1.2. <https://CRAN.R-project.org/package=taipan>.
- Kocmoud, C and D House (1998). A Constraint-based Approach to Constructing Continuous Cartograms. In: *Proc. Symp. Spatial Data Handling*, pp.236–246.
- Kraak, MJ (2017). “Cartographic Design”. English. In: *The International Encyclopedia of Geography: People, the Earth, Environment, and Technology*. United States: Wiley, pp. 1–16.
- Kreveld, M van and B Speckmann (2007). On rectangular cartograms. *Computational Geometry* **37**(3). Special Issue on the 20th European Workshop on Computational Geometry, 175–187.
- Kronenfeld, BJ and DWS Wong (2017). Visualizing Statistical Significance of Disease Clusters Using Cartograms. *International Journal of Health Geographics* **16**(1), 19.

- Levison, ME and W Haddon Jr (1965). The Area Adjusted Map. An Epidemiologic Device. *Public Health Reports* **80**, 55–59.
- Lucchesi, L and W C.K. (2017). Visualizing Uncertainty in Areal Data with Bivariate Choropleth Maps, Map Pixelation and Glyph Rotation. *Stat.*
- MacEachren, AM (1992). Visualizing Uncertain Information. *Cartographic Perspectives* (13), 10–19.
- Mackey, W. F. (2019). *absmapsdata: A catalogue of ready-to-use ASGS mapping data*. R package version 0.1.1.
- Madsen, R (2019). *Programming Design Systems*. <https://programmingdesignsystems.com/>.
- Matheron, G (1963). Principles of geostatistics. *Economic Geology* **58**, 1246–1266.
- McGranaghan, M (1993). A Cartographic View of Spatial Data Quality. *Cartographica: The International Journal for Geographic Information and Geovisualization* **30**(2-3), 8–19.
- Min Ouyang and P Revesz (2000). Algorithms for Cartogram Animation. In: *Proceedings 2000 International Database Engineering and Applications Symposium (Cat. No.PR00789)*, pp.231–235.
- Monmonier, M (2018). *How to Lie with Maps (Third Edition)*. University of Chicago Press. <https://books.google.com.au/books?id=MwdRDwAAQBAJ>.
- Monmonier, M (2005). Cartography: Distortions, World-views and Creative Solutions. *Progress in Human Geography* **29**(2), 217–224. eprint: <https://doi.org/10.1191/0309132505ph540pr>.
- Montanaro, D (2016). *NPR Battleground Map: Hillary Clinton Is Winning — And It's Not Close*. <https://www.npr.org/2016/10/18/498406765/npr-battleground-map-hillary-clinton-is-winning-and-its-not-close>.
- Moore, DA and TE Carpenter (1999). Spatial Analytical Methods and Geographic Information Systems: Use in Health Research and Epidemiology. *Epidemiologic Reviews* **21**(2), 143–161. eprint: <http://oup.prod.sis.lan/epirev/article-pdf/21/2/143/6727658/21-2-143.pdf>.
- Müller, K and H Wickham (2019). *tibble: Simple Data Frames*. R package version 2.1.3. <https://CRAN.R-project.org/package=tibble>.

- Neuwirth, E (2014). *RColorBrewer: ColorBrewer Palettes*. R package version 1.1-2. <https://CRAN.R-project.org/package=RColorBrewer>.
- Northern Ireland Cancer Registry (2011). *All-Ireland Cancer Atlas (1995-2007)*. <http://www.ncri.ie/publications/cancer-atlases>.
- Nusrat, S and SG Kobourov (2016). The State of the Art in Cartograms. *Computer Graphics Forum* 35(3), 619–642. arXiv: [1605.08485](https://arxiv.org/abs/1605.08485).
- Olson, JM (1976). Noncontiguous Area Cartograms. *The Professional Geographer* 28(4), 371–380. eprint: <https://doi.org/10.1111/j.0033-0124.1976.00371.x>.
- Pebesma, E (2018). Simple Features for R: Standardized Support for Spatial Vector Data. *The R Journal* 10(1), 439–446.
- Pebesma, E and B Graeler (2019). *gstat: Spatial and Spatio-Temporal Geostatistical Modelling, Prediction and Simulation*. R package version 2.0-3. <https://CRAN.R-project.org/package=gstat>.
- Pedersen, TL (2018). *The Grammar of Animation*. Keynote talk at useR!2018. <https://youtu.be/21ZWDrTukEs> (visited on 11/16/2018).
- Pedersen, TL and D Robinson (2019). *ggridge: A Grammar of Animated Graphics*. R package version 1.0.3. <https://CRAN.R-project.org/package=ggridge>.
- Pediatric Oncology Group of Ontario (2015). *Incidence Rate of Childhood Cancers, Atlas of Childhood Cancer in Ontario (1985-2004)*. https://www.pogo.ca/wp-content/uploads/2015/02/POGO_CC-Atlas-3-Incidence_Feb-2015.pdf (visited on 09/26/2019).
- Perin, C (2014). “Direct Manipulation for Information Visualization”. Theses. Université Paris Sud - Paris XI. <https://hal.inria.fr/tel-01096366>.
- Queensland Cancer Registry (2011). *The Atlas of Cancer in Queensland (1998 - 2007)*. <https://cancerqld.org.au/research/queensland-cancer-statistics/queensland-cancer-atlas/> (visited on 09/26/2019).
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>.
- Raisz, E (1963). Rectangular Statistical Cartograms of the World. *Journal of Geography* 35 (1), 8–10.

- Roberts, J (2019). "Communication of Statistical Uncertainty to Non-expert Audiences". MA thesis. Queensland University of Technology. <https://eprints.qut.edu.au/130786/>.
- Robinson, D and A Hayes (2019). *broom: Convert Statistical Analysis Objects into Tidy Tibbles.* R package version 0.5.3. <https://CRAN.R-project.org/package=broom>.
- Roy Chowdhury, N (2014). ""Explorations of the lineup protocol for visual inference: application to high dimension, low sample size problems and metrics to assess the quality"". <https://lib.dr.iastate.edu/etd/13988>.
- S., K, D Cook, and J Roberts (2019). Cancer Applications of Choropleth Maps, and the Potential of Cartograms and Alternative Map Displays. *Manuscript submitted for publication.*
- Skowronnek, A (2016). *Beyond Choropleth Maps – A Review of Techniques to Visualize Quantitative Areal Geodata.* https://alsino.io/static/papers/BeyondChoropleths_AlsinoSkowronnek.pdf.
- Statistics, AB of (2018). [https://www.abs.gov.au/websitedbs/D3310114.nsf/home/Australian+Statistical+Geography+Standard+\(ASGS\)](https://www.abs.gov.au/websitedbs/D3310114.nsf/home/Australian+Statistical+Geography+Standard+(ASGS)).
- Tobler, W (2004). Thirty Five Years of Computer Cartograms. *Annals of the Association of American Geographers* **94**(1), 58–73.
- Tufte, ER (2001). *The visual display of quantitative information.* Vol. 2. Graphics press Cheshire, CT.
- Tufte, ER (1990). *Envisioning Information.* Graphics Press.
- U.S. Department of Health and Human Services, Centers for Disease Control and Prevention and National Cancer Institute - Cancer Statistics Working Group (2019). *U.S. Cancer Statistics Data Visualizations Tool (data 1999-2016).* <http://www.cdc.gov/cancer/dataviz> (visited on 09/26/2019).
- Urbanek, S (2013). *png: Read and write PNG images.* R package version 0.1-7. <https://CRAN.R-project.org/package=png>.
- van der Walt, S. and Smith, N (2015). *mpl colormaps.* <https://bids.github.io/colormap/>.

- Van der Wel, FJ, RM Hootsmans, and F Ormeling (1994). "Visualization of Data Quality". In: *Modern Cartography Series*. Vol. 2. Elsevier, pp.313–331. <https://doi.org/10.1016/B978-0-08-042415-6.50023-5>.
- W., PL, DB Carr, and JB Pearson (2015). micromapST: Exploring and Communicating Geospatial Patterns in US State Data. *Journal of Statistical Software* **63**(3), 1–25.
- Walter, SD (2001). *Disease Mapping: A Historical Perspective*. Oxford University Press.
- Wickham, H, D Cook, H Hofmann, and A Buja (2010). Graphical inference for infovis. *IEEE Transactions on Visualization and Computer Graphics* **16**(6), 973–979.
- Wickham, H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, H, NR Chowdhury, D Cook, and H Hofmann (2018). *nullabor: Tools for Graphical Inference*. R package version 0.3.5. <https://CRAN.R-project.org/package=nullabor>.
- Wickham, H, R François, L Henry, and K Müller (2019). *dplyr: A Grammar of Data Manipulation*. R package version 0.8.3. <https://CRAN.R-project.org/package=dplyr>.
- Wilke, CO (2019). *cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'*. R package version 1.0.0. <https://CRAN.R-project.org/package=cowplot>.
- World Health Organization's International Agency for Research on Cancer (2018). *Globocan 2018: Estimated Cancer Incidence, Mortality and Prevalence*. <http://globocan.iarc.fr/Pages/Map.aspx> (visited on 09/26/2019).
- Xie, Y (2014). "knitr: A Comprehensive Tool for Reproducible Research in R". In: *Implementing Reproducible Computational Research*. Ed. by V Stodden, F Leisch, and RD Peng. ISBN 978-1466561595. Chapman and Hall/CRC. <http://www.crcpress.com/product/isbn/9781466561595>.
- Zhu, H (2019). *kableExtra: Construct Complex Table with 'kable' and Pipe Syntax*. R package version 1.1.0. <https://CRAN.R-project.org/package=kableExtra>.
- Zitner, A, R Yeip, and J Wolfe (2016). *Draw the 2016 Electoral College Map*. (<http://graphics.wsj.com/elections/2016/2016-electoral-college-map-predictions/>).

Appendix A

Appendix

A.1 Overall Performance

The detection rate is considered for each lineup. The detection rates for group A were less varied than the detection rates for the lineups seen by group B. Figure A.1 shows the distribution using a boxplot. This shows the median value for detection rate was extremely similar.

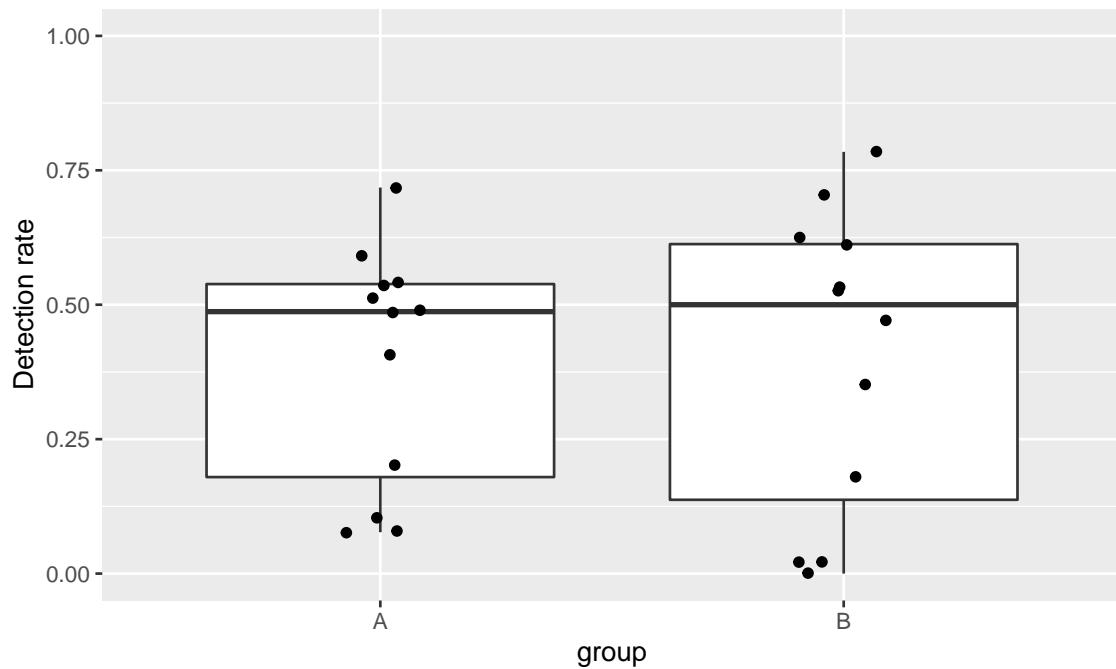


Figure A.1: Boxplots of the distribution of detection rates for each line up, separated by group.

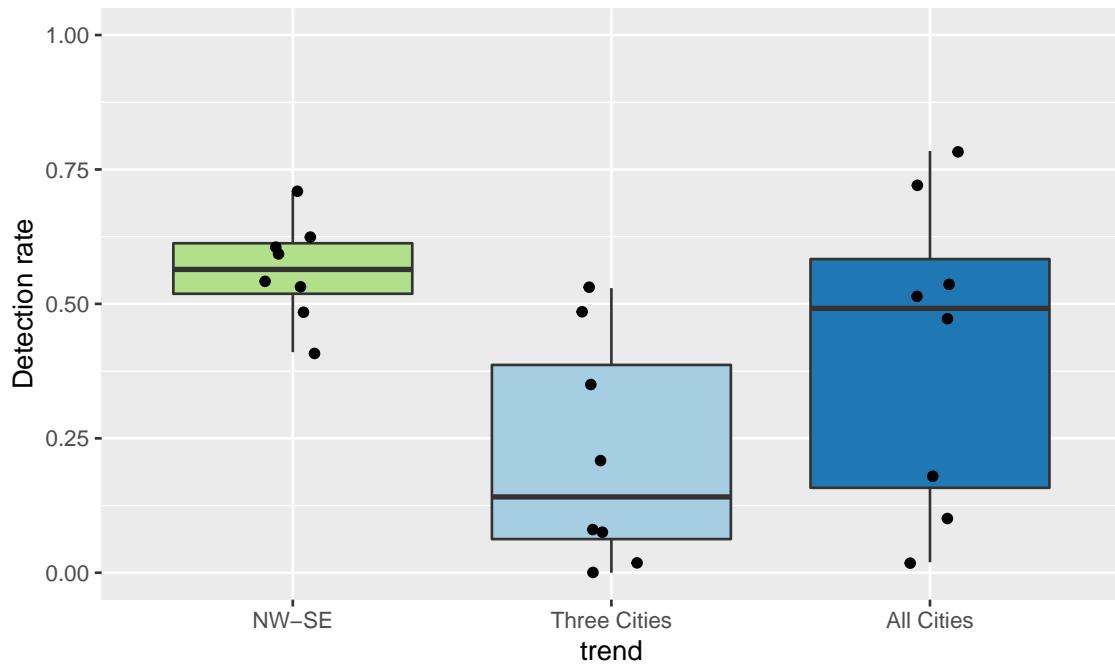


Figure A.2: Boxplots of the distribution of detection rates for each line up, separated by trend model.

The overall detection rate is considered for each trend model in Figure A.2. The detection rates for the NW-SE trend model was less varied than the detection rates for the Three Cities and All Cities trends. Figure A.1 shows the distribution using a boxplot. This shows the distributions of the rates do not overlap for NW-SE and Three Cities trends, the Three cities range was larger, but the median was much higher for the NW-SE trend. The All Cities trend model distribution overlaps with the NW-SE and All Cities trends.

The boxplots in Figure A.3 contrast the distribution of the detection rates for each type of display. The detection rates across the lineups was less varied for the hexagon display. There was a large difference in the medians for the types of displays. Without considering the relationship for each lineup, the hexagon lineup display allowed the participants to achieve higher detection rates.

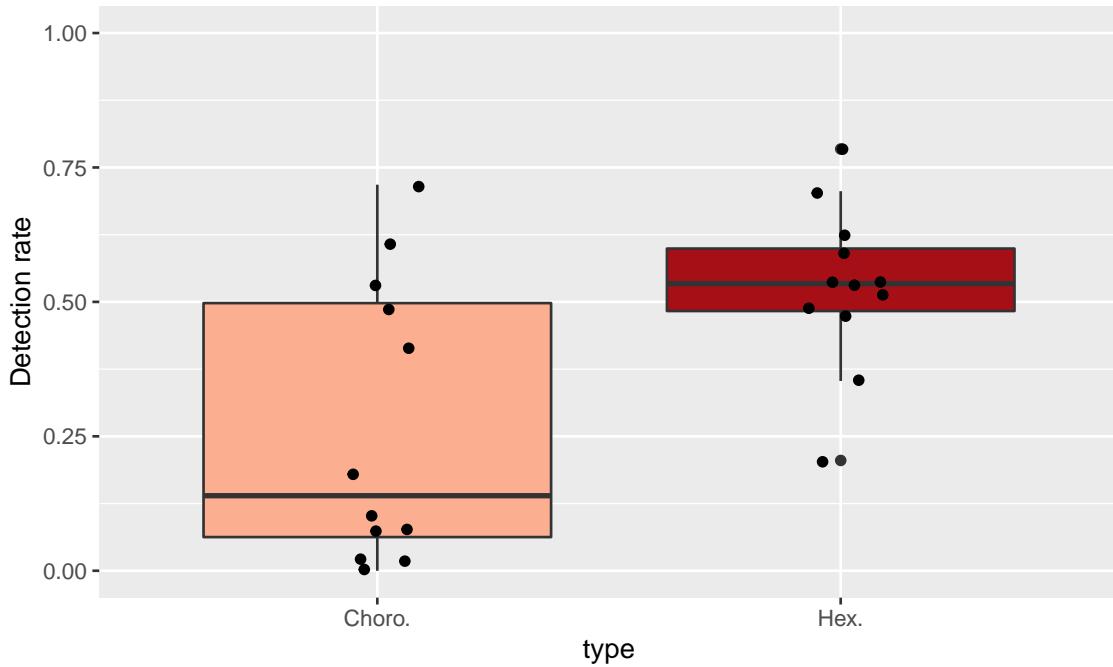


Figure A.3: Boxplots of the distribution of detection rates for each line up, separated by type of display.

A.2 Lineups

The choropleth map lineups were created using the Australian Statistical Areas at Level 3. Each lineup had twelve map displays, this gave participants the choice of any plot, and the choice to not provide a response. This non-response is indicated by 0. The choices made by participants are displayed in Figure A.4. The height of each orange lollipop indicates the proportion of participants that selected the map display of real data, they represent the correct choices. The green lollipops show the proportion of participants that selected the incorrect displays in each lineup.

The proportion of choices are also presented separately for each trend model in Table A.1, Table A.2, and Table A.3. The correct map display in lineups with a North West to South East trend was chosen correctly with much greater frequency. In the lineups of All Cities displays, participants were misled by the choropleth display, but not the hexagon display for all except (2). All of Three Cities displays, except (4), were detected in the hexagon display. All except one lineup had at least one participant select the correct map in the lineup as shown in Figure A.4.

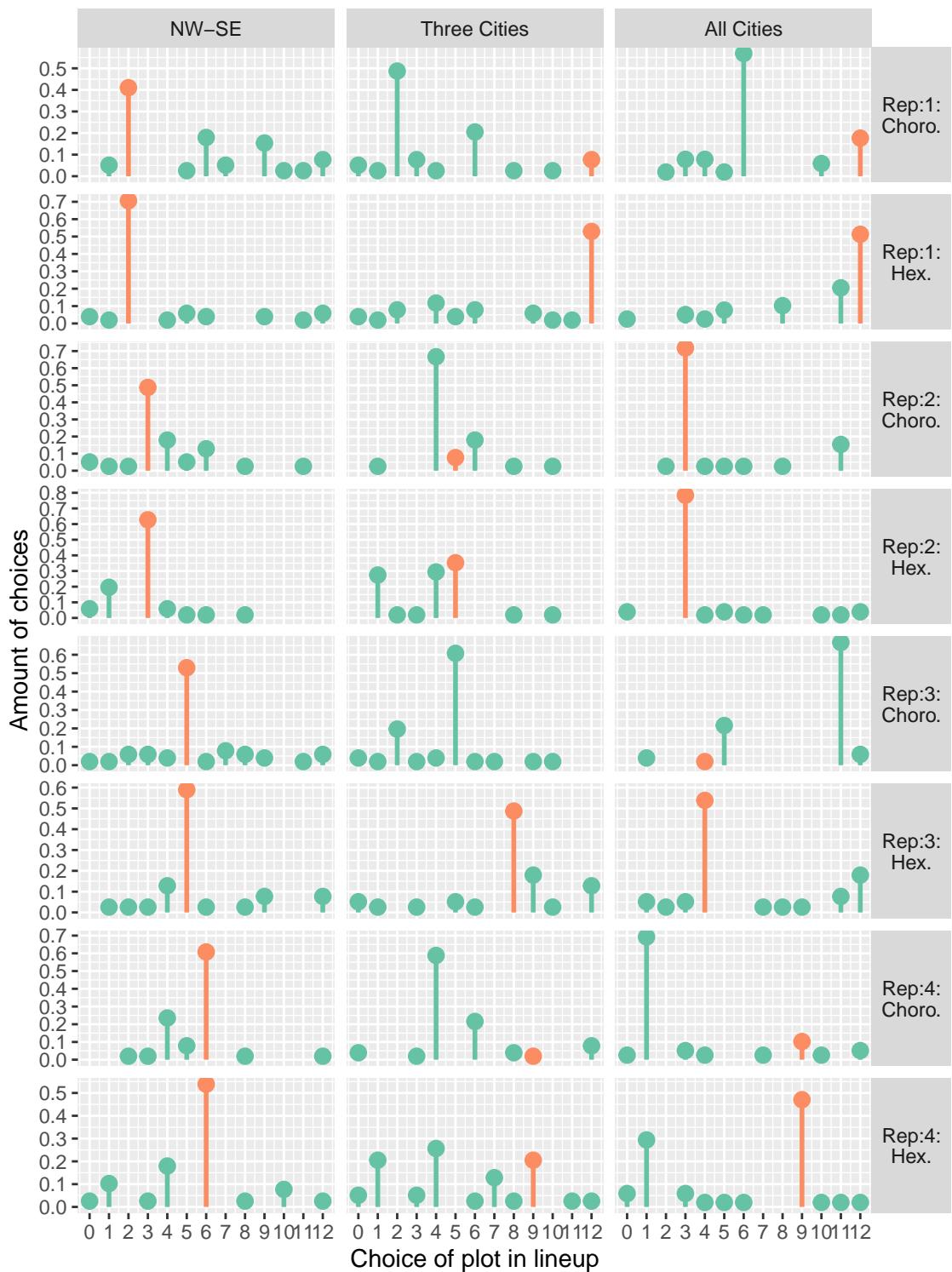


Figure A.4: Pin plots of the proportion of choices made for each lineup location. Each facet is associated with one lineup, the height of the points show the proportion of the participants that made each choice when considering each lineup. The points coloured orange show the map which contained a trend model, these are the correct choices. The numbers differentiate the replicates of each trend model and type of map display. Participants were able to select 0 to indicate they did not want to choose a map.

Table A.1: *The proportion of participants who selected each of the twelve map choices in each lineup for NW-SE displays.*

Rep	Type	0	1	2	3	4	5	6	7	8	9	10	11	12
1	Choro.	0.00	0.05	0.41	0.00	0.00	0.03	0.18	0.05	0.00	0.15	0.03	0.03	0.08
	Hex.	0.04	0.02	0.71	0.00	0.02	0.06	0.04	0.00	0.00	0.04	0.00	0.02	0.06
2	Choro.	0.05	0.03	0.03	0.49	0.18	0.05	0.13	0.00	0.03	0.00	0.00	0.03	0.00
	Hex.	0.06	0.20	0.00	0.63	0.06	0.02	0.02	0.00	0.02	0.00	0.00	0.00	0.00
3	Choro.	0.02	0.02	0.06	0.06	0.04	0.53	0.02	0.08	0.06	0.04	0.00	0.02	0.06
	Hex.	0.00	0.03	0.03	0.03	0.13	0.59	0.03	0.00	0.03	0.08	0.00	0.00	0.08
4	Choro.	0.00	0.00	0.02	0.02	0.24	0.08	0.61	0.00	0.02	0.00	0.00	0.00	0.02
	Hex.	0.03	0.10	0.00	0.03	0.18	0.00	0.54	0.00	0.03	0.00	0.08	0.00	0.03

Table A.2: *The proportion of participants who selected each of the twelve map choices in each lineup for Three Cities displays.*

Rep	Type	0	1	2	3	4	5	6	7	8	9	10	11	12
1	Choro.	0.05	0.03	0.49	0.08	0.03	0.00	0.21	0.00	0.03	0.00	0.03	0.00	0.08
	Hex.	0.04	0.02	0.08	0.00	0.12	0.04	0.08	0.00	0.00	0.06	0.02	0.02	0.53
2	Choro.	0.00	0.03	0.00	0.00	0.67	0.08	0.18	0.00	0.03	0.00	0.03	0.00	0.00
	Hex.	0.00	0.27	0.02	0.02	0.29	0.35	0.00	0.00	0.02	0.00	0.02	0.00	0.00
3	Choro.	0.04	0.02	0.20	0.02	0.04	0.61	0.02	0.02	0.00	0.02	0.02	0.00	0.00
	Hex.	0.05	0.03	0.00	0.03	0.00	0.05	0.03	0.00	0.49	0.18	0.03	0.00	0.13
4	Choro.	0.04	0.00	0.00	0.02	0.59	0.00	0.22	0.00	0.04	0.02	0.00	0.00	0.08
	Hex.	0.05	0.21	0.00	0.05	0.26	0.00	0.03	0.13	0.03	0.21	0.00	0.03	0.03

Table A.3: *The proportion of participants who selected each of the twelve map choices in each lineup for All Cities displays.*

Rep	Type	0	1	2	3	4	5	6	7	8	9	10	11	12
1	Choro.	0.00	0.00	0.02	0.08	0.08	0.02	0.57	0.00	0.00	0.00	0.06	0.00	0.18
	Hex.	0.03	0.00	0.00	0.05	0.03	0.08	0.00	0.00	0.10	0.00	0.00	0.21	0.51
2	Choro.	0.00	0.00	0.03	0.72	0.03	0.03	0.03	0.00	0.03	0.00	0.00	0.15	0.00
	Hex.	0.04	0.00	0.00	0.78	0.02	0.04	0.02	0.02	0.00	0.00	0.02	0.02	0.04
3	Choro.	0.00	0.04	0.00	0.00	0.02	0.22	0.00	0.00	0.00	0.00	0.00	0.67	0.06
	Hex.	0.00	0.05	0.03	0.05	0.54	0.00	0.00	0.03	0.03	0.03	0.00	0.08	0.18
4	Choro.	0.03	0.69	0.00	0.05	0.03	0.00	0.00	0.03	0.00	0.10	0.03	0.00	0.05
	Hex.	0.06	0.29	0.00	0.06	0.02	0.02	0.02	0.00	0.00	0.47	0.02	0.02	0.02

A.2.A All Cities

Replicate 1

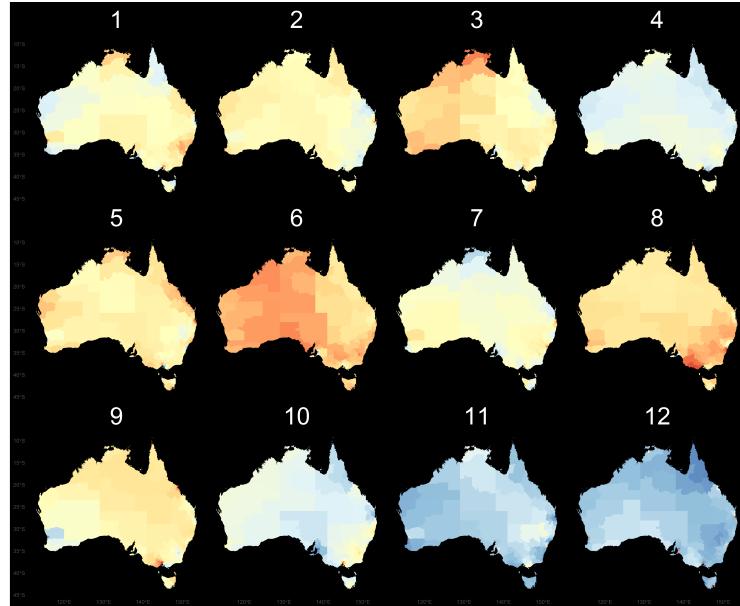


Figure A.5: A choropleth map lineup, location 12 contains a distribution that affects all capital cities.

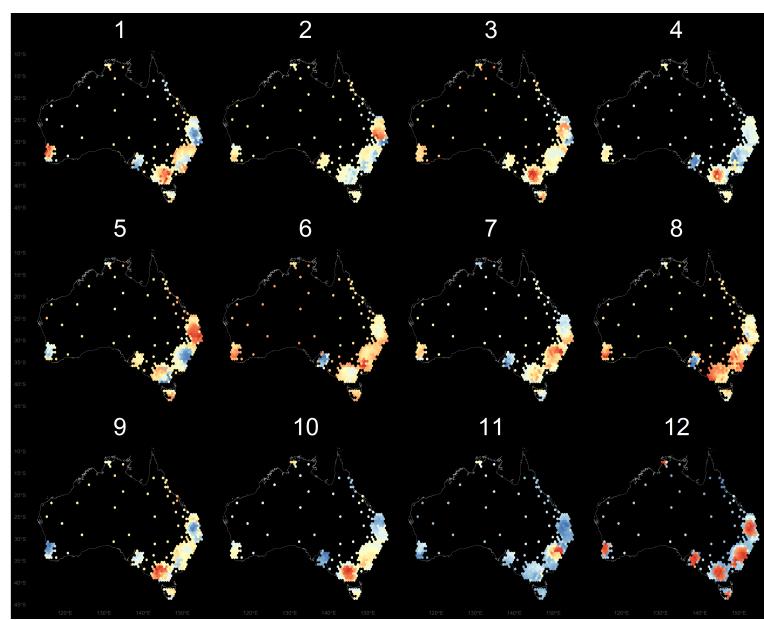


Figure A.6: A hexagon tile map lineup, location 12 contains a distribution that affects all capital cities.

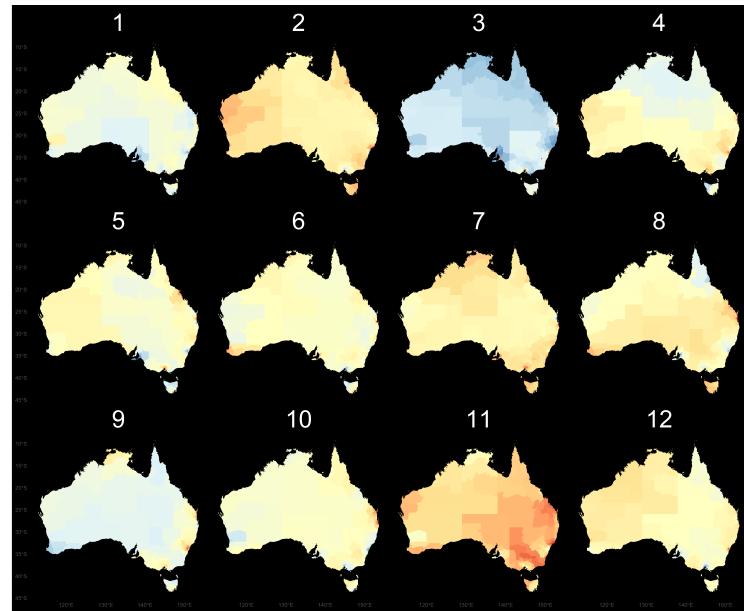
Replicate 2

Figure A.7: A choropleth map lineup, location 3 contains a distribution that affects all capital cities.

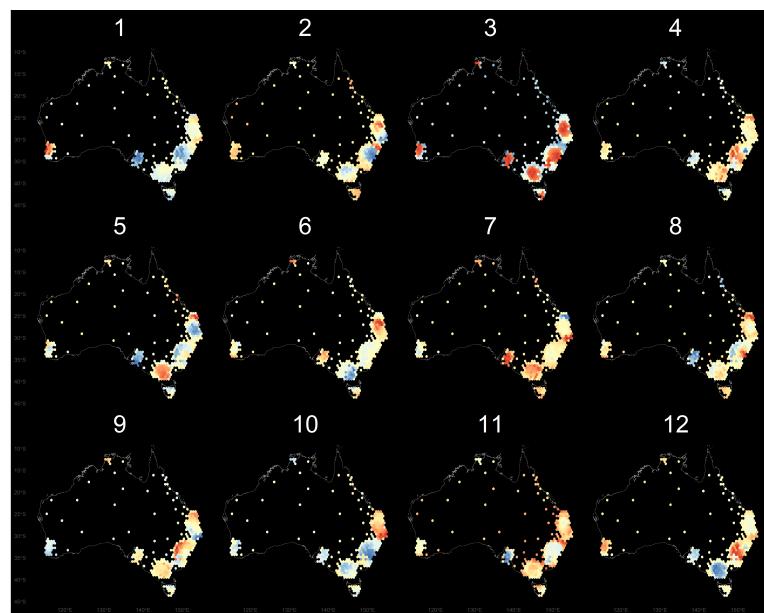


Figure A.8: A hexagon tile map lineup, location 3 contains a distribution that affects all capital cities.

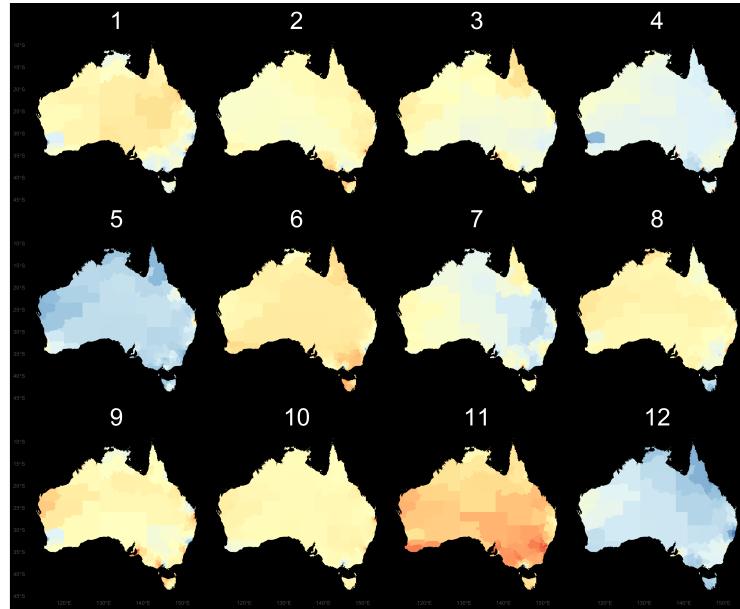
Replicate 3

Figure A.9: A choropleth map lineup, location 4 contains a distribution that affects all capital cities.

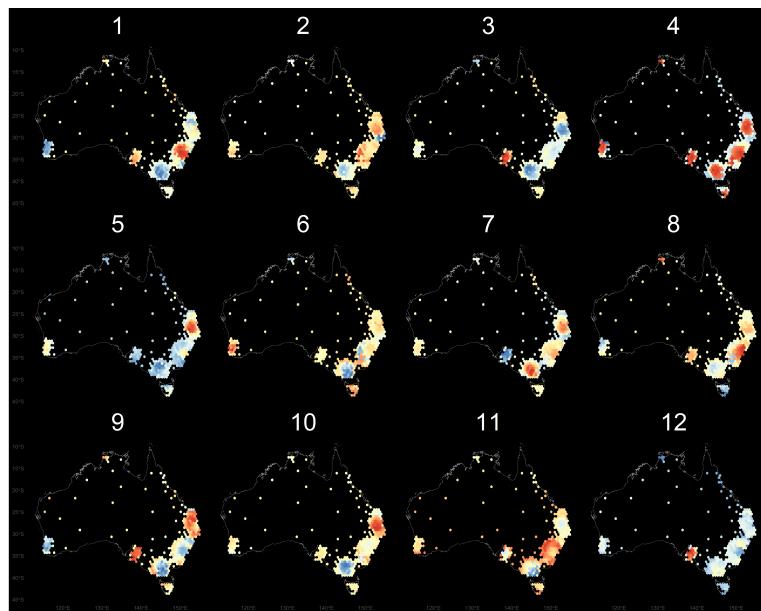


Figure A.10: A hexagon tile map lineup, location 4 contains a distribution that affects all capital cities.

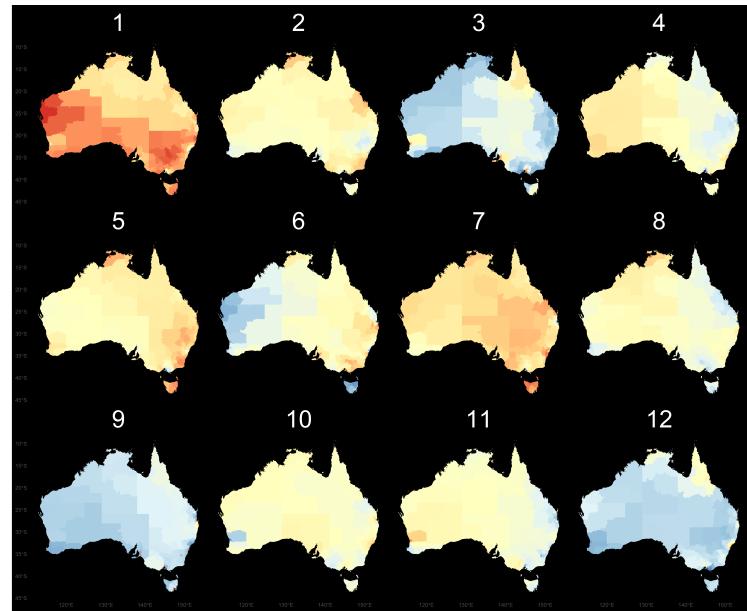
Replicate 4

Figure A.11: A choropleth map lineup, location 9 contains a distribution that affects all capital cities.

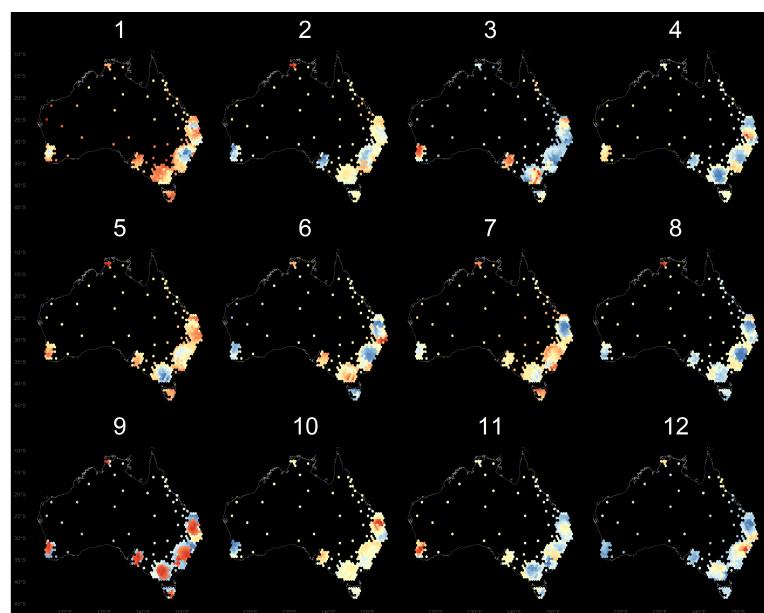


Figure A.12: A hexagon tile map lineup, location 9 contains a distribution that affects all capital cities.

A.2.B Three Cities

Replicate 1

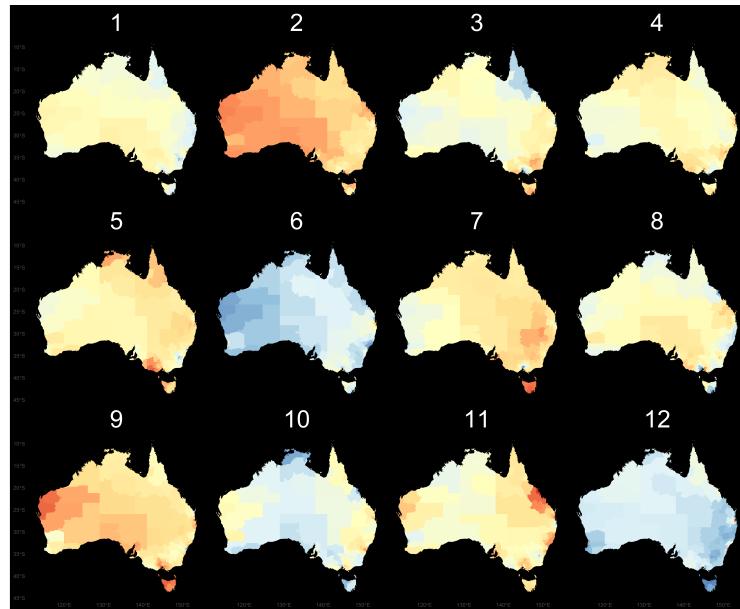


Figure A.13: A choropleth map lineup, location 12 contains a distribution that affects three of the Australian capital cities.

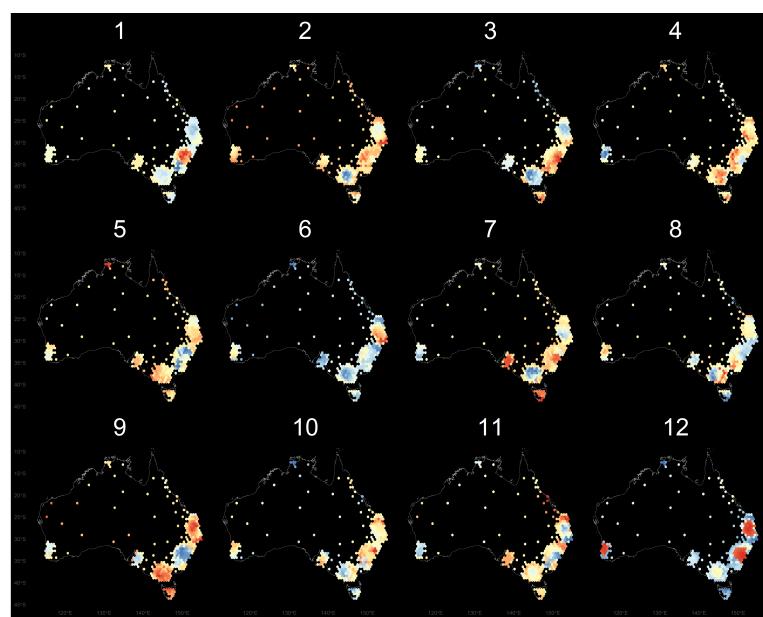


Figure A.14: A hexagon tile map lineup, location 12 contains a distribution that affects three of the Australian capital cities.

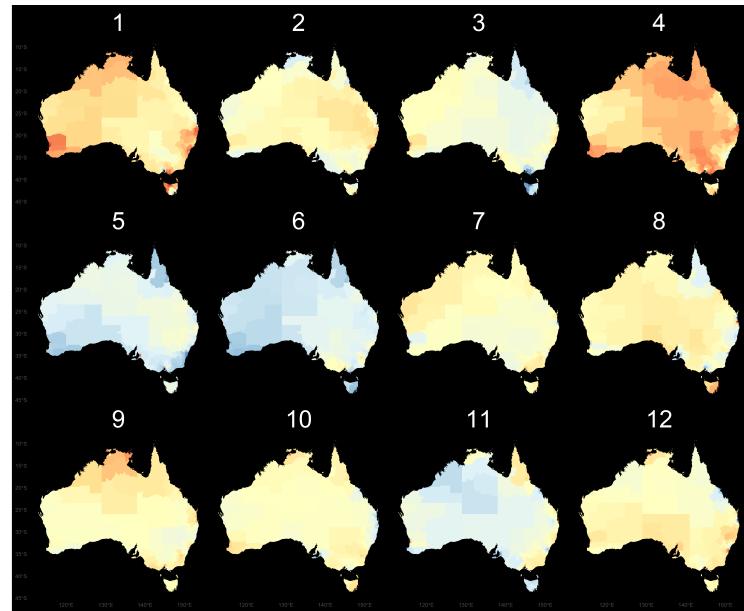
Replicate 2

Figure A.15: A choropleth map lineup, location 3 contains a distribution that affects three of the Australian capital cities.

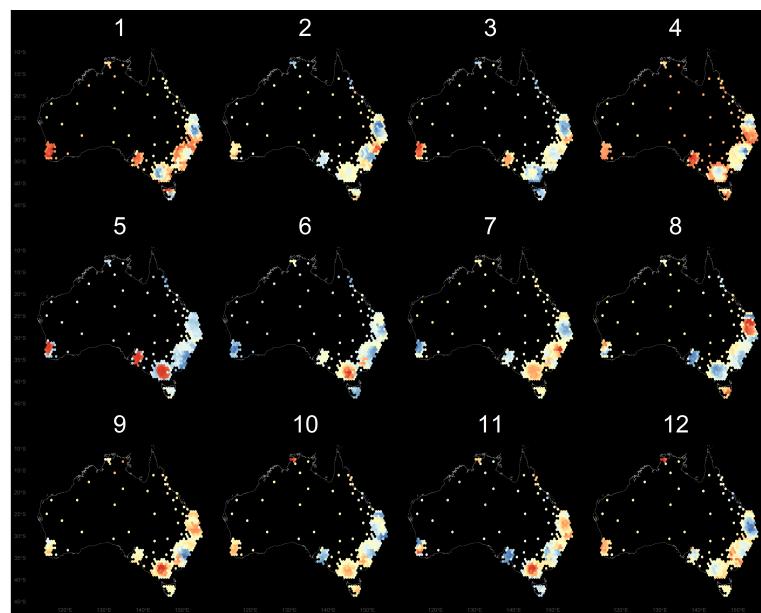


Figure A.16: A hexagon tile map lineup, location 3 contains a distribution that affects three of the Australian capital cities.

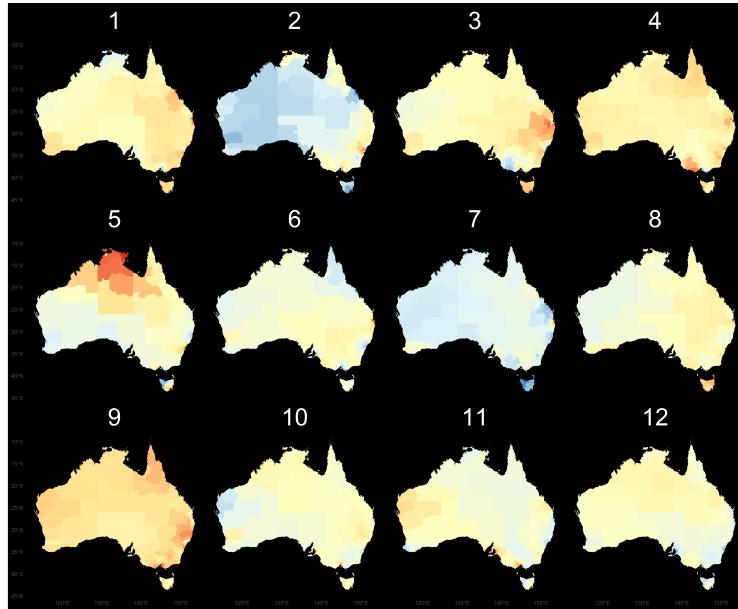
Replicate 3

Figure A.17: A choropleth map lineup, location 4 contains a distribution that affects three of the Australian capital cities.

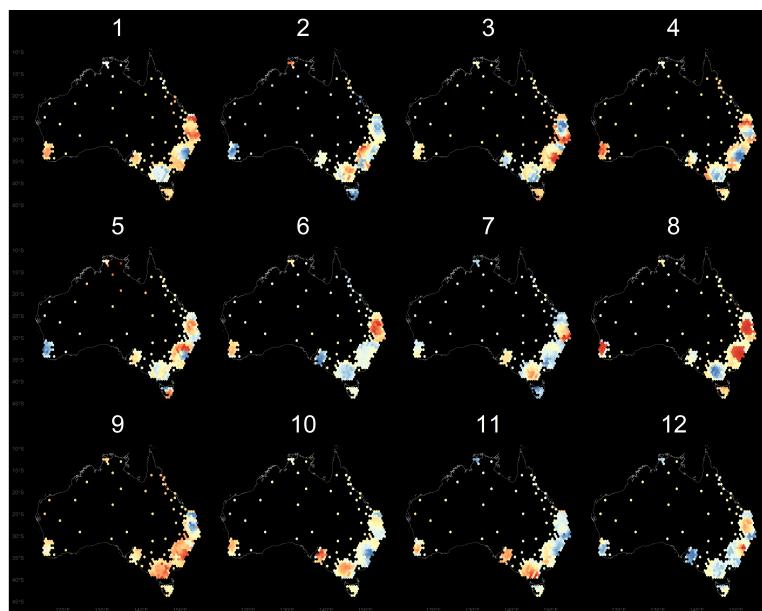


Figure A.18: A hexagon tile map lineup, location 4 contains a distribution that affects three of the Australian capital cities.

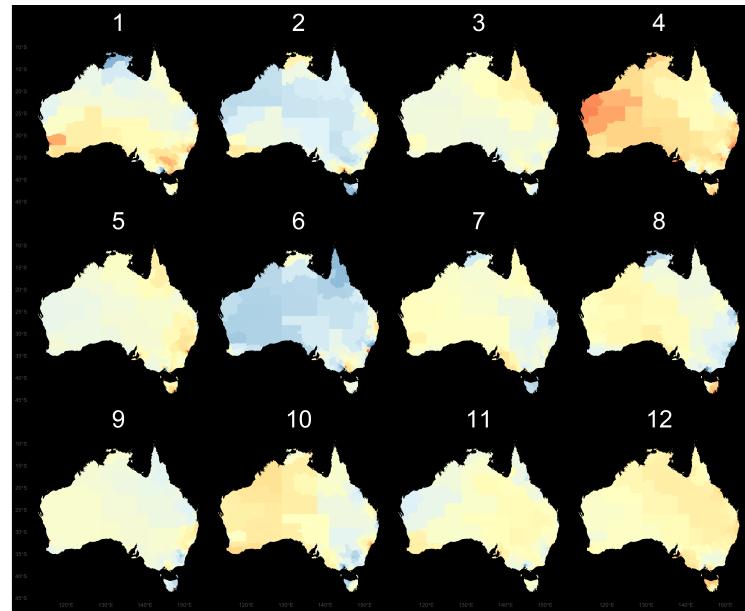
Replicate 4

Figure A.19: A choropleth map lineup, location 9 contains a distribution that affects three of the Australian capital cities.

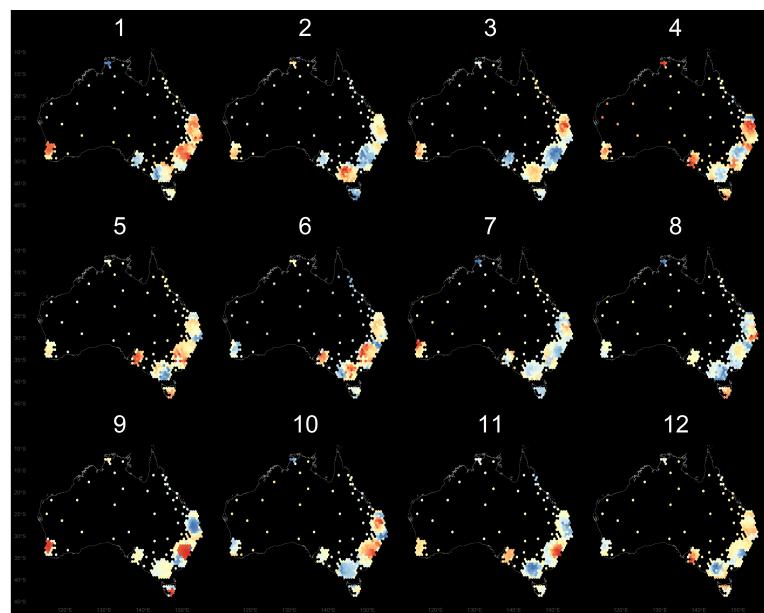


Figure A.20: A hexagon tile map lineup, location 9 contains a distribution that affects three of the Australian capital cities.

A.2.C North West to South East

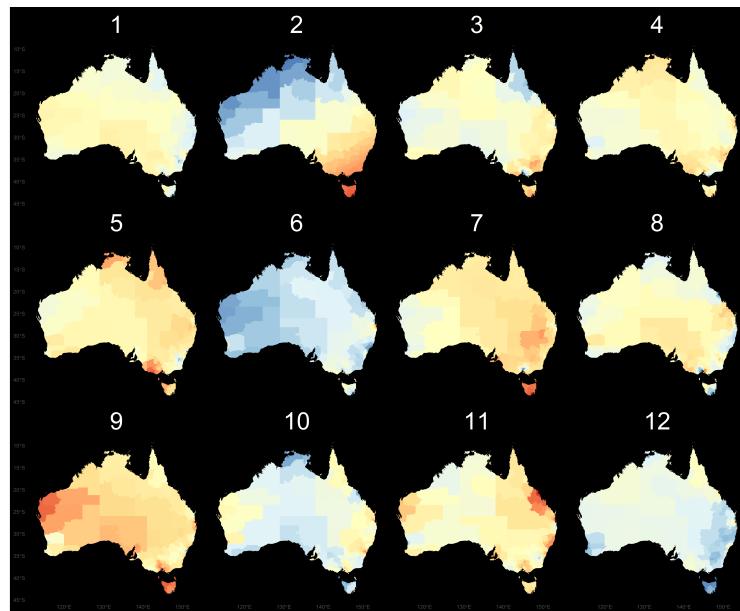
Replicate 1

Figure A.21: A choropleth map lineup, location 12 contains a distribution that affects all areas from North West to the South East.

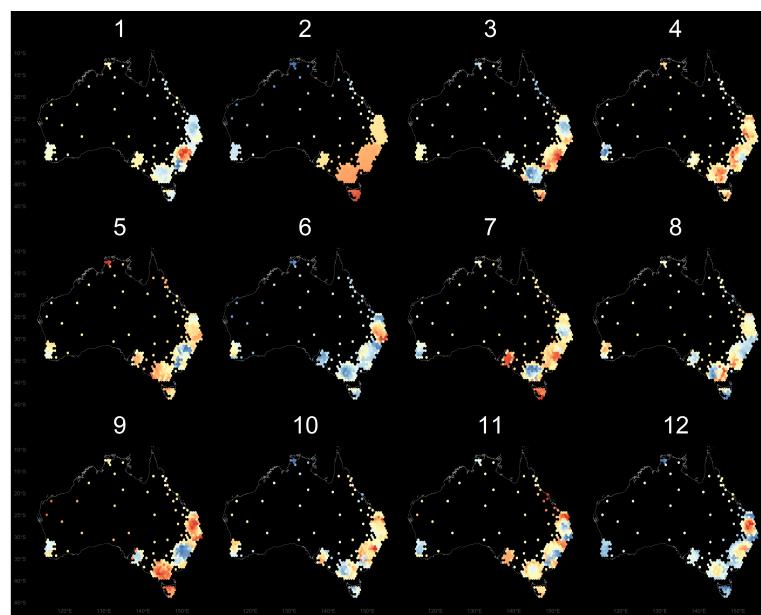


Figure A.22: A hexagon tile map lineup, location 12 contains a distribution that affects all areas from North West to the South East.

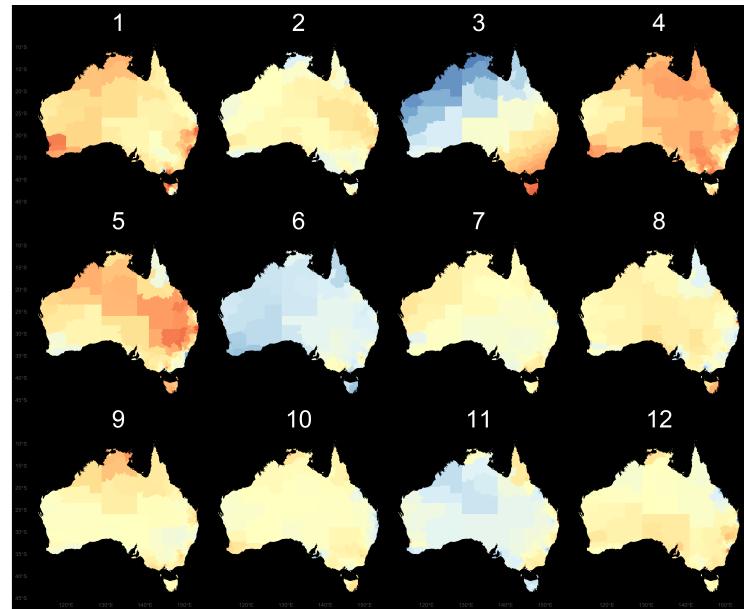
Replicate 2

Figure A.23: A choropleth map lineup, location 3 contains a distribution that affects all areas from North West to the South East.

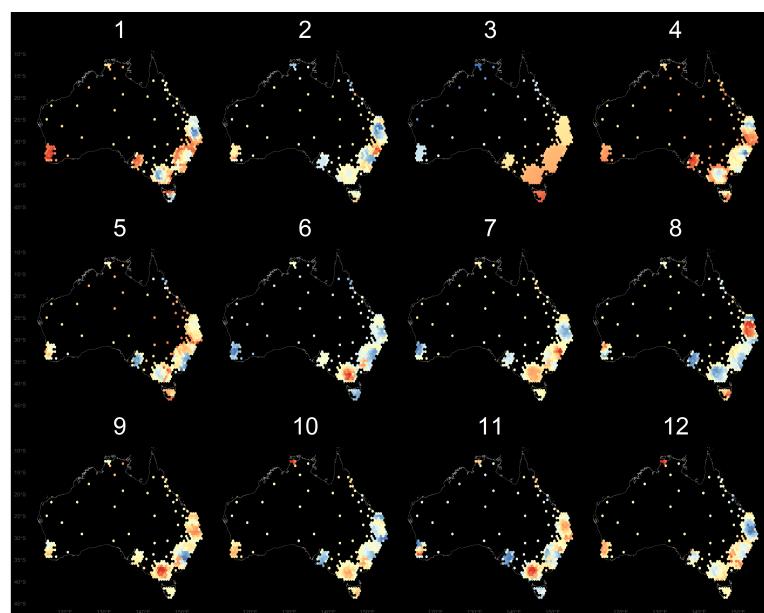


Figure A.24: A hexagon tile map lineup, location 3 contains a distribution that affects all areas from North West to the South East.

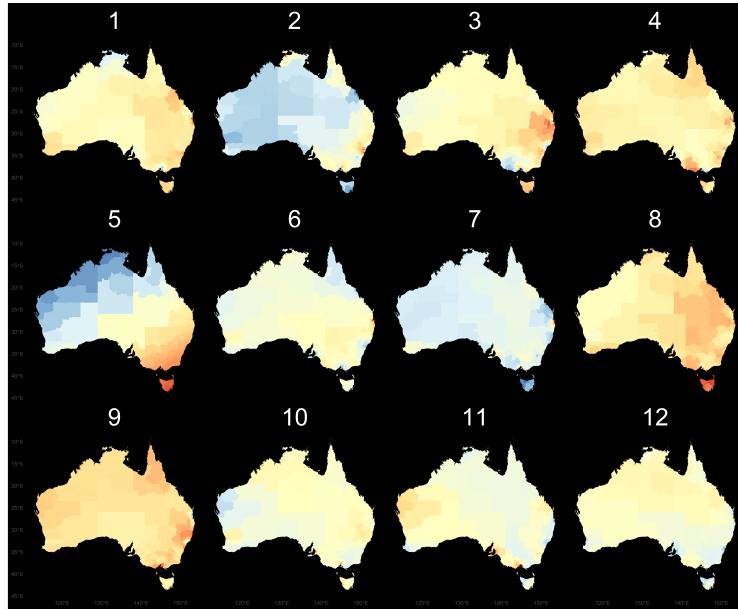
Replicate 3

Figure A.25: A choropleth map lineup, location 4 contains a distribution that affects all areas from North West to the South East.

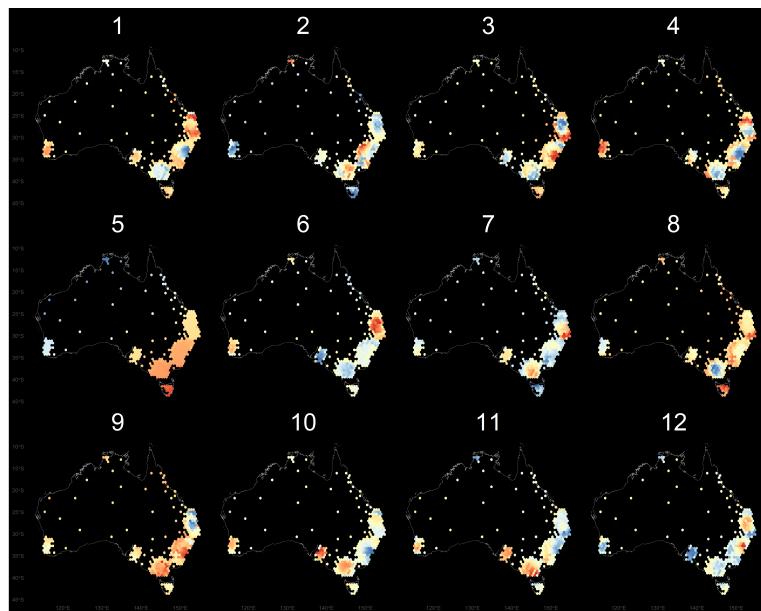


Figure A.26: A hexagon tile map lineup, location 4 contains a distribution that affects all areas from North West to the South East.

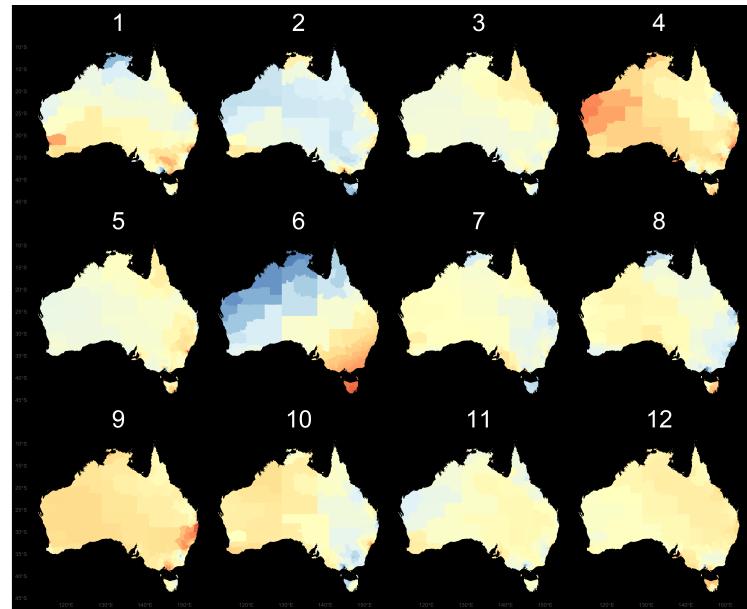
Replicate 4

Figure A.27: A choropleth map lineup, location 9 contains a distribution that affects all areas from North West to the South East.

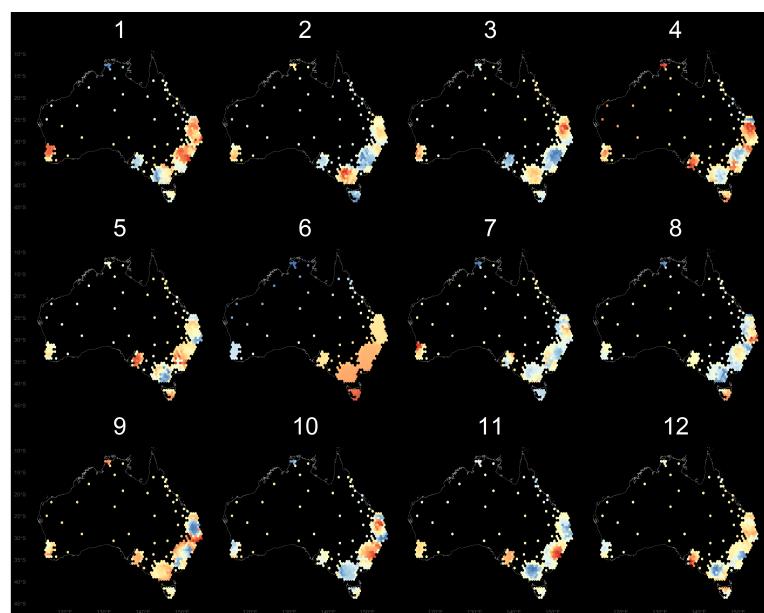


Figure A.28: A hexagon tile map lineup, location 9 contains a distribution that affects all areas from North West to the South East.

A.3 Experiment survey procedure

Participants were recruited via advertising on the Figure Eight crowdsourcing platform. Choosing the task from the list directed all potential participants to the page of instructions. This page contained written instructions and is shown in Figure A.29.

The screenshot shows a survey titled "Make A Choice Between Sets Of Australian Maps". At the top left is a "Instructions" button. Below it is a section titled "Overview" containing general information about the survey, participation, and rights. It also includes a statement about data handling and a note about the survey's purpose. The "Steps" section follows, detailing the collection of demographic information and the survey process. Finally, the "Rules Tips" section provides instructions for making selections.

Make A Choice Between Sets Of Australian Maps

Instructions ▾

Overview

Welcome! Thank you for considering participating in our survey.

Participation will involve completing a few test questions followed by the survey. It will take no more than 10 minutes of your time to complete the task.

To recognize your contribution should you choose to participate the research team is offering you a maximum of \$5.00, paid into your Figure-Eight account at the completion of the survey.

Your rights when participating in this research survey:

Your participation in this research project is entirely voluntary. If you agree to participate you do not have to complete any question(s) you are uncomfortable answering. If you do agree to participate you can withdraw from the research project during your participation without comment or penalty.

If you decide to withdraw, you can close the survey window at any time. This will not submit your results, it will also not take a record of your contributor ID, and consequently you will not be paid for your participation.

All comments and responses are anonymous. Your contributor ID will be converted into an anonymous unique identifier marking your responses. This data may be used by other researchers in the future. Any data collected as part of this research project will be stored securely as per QUT's Management of research data policy. All answers you provide during the survey will be available online, this will not contain any personally identifiable information. Data will be stored for a minimum of 5 years. If you would like to know more about this research, please email stephanie.kobakian@qut.edu.au.

Your decision to participate or not participate will in no way impact upon your current or future relationship with Queensland University of Technology or any associated external organizations. This survey is conducted as part of a research project as part of a Queensland University of Technology degree.

Steps

Demographic information will be collected in the first stage of the survey, this includes gender, age range and education level. You must give your consent to participate to continue to the survey questions.

The survey will include 12 pages, each containing 12 map displays. You will be asked to choose one map on each page, that is most different from the rest.

You will need to report your choice, the reason you selected it and how difficult it was to make your choice.

Rules Tips

You must make a selection for each page. If it is difficult to choose, try to make a selection and indicate your certainty about this decision is very low.

Figure A.29: The training lineups of choropleth maps.

A.3.A Training

The participants were trained using three displays. There were relatively simple lineups, they are displayed in Figure A.30 and Figure A.31.

Simple Examples

The sets of images below show three maps of Australia, either a choropleth map or a hexagon tile map. You will see similar maps in the survey. In a choropleth map, geographic regions are colored according to a numerical value. In a hexagon tile map, each hexagon is colored according to a numerical value.

Please look at these images and read the explanations prior to taking the survey.

Set 1:

In this set of three maps of Australia: Which map is most different?

Answer: MAP 1.
This map has many blue areas at the top, and many red areas at the bottom. Tasmania, the island on the bottom right has much more red than the rest of Australia. This is very different to the other maps, which do not show the same color trend from North to South.

Set 2:

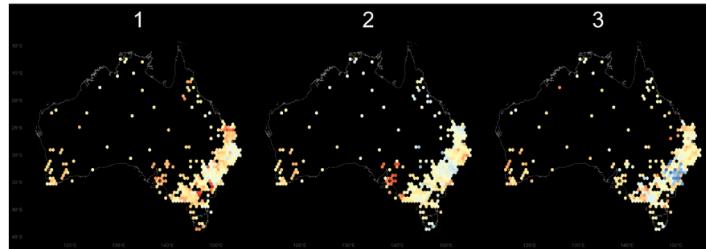
In this set of three maps of Australia: Which map is most different?

Answer: MAP 3.
This map has more red and orange areas around the coast.

Figure A.30: The training lineups of choropleth maps.

Set 3:

In this set of three maps of Australia: Which map is most different?

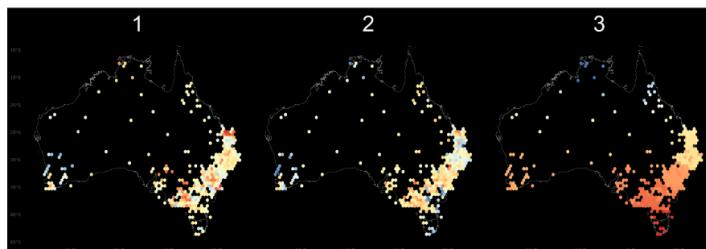


Answer: **MAP 2.**

This map has a group of red areas near each other. The other two maps have scattered colors.

Set 4:

In this set of three maps of Australia: Which map is most different?



Answer: **MAP 3.**

This map has some blue areas at the top, and many red areas at the bottom.

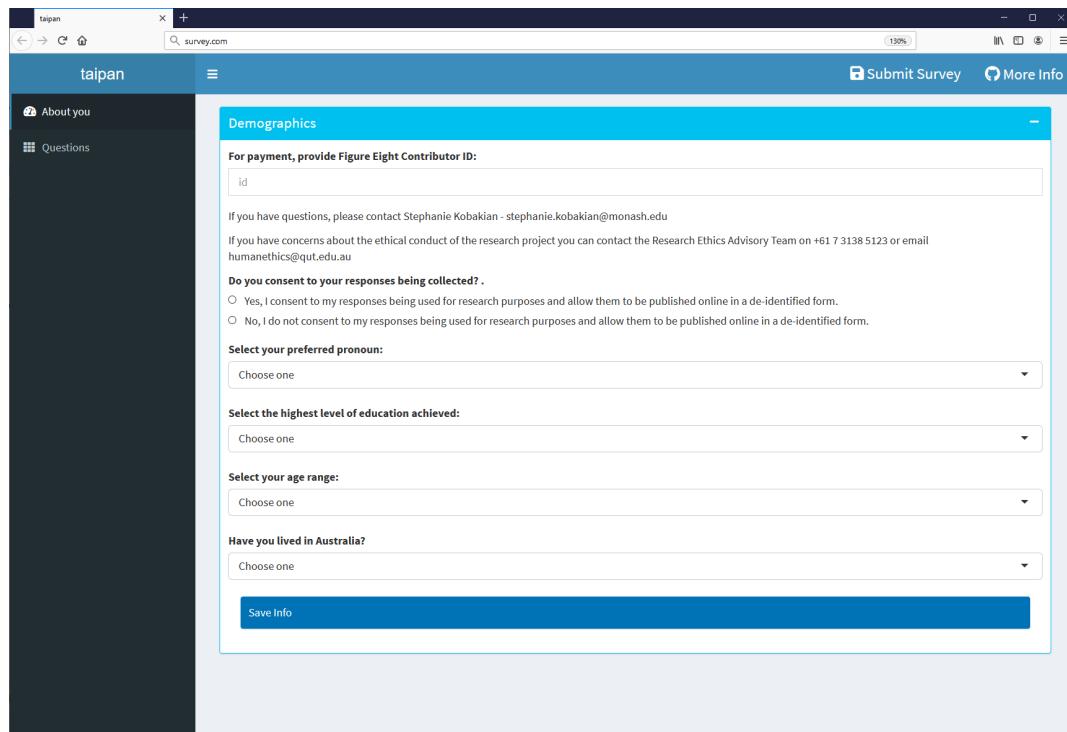
Figure A.31: *The training lineups of hexagon tile maps.*

A.3.B Survey application

The survey application was a shinydashboard we application, hosted on a website external to the Figure-Eight platform. The link to the survey was located at the bottom of the instructions and training page. Only participants who had read all of the instructions and seen the example image sets continued to the survey via the link. This page also contained a question that asked participants for a validation code. The participants unique validation code was generated upon them opening the web application. This code was released to participants when they had considered all twelve lineups and submitted their

responses to the `googlesheets` data set. Their validation codes were contained in the data set and associated with each of their responses.

The demographic and consent page of the `shinydashboard` web application are displayed in Figure A.32. Two example lineups are shown, one choropleth map lineup in Figure A.33 and one hexagon tile map lineup in Figure A.34.



The screenshot shows a web browser window titled 'taipan' with a URL of 'survey.com'. The main content area is titled 'Demographics'. It contains several form fields:

- 'For payment, provide Figure Eight Contributor ID:' with a text input field containing 'id'.
- A note: 'If you have questions, please contact Stephanie Kobakian - stephanie.kobakian@monash.edu'
- A note: 'If you have concerns about the ethical conduct of the research project you can contact the Research Ethics Advisory Team on +61 7 3138 5123 or email humanethics@qut.edu.au'
- 'Do you consent to your responses being collected?' with two radio button options: 'Yes, I consent to my responses being used for research purposes and allow them to be published online in a de-identified form.' and 'No, I do not consent to my responses being used for research purposes and allow them to be published online in a de-identified form.'
- 'Select your preferred pronoun:' with a dropdown menu showing 'Choose one'.
- 'Select the highest level of education achieved:' with a dropdown menu showing 'Choose one'.
- 'Select your age range:' with a dropdown menu showing 'Choose one'.
- 'Have you lived in Australia?' with a dropdown menu showing 'Choose one'.
- A blue 'Save Info' button at the bottom.

Figure A.32: The demographics questions tab of the `shinydashboard` survey application.

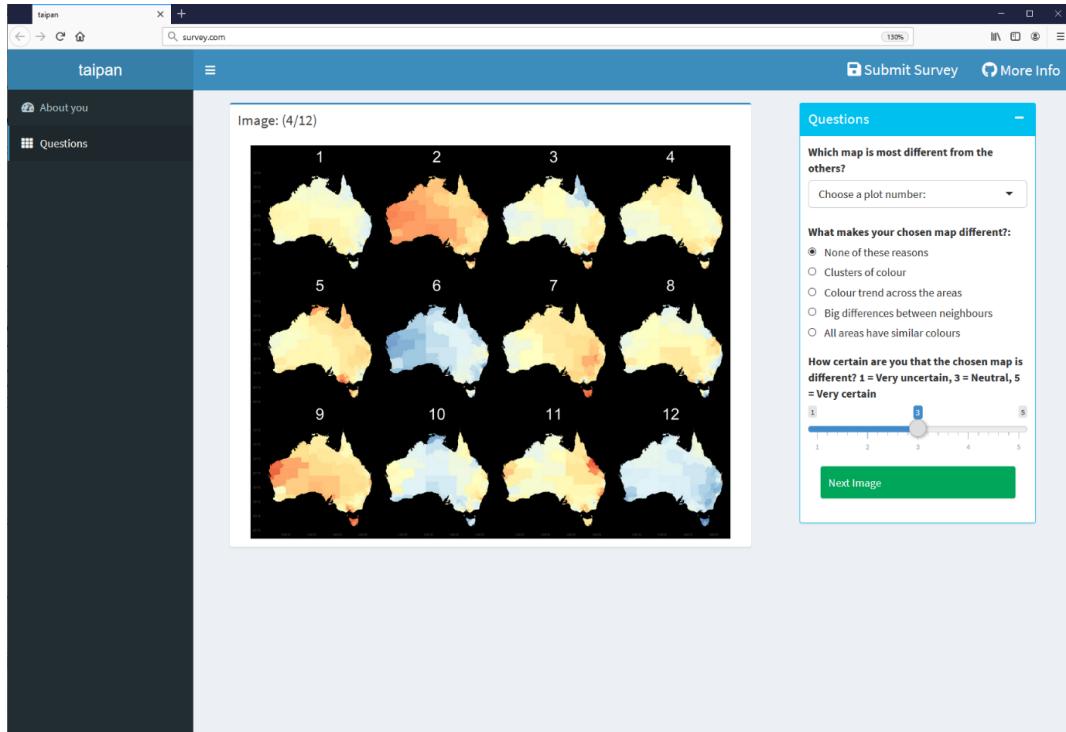


Figure A.33: An example of the choropleth map lineup shown in the survey tab of the shinydashboard app.

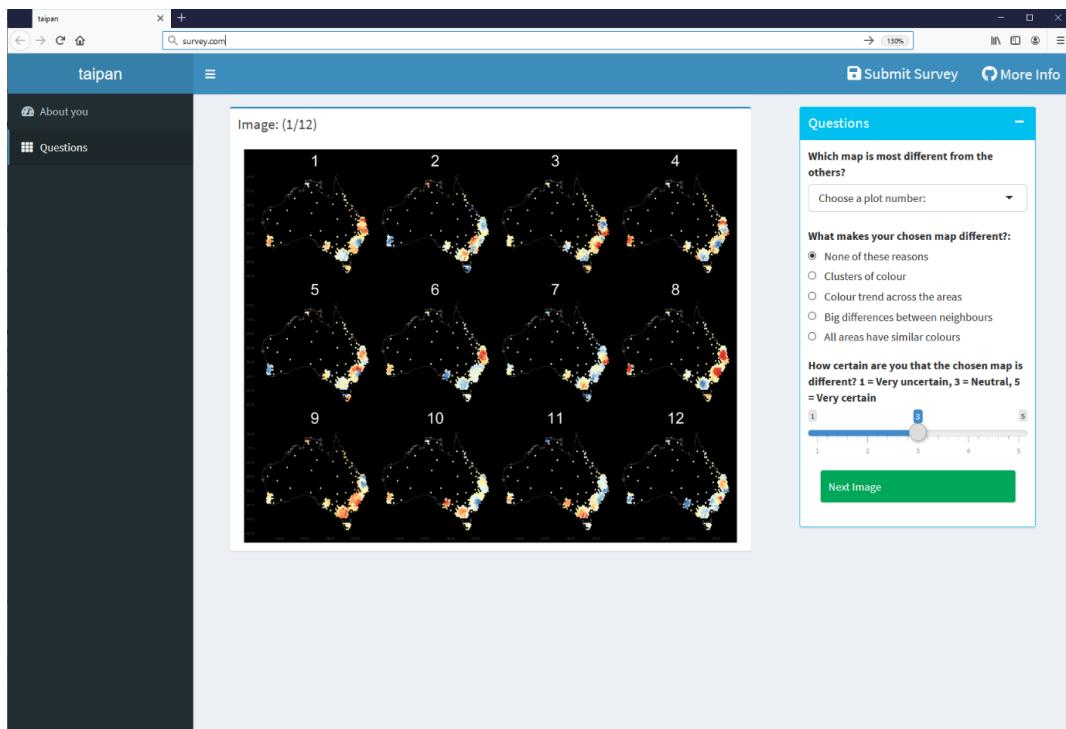


Figure A.34: An example of the hexagon tile map lineup shown in the survey tab of the shinydashboard app.

A.4 Ethics Approval

Assessing the effectiveness of different visualisation methods for Australian spatial data.

QUT Ethics Approval Number 1900000991

Research team

Principal Researcher:	Stephanie Kobakian	Masters Student
Associate Researchers:	Kerrie Mengersen	Principal Supervisor
	Earl Duncan	Associate Supervisor
	Dianne Cook	External Supervisor

Faculty of Science and Engineering
Queensland University of Technology (QUT)

Why is the study being conducted?

The purpose of this research project is to test the effectiveness of two types of spatial displays: a choropleth map, and a hexagon map, where each geographic region is represented by a hexagon. This will examine the use of different map styles in communicating a relationship between geographic areas. The purpose of these displays is to convey the spatial distribution of the disease occurrence, or incidence. This can mean detecting hot spots corresponding to outbreaks, spatial trends, for example, indicating occurrence is related to latitude or even rural vs urban differences. Effectiveness of the display will be measured by accurate and efficient perception of these patterns.

This research project is being undertaken as part of a Masters study for Stephanie Kobakian, a student at Queensland University of Technology. You are invited to participate in this research project because you have had experience answering surveys and participating in crowdsource activities.

What does participation involve?

Participation will involve completing a few test questions followed by a survey. Each survey item will contain a grid of maps, you will be asked the following question:

Which map is most different from the others?

Report your choice and the reason you selected it, and how difficult your decision was to make. It will take no more than 10 minutes of your time to complete the task.

Questions will include images similar to Figure 1 below:

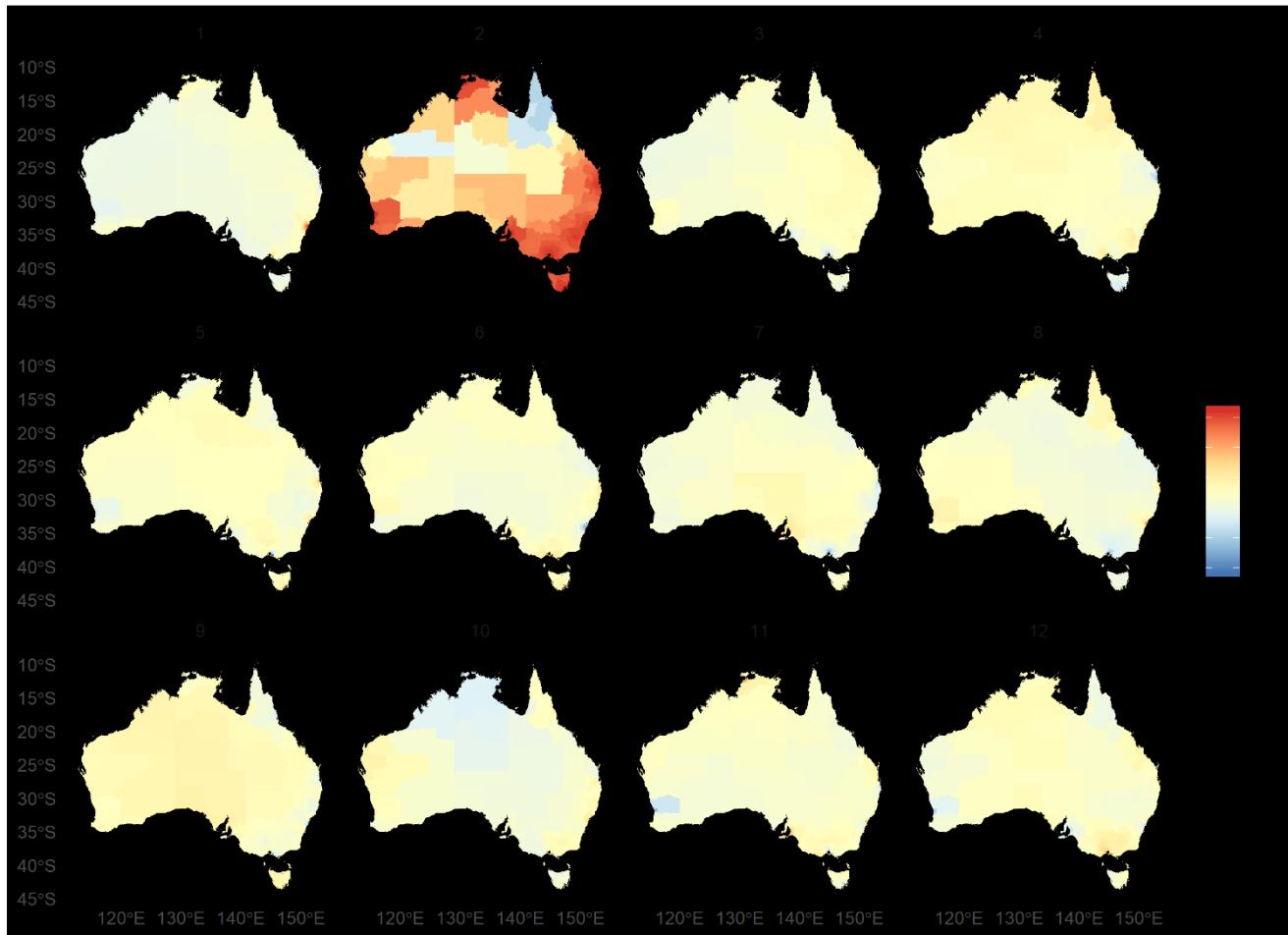


Figure 1. A lineup of geographic maps of Australia. Each sa3 has been coloured according to a simulated data set. Only one of these maps displays a spatial relationship, the rest are null plots, with colours shuffled between the areas.

Your participation in this research project is entirely voluntary. If you agree to participate you do not have to complete any question(s) you are uncomfortable answering. Your decision to participate or not participate will in no way impact upon your current or future relationship with QUT (for example your grades) or associated external organisation. If you do agree to participate you can withdraw from the research project during your participation without comment or penalty. However, as the survey does not request any personal identifying information, once it has been submitted it will not be possible to withdraw.

What are the possible benefits for me if I take part?

It is expected that this research project will directly benefit you as a paid member of the Figure-Eight platform. The outcomes of the research may also benefit researchers who have the options to consider the maps they use to communicate spatial information to the general public.

To recognise your contribution should you choose to participate the research team is offering you \$5.00 paid into your Figure-Eight account at the completion of the survey.

What are the possible risks for me if I take part?

Risks:

Monetary risk: if participants do not accurately provide their Contributor ID in our external survey we will not be able to confirm they participated, and provide the payment to their account. As the figure Eight platform encourages paying participants after the survey collection period has finished.

Privacy risk: data on contributors, such as location, channel, time, and IP address will be provided to researchers by the platform. This information will be held on the researcher's personal laptops.

Psychological risks of taking part in this survey involves a possible negative affective state such as anxiety as participants will be asked to evaluate maps that are very unfamiliar. There is also the potential risk of anxiety resulting from the colouring of red areas, this colour scheme is best for all colour blindness types except greyscale.

What about privacy and confidentiality?

All comments and responses are anonymous. It will only be possible to identify due to your contributor ID provided in the research, personal identifying information is not sought in any of the responses. It will not be possible to re-identify you using your contributor ID, but it will be removed and not stored in a public space after is received by the researchers. This data may be used by other researchers in the future.

Any data collected as part of this research project will be stored securely as per QUT's Management of research data policy. All answers you provide during the survey will be available online, this will not contain any personally identifiable information. Data will be stored for a minimum of 5 years. It will be available publicly at the web address: <https://github.com/srkobakian/experiment>.

The research project is funded by ACEMS and they will have access to the data obtained during the project as it will be publicly available.

How do I give my consent to participate?

The survey will ask if each participant gives their consent for their responses to be used.

The selection of the "yes" checkbox will allow continuation to the survey questions.

Submission of the completed survey is accepted as an indication of your consent to participate in this research project, you may withdraw by completing less than 50% of the questions, after checking the "yes" checkbox.

What if I have questions about the research project?

If you have any questions or require further information please contact one of the listed researchers:

Stephanie Kobakian	stephanie.kobakian@hdr.qut.edu.au	+61 433699797
Kerrie Mengersen	k.mengersen@qut.edu.au	+61 731382063
Earl Duncan	earl.duncan@qut.edu.au	+61 410874218
Dianne Cook	dicook@monash.edu	+61 399052608

What if I have a concern or complaint regarding the conduct of the research project?

QUT is committed to research integrity and the ethical conduct of research projects. If you wish to

discuss the study with someone not directly involved, particularly in relation to matters concerning policies, information or complaints about the conduct of the study or your rights as a participant, you may contact the QUT Research Ethics Advisory Team on +61 7 3138 5123 or email humanethics@qut.edu.au.

**Thank you for helping with this research project.
Please print this sheet for your information.**