

Spatial modelling of the electoral two-party preferred vote in Australia. A study of federal elections between 2001 and 2016 via the **eechidna** R package.

Jeremy Forbes, Di Cook & Rob Hyndman

2019-03-29

Contents

1	Introduction	3
2	Data collection, wrangling and imputation	5
2.1	Collecting the data	5
2.2	Joining Census and election data	5
3	Modelling	9
3.1	Data pre-processing	9
3.2	Model framework	10
4	Results	12
4.1	Spatial autoregressive parameter	12
4.2	Country-wide trend	12
4.3	Influential socio-demographics	13
4.4	A closer look at the residuals	15
5	Conclusion	19
6	Acknowledgements	20
7	Software	21

Chapter 1

Introduction

Australia has changed in many ways over the last two decades. Rising house prices, country-wide improvements in education, an ageing population, and a decline in religious affiliation, are just a few facets of the country's evolving socio-demographic characteristics. At the same time, political power has moved back and forth between the two major parties. In the 2007 and 2010 federal elections, the Australian Labor Party (Labor) was victorious, whereas the 2001, 2004, 2013 and 2016 elections were won by the Liberal National coalition (Liberal). The two-party preferred vote, a measure of support between these two parties, fluctuated between 47.3% and 53.5% (in favour of the Liberal party) over this period. This study explores how electoral characteristics relate to two-party preference, and whether their effects have changed over time. Electoral socio-demographics are derived from the Census, and vote counts are obtained from federal elections.

Joining these two data sources is problematic as there is an inherent asynchronicity in the two events. A Census is conducted by the Australian Bureau of Statistics (ABS) every five years, whereas a federal election (conducted by the Australian Electoral Commission (AEC)) usually occurs every three years. The first problem addressed is that of obtaining appropriate Census data for the 2004, 2007, 2010 and 2013 elections - election years in which a Census does not occur. The predominant approach in previous studies is to join voting outcomes to the nearest Census, without accounting for any temporal differences (see Davis and Stimson (1998), Stimson et al. (2006), Liao et al. (2009) and Stimson and Shyy (2009)). Furthermore, electoral boundaries change regularly, so spatial discrepancies also arise when matching electoral data. To obtain appropriate Census data for these four elections, electoral socio-demographics are imputed using a spatio-temporal imputation that combines areal interpolation (Goodchild et al., 1993) and linear time-interpolation. Collecting and wrangling the raw data, along with the imputation process, are detailed in section 2. All data and associated documentation relating to this procedure are available in the `eechidna` R package (Forbes et al., 2019), providing a resource for future analysis.

Previous work on modelling Australian federal elections have found that aggregate socio-demographics are relatively good predictors of voting outcomes. Forrest et al. (2001) does this using multiple regression of the Liberal and Labor primary vote for polling booths in the Farrer electorate in 1998. Stimson et al. (2006), Stimson and Shyy (2009) and Stimson and Shyy (2012) use principal component analysis of polling booths in the 2001, 2004 and 2007 elections respectively, also finding that socio-demographic characteristics of polling booths are linked to their two-party preferred vote. On the contrary, Stimson and Shyy (2009) models the polling booth swing vote (change in the two-party preferred vote) in the 2007 election, finding that little of swing vote can be explained by Census data. Instead of analyzing a single election in isolation, this paper employs a consistent model framework across six elections so that temporal changes in the effects of socio-demographics can be observed, where each federal elections is modelled with a cross-sectional data set. The use of a regression framework to examine these socio-political relationships over time is seemingly absent from previous Australian studies. It also appears that no study has attempted any type of statistical analysis of socio-demographics in conjunction with voter behaviour in Australia since 2007, making this paper distinctly different from those previous.

The cross-sectional data set for each election consists of the two-party preferred vote (response variable), and socio-demographic variables (explanatory variables) that characterise each electorate. To obtain these cross-sections, socio-demographic variables are first standardized, and then principal components are used to group variables into “factors”. To account for the inherent spatial structure of the data, a spatial error model is fit for each election. These steps are discussed in section 3. In section 4 inference is conducted on the models to see which effects are significant, how effects change over time and which electorates have abnormal voting behaviour.

Chapter 2

Data collection, wrangling and imputation

2.1 Collecting the data

The voting outcome of interest is the electoral two-party preferred vote, which is provided by the Australian Electoral Commission (AEC) for the 2001, 2004, 2007, 2010, 2013 and 2016 elections via the AEC Tally Room. The AEC divide Australia into 150 regions called electorates, with each corresponding to a single seat in the House of Representatives. Voting is compulsory in Australia, and each voter assigns a numbered preference to each available candidate in their electorate. The two-party preferred vote is determined by a tally of these preferences where, by convention, only the ranks of the Labor and Liberal candidates are considered. This is recorded as a percentage preference in favour of the Liberal party.

Socio-demographic variables are derived from the Census of Population and Housing (Census), which is a survey of every household in Australia, recording information such as age, gender, ethnicity, education level and income. There have been four Censuses in the 21st century, being that in 2001, 2006, 2011 and 2016. The Australian Bureau of Statistics (ABS) conducts the Census and publishes aggregated information. The ABS approximation of electorates at the time of the Census is chosen. From this aggregate information, 67 socio-demographic variables are computed for each of the electorates.

Raw data is sourced online from the AEC and ABS websites in .csv and .xlsx files. The format of these files change over the years, making extracting the appropriate information a big task. The functions available in the `dplyr` (Wickham et al., 2019b) and `readxl` (Wickham et al., 2019a) R packages are very useful, as they provide fast consistent tools for data manipulation and functions to import .xlsx files (respectively). The 2001 and 2006 Census data are however published in a format where each electorate has a separate document, making it difficult to use the `dplyr` tools and instead cells have to be selected from each individual file to construct the desired variables. All scripts required for the data wrangling process can be found in the github repository for the `eechidna` R package (Forbes et al., 2019), along with the raw data. The `eechidna` package makes this study entirely reproducible and provides a resource to help wrangle data for future Censuses and elections, when they become available.

2.2 Joining Census and election data

2.2.1 Differences between Census and election data

Between 2001 and 2016 there were six elections and four Censuses (see Figure 2.1). Electoral boundaries are redistributed regularly by the AEC, meaning that only in the years where both a Census and election occur

will the boundaries match - the case for the 2001 and 2016 election. Therefore, for the four elections between 2004 and 2013, both temporal and spatial differences in electorates need to be accounted for when joining the electoral two-party preferred vote with Census data. For these elections a spatio-temporal imputation method is employed to obtain electoral socio-demographics. This method uses Census information from both before and after the election of interest.

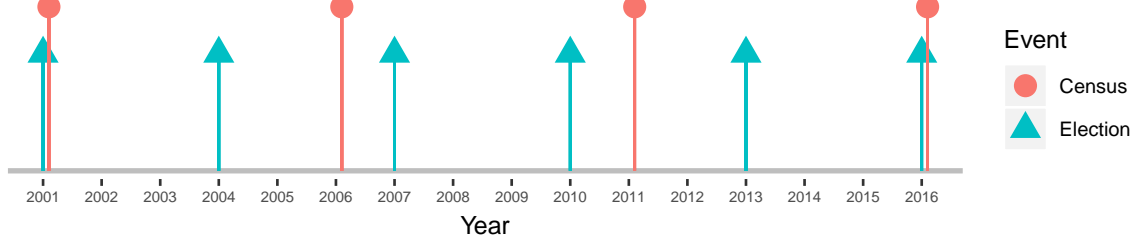


Figure 2.1: Timeline of Australian elections and Censuses. They do not always occur in the same year.

2.2.2 Spatio-temporal imputation

To account for spatial differences, the piece-wise approximation method in Goodchild et al. (1993) is adopted. Consider a map of source zones $s = 1, \dots, S$, for which socio-demographic information is available, and a set of target zones $t = 1, \dots, T$ for which information is to be imputed. In this context the map of electoral boundaries at the time of a Census would be the source zones, and the boundaries at the time of the election would be the target zones. Denote the area of intersection between source zone s and target zone t as $A_{s,t}$, the population of the source zone s as U_s , and the population of intersection between source zone s and target zone t as $P_{s,t}$.

Compute each $A_{s,t}$ and estimate population of the intersection:

$$\hat{P}_{s,t} = \frac{U_s * A_{s,t}}{\sum_{t=1}^T A_{s,t}}$$

This assumes that populations are uniformly distributed within each source zone.

In order to calculate socio-demographic information for each of the target zones, a weighted average is taken using the estimated population as weights. Denote a given Census variable for the target zone C_t , and the same Census variable for the source zone D_s :

$$\hat{C}_t = \frac{\sum_{s=1}^S D_s * \hat{P}_{s,t}}{\sum_{s=1}^S \hat{P}_{s,t}}$$

This assumes that each individual in a source zone assumes the aggregate characteristics of the zone.

Applying this to each of the target zones addresses the spatial component, as it imputes the required socio-demographic for the desired electoral boundaries. However these are applicable at the time of the Census (source year) and are not yet appropriate for the election (target year).

Denote year y , with a Census falling on y_1 and y_3 , and an election on year y_2 , and add this subscript to the Census variable estimate, $\hat{C}_{t,y}$. To account for temporal changes, linear interpolation is used between Census years to get the final estimate of a Census variable for the target zone in the election year y_2 . This assumes that population evolves in a linear manner over time.

$$\hat{C}_{t,y_2} = \frac{y_3 - y_2}{y_3 - y_1} * \hat{C}_{t,y_1} + \frac{y_2 - y_1}{y_3 - y_1} * \hat{C}_{t,y_3}$$

2.2.3 Applied

Publically available Census data is aggregated and there are different resolutions accessible, ranging from SA1 (over 50,000 zones) to electoral divisions (150 zones). For this study electoral divisions are used as source zones, and this imputation method is applied for each of the 2004, 2007, 2010 and 2013 elections. To demonstrate its functionality, consider the imputation of socio-demographic variables for the electorate of Hume in New South Wales (NSW), at the time of the 2013 federal election. Figure 2.2 shows this region amongst other NSW electorates.

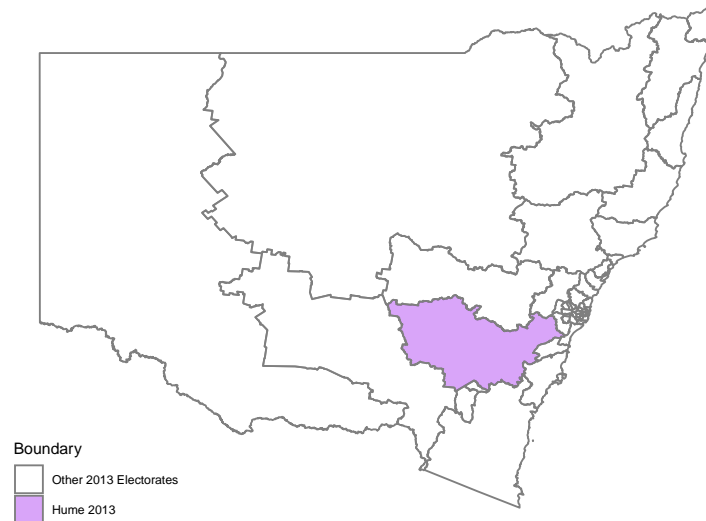


Figure 2.2: Some of the electoral boundaries in NSW for 2013, with the electoral boundary for Hume, shown in purple.

The Censuses neighbouring the 2013 election are those in 2011 and 2016, and the Hume boundary is changed, as seen by plotting the Hume boundary (purple) in the 2013 election over the divisions in 2016 (see Figure 2.3).

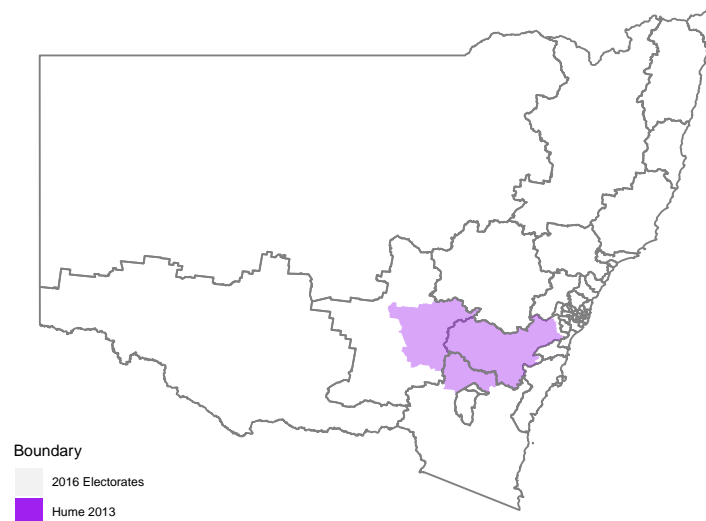


Figure 2.3: Census division boundaries in NSW for 2016, with the 2013 electoral boundary for Hume, shown in purple. The purple region is not contained within a single Census division.

There are many electorates in 2016 that intersect with the purple region (Hume boundary for 2013), these

include the divisions of Riverina, Eden-Monaro and Hume, along with smaller intersecting areas with Fenner, Calare, Gilmore and Whitlam. To impute Census information for this purple region, calculate the percentage of each 2016 electorate that intersects with the purple region, which is then used to estimate intersection populations $\hat{P}_{s,t}$.

Electorate (2016)	Percentage	Population in Electorate	Estimated Population Allocated to Purple Region: $\hat{P}_{s,t}$
HUME	96.54%	150643	145427
RIVERINA	25.11%	155793	39117
EDEN-MONARO	11.09%	147532	16358
CANBERRA	0.28%	196037	548
FENNER	0.23%	202955	474
WHITLAM	0.06%	152280	92
GILMORE	0.06%	150436	86
CALARE	0.01%	161298	21

Now consider the socio-demographic *AusCitizen* - the proportion of people in the region who are Australian citizens.

DivisionNm	AusCitizen (%): D_s	Estimated Population Allocated to Purple Region: $\hat{P}_{s,t}$
HUME	90.02	145427
RIVERINA	89.11	39117
EDEN-MONARO	88.00	16358
CANBERRA	85.48	548
FENNER	83.64	474
WHITLAM	89.52	92
GILMORE	89.03	86
CALARE	87.56	21

Then taking a weighted average of *AusCitizen* using the estimated population as weights yields $\hat{C}_{Hume,2016} = 89.65\%$. Repeating this process using the 2011 Census and electoral boundaries yields $\hat{C}_{Hume,2011} = 91.00\%$

Finally, linearly interpolate between 2011 and 2016 to arrive at the 2013 estimate:

$$\begin{aligned}
 \hat{C}_{Hume,2013} &= \frac{3}{5} \cdot \hat{C}_{Hume,2011} + \frac{2}{5} \cdot \hat{C}_{Hume,2016} \\
 &= \frac{3}{5} \cdot 91.00\% + \frac{2}{5} \cdot 89.65\% \\
 &= 90.46\%
 \end{aligned}$$

This is done for each of the socio-demographic variables, and repeated each of the 2013 electorates.

Chapter 3

Modelling

3.1 Data pre-processing

With socio-demographic information now available for each electorate, each election is joined to the data corresponding with its two-party preferred vote. Socio-demographic variables within each election year are standardized to have mean zero and variance one, to adjust for changing variable scales. For example, inflation-adjusted median rental prices increased across almost all electorates, with median rent of 200 dollars per week placing an electorate in the 90th percentile in 2001, but only the 30th percentile in 2016.

3.1.1 Dimension reduction

With only $N = 150$ observations (electorates) in each election and $p = 65$ socio-demographic variables in each cross-section, any model using all variables would face serious problems with multi-collinearity and overfitting, likely leading to erroneous conclusions regarding variable significance. Therefore a form of dimension reduction is adopted before models are fit.

Socio-demographic variables¹ that represent similar information are combined into “factors” using principal component analysis (PCA). The scree plots of the principal components for each election all level off after four components, and the loadings of these four components are similar across the elections. Principal components are then computed on the combined set of socio-demographics across all six elections. A factor is created by combining several variables all have large loadings in a particular component and when there is an intuitive reason as to why these variables could represent common information. A loading with magnitude greater than 0.15 is considered large. After computing these sums, each factor is again standardized to have mean zero and variance one, within each election.

Consider the **Incomes** factor as an illustration. Independent of principal components, we may suspect that median personal income, median household income and median family income are providing similar information about the financial wellbeing of an electorate. Their loadings in the first principal component are large (0.19, 0.21 and 0.22 respectively), which provides the evidence needed to combine these variables into a single factor, which is called **Incomes**.

This process reduces the predictor set to $p = 30$.

¹A preliminary step involved removing all age bands, because age is represented by median age, and to remove variables relating to particular denominations of Christianity.

3.2 Model framework

An identical model specification is used across the six elections, with each election modelled separately. This allows for the socio-demographic effects to be estimated separately for each year, allowing for interpretation of temporal changes in these effects. This is preferable over a single longitudinal model because it avoids any concerns of undue bias stemming from an incorrectly imposed time-varying restriction on any variable. Without such restrictions, a pooled cross-sectional model does not yield any distinct advantage over separate cross-sections. The panel approach is avoided because of how frequently electoral boundaries change, meaning that electorates that have the same name across elections are not guaranteed to represent the same geographical region. Therefore any fixed or random effects models would be difficult to estimate without implementing consistent boundaries, which would require further imputation.

For each cross-section, let the response variable be the two-party preferred vote in favour of the Liberal party, denoted Y , with $Y = 70$ representing a 70% preference for Liberal, 30% for Labor. Although Y lies in the interval $(0, 100)$, observed values are never very close to 0 or 100 (minimum 24.05% and maximum 74.90%), so there is no need to impose the constraint of $Y \in [0, 100]$. Furthermore, the response is found to be spatially correlated in each election (Moran's I test, $p \leq 7 \cdot 10^{-15}$). This is expected, as electorates are aggregate spatial units, and hence the spatial structure of the data must be modelled appropriately.

The spatial error model (Anselin, 1988) is chosen because it captures spatial heterogeneity by incorporating a spatially structured random effect vector (LeSage et al., 2009). In this context, the random effect can be thought of as capturing the unobserved political climate in each electorate, where the climate is correlated with the climate in neighbouring electorates. This functions under the assumption that the climate is independent of electoral socio-demographics, and that an electorate is equally correlated with any electorate that shares a part of its boundary. Spatial weights are calculated in accordance with these assumptions. The spatial error model is specified as follows:

Let ρ be spatial autoregressive coefficient, \mathbf{v} be a spherical error term, \mathbf{W} be a matrix of spatial weights (containing information about the neighbouring regions), \mathbf{X} be a matrix of socio-demographic covariates, $\boldsymbol{\beta}$ be a vector of regression coefficients and \mathbf{a} be a spatially structured random effect vector.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{a}$$

and

$$\mathbf{a} = \rho\mathbf{W}\mathbf{a} + \mathbf{v}$$

where

$$\mathbf{v} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

.

so it can be written

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \rho\mathbf{W})^{-1}\mathbf{v}$$

Estimation is done using feasible generalized least squares.

Table 3.1 details the resultant estimated model coefficients and their estimated standard errors for each of the six elections. These are interpreted in the next section.

Table 3.1: Estimated model for each of the six elections.

	<i>Dependent variable:</i>					
	Two-party preferred vote in favor of the Liberal party					
	2001 (1)	2004 (2)	2007 (3)	2010 (4)	2013 (5)	2016 (6)
ρ	0.46*** (0.15)	0.29* (0.17)	0.24 (0.17)	0.19 (0.16)	0.27* (0.16)	0.50*** (0.17)
AusCitizen	-3.13 (2.26)	-2.64 (2.43)	-2.53 (2.34)	-0.08 (2.79)	-3.40 (2.76)	-1.80 (2.71)
Born_Asia	2.22 (2.18)	-0.95 (2.44)	-1.60 (2.19)	-6.83** (2.73)	-3.03 (2.71)	-0.55 (2.17)
Born_MidEast	-1.15 (1.07)	-1.59 (1.20)	-2.01* (1.11)	-2.03 (1.27)	-0.92 (1.24)	-1.44 (1.13)
Born_SE_Europe	-3.21** (1.42)	-4.24*** (1.46)	-3.61*** (1.02)	-4.14*** (1.19)	-3.69*** (1.07)	-2.72*** (0.97)
Born_UK	0.25 (1.00)	-0.07 (0.98)	0.34 (0.90)	0.56 (1.07)	-0.09 (1.04)	-1.32 (1.04)
BornElsewhere	-5.04 (3.30)	-4.91 (3.68)	-4.13 (3.38)	2.35 (4.23)	-5.23 (4.15)	-4.14 (3.97)
Buddhism	-0.49 (1.39)	-0.17 (1.61)	-1.37 (1.61)	-0.83 (1.80)	-0.12 (1.68)	-1.60 (1.56)
Christianity	-2.48 (1.73)	-1.23 (1.85)	0.38 (1.83)	0.50 (1.99)	2.41 (1.85)	1.68 (1.78)
CurrentlyStudying	-2.19** (0.99)	-0.13 (1.13)	2.06* (1.17)	2.12* (1.25)	1.15 (1.26)	-0.16 (1.18)
DeFacto	-6.44*** (1.87)	-5.37** (2.48)	-6.43*** (2.31)	-8.07*** (3.06)	-6.56** (3.11)	-8.53*** (2.83)
DiffAddress	3.88*** (0.94)	5.06*** (1.12)	4.22*** (0.99)	5.57*** (1.76)	3.53* (1.91)	5.67*** (1.60)
Distributive	1.27 (1.12)	2.01* (1.21)	1.36 (1.13)	1.57 (1.34)	2.10* (1.27)	1.20 (1.21)
Education	1.08 (2.38)	0.52 (3.12)	-5.52* (3.27)	-4.08 (3.95)	-4.44 (3.78)	-7.07** (3.55)
Extractive	4.83*** (1.48)	5.45*** (1.42)	5.37*** (1.36)	7.31*** (1.56)	6.71*** (1.47)	7.43*** (1.39)
FamHouseSize	-0.16 (2.19)	0.87 (2.72)	-2.40 (2.69)	-2.53 (3.25)	-3.26 (3.28)	-2.91 (2.90)
Incomes	4.36** (1.77)	5.03* (2.66)	9.45*** (2.75)	7.09** (3.25)	7.97*** (2.92)	12.20*** (2.75)
Indigenous	2.91* (1.68)	1.97 (1.95)	2.48 (1.75)	2.84 (2.16)	0.67 (2.14)	-0.05 (2.00)
Islam	-0.92 (1.22)	-0.97 (1.36)	-0.54 (1.27)	-2.50 (1.52)	-0.82 (1.42)	-0.95 (1.34)
Judaism	1.88* (1.05)	1.78 (1.13)	2.66*** (1.01)	1.97* (1.15)	2.74** (1.10)	1.65* (1.00)
ManagerAdminClericalSales	2.06*** (0.71)	3.32*** (0.93)	6.00*** (0.90)	5.47*** (1.08)	5.04*** (1.03)	5.78*** (1.06)
Married	0.44 (2.31)	0.11 (2.96)	-1.22 (2.83)	-0.22 (3.15)	0.91 (3.03)	-2.34 (2.81)
MedianAge	2.32* (1.32)	4.96*** (1.65)	3.66** (1.81)	4.00* (2.26)	2.30 (2.08)	2.87 (1.79)
NoReligion	-1.57 (1.59)	-0.92 (1.71)	0.56 (1.73)	-0.30 (1.92)	1.02 (1.94)	1.31 (2.04)
OneParent_House	-1.73 (1.36)	-0.45 (1.59)	-0.75 (1.49)	-1.46 (1.69)	-0.77 (1.57)	-0.74 (1.47)
OtherLanguageHome	-0.44 (3.22)	5.92 (4.16)	9.98** (3.91)	11.24** (4.76)	9.00* (4.66)	9.84** (4.44)
PropertyOwned	-0.46 (1.37)	-0.53 (1.50)	0.67 (1.43)	-0.94 (1.76)	-0.48 (1.67)	1.41 (1.50)
RentLoanPrice	-1.57	-3.32*	-4.01**	-0.97	-0.70	-0.89

Chapter 4

Results

4.1 Spatial autoregressive parameter

The spatial autoregressive coefficient ρ is positive and significant in only the 2001 and 2016 elections (Figure 4.1), meaning that in these elections, an electorate's political climate was affected by the attitudes of its neighbours. Conversely, in the other four elections, the spatial effect weakens to become insignificant. In these years, it appears that the spatial component does not explain anything over and above the electoral socio-demographics, meaning electorates voted effectively independently.

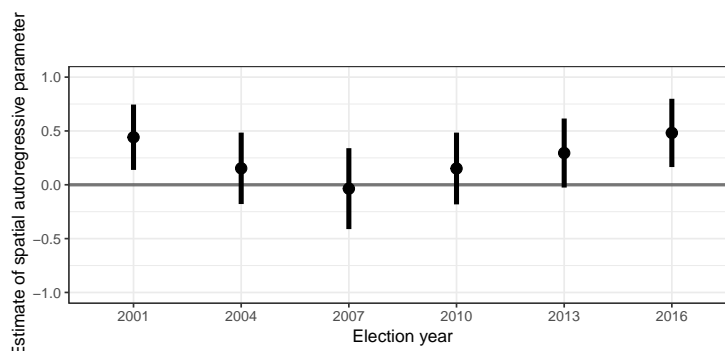


Figure 4.1: Estimates of the spatial autoregressive parameter for each of the six elections, with a 95% confidence interval. Only in 2001 and 2016 is there a significant spatial component

4.2 Country-wide trend

Since all socio-demographics have been standardized to have a mean of zero and a variance of one, the intercept in each model can be interpreted as the estimated two-party preferred vote for an electorate with mean characteristics¹. The baseline of party preference has varied over the elections, with the biggest swing occurring in the 2007 election where the mean electorate shifted more than five percentage points in favour of the Labor party.

¹Mean of all variables aside from Judaism, Indigenous, Islam and Buddhism, where it assumes the mean of the log value.

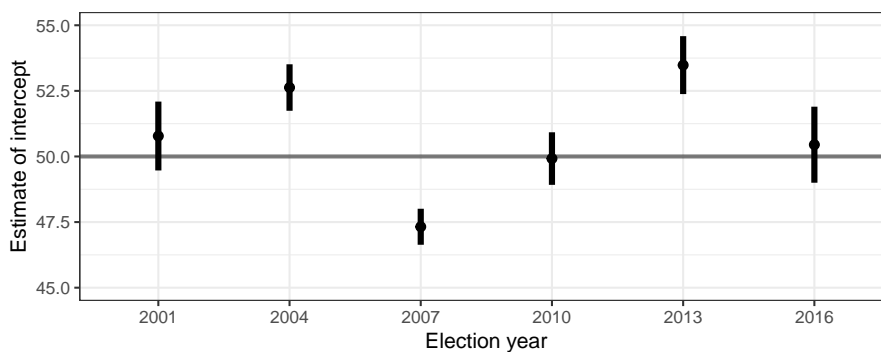


Figure 4.2: Estimated intercept for each election, which represents the two-party preferred vote for an electorate with mean characteristics.

4.3 Influential socio-demographics

To investigate the socio-demographics that have a strong effect on the two-party preferred vote, partial residual plots are used. These show the direction, size and significance of an estimated effect - the slope of the prediction line matches the estimated coefficient, and the shaded region represents a 95% confidence band. If a horizontal line can be drawn through the confidence band, then the effect is insignificant. Plots for each election are faceted to compare the effects over time in Figure 4.3 and Figure 4.4.

It is important here to note the ecological fallacy - insights are being drawn at the electorate level, and cannot be inferred for another disaggregate level (e.g. individual voters).

4.3.1 Income and unemployment

Typically the Labor party campaigns on more progressive policies, which often include tax reform that adversely affects higher income earners, and more generous social assistance programs. Perhaps due to these policies, higher income electorates have been more likely to support the Liberal party, as the **Incomes** factor has a positive effect on Liberal preference (see row 1 in Figure 4.3). This effect is significant in every election aside from 2004, in which it is only marginally insignificant ($p = 0.0613$). Unemployment however, has not been as influential. In 2001 and 2004, electorates with higher unemployment were more likely to support Labor, but over time this has shifted towards support for the Liberal party, culminating in a significantly positive effect in 2016.

4.3.2 Industry and type of work

Electorates with higher proportions of workers in mining, gas, water, agriculture, waste and electricity (grouped as **Extractive** industries) are consistently linked with higher support for the Liberal party, with the magnitude of this effect slightly increasing over the years (see row 3 in Figure 4.3). This is unsurprising, as the Liberal party has close ties with these traditional energy industries, and typically present policies to reduce taxation on energy production. Furthermore, electorates with more workers in construction or manufacturing industries (**Transformative**) are also more likely to support the Liberal party (see row 4 in Figure 4.3).

Similarly, workers in managerial, administrative, clerical and sales roles (**ManagerAdminClericalSales**) is also a significant predictor of two-party preference vote across all six elections, with a higher proportion of people working these jobs increasing Liberal support. The magnitude of this effect has also increased over the years.

4.3.3 Household mobility

In each of the six elections, electorates with a higher proportion of people that have recently moved house (meaning in last five years) were more likely to support the Liberal party, although this effect was marginally insignificant in 2013 (see row 6 in Figure 4.3). Having controlled for characteristics of house ownership and rental prices (via the factors `PropertyOwned` and `RentLoan` respectively), this effect is somewhat surprising.

4.3.4 Relationships

De facto relationships, but not marriages, are found to be an important (and significant) predictor of the two-party preferred vote in all six elections, with more de facto relationships associated with higher support for the Labor party. Marriages however, are insignificant.

4.3.5 Age

Older regions are often believed to be more conservative, and it can be seen that electorates with a higher median age have been more likely to support the Liberal party - although this effect is only significant in 2007 and 2010 (see row 2 in Figure 4.4).

4.3.6 Education

Since 2007, electorates with higher education levels have been associated with supporting the Labor party, although this effect is only significant in 2016. Prior to 2007, education had an almost zero effect (see row 3 in 4.3).

4.3.7 Diversity

Larger migrant populations from Asia, the Middle East, South-Eastern Europe, the United Kingdom and elsewhere, are either associated with Labor support, or have no effect. Of these areas, only South-Eastern European populations were significant in each election and Asian migrants were a significant in 2010. Speaking other language (aside from English) however, have a far stronger effect, as observed through `OtherLanguageHome`. Electorates with more diverse speech were more likely to support the Liberal party from 2004 onwards, with this effect being significant in 2007, 2010 and 2016. Furthermore, of the variables relating to religion, only Judaism shows a consistent effect, with electorates with larger Jewish populations more likely to vote Liberal.

4.3.8 Partial residual plots

4.3.9 A note on similar variables

Many of the Census variables represent similar information, which is why factors were created and some variables were removed. However, there still remain some variables which are closely related. For example, electorate income levels (via `Incomes`) is likely to be related to electoral unemployment and labor force participation (via `Unemployment`). In 2001, the coefficient estimate for `Unemployment` negative but not significant, and the `Incomes` variable is significant. If the `Incomes` variable is removed from the model, `Unemployment` absorbs the negative effect, becoming significant ($p = 0.0056$).

4.4 A closer look at the residuals

4.4.1 Residuals by state

It is often hypothesized that states have systematic differences that cause their electorates to vote differently. Boxplots of residuals grouped by state reveal that only Tasmania, the Australian Capital Territory and the Northern Territory appear to have a state-specific effect that is not captured by the models. Tasmania and the Australian Capital Territory appear to have a bias towards Labor, whereas the Northern Territory has one towards Liberal. There are few electorates in these states (five, two and two respectively), so this might be due to incumbent effects rather than an actual state-specific bias.

4.4.2 Outlier electorates

Based on the distribution of the cook's distance values, a cook's distance greater than 0.1 is considered to be influential and a potential outlier. The electorate of Sydney (NSW) has a large cook's distance from 2001 to 2013, due to its diverse population (language, birthplace and religion), high number of defacto relationships, high income, high household mobility and small amount of workers in extractive and transformative jobs. It has remained a strong supporter of the Labor party and Liberal vote is severely overpredicted by the model, making it an outlier. Nearby in metropolitan NSW, the electorate of Wentworth is found to be an outlier in all but the 2007 election. Although historically Liberal, its two-party vote jumped by over 10 percentage points in 2010 without experiencing any notable changes in its socio-demographic makeup - implying that this may be the direct effect of its Liberal member, Malcolm Turnbull, becoming the leader of the Liberal party. Liberal support in Wentworth is underpredicted by the model in each year, and more so with Turnbull as Liberal leader.

Lingiari, an electorate taking up almost all of the Northern Territory, is an outlier in the 2001-2007 elections due to its large Indigenous population, young age profile and low rates of property ownership. Fowler (NSW) has a diverse population with a high proportion of migrants, many Buddhists and Muslims, and has very strong Labor support, making it influential in 2001, 2004 and 2010. Other electorates with large cook's distance are Barton (NSW) and Leichhardt (QLD) in 2016, and Canberra (ACT) in 2007.

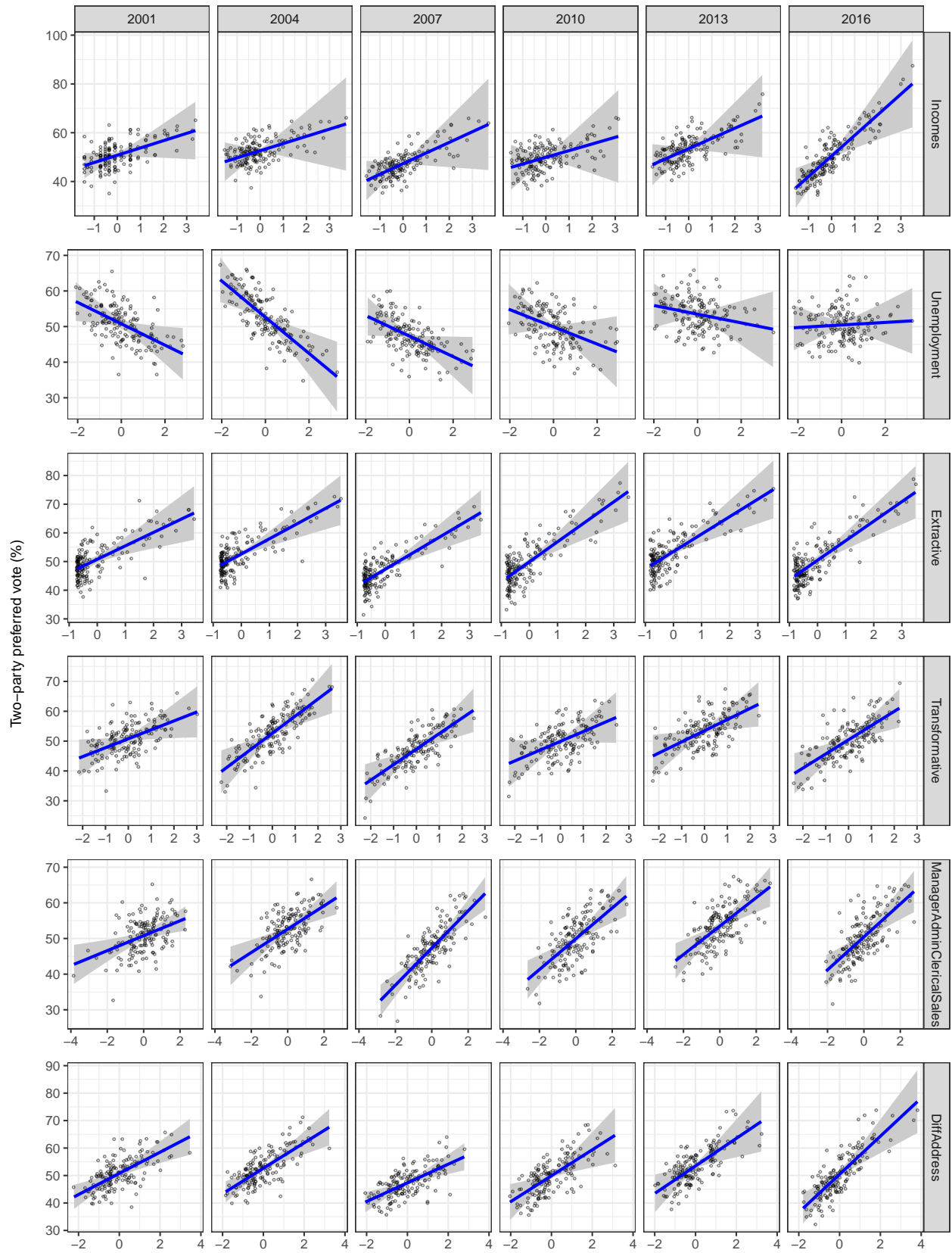


Figure 4.3: Partial residual plots for: income and unemployment, industry and type of work, and household mobility.

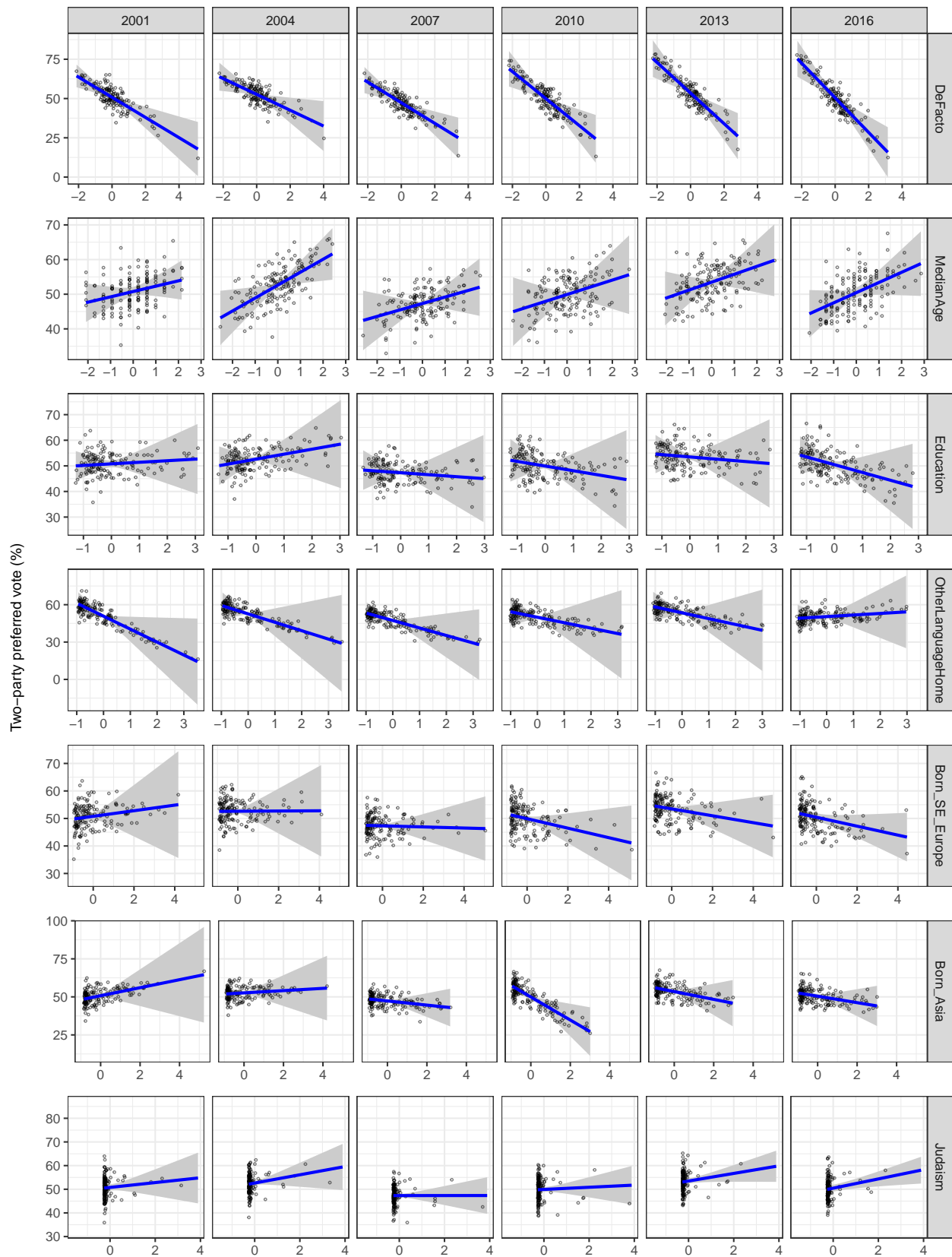


Figure 4.4: Partial residual plots for: relationships and age, education and diversity.

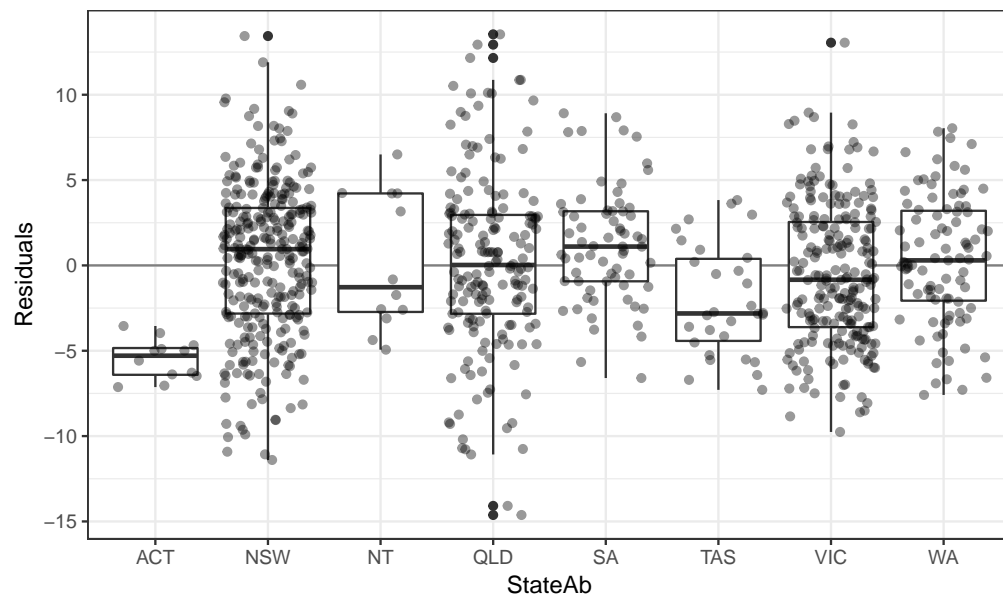


Figure 4.5: Boxplot of residuals by state with jittered points.

Chapter 5

Conclusion

Across the past six elections, most of the socio-demographics that drive the electoral two-party preferred vote have remained steady, whilst a few (typically weaker) effects have varied over time. Industry and type of work are particularly influential, with energy-related and manufacturing/construction jobs, as well as administrative roles being strongly linked with the Liberal party in all elections. Incomes have had a similarly consistent effect, with higher income areas supporting Liberal. Higher levels of unemployment have shifted from weak association with Labor to a significant Liberal effect over the years, and higher education levels have been associated with Labor since 2007 (although only significant in 2016). It is also found that electorates with higher household mobility support Liberal, birthplace diversity favours Labor and more de facto relationships align with Labor preference - although marriages, family and household sizes have no material influence. Furthermore, the neighbourhood (spatial) effects were found to be positive in all elections, although only significant in 2001 and 2016, meaning that in the 2004-2013 elections, electorates effectively voted independently.

The findings in this paper complement the existing literature by modelling temporal trends, which as far as the authors are aware, has not been done previously for Australian elections using a regression framework. It is also the first study to model any Australian election since 2010 using Census information. Additionally, a key contribution of this research is the spatio-temporal imputation, which is applied to obtain Census data for the 2004, 2007, 2010 and 2013 elections. All data sets used in this study have been contributed to the `eechidna` R package, which provides a rich, accessible data resource for future Australian electoral analysis.

Chapter 6

Acknowledgements

This paper was produced using `RMarkdown` (Allaire et al., 2019) and `knitr` (Xie, 2015). All corresponding code for this paper can be found in the github repository (<https://github.com/jforbes14/journal-writing>), and the data used is available in the `eechidna` package (Forbes et al., 2019). All raw data was obtained from the Australian Electoral Commission, the Australian Bureau of Statistics and the Australian Government.

Chapter 7

Software

All election and Census data sets, along with electoral maps and more, are available in the **eechidna** R package, which can be downloaded from CRAN. **eechidna** (Exploring Election and Census Highly Informative Data Nationally for Australia) makes it easy to look at the data from the Australian Federal elections and Censuses that occurred between 2001 and 2016. This study contributed a large revision to the **eechidna** package, which included the addition of election and Census data for 2001-2010, vote outcomes for polling booths and imputed Census data for election years. For more details on using **eechidna**, please see the articles (vignettes) on the github page (<https://ropenscilabs.github.io/eechidna/>).

The authors would like to give a sincere thanks to Anthony Ebert, Heike Hofmann, Thomas Lumley, Ben Marwick, Carson Sievert, Mingzhu Sun, Dilini Talagala, Nicholas Tierney, Nathaniel Tomasetti, Earo Wang and Fang Zhou, all of whom have contributed to the **eechidna** package.

Bibliography

- Allaire, J., Xie, Y., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., Wickham, H., Cheng, J., Chang, W., and Iannone, R. (2019). *rmarkdown: Dynamic Documents for R*. R package version 1.12.
- Anselin, L. (1988). *Spatial econometrics: methods and models*, volume 4. Springer Science & Business Media.
- Davis, R. and Stimson, R. (1998). Disillusionment and disenchantment at the fringe: explaining the geography of the one nation party vote at the queensland election. *People and place*, 6(3):69–82.
- Forbes, J., Cook, D., Ebert, A., Hofmann, H., Hyndman, R., Lumley, T., Marwick, B., Sievert, C., Sun, M., Talagala, D., Tierney, N., Tomasetti, N., Wang, E., and Zhou, F. (2019). *eechidna: Exploring Election and Census Highly Informative Data Nationally for Australia*. R package version 1.3.0.
- Forrest, J., Alston, M., Medlin, C., and Amri, S. (2001). Voter behaviour in rural areas: a study of the farrer electoral division in southern New South Wales at the 1998 federal election. *Australian geographical studies*, 39(2):167–182.
- Goodchild, M. F., Anselin, L., and Deichmann, U. (1993). A framework for the areal interpolation of socioeconomic data. *Environment and planning A*, 25(3):383–397.
- LeSage, J., Kelley Pace, R., and Pace, R. K. (2009). *Introduction to spatial econometrics*. Chapman and Hall/CRC.
- Liao, E., Shyy, T., and Stimson, R. (2009). Developing a web-based e-research facility for socio-spatial analysis to investigate relationships between voting patterns and local population characteristics. *Journal of spatial science*, 54(2):63–88.
- Stimson, R., McCrea, R., and Shyy, T. (2006). Spatially disaggregated modelling of voting outcomes and socio-economic characteristics at the 2001 australian federal election. *Geographical research*, 44(3):242–254.
- Stimson, R. and Shyy, T. (2012). And now for something different: modelling socio-political landscapes. *Ann Reg Sci*, 50:623–643.
- Stimson, R. and Shyy, T.-K. (2009). A socio-spatial analysis of voting for political parties at the 2007 federal election. *People and place*, 17(1):39–54.
- Wickham, H., Bryan, J., Kalicinski, M., Valery, K., Leitenne, C., Colbert, B., Hoerl, D., and Miller, E. (2019a). *readxl: Read Excel Files*. R package version 1.3.1.
- Wickham, H., François, R., Henry, L., and Müller, K. (2019b). *dplyr: A Grammar of Data Manipulation*. R package version 0.8.0.1.
- Xie, Y. (2015). *Dynamic Documents with R and knitr*. Chapman and Hall/CRC, Boca Raton, Florida, 2nd edition. ISBN 978-1498716963.