

# How who we are affects how we vote

*Rob J Hyndman and Di Cook*

## 1 Introduction

We often hear stereotypes about voting patterns — people in their 20s without kids are more likely to be left wing, migrants are more conservative, wealthier people tend to favour the conservative parties, and so on. So we thought we would test these ideas by matching census data to election data, to see if we can identify what are the socio-demographic characteristics of an electorate that are most closely related to how they vote.

## 2 Previous intro

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 socio-demographic characteristics relate to two-party preference, and whether their effects have changed over time.

Data on electoral socio-demographics are derived from the Australian Census, and vote counts are obtained from Australian federal elections. Joining these two data sources is problematic as there is an inherent asynchronicity in the two types of events. A Census is conducted by the Australian Bureau of Statistics (ABS) every five years, whereas federal elections, conducted by the Australian Electoral Commission (AEC), usually occur every three years or so. The first problem addressed is that of constructing 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 was to join voting outcomes to the nearest Census, without accounting for any temporal differences (see Davis and R. Stimson 1998; R. Stimson, McCrea, and T. Shyy 2006; Liao, T. Shyy, and R. Stimson 2009; Robert Stimson and T.-K. Shyy 2009). Furthermore, electoral boundaries change regularly, so spatial discrepancies also arise when matching with electoral data. To obtain appropriate “Census-like” data for these four elections, electoral socio-demographics are constructed using a spatio-temporal imputation that combines areal interpolation (Goodchild, Anselin, and Deichmann 1993) and linear time-interpolation. Collecting and wrangling the raw data, along with the imputation process, are detailed in Section 3. 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 has found that aggregate socio-demographics are relatively good predictors of voting outcomes. Forrest et al. (2001) used multiple regression to model the Liberal and Labor primary vote for polling booths in the Farrer electorate in 1998 as a function of Census variables from 1996. R. Stimson, McCrea, and T. Shyy (2006), Robert Stimson and T.-K. Shyy (2009) and R. Stimson and T. Shyy (2012) used 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. In contrast, Robert Stimson and T.-K. Shyy (2009) models the polling booth swing vote (change in the two-party preferred vote) in the 2007 election, finding that little of the 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. Each federal election is modelled with a cross-sectional dataset. The cross-sectional dataset for each election used here consists of the two-party preferred vote (as the response variable), and a set of common socio-demographic variables (as

the explanatory variables) that characterize each electorate. To prepare these datasets, socio-demographic variables are first standardized, and then a principal component analysis is used to group variables into “factors”. To account for the inherent spatial structure of the data, a spatial error model is then estimated for each election.

The paper is organised as follows. Section 3 describes the data collection, joining and cleaning. These pre-processing steps and model details are discussed in Section 4. Section 5 describes the inference conducted to determine significance of effects and how these change over time. Section 6 summarises the work. Two supplementary sections document the contributions of others to this work and the software.

## 3 Data collection, wrangling and imputation

### 3.1 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 divides 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 were derived from the Australian 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 so far in the 21st century, conducted in 2001, 2006, 2011 and 2016. The Australian Bureau of Statistics (ABS) conducts the Census and publishes aggregated information. The ABS uses electoral boundaries as defined by the AEC at the time of each Census, which may not match those in place at the subsequent and previous elections. From the available Census information aggregated at the electorate level, 65 socio-demographic variables were defined for each of the electorates to be used in the analysis.

Raw data was sourced online from the AEC and ABS websites in `.csv` and `.xlsx` files. The formats of these files differ over the years, making extracting the appropriate information a big task. The functions available in the `dplyr` (Wickham, François, et al. 2019) and `readxl` (Wickham, Bryan, et al. 2019) R packages are particularly useful, as they provide fast consistent tools for data manipulation and functions to import `.xlsx` files (respectively). The 2001 and 2006 Census data are published in a format where the information for each electorate is held in a separate document making it difficult to use the `dplyr` tools. 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.

### 3.2 Joining Census and election data

#### 3.2.1 Differences between Census and election data

Between 2001 and 2016 there were six elections and four Censuses (see Figure 1). Electoral boundaries are redistributed regularly by the AEC, meaning that only in the years where both a Census and election occur are all boundaries likely to match — the case for the 2001 and 2016 elections. 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 was employed to obtain electoral socio-demographics. This method uses Census information from both before and after the election of interest.

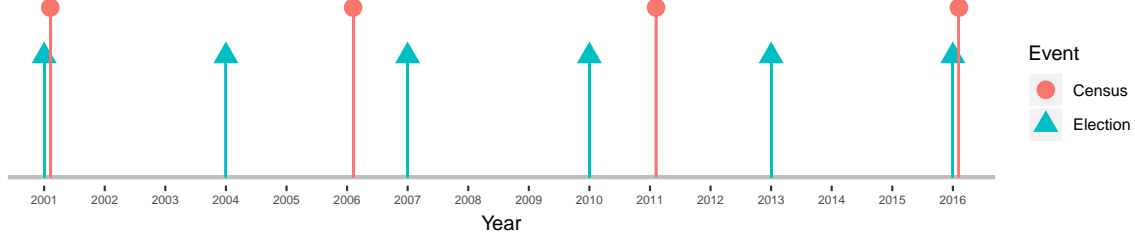


Figure 1: Timeline of Australian elections and Censuses. They do not always occur in the same year.

### 3.2.2 Spatio-temporal imputation

For each election, neighbouring Census information has to be combined in some way so that it represents the boundaries in place at the time of the election. This is done by taking the electoral boundaries and imputing the corresponding socio-demographic characteristics for each of neighbouring Censuses, thereby addressing the spatial aspect. Next, to deal with the temporal component, characteristics at the time of the election are constructed using linear interpolation between the spatially imputed neighbouring Census variables.

To account for spatial differences, the piecewise approximation method in Goodchild, Anselin, and Deichmann (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. This is described in the context of a single election, and a single neighbouring Census.

Let the map of electoral boundaries at the time of a Census define the source zones, and let the boundaries at the time of the election be the target zones. Denote the area of intersection between a source zone  $s$  and a target zone  $t$  as  $A_{s,t}$ . Additionally, let the population of the source zone  $s$  be  $U_s$  and the population of intersection between source zone  $s$  and target zone  $t$  be  $P_{s,t}$ . The estimated population of intersection is given by

$$\hat{P}_{s,t} = \frac{U_s * A_{s,t}}{\sum_{t=1}^T A_{s,t}}, \quad \text{for all } s = 1, \dots, S \text{ and } t = 1, \dots, T.$$

Note that this estimator implicitly 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 intersection populations as weights. Denote a given Census variable for the target zone by  $C_t$ , and the same Census variable for the source zone as  $D_s$ . Then, estimate  $C_t$  using

$$\hat{C}_t = \frac{\sum_{s=1}^S D_s * \hat{P}_{s,t}}{\sum_{s=1}^S \hat{P}_{s,t}}, \quad \text{for each } t = 1, \dots, T.$$

This concludes the spatial imputation of the socio-demographic characteristics for one target zone (a single electoral boundary), at the time of only one of the neighbouring Censuses. This process is repeated for all of the target zones, and then for the other neighbouring Census.

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. Let  $y_1$  be the year of the Census preceding an election, let  $y_2$  be the year of the election, and  $y_3$  be the year of the Census that follows. Add this year subscript to the Census variable estimate  $\hat{C}_t$ , resulting in  $\hat{C}_{t,y}$ . Linear interpolating between these Census years results an imputed value for the election year, given by

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

Implicitly this assumes that population characteristics change in a linear manner over time.

### An illustration of the spatio-temporal imputation

Census data is publicly available at different levels of aggregation, ranging from SA1 (over 50,000 zones) to electoral divisions (150 zones). For this study, electoral divisions are used as source zones, and the imputation

method is applied to produce socio-demographic variables for each of the 2004, 2007, 2010 and 2013 elections. As mentioned earlier, there is no need to impute socio-demographic variables for the 2001 and 2016 elections. To illustrate the method, 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. The boundaries shown in Figure 2 define all target zones in NSW for 2013, with the target zone of interest (Hume) shaded purple.

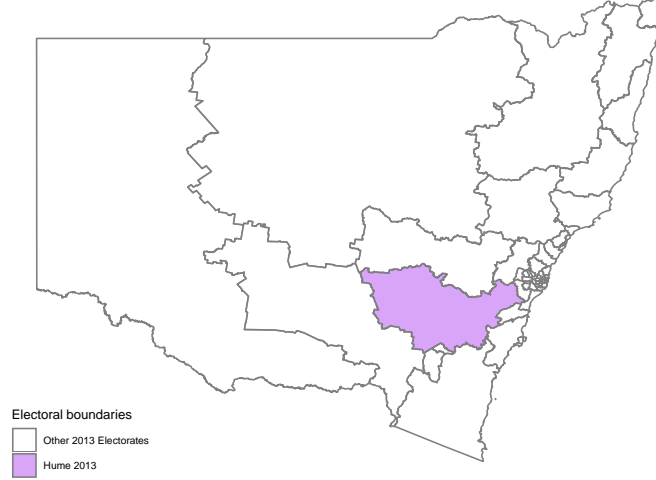


Figure 2: Some of the electoral boundaries in NSW for 2013, with the electoral boundary for Hume shown in purple.

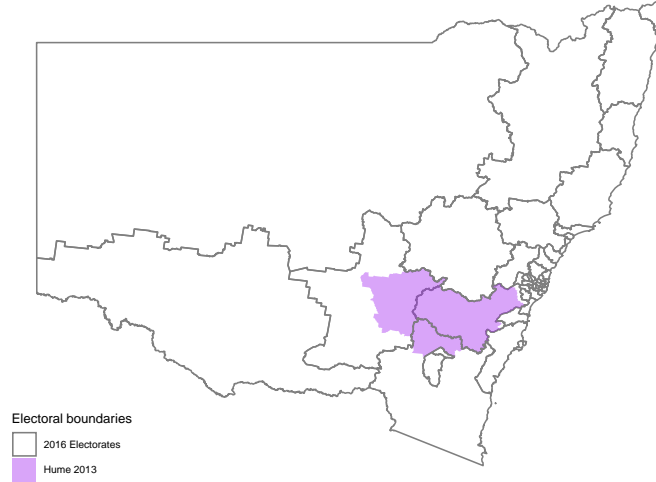


Figure 3: Census division boundaries in NSW for 2016, with the 2013 electoral boundary for Hume, shown in purple. The purple region is not contained within a single Census division.

The corresponding source zones from the 2016 Census are shown in Figure 3. As can be seen, the Hume boundary from the 2013 election (shaded purple) does not perfectly match any of the source zones.

There are many source zones from the 2016 Census that intersect with this purple region, including the divisions of Riverina, Eden-Monaro and Hume, along with smaller intersecting areas with Fenner, Calare, Gilmore and Whitlam. The proportion of each source zone that overlaps with the purple region is calculated, and used to obtain the intersecting populations  $\hat{P}_{s,t}$ .

Now consider the socio-demographic variable *AusCitizen*, the proportion of people in the region who are Australian citizens. A weighted average of *AusCitizen*, with the allocated population from each source zone as

Table 1: Population from each intersecting Census source zone (2016) that is allocated to the target zone (purple region - Hume electoral boundary in 2013), and the corresponding proportion of Australian citizens in each of these source zones.

| Source zone (2016) | Proportion | Source zone population | Population allocated to purple region: $\hat{P}_{s,t}$ | AusCitizen (%): $D_s$ |
|--------------------|------------|------------------------|--|-----------------------|
| Hume               | 0.9654     | 150643                 | 145427   | 90.0168               |
| Riverina           | 0.2511     | 155793                 | 39117  | 89.1144               |
| Eden-Monaro        | 0.1109     | 147532                 | 16358  | 87.9999               |
| Canberra           | 0.0028     | 196037                 | 548  | 85.4793               |
| Fenner             | 0.0023     | 202955                 | 474  | 83.6432               |
| Whitlam            | 0.0006     | 152280                 | 92   | 89.5173               |
| Gilmore            | 0.0006     | 150436                 | 86   | 89.0266               |
| Calare             | 0.0001     | 161298                 | 21   | 87.5603               |

weights, yields  $\hat{C}_{\text{Hume},2016} = 89.65\%$ . Repeating this process using the 2011 Census yields  $\hat{C}_{\text{Hume},2011} = 91.00\%$ . Finally, linear interpolation between 2011 and 2016 yields the 2013 estimate:

$$\hat{C}_{\text{Hume},2013} = \frac{3}{5}\hat{C}_{\text{Hume},2011} + \frac{2}{5}\hat{C}_{\text{Hume},2016} = 90.46\%.$$

This is done for each of the 65 socio-demographic variables, and is repeated for each of the 149 remaining target zones corresponding with 2013 electorates.

## 4 Modelling

Following this process, electoral socio-demographic variables are available for each of the six elections and can be joined with their corresponding two-party preferred votes. Before choosing an appropriate model, two issues with the socio-demographic variables need to be addressed. First, variable scales change over the years, making it important to standardize variables. Second, many variables represent similar information and where appropriate, should be combined in some way. Principal component analysis is used to identify variables covary, from which intuitive groupings are selected to be combined into a single variable. This also reduces the dimension of the data. After these steps, a model specification is chosen.

### 4.1 Standardizing variables

Many of the socio-demographic variables have changing scales over the years. 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. In order for socio-demographic effects to be comparable across years, all explanatory variables are standardized to have mean zero and variance one within each election year. By standardizing, each variable is reported as a relative measure compared to all other electorates in the same year.

### 4.2 Creating factors

There are only  $N = 150$  observations (electorates) in each election and  $p = 65$  socio-demographic variables in each cross-section, with many variables represent similar information about an electorate. Any model that uses all variables would face serious problems with multi-collinearity and over-fitting, which would likely lead to erroneous conclusions regarding variable significance. To address this, groups of variables that represent similar information are combined into “factors”<sup>1</sup>.

A factor is created from a group of variables if there is an intuitive reason as to why they should represent similar information and if there is evidence to suggest that they covary. For example, a potential group would be variables relating to electoral incomes — median family, household and personal incomes. To determine

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

which variables covary, principal component analysis is used on a combined dataset of socio-demographic variables from all six elections. It is appropriate to compute principal components in this way because when computed separately for each election, scree plots level off after four components and the loadings of the first four components are similar across the elections.

Only the first four principal components from the combined dataset are considered, as the scree plot corresponding to the combined dataset levels off after the fourth component. Variables that have a large loading in a particular component are deemed to covary, with a loading with magnitude greater than 0.15 being considered large. Six factors are created using this criteria. These are: **Incomes** (median personal, household and family incomes); **Unemployment** (unemployment and labor force participation rates); **PropertyOwned** (rates of housing ownership, mortgages, renting and government housing); **RentLoanPrice** (median rental and loan repayments); **FamHouseSize** (average household size, ratio of people to families and household makeup (single person, couple with kids and couple without kids); and **Education** (high school and university qualifications, jobs requiring higher levels of education as well as vocational course completions and jobs that do not require higher education levels, such as laborer or tradesperson). For each of these groups, variables with positive loadings are added and those with negative loadings are subtracted to create a factor. After computing these sums, each factor is standardized to have mean zero and variance one, within each election.

There are  $p = 30$  variables in the resultant predictor set, with all of these used in the regression for each election.

### 4.3 Regression incorporating spatially dependent errors

An identical model specification is used for each of the six elections, with each election modelled separately. This allows the socio-demographic effects to be estimated separately for each election year, facilitating analysis of temporal changes in variable effects. This approach is preferable to using a single longitudinal model because it avoids any concerns of undue bias stemming from incorrectly imposed time-varying restrictions 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, noting that electorates with 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  $\mathbf{y}$  be the vector two-party preferred vote in favour of the Liberal party; for example,  $y_i = 70$  represents a 70% preference for Liberal, 30% for Labor, in electorate  $i$ . Although  $y_i$  lies in the interval  $(0, 100)$ , observed values are never close to 0 or 100 (minimum 24.05% and maximum 74.90%), so there is no need to formally impose the constraint of  $y_i \in [0, 100]$ . Furthermore, the responses are found to be spatially correlated in each election (Moran's I test,  $p \leq 7 \cdot 10^{-15}$ ). This is not surprising as electorates are aggregate spatial units, and hence the spatial structure of the data must be modelled appropriately.

The spatial error model (Luc Anselin 1988) is chosen because it captures spatial heterogeneity by incorporating a spatially structured random effect vector (LeSage, Kelley Pace, and Pace 2009). In this context, the random effect can be thought of as capturing the unobserved political climate in each electorate, where this climate is correlated with the climate in neighbouring electorates, under the assumption that the climate is independent of electoral socio-demographics.

Spatial weights are calculated in accordance with the assumption that an electorate is equally correlated with any electorate that shares a part of its boundary. Let  $\rho$  be the 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  $\mathbf{v} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ , and hence

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

Table 2: Estimated spatial error model parameters (standard errors) for each of the six election years.

|                               | 2001               | 2004               | 2007               | 2010               | 2013               | 2016               |
|-------------------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| $\rho$                        | 0.46***<br>(0.15)  | 0.29*<br>(0.17)    | 0.24<br>(0.17)     | 0.19<br>(0.16)     | 0.27*<br>(0.16)    | 0.50***<br>(0.17)  |
| AusCitizen                    | -3.13<br>(2.26)    | -2.64<br>(2.43)    | -2.53<br>(2.34)    | -0.08<br>(2.79)    | -3.40<br>(2.76)    | -1.80<br>(2.71)    |
| BornAsia                      | 2.22<br>(2.18)     | -0.95<br>(2.44)    | -1.60<br>(2.19)    | -6.83**<br>(2.73)  | -3.03<br>(2.71)    | -0.55<br>(2.17)    |
| Born_MidEast                  | -1.15<br>(1.07)    | -1.59<br>(1.20)    | -2.01*<br>(1.11)   | -2.03<br>(1.27)    | -0.92<br>(1.24)    | -1.44<br>(1.13)    |
| BornSEEuro                    | -3.21**<br>(1.42)  | -4.24***<br>(1.46) | -3.61***<br>(1.02) | -4.14***<br>(1.19) | -3.69***<br>(1.07) | -2.72***<br>(0.97) |
| Born_UK                       | 0.25<br>(1.00)     | -0.07<br>(0.98)    | 0.34<br>(0.90)     | 0.56<br>(1.07)     | -0.09<br>(1.04)    | -1.32<br>(1.04)    |
| BornElsewhere                 | -5.04<br>(3.30)    | -4.91<br>(3.68)    | -4.13<br>(3.38)    | 2.35<br>(4.23)     | -5.23<br>(4.15)    | -4.14<br>(3.97)    |
| Buddhism                      | -0.49<br>(1.39)    | -0.17<br>(1.61)    | -1.37<br>(1.61)    | -0.83<br>(1.80)    | -0.12<br>(1.68)    | -1.60<br>(1.56)    |
| Christianity                  | -2.48<br>(1.73)    | -1.23<br>(1.85)    | 0.38<br>(1.83)     | 0.50<br>(1.99)     | 2.41<br>(1.85)     | 1.68<br>(1.78)     |
| CurrentlyStudying             | -2.19**<br>(0.99)  | -0.13<br>(1.13)    | 2.06*<br>(1.17)    | 2.12*<br>(1.25)    | 1.15<br>(1.26)     | -0.16<br>(1.18)    |
| DeFacto                       | -6.44***<br>(1.87) | -5.37**<br>(2.48)  | -6.43***<br>(2.31) | -8.07***<br>(3.06) | -6.56**<br>(3.11)  | -8.53***<br>(2.83) |
| DiffAddress                   | 3.88***<br>(0.94)  | 5.06***<br>(1.12)  | 4.22***<br>(0.99)  | 5.57***<br>(1.76)  | 3.53*<br>(1.91)    | 5.67***<br>(1.60)  |
| Distributive                  | 1.27<br>(1.12)     | 2.01*<br>(1.21)    | 1.36<br>(1.13)     | 1.57<br>(1.34)     | 2.10*<br>(1.27)    | 1.20<br>(1.21)     |
| Education                     | 1.08<br>(2.38)     | 0.52<br>(3.12)     | -5.52*<br>(3.27)   | -4.08<br>(3.95)    | -4.44<br>(3.78)    | -7.07**<br>(3.55)  |
| Extractive                    | 4.83***<br>(1.48)  | 5.45***<br>(1.42)  | 5.37***<br>(1.36)  | 7.31***<br>(1.56)  | 6.71***<br>(1.47)  | 7.43***<br>(1.39)  |
| FamHouseSize                  | -0.16<br>(2.19)    | 0.87<br>(2.72)     | -2.40<br>(2.69)    | -2.53<br>(3.25)    | -3.26<br>(3.28)    | -2.91<br>(2.90)    |
| Incomes                       | 4.36**<br>(1.77)   | 5.03*<br>(2.66)    | 9.45***<br>(2.75)  | 7.09**<br>(3.25)   | 7.97***<br>(2.92)  | 12.20***<br>(2.75) |
| Indigenous                    | 2.91*<br>(1.68)    | 1.97<br>(1.95)     | 2.48<br>(1.75)     | 2.84<br>(2.16)     | 0.67<br>(2.14)     | -0.05<br>(2.00)    |
| Islam                         | -0.92<br>(1.22)    | -0.97<br>(1.36)    | -0.54<br>(1.27)    | -2.50<br>(1.52)    | -0.82<br>(1.42)    | -0.95<br>(1.34)    |
| Judaism                       | 1.88*<br>(1.05)    | 1.78<br>(1.13)     | 2.66***<br>(1.01)  | 1.97*<br>(1.15)    | 2.74**<br>(1.10)   | 1.65*<br>(1.00)    |
| ManagerAdmin                  | 2.06***<br>(0.71)  | 3.32***<br>(0.93)  | 6.00***<br>(0.90)  | 5.47***<br>(1.08)  | 5.04***<br>(1.03)  | 5.78***<br>(1.06)  |
| Married                       | 0.44<br>(2.31)     | 0.11<br>(2.96)     | -1.22<br>(2.83)    | -0.22<br>(3.15)    | 0.91<br>(3.03)     | -2.34<br>(2.81)    |
| MedianAge                     | 2.32*<br>(1.32)    | 4.96***<br>(1.65)  | 3.66**<br>(1.81)   | 4.00*<br>(2.26)    | 2.30<br>(2.08)     | 2.87<br>(1.79)     |
| NoReligion                    | -1.57<br>(1.59)    | -0.92<br>(1.71)    | 0.56<br>(1.73)     | -0.30<br>(1.92)    | 1.02<br>(1.94)     | 1.31<br>(2.04)     |
| OneParentHouse                | -1.73<br>(1.36)    | -0.45<br>(1.59)    | -0.75<br>(1.49)    | -1.46<br>(1.69)    | -0.77<br>(1.57)    | -0.74<br>(1.47)    |
| OtherLanguage                 | -0.44<br>(3.22)    | 5.92<br>(4.16)     | 9.98**<br>(3.91)   | 11.24**<br>(4.76)  | 9.00*<br>(4.66)    | 9.84**<br>(4.44)   |
| PropertyOwned                 | -0.46<br>(1.37)    | -0.53<br>(1.50)    | 0.67<br>(1.43)     | -0.94<br>(1.76)    | -0.48<br>(1.67)    | 1.41<br>(1.50)     |
| RentLoanPrice                 | -1.57<br>(1.49)    | -3.32*<br>(1.76)   | -4.01**<br>(1.67)  | -0.97<br>(2.07)    | -0.70<br>(2.07)    | -0.89<br>(2.16)    |
| SocialServ                    | 2.51*<br>(1.33)    | 1.65<br>(1.41)     | 2.47*<br>(1.29)    | 2.53*<br>(1.47)    | 2.35*<br>(1.32)    | 4.45***<br>(1.19)  |
| Transformative                | 3.24**<br>(1.55)   | 4.73***<br>(1.78)  | 4.84***<br>(1.74)  | 4.46**<br>(1.98)   | 3.56**<br>(1.78)   | 4.58***<br>(1.53)  |
| Unemployment                  | -2.45*<br>(1.40)   | -3.07*<br>(1.63)   | 0.29<br>(1.51)     | 0.08<br>(1.76)     | 1.67<br>(1.51)     | 2.79**<br>(1.37)   |
| Constant                      | 50.81***<br>(0.71) | 52.60***<br>(0.58) | 47.31***<br>(0.52) | 49.93***<br>(0.57) | 53.51***<br>(0.58) | 50.49***<br>(0.80) |
| Observations                  | 150                | 150                | 150                | 150                | 150                | 150                |
| Residual Standard Error (GLS) | 4.69               | 5.04               | 4.79               | 5.63               | 5.18               | 4.88               |

Note:

\* p&lt;0.1; \*\* p&lt;0.05; \*\*\* p&lt;0.01

Estimation of the above spatial error model is undertaken using feasible generalized least squares.

Table 2 details the estimated model coefficients and their estimated standard errors, for each of the six elections. An interpretation of these estimated values is provided in the next section.

## 5 Results

### 5.1 Spatial autoregressive parameter

The spatial autoregressive coefficient  $\rho$  is positive and significant in only the 2001 and 2016 elections (Figure 4), meaning that in these elections, the political climate of an electorate appears to be 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 not already explained by the electoral socio-demographics, meaning electorates effectively voted independently.

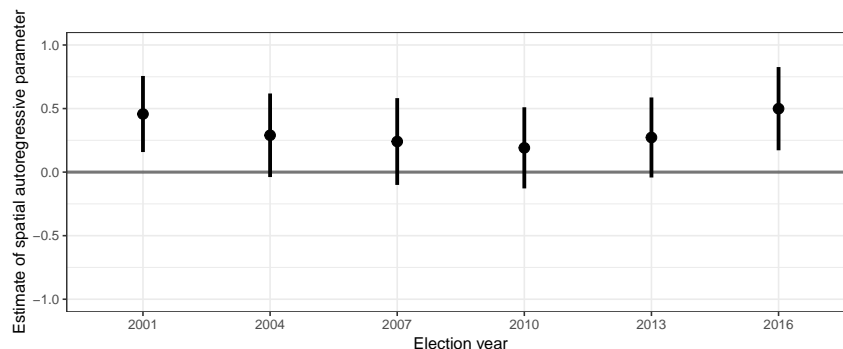


Figure 4: Estimates of the spatial autoregressive parameter for each of the six elections, reported with their individual 95% confidence intervals. Only in 2001 and 2016 is there a significant spatial component.

### 5.2 Country-wide trend

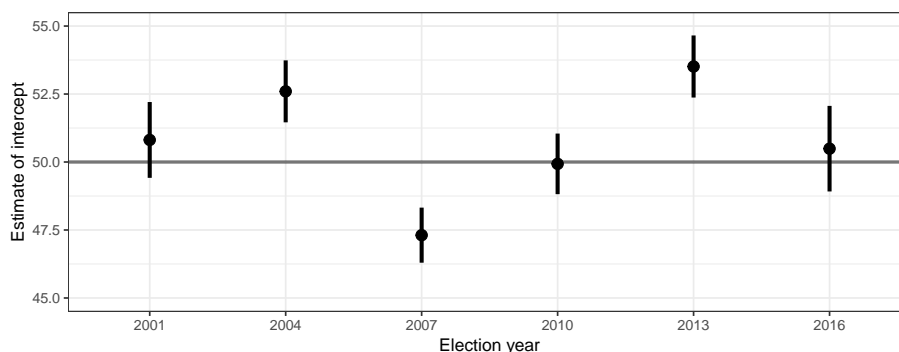


Figure 5: Estimated intercept for each election, which represents the two-party preferred vote for an electorate with mean characteristics.

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 characteristics<sup>2</sup>. Figure 5 shows that 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.

<sup>2</sup>Mean of all variables aside from Judaism, Indigenous, Islam and Buddhism, where it assumes the mean of the log value.



### 5.3 Influential socio-demographics

To investigate the socio-demographics that have a strong effect on the two-party preferred vote, partial residual plots are used and shown in Figures ?? and ?. The partial residuals are the residuals from the fitted model with the estimated effect an individual variable added. 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, computed using the method in Breheny and Burchett (2017). If a horizontal line can be drawn through the confidence band, then the effect is insignificant. The estimated intercept is also added to the partial residuals for interpretability. Plots for each election are faceted to compare the effects over time in Figures ?? and ?. Only socio-demographics that have a significant effect in at least one election are displayed in Figures ?? and ?.

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).

#### 5.3.1 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 it is due to these policies that higher income electorates appear more likely to support the Liberal party, as the **Incomes** factor has a positive effect on Liberal preference (see row 1 in Figure ?). This effect is significant in every election aside from 2004, where it is only marginally insignificant ( $p = 0.0613$ ). Unemployment however, is not as influential. In 2001 and 2004, electorates with higher unemployment align with Labor, but over time this shifts towards support for the Liberal party, culminating in a significantly positive effect in 2016.

#### 5.3.2 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 ?). 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 ?).

Similarly, the proportion of workers in managerial, administrative, clerical and sales roles (**ManagerAdmin**) 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 also seems to increase over the years.

#### 5.3.3 Household mobility

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

#### 5.3.4 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. The proportion of individuals who are married however, is insignificant (not shown).

### 5.3.5 Age

Regions comprising more older people are often believed to be more conservative, and indeed it found that electorates with a higher median age are more likely to support the Liberal party — although this effect is significant only in 2007 and 2010 (see row 2 in Figure ??).

### 5.3.6 Education

Since 2007, electorates with higher education levels are associated with supporting the Labor party, although this effect is significant only in 2016. Before 2007, education has an almost zero effect (see row 3 in Figure ??).

### 5.3.7 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 appear significant in each election, with the proportion of Asian migrants also being significant in 2010. Speaking other languages (aside from English) however, appears to have a far stronger effect, as observed through the `OtherLanguage` variable. Electorates with more diverse speech are associated with higher support for 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 relatively large Jewish populations more likely to vote Liberal.

### 5.3.8 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, some variables remain which are closely related. For example, an electorate's income level (via `Incomes`) is likely to be related to electoral unemployment and labor force participation (via `Unemployment`). In 2001, the coefficient estimate for `Unemployment` is negative but not significant, whilst the `Incomes` variable is significant. If the `Incomes` variable is removed from the model in 2001, `Unemployment` absorbs the negative effect, becoming significant ( $p = 0.0056$ ).

## 5.4 A closer look at the residuals

### 5.4.1 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 (Figure 6) reveal that the data reflects this – there appears to be a state-specific effect not captured by the models. Tasmania and the Australian Capital Territory appear to have a bias towards Labor, whereas the Northern Territory tends towards voting Liberal. However, there are relatively few electorates in each of these states (five, two and two respectively), so this apparent result may be due to incumbent effects rather than an actual state-specific bias.

### 5.4.2 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 the 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.

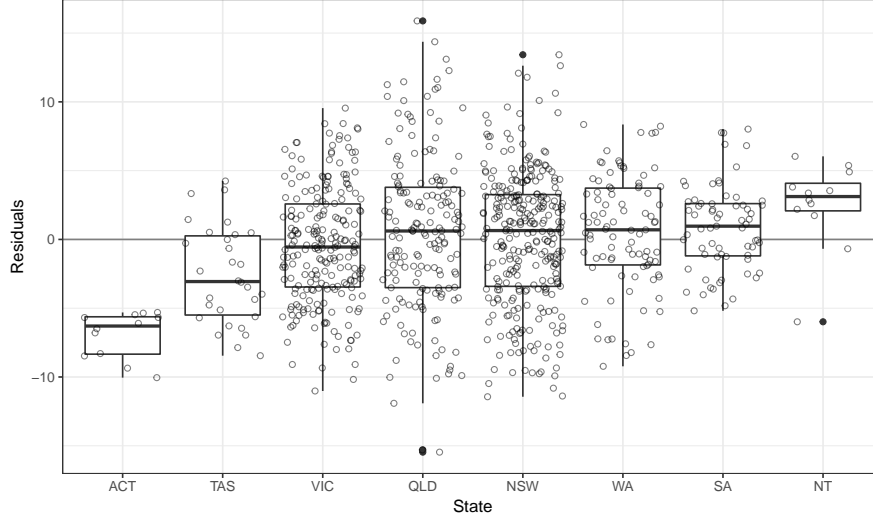


Figure 6: Boxplot of residuals by state with jittered points. States ordered by median residual. A state-specific bias not captured by the model is evident.

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 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.

## 6 Conclusion

This paper explores the effects of electoral socio-demographic characteristics on the two-party preferred vote in the 2001–2016 elections, using information from the corresponding Australian federal elections and Censuses. As a Census does not always occur in the same year as an election, Census data for the 2004–2013 elections are generated by employing a method of spatio-temporal imputation. This imputes electoral socio-demographics for the electoral boundaries in place at the time of the election — an approach that is distinctly different from previous work on modelling election outcomes, where Census and election data are typically joined without addressing their temporal differences. Before estimating a model, these socio-demographic variables are standardized (to adjust for changing variable scales) and many variables (representing similar information) are combined into factors, resulting in a reduced predictor set. A spatial error model is then estimated for each election, accounting for the inherent spatial structure of the data.

Across the past six elections, most of the socio-demographics that drive the electoral two-party preferred vote are found to remain steady, whilst a few (typically weaker) effects vary 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 a similarly consistent effect, with higher income areas supporting Liberal. Higher levels of unemployment shift from weak association with Labor to a significant Liberal effect over the years, and higher education levels are associated with Labor from 2007 (although significant only 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 are found to be positive in all elections, although significant only 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 wrangling of the raw data and imputed data sets for the 2004, 2007, 2010 and 2013 elections, which have been contributed to the **eechidna** R package — providing a rich, accessible data resource for future Australian electoral analysis.

## 7 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 [github.com/jforbes14/eechidna-paper](https://github.com/jforbes14/eechidna-paper), 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.

## 8 Software

All election and Census datasets, along with electoral maps and more, are available in the **eechidna** (Exploring Election and Census Highly Informative Data Nationally for Australia) R package, which can be downloaded from CRAN. The **eechidna** package 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, voting 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 [ropenscilabs.github.io/eechidna/](https://ropenscilabs.github.io/eechidna/).

The authors would like to sincerely thank 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.

## References

- Allaire, JJ et al. (2019). *rmarkdown: Dynamic Documents for R*. R package version 1.12. URL: <https://rmarkdown.rstudio.com>.
- Anselin, Luc (1988). *Spatial econometrics: methods and models*. Vol. 4. Springer Science & Business Media.
- Breheny, Patrick and Woodrow Burchett (2017). “Visualization of regression models using visreg”. In: *The R Journal* 9.2, pp. 56–71. URL: <https://journal.r-project.org/archive/2017/RJ-2017-046/index.html>.
- Davis, R. and R. Stimson (1998). “Disillusionment and disenchantment at the fringe: explaining the geography of the One Nation party vote at the Queensland election”. In: *People and place* 6.3, pp. 69–82.
- Forbes, Jeremy et al. (2019). *eechidna: Exploring Election and Census Highly Informative Data Nationally for Australia*. R package version 1.3.0. URL: <https://CRAN.R-project.org/package=eechidna>.
- Forrest, James et al. (2001). “Voter behaviour in rural areas: a study of the Farrer electoral division in southern New South Wales at the 1998 federal election”. In: *Australian Geographical Studies* 39.2, pp. 167–182.
- Goodchild, M F, L Anselin, and U Deichmann (1993). “A framework for the areal interpolation of socioeconomic data”. In: *Environment and Planning A* 25.3, pp. 383–397.
- LeSage, James, R Kelley Pace, and Robert Kelley Pace (2009). *Introduction to spatial econometrics*. Chapman and Hall/CRC, pp. 50–52.
- Liao, E., T. Shyy, and R. Stimson (2009). “Developing a web-based e-research facility for socio-spatial analysis to investigate relationships between voting patterns and local population characteristics”. In: *Journal of Spatial Science* 54.2, pp. 63–88.
- Stimson, R., R. McCrea, and T. Shyy (2006). “Spatially disaggregated modelling of voting outcomes and socio-economic characteristics at the 2001 Australian federal election”. In: *Geographical Research* 44.3, pp. 242–254.
- Stimson, R. and T. Shyy (2012). “And now for something different: modelling socio-political landscapes”. In: *Annals of Regional Science* 50, pp. 623–643.
- Stimson, Robert and Tung-Kai Shyy (2009). “A socio-spatial analysis of voting for political parties at the 2007 federal election”. In: *People and Place* 17.1, pp. 39–54.

Wickham, Hadley, Jennifer Bryan, et al. (2019). *readxl: Read Excel Files*. R package version 1.3.1. URL: <https://CRAN.R-project.org/package=readxl>.

Wickham, Hadley, Romain François, et al. (2019). *dplyr: A Grammar of Data Manipulation*. R package version 0.8.0.1. URL: <https://CRAN.R-project.org/package=dplyr>.

Xie, Yihui (2015). *Dynamic Documents with R and knitr*. 2nd ed. Boca Raton, Florida: Chapman and Hall/CRC. URL: <http://yihui.name/knitr/>.