

A Clustering Approach for Characterizing Voter Types: An Application to High-Dimensional Ballot and Survey Data

Shiro Kuriwaki (Harvard University)

Paper URL: <https://osf.io/v3rhz>

Author Homepage: <https://www.shirokuriwaki.com>

July 2020

Motivation for Clustering

Social scientists and political journalists want to group the electorate into clusters:

- by vote choice: e.g. Swing voters
- by demographics: e.g. Soccer moms

But vote data hard to summarize due to large number of combinations

Data on Vote Choice

- ✓ Individual-level
- ✓ Unordered outcomes (e.g. "D", "R", "Green", "Abstain")
- ✓ High-dimensional (e.g., Congress, Governor, Sheriff)

Existing approaches

- Focus on one pair of offices at a time
- Show all pair-wise correlations
- Impose an ideal point model

EM Algorithm for Clustering Categorical Data

Suppose voter i belongs to an unobserved cluster $Z_i \in \{1, \dots, K\}$

We only observe their vote $Y_{ij} \in \{\text{abstain}, \text{split}, \text{straight}\}$ in office j .

Quantities of interest:

- How large is cluster k ?

$$\pi_k \equiv \Pr(Z_i = k)$$

- How likely is cluster k to split in office j ?

$$\mu_{(k,j,\text{split})} \equiv \Pr(Y_{ij} = \text{split} | Z_i = k)$$

- Given their vote pattern, is voter i in cluster k ?

$$\zeta_{ik} = \Pr(Z_i = k | Y_i, \pi, \mu)$$

Model of vote choice: given cluster, offices are independent

$$\text{Likelihood (Voter in cluster } k \text{ votes } \ell \text{ in office } j) = \prod_{j=1}^J \prod_{\ell} (\mu_{kj\ell})^{1(Y_{ij}=\ell)}$$

where $\ell \in \{\text{abstain}, \text{split}, \text{straight}\}$

Estimation Challenge:

Find values of μ, π that are most likely given the model and data

... with a multinomial logit IIA assumption for varying choice sets due to **uncontested** races

$$\Pr(Y_{ij} = \ell | Z_i = k) = \frac{\exp(\psi_{kj\ell})}{\sum_{\ell' \in \mathcal{Y}_{ij}} \exp(\psi_{kj\ell'})}$$

... with allowing covariates to predict cluster classification

$$\pi_{ik} = \frac{\exp(\mathbf{X}_i^\top \boldsymbol{\gamma}_k)}{\sum_{k'=1}^K \exp(\mathbf{X}_i^\top \boldsymbol{\gamma}_{k'})}$$

where \mathbf{X} is a $N \times P$ matrix of voter-level covariates.

Expectation Maximization (EM) algorithm (via new package, clusterCVR)

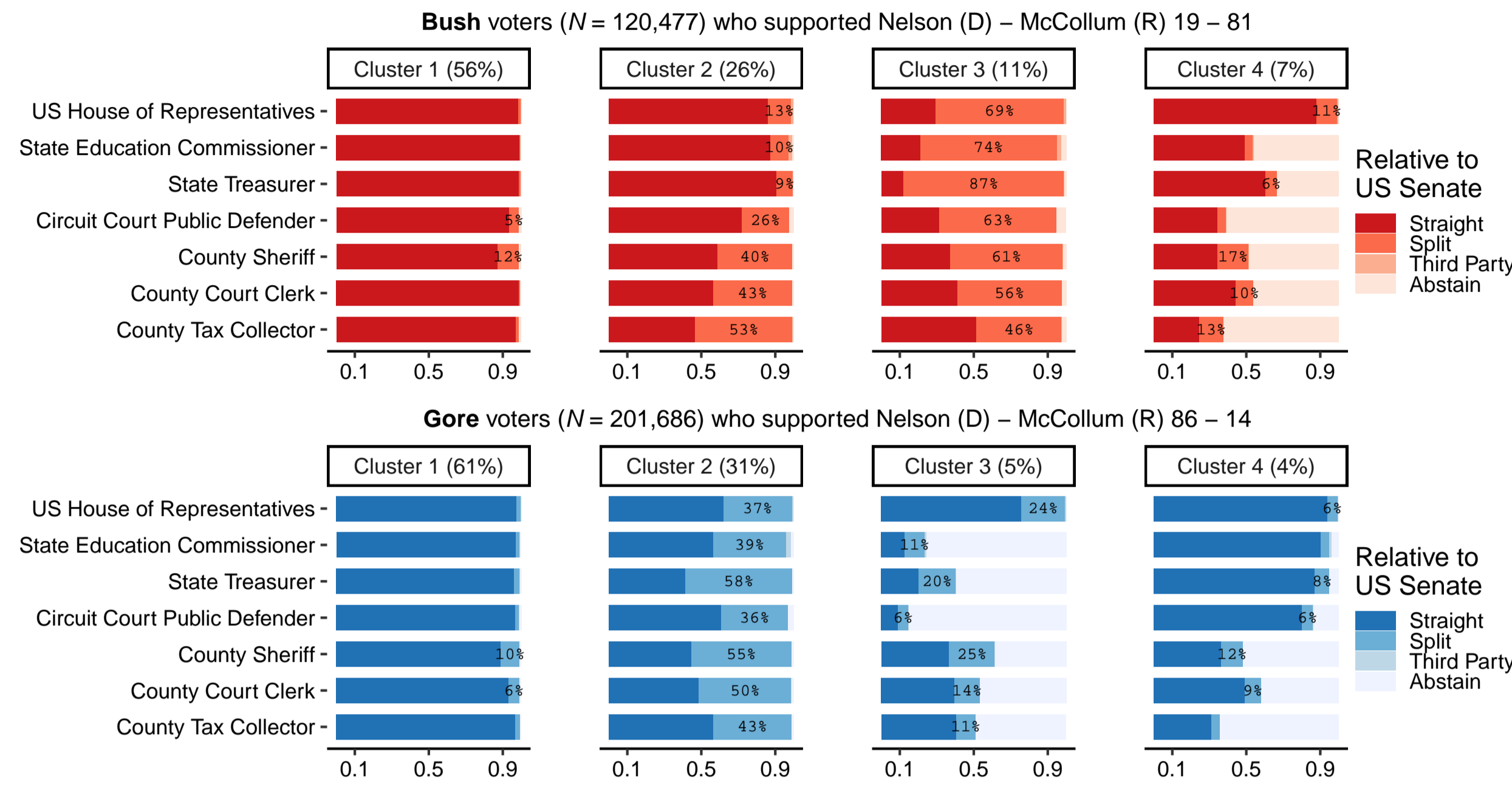
Guess $Z_i \rightarrow$ update parameters by MLE \rightarrow update $E[Z_i = k | Y_i], \dots, \odot$ until convergence.

Application 1: Cast Vote Records (CVRs)

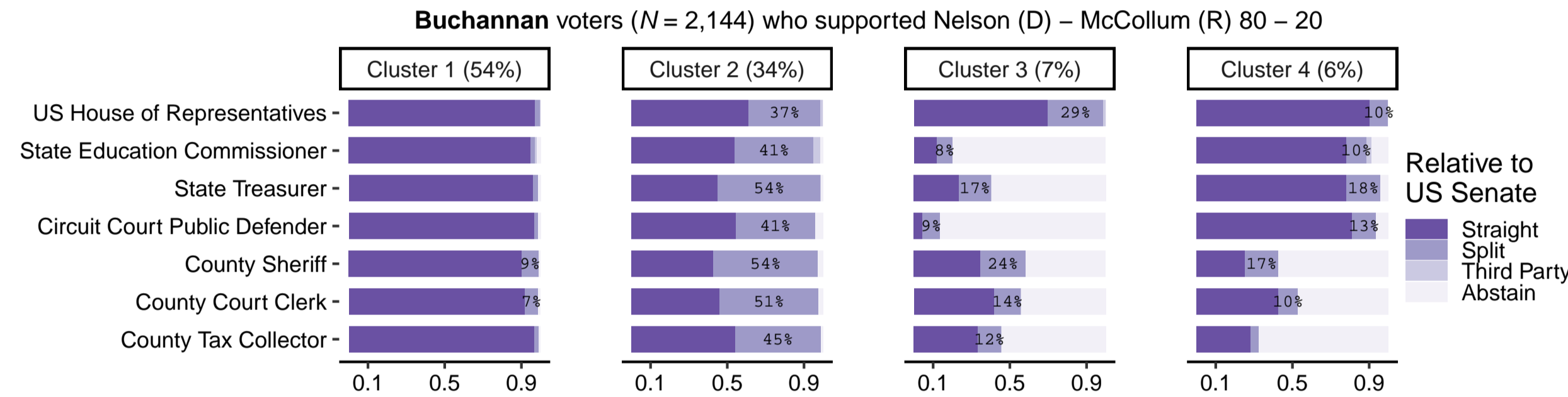
Re-analyze Herron and Lewis (2007) ballot data in Florida's 2000 Election.

Data: Palm Beach County, recoded "Straight party", "Split ticket", "Third party", or "Abstain"

Cluster analysis of CVRs show more ticket splitting down-ballot



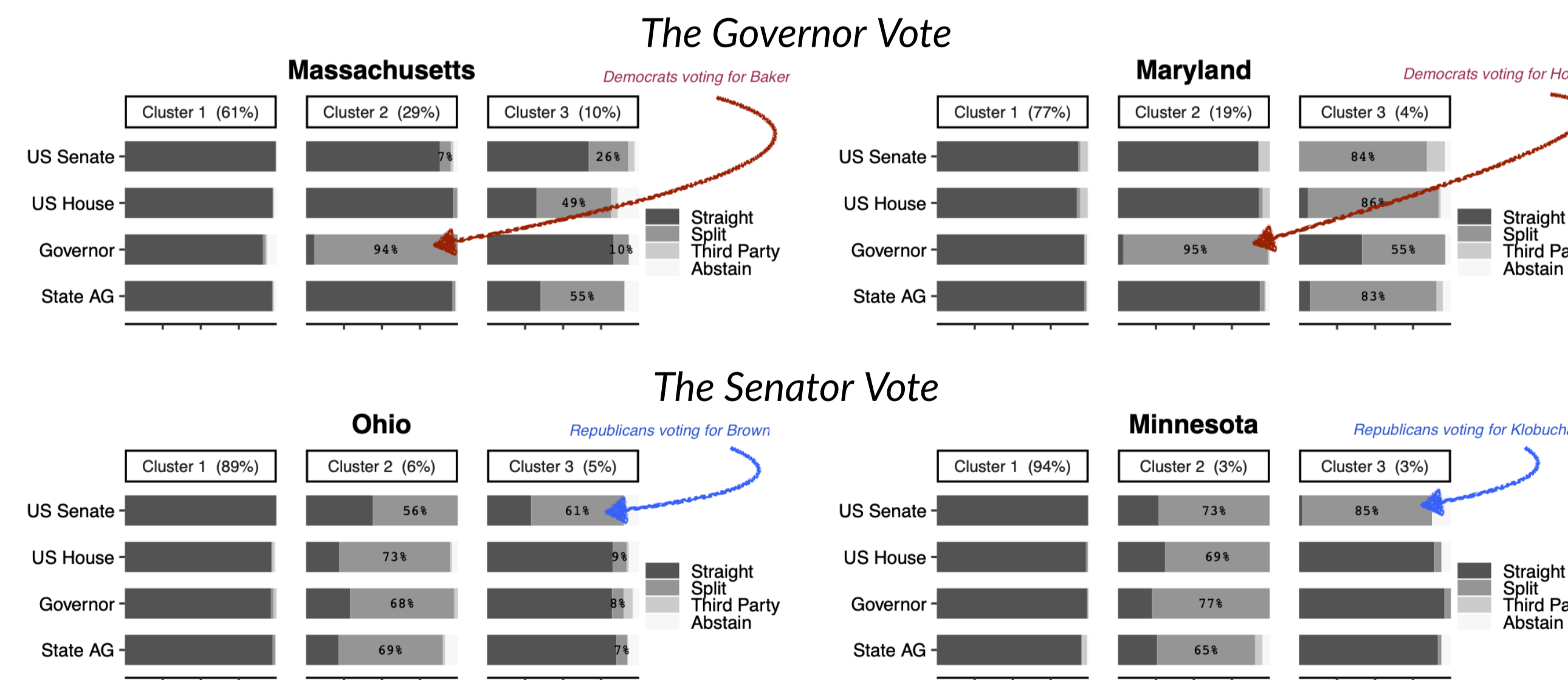
"Buchanan voters" on the infamous Butterfly ballot look awfully like Gore voters



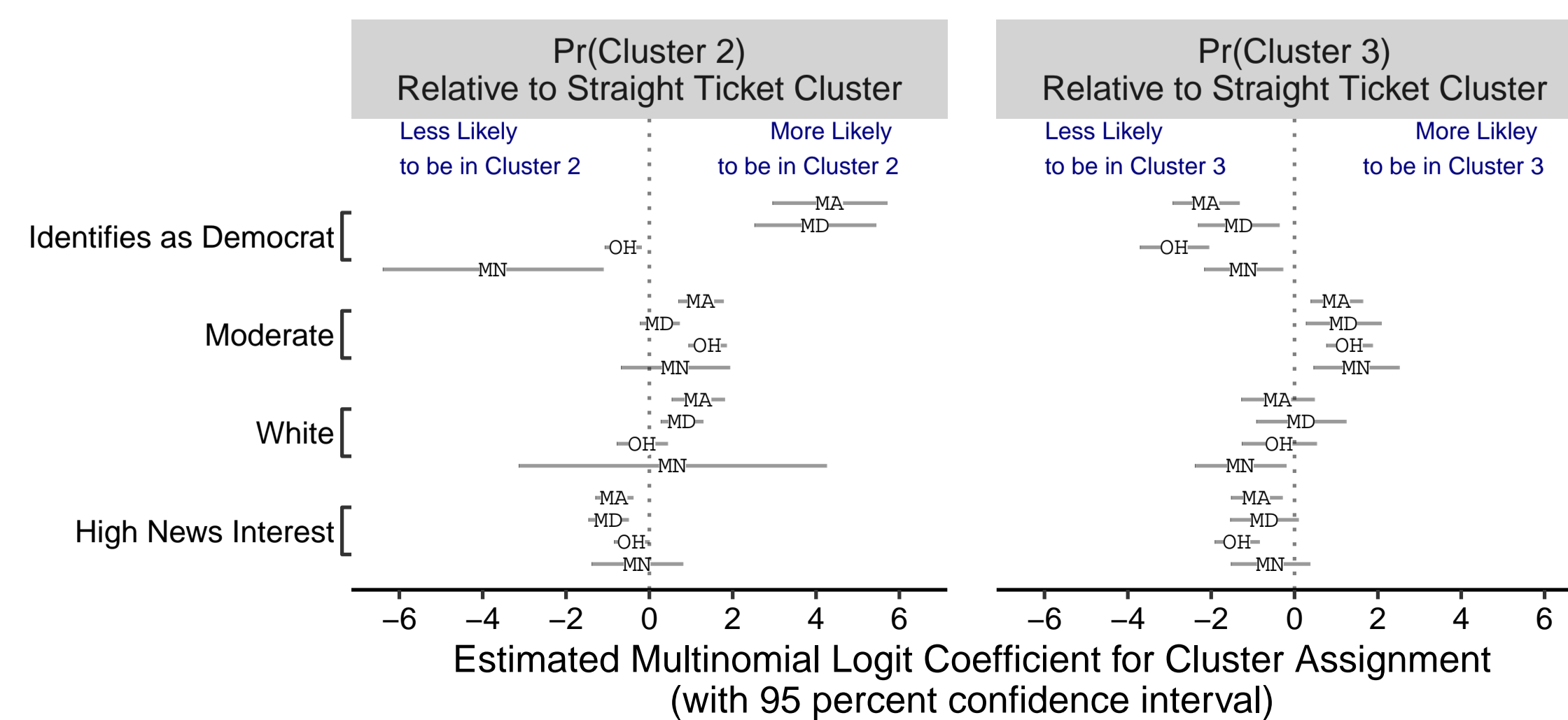
Application 2: Multi-state Survey Data

Data: 2018 CCES in states with US Senate, US House, Governor, and State Attorney General

Although most voters vote the party line, popular Governors and Senators have a personal vote



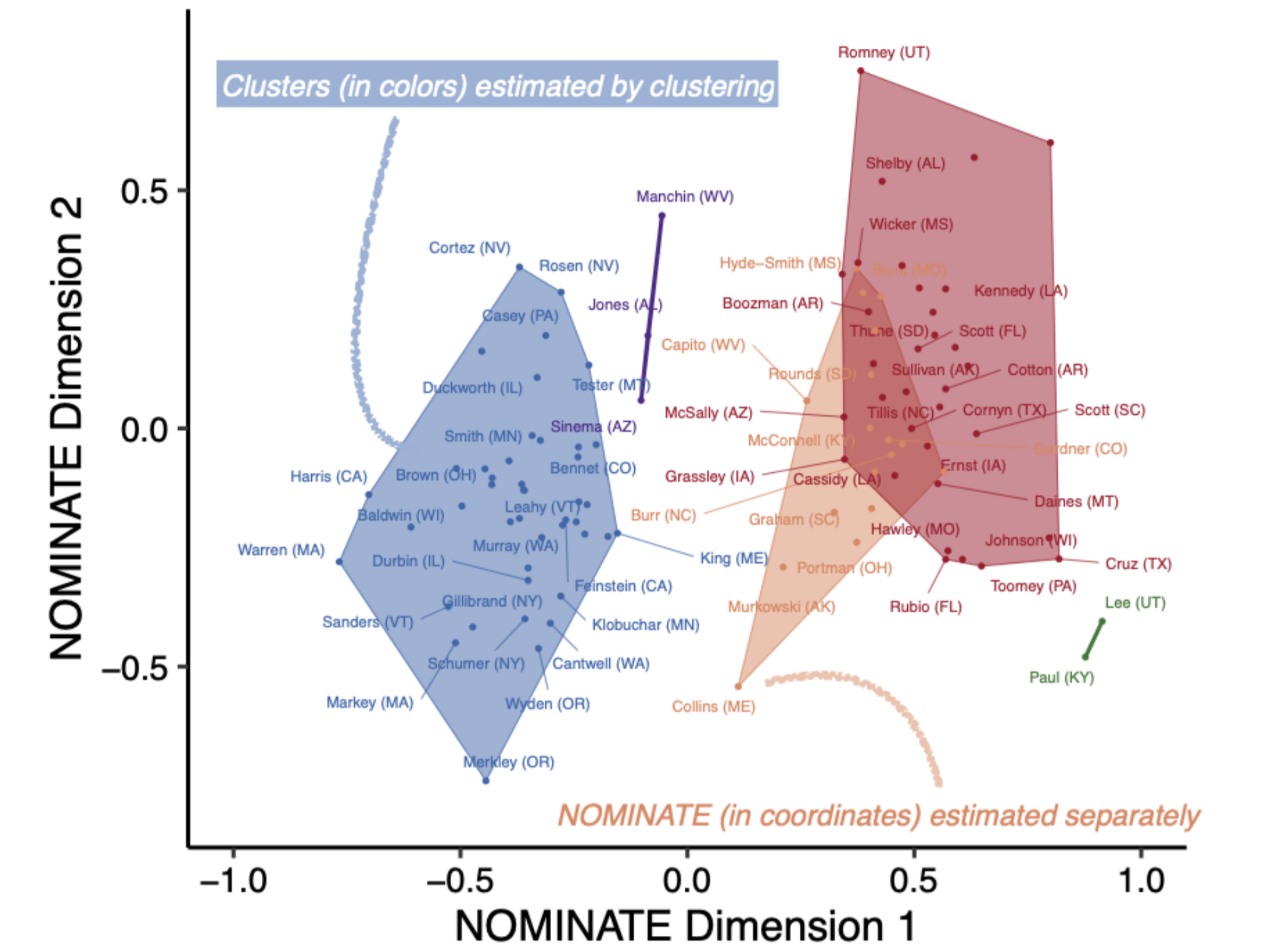
Incorporate survey covariates to cluster covariates



Comparison with Existing IRT Methods

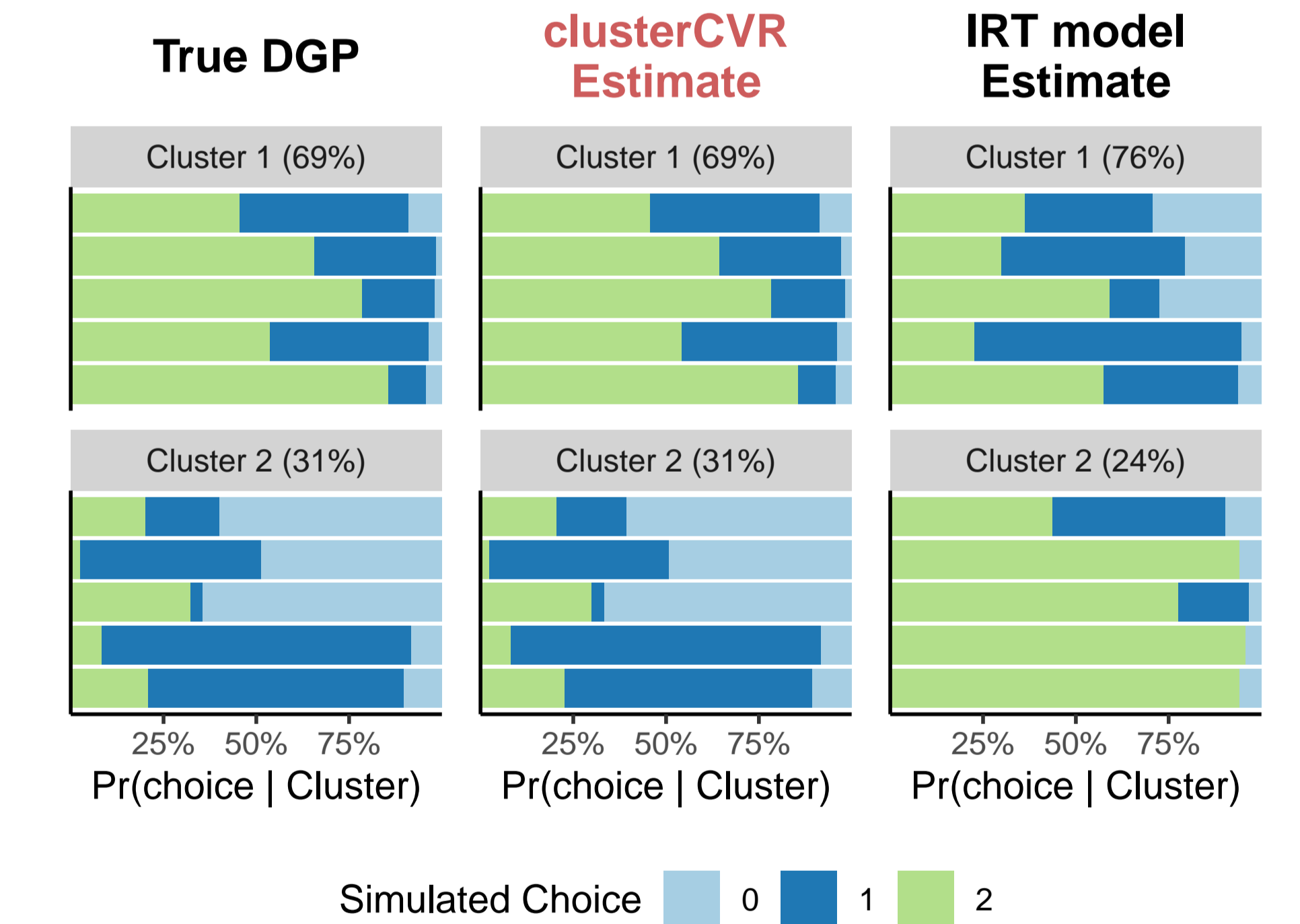
Most existing work summarizes high-dimensional vote choice data with ideal points (IRT).

Clustering recovers reasonable clusters in rollcall votes, similar to NOMINATE



Source: Voteview. Covers 406 votes in 2019. clusterCVR estimated with 5 clusters.

In simulated data, clustering model recovers true DGP more reliably



Run IRT with

- One dimensional emIRT on binary data (with abstain)
- Then cluster ideal points by k-means.
- Characterize choice probabilities by sample means (as in clustering)

Cautions on Interpretation

- Interpretation depends on number of clusters users to pick
- Interpretation and labelling requires substantive knowledge

Selected References

- Kuriwaki, Shiro (2020), "Party Loyalty on the Long Ballot: Is Ticket Splitting More Prevalent in State and Local Elections?" ([10.31235/osf.io/bvgz3](https://doi.org/10.31235/osf.io/bvgz3))
- Linzer, Drew, and Jeff Lewis (2011), "poLCA: An R Package for Polytomous Variable Latent Class Analysis". JSS
- Yamauchi, Soichiro (2020). emLogit: A ECM algorithm for the Multinomial Logit Model.