

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.

It is also interesting to see how this might have changed over time. For example, if wealth was a good predictor of voting patterns in the 2001 election, is it still a good predictor in 2019? 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 some of the ways we have changed. At the same time, political power has moved back and forth between the two major parties. How much can we attribute changes in political power to changes in who we are?

2 Census and electoral data

The Census provides data on electoral socio-demographics, and vote counts in each electorate can be obtained from Australian federal elections. However, joining these two data sources is difficult because the Censuses are not held at the same time as the elections. Between 2001 and 2016 there were six elections and four Censuses, as shown in the timeline below.

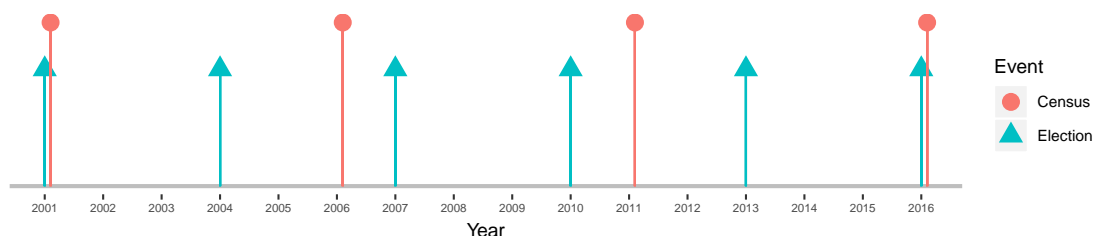


Figure 1: *Timeline of Australian elections and Censuses. They do not always occur in the same year.*

Not only can an electorate change between the last Census and an election, but even the electorate boundaries can change. Some electorates can disappear altogether and

new electorates can arise. 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. So we first had to estimate what the socio-demographic characteristics of an electorate would have been at the time of each election using a complicated method of interpolation over time and geography. This method uses Census information from both before and after the election of interest, and information from neighbouring electorates when boundaries have changed.

A simple way to measure voting patterns is to consider the two-party preferred (2PP) vote, which is based on the tally of preferences for the Labour and Liberal candidates, ignoring all other candidates. By convention, this is recorded as a percentage preference in favour of the Liberal party – so a 2PP value of 45% indicates that 45% of voters preferred a Liberal candidate and 55% of voters preferred a Labour candidate.

We consider how various socio-demographic variables obtained from Census data can be used to explain the 2PP values for each of the 150 electorates in each of the federal elections between 2001 and 2016.

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.

3 Modelling

There are dozens of socio-demographic variables available in each Census, with many variables representing similar information about an electorate. To reduce the numbers of variables, we combine variables with similar information. For example, a potential group would be variables relating to electoral incomes — median family, household and personal incomes. To determine 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); FamHouse-Size (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.

3.1 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 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 (Anselin 1988) is chosen because it captures spatial hetero-

geneity by incorporating a spatially structured random effect vector (LeSage, Kelley Pace & 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 ρ be the spatial autoregressive coefficient, v be a spherical error term, W be a matrix of spatial weights (containing information about the neighbouring regions), X be a matrix of socio-demographic covariates, β be a vector of regression coefficients and a be a spatially structured random effect vector.

$$y = X\beta + a,$$

and

$$a = \rho W a + v,$$

where $v \sim N(0, \sigma^2 I_n)$, and hence

$$y = X\beta + (I_n - \rho W)^{-1}v.$$

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

Table 1 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.

4 Results

4.1 Spatial autoregressive parameter

The spatial autoregressive coefficient ρ is positive and significant in only the 2001 and 2016 elections (Figure 2), 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.

Table 1: *Estimated spatial error model parameters (standard errors) for each of the six election years.*

	2001	2004	2007	2010	2013	2016
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)
BornAsia	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)
BornSEEUro	-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)
ManagerAdmin	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)
OneParentHouse	-1.73 (1.36)	-0.45 (1.59)	-0.75 (1.49)	-1.46 (1.69)	-0.77 (1.57)	-0.74 (1.47)
OtherLanguage	-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 (1.49)	-3.32* (1.76)	-4.01** (1.67)	-0.97 (2.07)	-0.70 (2.07)	-0.89 (2.16)
SocialServ	2.51* (1.33)	1.65 (1.41)	2.47* (1.29)	2.53* (1.47)	2.35* (1.32)	4.45*** (1.19)
Transformative	3.24** (1.55)	4.73*** (1.78)	4.84*** (1.74)	4.46** (1.98)	3.56** (1.78)	4.58*** (1.53)
Unemployment	-2.45* (1.40)	-3.07* (1.63)	0.29 (1.51)	0.08 (1.76)	1.67 (1.51)	2.79** (1.37)
Constant	50.81*** (0.71)	52.60*** (0.58)	47.31*** (0.52)	49.93*** (0.57)	53.51*** (0.58)	50.49*** (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<0.1; **p<0.05; ***p<0.01

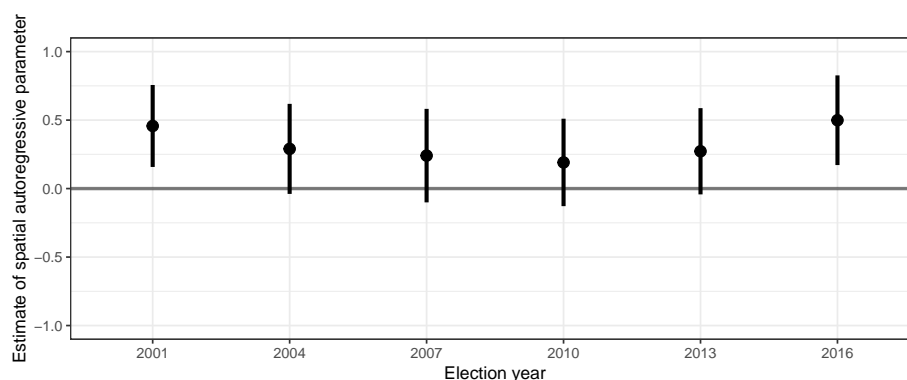


Figure 2: Estimates of the spatial autoregressive parameter for each of the six elections, reported with their individual 95% confidence intervals. Only in 2001 and 2016 is there a significant spatial component.

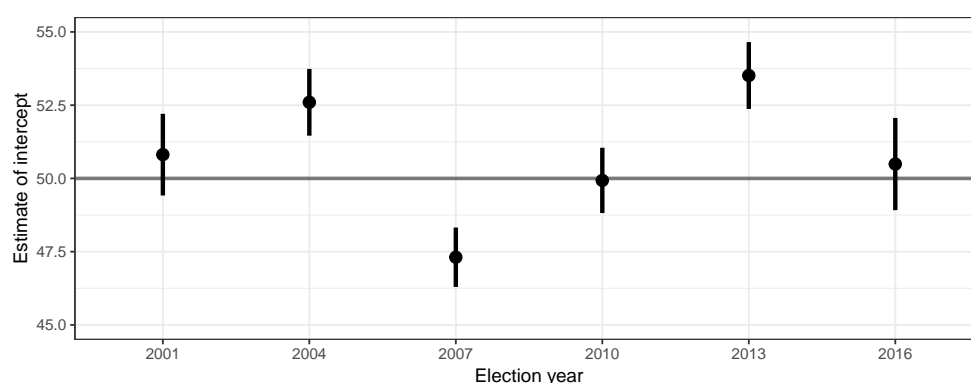


Figure 3: Estimated intercept for each election, which represents the two-party preferred vote for an electorate with mean characteristics.

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. Figure 3 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.

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 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

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

region represents a 95% confidence band, computed using the method in Breheny & 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).

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 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.

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 ?). 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.

4.3.3 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.

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

4.3.5 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 ??).

4.3.6 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 ??).

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 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.

4.3.8 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 ($p = 0.0056$).

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 (Figure 4) 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.

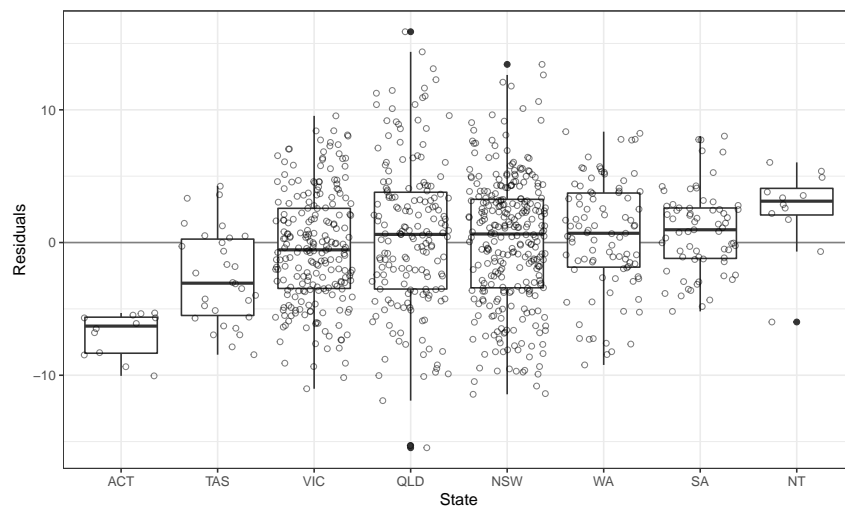


Figure 4: Boxplot of residuals by state with jittered points. States ordered by median residual. A state-specific bias not captured by the model is evident.

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 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.

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.

5 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.

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 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.

7 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/.

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, J, Y Xie, J McPherson, J Luraschi, K Ushey, A Atkins, H Wickham, J Cheng, W Chang & R Iannone (2019). *rmarkdown: Dynamic Documents for R*. R package version 1.12. <https://rmarkdown.rstudio.com>.
- Anselin, L (1988). *Spatial econometrics: methods and models*. Vol. 4. Springer Science & Business Media.
- Breheny, P & W Burchett (2017). Visualization of regression models using visreg. *The R Journal* **9**(2), 56–71. <https://journal.r-project.org/archive/2017/RJ-2017-046/index.html>.
- Forbes, J, D Cook, A Ebert, H Hofmann, RJ Hyndman, T Lumley, B Marwick, C Sievert, M Sun, D Talagala, N Tierney, N Tomasetti, E Wang & F Zhou (2019). *eechidna: Exploring Election and Census Highly Informative Data Nationally for Australia*. R package version 1.3.0. <https://CRAN.R-project.org/package=eechidna>.
- LeSage, J, R Kelley Pace & RK Pace (2009). *Introduction to spatial econometrics*. Chapman and Hall/CRC, pp. 50–52.
- Xie, Y (2015). *Dynamic Documents with R and knitr*. 2nd ed. Boca Raton, Florida: Chapman and Hall/CRC. <http://yihui.name/knitr/>.