

# The Racial Voting Power Gap: Analyzing Racial Gerrymandering Through Solving Ecological Inference with The Discrete Voter Model

A THESIS PRESENTED  
BY  
HAKEEM OLAKUNLE ISA ANGULU  
TO  
THE DEPARTMENT OF COMPUTER SCIENCE AND STATISTICS  
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF  
BACHELOR OF ARTS  
IN THE SUBJECT OF  
COMPUTER SCIENCE AND STATISTICS  
HARVARD UNIVERSITY  
CAMBRIDGE, MASSACHUSETTS  
MAY 2020

©2020 – HAKEEM OLAKUNLE ISA ANGULU  
ALL RIGHTS RESERVED.

# The Racial Voting Power Gap: Analyzing Racial Gerrymandering Through Solving Ecological Inference with The Discrete Voter Model

## ABSTRACT

Section 2 of the Voting Rights Act of 1965 (VRA)<sup>vra</sup> prohibits voting practices or procedures that discriminate based on race, color, or membership in a language minority group, and is often cited by plaintiffs seeking to challenge racially-gerrymandered districts in court.

In 1986, with *Thornburg v. Gingles*<sup>tho</sup>, the Supreme Court held that in order for a plaintiff to prevail on a section 2 claim, they must show that:

1. the racial or language minority group is sufficiently numerous and compact to form a majority in a single-member district
2. that group is politically cohesive
3. and the majority votes sufficiently as a bloc to enable it to defeat the minority's preferred candidate

All three conditions are notoriously hard to show, given the lack of data on how people vote by race.

In the 1990s and early 2000s, Professor Gary King's ecological inference method tackled the second condition: racially polarized voting, or racial political cohesion. His technique became the standard

technique for analyzing racial polarization in elections by inferring individual behavior from group-level data. However, for more than 2 racial groups or candidates, that method hits computational bottlenecks.

A new method of solving the ecological inference problem, using a mixture of contemporary statistical computing techniques, is demonstrated here. It can be used for multiple racial groups and candidates, and is shown to work well on randomly-generated mock election data.

# Contents

0	INTRODUCTION	I
1	BACKGROUND AND LITERATURE REVIEW	6
1.1	Racial Gerrymandering and the Courts . . . . .	6
1.2	Ecological Inference . . . . .	11
1.3	Discretization . . . . .	15
2	METHODS	22
2.1	phc . . . . .	23
2.2	expec_votes . . . . .	29
2.3	prob_votes . . . . .	30
2.4	elect . . . . .	33
2.5	dvm . . . . .	34
2.6	Evaluation . . . . .	38
3	RESULTS	40
4	DISCUSSION	46
5	CONCLUSION	49
	APPENDIX A APPENDIX	53
	REFERENCES	64

# List of Tables

1.1	A $3 \times 3$ Table of Voting By Race, using Counts of Voting Age Population . . . . .	12
1.2	Quantities in King's EI . . . . .	13
1.3	Parameters in King's EI . . . . .	14
2.1	Partitions of 10 votes among 3 demographic groups . . . . .	31
3.1	Demographic Distributions Tested in the $2 \times 2$ Case . . . . .	41
3.2	Demographic Voting Patterns Tested in the $2 \times 2$ Case . . . . .	41
3.3	Demographic Distributions Tested in the $3 \times 2$ Case . . . . .	42
3.4	Demographic Voting Patterns Tested in the $3 \times 2$ Case . . . . .	42
3.5	Accuracy and Runtime of the Discrete Voter Model (DVM) and King's EI (KEI) . .	43

# List of Figures

1.1	The Creation of Texas' 25th Congressional District (orange) in 2003 <sup>II</sup> . . . . .	11
1.2	Kernel Density Plot of the Beta(0.1, 0.1) Distribution . . . . .	16
1.3	A Bivariate Multimodal Distribution . . . . .	17
2.1	Example of a Rank 3, Granularity 10 Probabilistic Hypercube . . . . .	26
2.2	Example of a Rank 2, Granularity 10 Probabilistic Hypercube . . . . .	27
2.3	The Probability Distribution within of a Rank 2, Granularity 10 Probabilistic Hypercube . . . . .	28
3.1	Distribution of Probability over the 2D Hypercube for a Model in the $2 \times 2$ Case .	44
3.2	Distribution of Probability over the 3D Hypercube for a Model in the $3 \times 2$ Case .	45

THIS IS THE DEDICATION.



# Acknowledgments

I would like to express my deep gratitude to Professor Latanya Sweeney, Jinyan Zang, and the students of the Fall 2019 iteration of GOV<sub>93B</sub> at Harvard University. They provided me with patient guidance, enthusiastic encouragement and useful critiques of the scope and design of this research work. In addition, I would like to thank Professor Moon Duchin of Tufts University for her guidance of this research and my overall thesis work.



## Introduction

THE VOTING RIGHTS ACT of 1965 (VRA)<sup>vra</sup> was the result of decades of activism and advocacy around unencumbered suffrage in the United States. Specifically, this act was meant to “enforce the fifteenth amendment to the constitution,” which was ratified 95 years prior, and came after African Americans in the South protested the tremendous obstacles to voting they encountered, including poll

taxes, literacy tests, harassment, intimidation, physical violence, and other systemic and infrastructural facets of the practical application of the right to vote. The United States Department of Justice has cited it as “the single most effective piece of civil rights legislation ever passed by Congress.”<sup>7</sup>

Section 2 of that act specifically prohibits any jurisdiction in the nation from implementing any “voting qualification or prerequisite to voting or standard, practice, or procedure...which results in a denial or abridgment of the right of any citizen of the United States to vote on account of race or color.” Today an often challenged voting procedure is the drawing of political districts themselves, on the grounds that a political tool, called gerrymandering, dilutes the effectiveness of the vote.

Gerrymandering is the process of manipulating the boundaries of electoral districts to establish an unfair political advantage for some group of people. This process occurs at all levels of government in the United States of America, and regularly occludes the representative characteristic of the democracy.

Gerrymandering is the abusive version of the more general *redistricting*, which is simply the process by which elected and appointed officials redraw electoral districts. In order to guarantee fair representation, electoral districts ought to be drawn based on data about the population and its distribution in space. There is a consensus that “fairness” requires population balance, hence redistricting is done at least every ten years in the United States, immediately after a census. When redistricting is corrupted, it can become a powerful political tool for diluting or magnifying a demographic group’s voting power in elections. When that demographic group is a race, gerrymandering becomes addressable by the VRA.

Quantifying the effects of voting standards, like the electoral map, on racial groups is notoriously difficult. Both Congress and the Supreme Court have attempted to set aside rules to prove that voting

standards that sound neutral have racially disparate effects.

In *Thornburg v. Gingles* (1986)<sup>tho</sup>, the Court established that plaintiffs needed to show that:

1. the racial or language minority group is sufficiently numerous and compact to form a majority in a single-member district
2. that group is politically cohesive
3. and the majority votes sufficiently as a bloc to enable it to defeat the minority's preferred candidate

This amounts to two kinds of evidence, now known colloquially as the Gingles Test, needed to successfully prove a violation of the Voting Rights Act:

1. the geographic conditions to draw an effective district
2. racial polarization that blocks the will of minority voters

The racially polarized voting test requires showing that a marginalized racial group is *politically cohesive*, and that the “majority votes sufficiently as a bloc to enable it to usually defeat the minority’s preferred candidate.” Both of these are uniquely hard to show given the secret ballot; that is, all voter choices in an election are anonymous. Thus, it is impossible to access the ground truth about how racial groups in the majority and minority vote. This information is necessary to challenging an electoral map that one believes is gerrymandered.

To fill this gap, statistical methods, most notably King’s Ecological Inference (EI)<sup>18</sup>, exist to infer that information from publicly available data.

EI was pioneered in the 1990s and 2000s by Harvard professor Gary King and Columbia professor Andrew Gelman, and it became the standard technique for attempting to analyze racial polarization

in elections in service of mounting VRA litigation. As King himself billed EI, it is “the process of learning about discrete individual-level behavior by analyzing data on groups”<sup>18</sup>.

From only the aggregate vote counts and aggregate racial distribution, EI produces inferred candidate preferences for each racial group. This is important because the courts that litigate VRA cases typically do not accept data from non-public sources. Often, this means that the only data that plaintiffs have to bring forward are Census data and election results.

One method of solving the Ecological Inference problem is King’s own Ecological Inference method. The simplest versions of this method use Normal distributions, and the most complex hinge on Binomial-Beta hierarchical models for Bayesian inference. As a point of departure, this paper explores a new way of providing data to the hierarchical model pipeline: by discretizing the whole parameter space of possible vote outcomes. With this discretized space, one can replace the limited Beta class of distributions with a much more flexible description of voter behavior. This new method for solving the Ecological Inference problem, the *discrete voter model* (DVM), was not computationally tractable at the full scale of a U.S. state in the 1990s and early 2000s, but is now, due to advances in algorithm design and computing power.

While gerrymandering is difficult to prove, it is clear that it is one of the most pressing contemporary voting rights problems. The nation’s electorate, and thus the integrity of the democracy, has to constantly contend with confusing and dilutive maps that are drawn by people, parties, and agencies with specific agendas. It is important that the methods for analyzing electoral maps are constantly updated and improved with new knowledge and computing power. Additionally, it is critical that organizers, potential plaintiffs, and members of various affected communities all have the ability to quickly and

intuitively synthesize the work of mathematicians, statisticians, geographers, political scientists, and other key stakeholders towards resolving this longstanding and unavoidably complex problem.

*tbd*

tbd

# 1

## Background and Literature Review

### 1.1 RACIAL GERRYMANDERING AND THE COURTS

GERRYMANDERING is the process of manipulating the boundaries of electoral districts to establish an unfair political advantage for some group of people. Racial gerrymandering has been affirmed by the court as a direct violation of constitutional rights, with the Voting Rights Act (VRA) used to

further bolster those rights. A racial gerrymander constitutes a “systemic and infrastructural facet of the practical application of the right to vote,”<sup>15</sup> and plaintiffs often invoke the VRA to force officials to redistrict.

Racially polarized voting, defined as different racial groups voting for different candidates in an election, has been central to redistricting law since the 1982 amendments to section 2 of the Voting Rights Act.<sup>15</sup> Two kinds of evidence, now known colloquially as the Gingles Test, after the 1986 *Thornburg v. Gingles*<sup>tho</sup> case, further codified that in order to prove a violation of the Voting Rights Act one needed to show that:

1. the racial or language minority group is sufficiently numerous and compact to form a majority in a single-member district
2. that group is politically cohesive
3. and the majority votes sufficiently as a bloc to enable it to defeat the minority’s preferred candidate

The final two requirements require showing racially polarized voting. “Political cohesion” is often interpreted as racial groups voting for one candidate, and that racial voting information is necessary to determine if racial minority group’s candidates are often defeated by the majority racial group’s candidates.

This means, practically, that in order to win a lawsuit plaintiffs must prove that voters in districts vote along racial lines. While not every instance of racially polarized voting occurs because of a gerrymandered electoral map, proof of racially polarized voting is often the first step in successfully litigating a racial gerrymandering case as a violation of the VRA.



### 1.1.1 ESTIMATING RACIALLY POLARIZED VOTING

ALTHOUGH racially polarized voting is critically important to racial gerrymandering litigation, it is quite difficult to prove. The *secret ballot*, an essential part of the United States' democracy, ensures that one's vote is anonymous. Hence, any data collection around voting patterns that separates along any axis of identity, is impossible, and all information about those voting patterns must be inferred.

The *Gingles* decision noted that two techniques were “standard in the literature for the analysis of racially polarized voting”<sup>tho</sup>: bivariate ecological regression and homogeneous precincts.

Dozens of court cases, from the 1980s to now, have relied on either of the above methods for litigating racial gerrymandering cases<sup>15</sup>. However, some appellate courts and judges have been hostile to them. In *Lewis v. Alamance County* (1996)<sup>Lew</sup>, the court said that the critical assumption of regression “runs counter to common sense.” In *Holder v. Hall* (1994)<sup>Hol</sup>, said that “bivariate regression analysis...does not directly control for...factors [other than race].”

These criticisms were often ignored by peer and other courts because they did not posit new ways of deciding cases without using regression. Chief Justice Roberts chimed into the debate as well, saying “At trials, assumptions and assertions give way to facts. In voting rights cases, that is typically done through regression analyses of past voting records.”<sup>LUL</sup>

However, the problem has gotten increasingly more complex. Jim Greiner, Professor of Law at Harvard Law School, stated in his paper *Ecological Inference in Voting Rights Act Disputes: Where Are We Now, and Where Do We Want To Be?*<sup>15</sup>

According to some scholars, the patterns themselves have shifted in that a certain percentage of white citizens in some jurisdictions [like Chicago] have become willing to vote for minority-preferred candidates. This shift suggests that the days in which “documenting racially polarized voting is generally...just beating the obvious” may be gone. Racial bloc voting, often a hotly contested issue in an “ordinary” redistricting dispute, may now become even more so, and as a result, using the best methodology is more critical.<sup>15</sup>

Greiner’s statement is a good introduction to the future of redistricting. He notes not only that the electorate’s voting patterns are becoming more complex, but also that to ensure that the courts are equipped with the correct tools to tackle the complexity, the methodologies supplied by and to expert witnesses must be continually improved.

The ecological regression and homogeneous precincts methods had several disadvantages. Ecological regression regularly produces impossible estimates of support, like that –400% of Hispanic voters in a district supported a candidate. The analysis of homogeneous precincts depends on the assumption that white people living in an all-white precinct vote similarly to white people living in a more racially-diverse precinct, which has been shown empirically, and can be thought intuitively, to be false.<sup>15</sup> In addition, neither method generalizes well past 2 racial groups and 2 candidates in an election.

The dominant racial bifurcation in both legal and lay discourse is that between Black and white citizens. Thus, much of the litigation of racial gerrymandering and VRA cases, along with the methods used, have focused on these two groups. However, many contemporary cases explore the more complex racial dynamics that exist in our more, and increasingly, diverse electorate.

Take, for example, a recent case that used ecological regression: *League of United Latin American Citizens v. Perry*, 126 S. Ct. 2594 (2006) (LULAC). In this case, Latinx plaintiffs, a mixture of private

citizens and representatives of the Mexican American Legal Defense and Educational Fund, argued that Texas' congressional districts were racially gerrymandered to negatively affect Latinx citizens in Texas. With the help of a dubiously-adjusted ecological regression model (to be used for a larger combination of racial groups), the court decided to throw out District 23 within a new plan proposed in 2003, ordering a remedial redistricting.<sup>LUL</sup> The court also stated that it need not rule that District 25 was a racial gerrymander, given that the remedial changes to District 23 would necessarily affect District 25.

Although not verbally verified by the court, District 25 is widely believed to have been drawn to pack Latinx voters into a district and dilute their voting power, so the redrawing of it and District 23 was a win for that section of the electorate. The change had incredibly practical, and almost immediate, effects. In the year of the remedial redistricting, the affected districts held new primary elections. Henry Bonilla, the Republican representative for the 23rd congressional district, was defeated by Democrat Ciro Rodriguez, a favorite of the Latinx electorate.

Figure L.I shows the creation of District 25 after the 2003 redistricting in question. Its shape is reminiscent of popular gerrymanders.

In that decision, the judges referred to a previous one from 1993, where the court stated:

In countless areas of the law weighty legal conclusions frequently rest on methodologies that would make scientists blush. The use of such blunt instruments in examining complex phenomena and corresponding reliance on inference owes not so much to a lack of technical sophistication among judges, although this is often true, but to an awareness that greater certitude frequently may be purchased only at the expense of other values.<sup>cle</sup>

Here, the court explicitly acknowledges the disadvantages of the methods used to infer racially po-

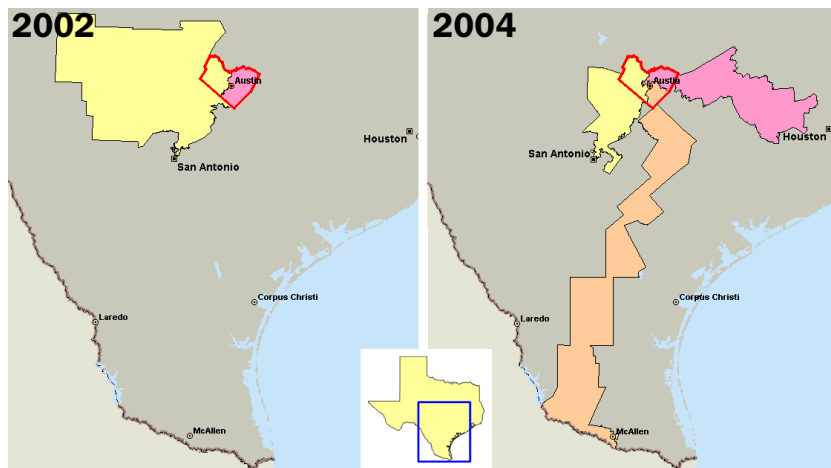


Figure 1.1: The Creation of Texas' 25th Congressional District (orange) in 2003<sup>11</sup>

larized voting. These disadvantages, the criticisms lain on ecological regression and homogeneous precincts by some courts, and improvements in theory inspired a burst of energy and innovation in the academy in the late 1990s and early 2000s. At the center of this burst is King's EI.

## 1.2 ECOLOGICAL INFERENCE

ECOLOGICAL INFERENCE (EI) is the process of using aggregate (or, ecological) data to infer discrete individual-level relationships of interests when individual-level data are not available.<sup>18</sup> This is distinct from the *ecological regression* discussed above, which is the statistical method of running regressions on aggregates and interpreting those regressions as predictive relations on the level of individual units.<sup>14</sup>

In a world with two categories, groups  $\{1, 2, \dots\}$  and groups  $\{A, B, \dots\}$ , both methods seek to answer the question “what percentage of members in group 1 are also in group A?” This is precisely the answer that plaintiffs need in racial gerrymandering cases, where groups  $\{1, 2, \dots\}$  are races and

**Table 1.1:** A  $3 \times 3$  Table of Voting By Race, using Counts of Voting Age Population

	Cand. 1	Cand. 2	$\Sigma$
Race 1	-	-	Race 1 Pop.
Race 2	-	-	Race 2 Pop.
$\Sigma$	Cand. 1 Votes	Cand. 2 Votes	Precinct Pop.

groups  $\{A, B, \dots\}$  are candidates.

Greiner provides an alternate description of EI that fits well within mathematical models.<sup>15</sup> EI can also be understood as the attempt to predict internal cell values of a set of contingency tables, like Table 1.1, when only the column and row totals (the margins of the table) are observed.

Each precinct has its own table, with an election's vote counts as the column totals and the voting age population, or demographic split of the district, often from the Census, as the row totals. The internal cells, how many people of each racial group that voted for each candidate, are protected by the secret ballot. Those internal cells are exactly the evidence needed in racial gerrymandering litigation, and their values have been inferred by ecological regression and homogeneous precinct analysis in the past.

Theoretically, these tables can be expanded to be as large as necessary:  $R \times C$ , where  $R$  corresponds to the number of rows (demographic groups or races) and  $C$  corresponds to the number of columns (candidates). This would be ideal given that the diversification of the electorate has created more politically powerful demographic groups, and that there are frequently more than 2 candidates in elections, especially in local cases (where many infringements of the VRA are suspected to occur).

Professor Gary King developed and proposed his own method to solve this problem in the late

1990s. It has been used in a few court cases, mainly for the advantage that unlike ecological regression, this method will never produce impossible estimates. It is also understood to be more intuitive than regression, addressing the court’s seldom complaint.

### 1.2.1 KING’S ECOLOGICAL INFERENCE

PROFESSOR KING introduced his method to solve the ecological inference problem in his 1997 book *A Solution to the Ecological Inference Problem: Reconstructing Individual Behavior from Aggregate Data*<sup>16</sup>. His method, which he expands upon with several papers from 1997 to the late 2000s, uses Bayesian inference and bounds methods to infer the values of the empty cells in Table 1.1.

Using Table 1.1 as a reference, Table 1.2 describes the quantities that King used to describe this problem<sup>18</sup>.

**Table 1.2:** Quantities in King’s EI

Quantity	Description
Observed	
$X_i$	fraction of population $i$ in race 1
$T_i$	fraction of population $i$ who voted for candidate 1.
$N_i$	size of population $i$
$T'_i$	number of people in population $i$ in race 1
Unobserved	
$b1_i$	fraction of members of race 1 who voted for candidate A
$b2_i$	fraction of members of race 2 who voted for candidate A

King’s method then uses the observed quantities, and a few parameters, to find the unobserved quantities and answer the ecological inference question. The parameters used are described in Table

### 1.3.

**Table 1.3:** Parameters in King’s EI

Parameter	Description
$c_1, d_1$	hyperparameters for the Beta distribution for $b_1$
$c_2, d_2$	hyperparameters for the Beta distribution for $b_2$ .
$\theta$	the expected fraction of the population of $i$ registered to vote $i$
$\lambda$	the fixed hyperparameter for the distributions of $c_1, c_2, d_1, d_2$

The model is a hierarchical Bayesian inference model that proceeds as follows:

1.  $c_1, c_2, d_1, d_2 \sim \text{Exponential}(\lambda)$
2.  $b_{1i}|c_1, d_1 \sim \text{Beta}(c_1, d_1)$
3.  $b_{2i}|c_2, d_2 \sim \text{Beta}(c_2, d_2)$
4.  $T'_i|b_{1i}, b_{2i}, X_i \sim \text{Binomial}(N_i, \theta_i)$

The Python implementation of King’s EI is given in Listing A.1. Using the observed quantities, King’s method produces a model for the unobserved quantities of interest, which can then be used in litigation as described above.

However, in King’s original paper, the bulk of his examples are for  $2 \times 2$  tables, and proposed extensions of that method beyond that bound have been dubious.

In King’s second iteration of his method and accompanying software, *EI2*, King implements a “Multiple Imputation Approach” to try to handle  $2 \times 3$  tables. In essence, this approach collapses the  $2 \times 3$  table into multiple  $2 \times 2$  subtables, then applies the original method to each of those subtables.<sup>17</sup> Statisticians and scholars, like Jim Greiner, have noted that this method as applied is statistically invalid, as it does not, and cannot, fully adhere to combinatorial rules that the pioneer of multiple im-

putation, Professor Donald Rubin of Harvard University, notes in his writings.<sup>23</sup>

The problems increase if the size of the table increases to  $3 \times 3$  or higher. Professor Karen Farree of the University of California, San Diego investigated King's EI applied to tables of that size, and found that the method produced biased estimates and relied on inappropriate modeling assumptions.<sup>13</sup>

Hence, although King's method has proved quite useful for some cases of ecological inference, with an increasingly diverse electorate, and complex multi-candidate elections, a method that generalizes well to an  $R \times C$  table, using the same data available to previous methods, is necessary.

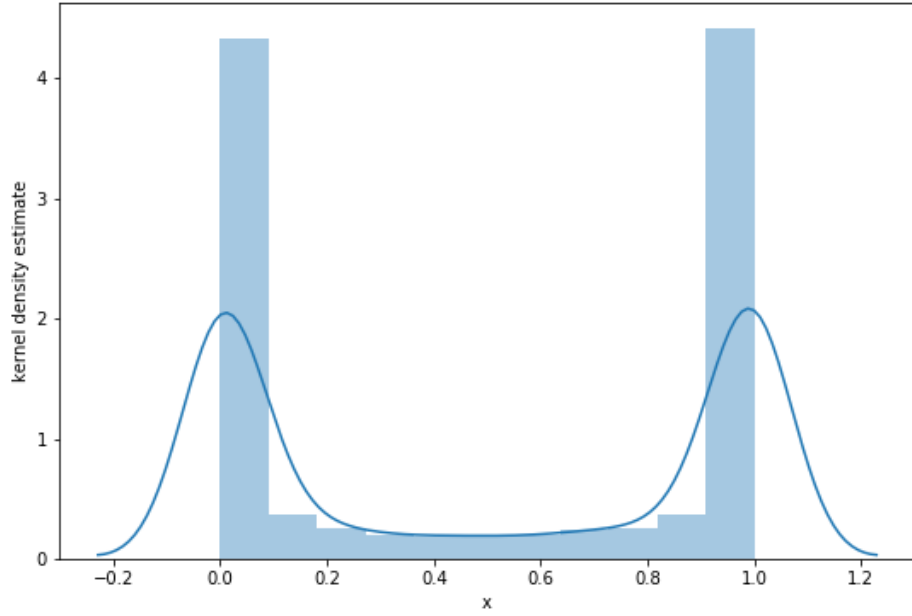
### 1.3 DISCRETIZATION

KING'S METHOD assumes the quantities of interest, the demographic voting patterns of a district, can be modeled completely with Beta distributions. The hierarchical model allows some freedom to those distributions, generating their parameters from Exponential distributions.

Beta distributions are known to be approximately bimodal if both shape parameters are below 1. However, Beta distributions can never be *fully* bimodal, which is the expectation for racially polarized voting: that two different groups vote quite differently. Figure 1.2 is an example of a very, but not fully, bimodal Beta distribution.

Perhaps more importantly, these distributions can never be multimodal (with more than 2 modes), which is the case that is necessary for multiple racial groups and candidates. This is where discretization becomes quite useful. Figure 1.3 is an example of a true multimodal distribution over two variables.



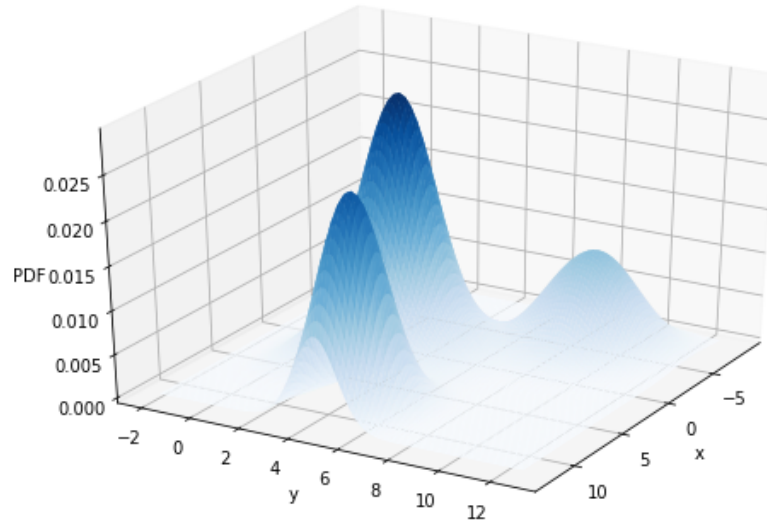


**Figure 1.2:** Kernel Density Plot of the Beta(0.1, 0.1) Distribution

A distribution like this could be created easily with Neural Networks, but that architecture requires tens of thousands or millions of data points to be reliable, where a “data point” would probably be an electoral result. The problem of ecological inference, especially in local elections, is simply too small for those methods.

However, by specifying very little about the space (even less than something like a Beta distribution), then allowing the space to take the shape of the available data, more complex, interesting, and possibly accurate distributions can be found. This is called the *Discrete Voter Model (DVM)*.

For a given candidate, each state of this space is a hypercube (the  $n$ -dimensional abstraction of a 2-D square/grid or a 3-D cube) of values, where each dimension corresponds to a demographic group



**Figure 1.3:** A Bivariate Multimodal Distribution

of interest. Each cell of the hypercube has a probability, where the sum of the probabilities of all cells is 1. Each cell also uses its position (the values on each axis) in the full hypercube to represent the demographic voting patterns of a given precinct.

For example, in the 2-dimensional case, for a hypercube with granularity 10 (each axis goes from 0 to 10), let the cell at (2, 5) have an internal probability of 0.4. This corresponds to there being a 40% chance that this candidate's demographic voting pattern in this precinct is:

- 20% of group 1 votes for them
- 50% of group 2 votes for them

Hence, a single hypercube represents the distribution of possible demographic voting patterns in a

precinct. Therefore, the distribution of hypercubes represents possibilities for all precincts in a map, and for an election.

The procedure by which the state space of hypercubes takes the shape of the data is Markov chain Monte Carlo.

### 1.3.1 MARKOV CHAIN MONTE CARLO

MARKOV CHAIN MONTE CARLO (MCMC) is a class of algorithms that sample from probability distributions that are hard to sample from, typically called *target distributions*.<sup>10</sup> With many distributions, like the Normal, Beta, and Exponential mentioned above, obtaining a sample or specifying those distributions is simple and fast. There exist many packages for Python, R, and other programming languages that generate samples instantly with closed form methods.

The target distribution for the Discrete Voter Model, that over the space of hypercubes that could have possibly resulted in the outcome of an election, is not well-specified, and differs not only by election, but also by precinct. Attempting to apply a well-known distribution requires many approximations and assumptions. In lieu of a well-known distribution, the Discrete Voter Model uses discretization and MCMC to find the estimates of the quantities of interest.

All MCMC algorithms are iterative, and the samples converge to the target distribution over time, so the understanding of that target distribution grows with each iteration. All MCMC algorithms construct Markov chains, which means that each sample is generated from, and only depends on, the prior sample in the process, establishing the Markov relationship.<sup>20</sup> At each step, a candidate for the

next sample generated according to some transition rules, determined by a *transition kernel*. That candidate is either accepted or rejected based off of some acceptance probability function. More precisely:

Let the chain be represented by  $X_0, X_1, \dots, X_t$ , where  $X_i$  is a step in the chain at iteration or time  $i$ .

Let the target distribution be  $\pi(\cdot)$ , which has an unnormalized density  $\pi_u$ .

Let there be some proposal distribution,  $Q(\cdot)$ , which generates new candidates at every step of the chain. Therefore, as a function  $q(x, \cdot)$ , it generates  $Y_{i+1}$  as a candidate for  $X_{i+1}$ , given  $X_i$ .

Let there be an acceptance function:

$$\alpha(x, y) = \min \left[ 1, \frac{\pi_u(y) \cdot q(y, x)}{\pi_u(x) \cdot q(x, y)} \right]$$

The candidate state  $Y_{i+1}$  is accepted with the probability  $\alpha(X_i, Y_{i+1})$ . If the state is accepted,  $X_{i+1} \leftarrow Y_{i+1}$ . If not,  $X_{i+1} \leftarrow X_i$ .

The differences in transition kernels,  $Q(x, \cdot)$ , and acceptance probability functions,  $\alpha(x, y)$ , are they key delineators of different types of MCMC algorithms. The Discrete Voter Model uses two types of transition kernels: *Random Walk Metropolis* and *Hamiltonian Monte Carlo*.

### 1.3.2 RANDOM WALK METROPOLIS

ONE MARKOV CHAIN MONTE CARLO METHOD which uses the Random Walk Metropolis transition kernel is the Metropolis-Hastings algorithm. Metropolis-Hastings is one of the simplest MCMC methods, but is still quite reliable and widely used for many applications.

This transition kernel requires that  $q(x, y) = q(y - x)$ . This is one way to ensure that the chain is reversible and detail balanced, necessary conditions for the chain to converge to the target distribution.<sup>22</sup> Examples of such distributions are  $Q(x, \cdot) = N(x, \sigma^2)$  and  $Q(x, \cdot) = \text{Uniform}(x - 1, x + 1)$ .

The Discrete Voter Mode’s proposal randomly selects two cells in the probability hypercube, and adds a small value,  $\varepsilon$ , to one and subtracts it from the other. This is a reversible chain – one could get theoretically get back to the probability hypercube that was changed later in the chain – which is essential for convergence.<sup>10</sup>

### 1.3.3 HAMILTONIAN MONTE CARLO

HAMILTONIAN MONTE CARLO (HMC) differs from Random Walk Metropolis in that its proposal distribution adapts to the current state and its position in the parameter space.<sup>9</sup>

HMC uses the gradient of the target distribution’s density, and adapts the proposal distribution toward that direction. This helps combat some of the inefficiencies of Random Walk Metropolis, like its tendency to explore parts of the state space where there is very little probability density. Hence, fewer samples have to be created to approximate the target distribution.

The “Hamiltonian” in “Hamiltonian Monte Carlo” refers to the inspiration of this kernel: Hamiltonian mechanics, a reformation of Newtonian mechanics that reconsiders the basis of time evolution of a system.<sup>12</sup> This kernel generates some *potential function/energy* from the target density, then uses the curvature of that function, some random seed, and the current state to determine the candidate for the next state.<sup>21</sup>

The primary benefit of HMC to the Discrete Voter Model is that it can theoretically produce a better sample of the target distribution with fewer steps in the algorithm. In addition, HMC has been shown to work well on state spaces with many parameters, like the potentially large probability hypercubes in this model.<sup>21</sup>

#### 1.3.4 ROBUSTNESS

THE MARKOV CHAIN CENTRAL LIMIT THEOREM guarantees convergence and accurate sampling with the correct specifications.<sup>10</sup> Furthermore, for both the kernels above, chains are robust to their assumptions – they often converge to reliable samples in enough time, even when assumptions (like a perfectly ergodic chain) are broken.

The main disadvantages of using Markov chains for such complex state spaces are that it is difficult to assess convergence, and the methods themselves are computationally expensive. However, given the increasing availability of computing power, including the prevalence of GPU-based computing and the maturity of programming packages like Stan, PyMC3, and TensorFlow Probability, these methods can be more widely applied.

*This is some random quote to start off the chapter.*

Firstname lastname

# 2

## Methods

The Discrete Voter Model is implemented in Python 3, and the accompanying code is found in the repository noted in [Appendix A](#). A set of 5 modules supports the Discrete Voter Model, with each implementing an aspect of the statistical framework. Those modules are:

1. `phc`: creates probabilistic hypercubes
2. `exp_votes`: finds the expectation of votes for a candidate

3. `prob_votes`: finds the probability of an electoral outcome
4. `elect`: aggregates election and electorate data
5. `dvm`: executes the Markov chain on a state space of hypercubes

The modules work together to run the Discrete Voter Model on an `Election` object to generate a distribution of probabilistic hypercubes for the election. Each probabilistic hypercube encodes a probability distribution within a potential precinct of the district, and is used to extract the likely demographic voting probabilities.

These demographic voting probabilities are the goal of the ecological inference, and provide information to answer whether voting in an election is polarized by race.

## 2.1 PHC

THE DISCRETE VOTER MODEL introduces the probabilistic hypercube (PHC) to represent the distribution of possible demographic voting probabilities in a precinct, for a candidate.

A probabilistic hypercube (PHC) is defined as a rank  $n$  tensor (matrix, multidimensional array, or holor) in  $\mathbb{R}^n$  space with the following characteristics:

1. the component values of the cells sum to 1
2. each dimension corresponds to a demographic group in an electorate
3. each cell's position, or index, represents the demographic voting probabilities of a precinct
4. each cell's component value represents the inferred probability of that cell representing the true



demographic voting probability of a precinct

The function `make_phc` initializes a PHC with either random or uniform cell components. The rank of the PHC corresponds to the number of demographic groups being analyzed by the model. The size of each dimension is the *granularity* of the discretization, and thus the model.

For example, a PHC with granularity 100 meant to represent 3 demographic groups will have the shape  $(100, 100, 100)$ , and can be thought of as a  $100 \times 100 \times 100$  matrix.

The demographic voting probabilities (DVP), the collection of  $b_i$  from King's model specification, are derived for a cell directly from that cell's position within the PHC. The higher a cell is along a dimension of the tensor, the higher the probability that members of that dimension's corresponding demographic group will vote for a candidate.

For example, in a precinct with 2 demographic groups, and in a PHC with granularity 10, a cell at the position  $10, 0$  represents  $\frac{10}{10} = 100\%$  of demographic group 1 voting for the candidate, and  $\frac{0}{10} = 0\%$  of demographic group 2 voting for the candidate.

In general, in some PHC with granularity  $d$  representing  $n$  demographic groups, for some cell at index  $\{i_0, i_1, i_2, \dots, i_n\} = \{i_j\}_{j=0}^n$ :

$$\text{DVP} = \left\{ \frac{i_j}{d} \right\}_{j=0}^n$$

When given a precinct, candidate, and some PHC, one can determine the demographic voting probability of the precinct for the candidate by sampling a cell, or set of cells, from the PHC randomly, using the cell's component values as probabilities. The higher the component value of a cell, the more

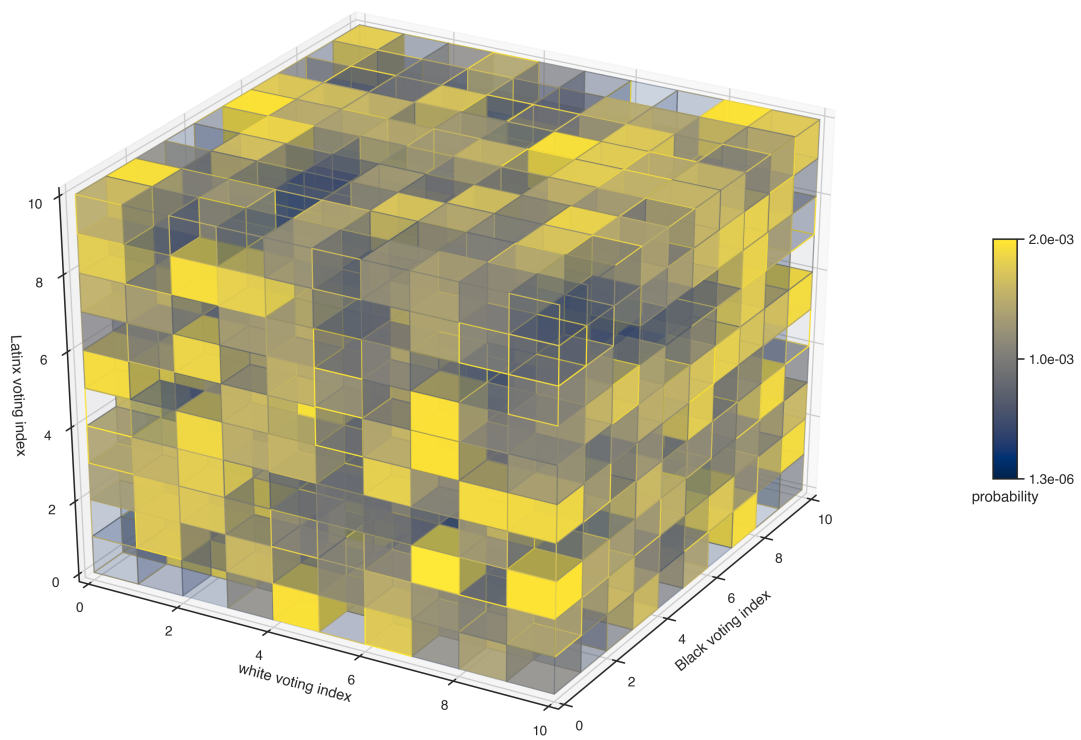
likely it is to be chosen to be the inferred demographic voting probability of the given precinct and candidate.

The granularity of the discretization is directly related to how precise the demographic voting probabilities are. In general, the higher the granularity of a PHC, the more precise its predictions may be. However, higher granularity PHCs have many more cells, and thus parameters, making sampling time/computing intensive.

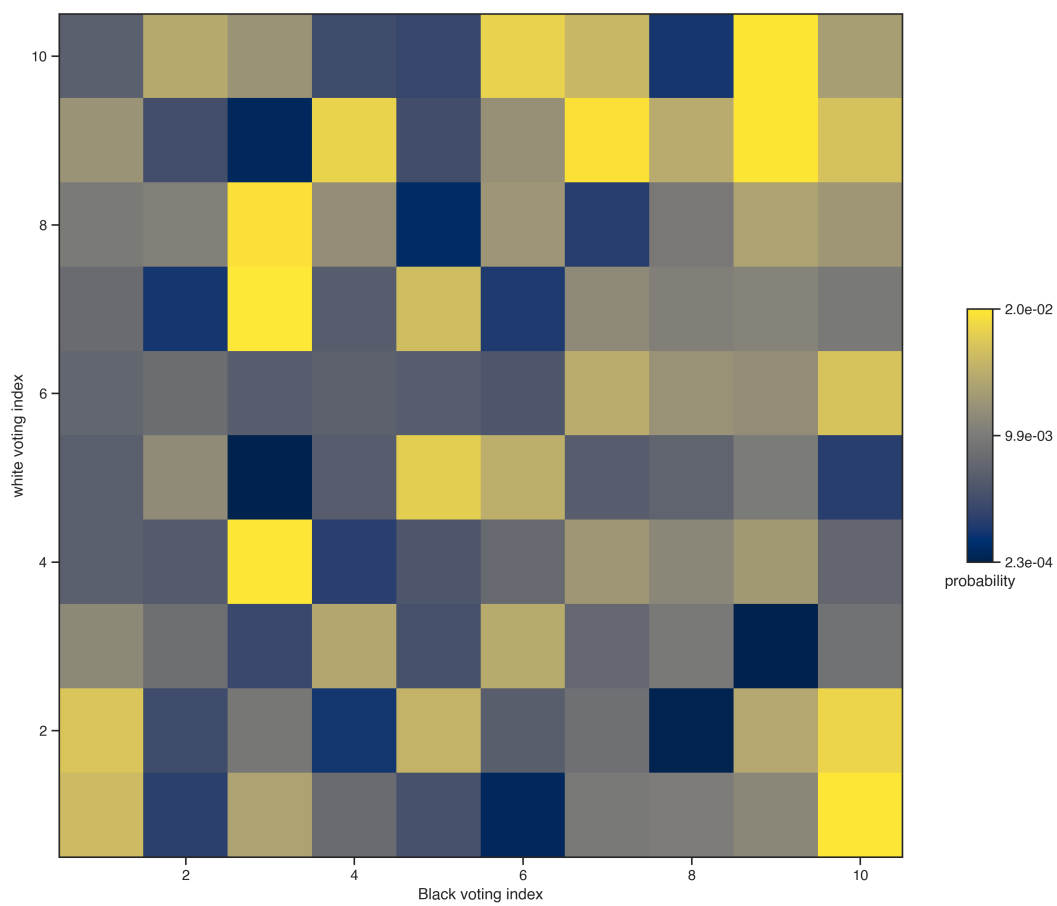
Figure 2.1 is a rendering of a rank 3, granularity 10 PHC, with the dimensions corresponding to the Latinx, white, and Black members of an electorate. Figure 2.2 is a similar rendering, but of a rank 2, granularity 10 PHC, with the dimensions corresponding to the white and Black members of an electorate. In each, the cells are colored according to their component values, or probabilities – the more yellow and opaque a cell, the higher the probability is in that cell.

The rank 2 case can also be visualized as a 3-dimensional plot, where the third dimension is the probability associated with the cells of the PHC. Figure 2.3 illustrates this.

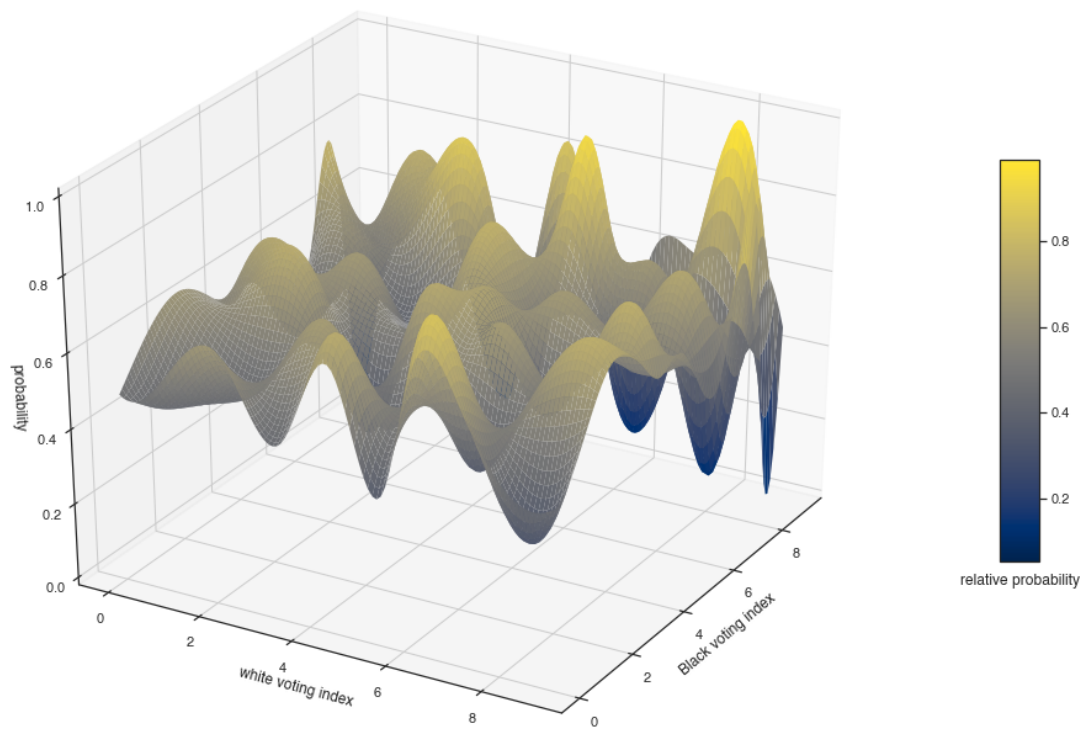
Higher rank PHCs are difficult to visualize in similar ways, but the representational and intrinsic characteristics of the object remain the same. This is critical to the Discrete Voter Model's ability to generalize to arbitrarily diverse electorates, with arbitrarily precise demographic voting probabilities.



**Figure 2.1:** Example of a Rank 3, Granularity 10 Probabilistic Hypercube



**Figure 2.2:** Example of a Rank 2, Granularity 10 Probabilistic Hypercube



**Figure 2.3:** The Probability Distribution within of a Rank 2, Granularity 10 Probabilistic Hypercube

## 2.2 EXPEC\_VOTES

TO BE ABLE TO USE PROBABILITY HYPERCUBES as states in a Markov chain Monte Carlo method, the Discrete Voter Model employs two methods of determining how “good” each PHC is. This determination, or scoring, is done with the modules `expec_votes` and `prob_votes`.

`expec_votes` calculates the expectation for the vote outcome from a PHC, and compares that outcome to the observed electoral outcome. The closer they are, the higher the probability that this PHC represents a probable set of demographic voting probabilities for the given precinct.

Let there be some PHC with granularity  $d$  representing  $n$  demographic groups, for some cell at index  $\{i_0, i_1, i_2, \dots, i_n\} = \{i_j\}_{j=0}^n$ :

Let  $i = 0, 1, \dots, m$  be the cells in the hypercube, where  $m = d^n$ .

Let  $p_i$  be the probability of being in cell  $i$ .

Let  $j = 0, 1, \dots, n$  be the demographic groups.

Let  $q_j$  be the probability of demographic group  $j$  to vote for the given candidate.  $q_{j=0}^n =$  the demographic voting probabilities from Section 2.1.

Let  $s_j$  be the population of demographic group  $j$  in the precinct.

Let  $v$  be the observed number of votes cast for the candidate in the precinct.

The expectation of a PHC is thus given by Equation 2.1.

$$\mathbb{E}(\text{PHC}) = \sum_{i=0}^m p_i \left( \sum_{j=0}^n q_j \cdot s_j \right) \quad (2.1)$$

In addition to calculating the expectation of a PHC, the `expes_votes` module compares that expectation to the observed electoral vote outcome in a precinct. It calculates the  $L1$ -norm of the difference between the two, then uses the function in Equation 2.2 to convert it to a representation of the log-probability that the PHC has the same expectation as the vote outcome.

Let  $x$  be the  $L1$ -norm of the difference between the expectation of a PHC and the electoral vote outcome,  $v$ :

$$x = \|\mathbb{E}(\text{PHC}) - v\|_1$$

$$f(x) = \log \left( \frac{1}{x + 1} \right) \quad (2.2)$$

The more negative the value of  $f(x)$ , the less likely the PHC is to be able to generate the outcome of the election and be accepted in the Markov chain.

### 2.3 PROB\_VOTES

THE DISCRETE VOTER MODEL also uses the `prob_votes` module to score a probability hypercube. `prob_votes` directly calculates the probability that a given hypercube produced the observed election

outcome.

With the same notation as `expes_votes` in Section 2.2, additionally:

Let  $a_j$  be the number of people in demographic group  $j$  that voted for the given candidate.

Equation 2.3 is the probability that a given hypercube produced a given election outcome.

$$\mathbb{P}(\text{PHC} \rightarrow v) = \sum_{\forall a_j \text{ s.t. } \sum_{j \in [0, n]} a_j = v} \left( \prod_{j=0}^n q_j^{a_j} \cdot (1 - q_j)^{s_j - a_j} \right) \cdot \frac{v!}{\prod_{a_j} a_j!} \quad (2.3)$$

Equation 2.3 uses the Binomial distribution to calculate the probability of members of different demographic groups voting together in a way to produce the desired outcome. This is extensible to any number of demographic groups and can be repeated for any number of candidates.

A crucial step in this process is generating sets of  $\{a_j\}_{j=0}^n$  such that:

$$\sum_{\forall a_j \text{ s.t. } \sum_{j \in [0, n]} a_j = v}$$

These sets are all possible partitions of the observed electoral outcome into the demographic groups.

For example, if 10 people voted for a candidate, and there are three demographic groups, 3 possible partitions for the people who voted for the candidate are shown in Table 2.1.

**Table 2.1:** Partitions of 10 votes among 3 demographic groups

Partition	Group 1 Votes	Group 2 Votes	Group 3 Votes
1	4	3	3
2	2	0	8
3	2	4	4



Integer partitioning can be computationally expensive. The time complexity of most algorithms for it is  $O(v(\log v))$ , where  $v$  is the vote outcome of an election. However, the Discrete Voter Model limits the size of partitions to be the number of demographic groups in the electorate,  $n$ . While DVM generalizes well to a large number of demographic groups,  $n$  is unlikely to get large enough to make this integer partitioning step a time bottleneck. With that said, most algorithms for this process have space complexity  $O(n^2)$ . While limiting the size of partitions does reduce the space necessary, to further mitigate that use of space, DVM creates the set of integer partitions as a Python generator. Hence, values are treated like a stream, where computation is delayed until a partition is necessary, and the set of partitions is treated like a regular iterator. Finally, since integer partitions are used several times in the computation, and only depend on the vote outcome and the shape of the PHC, not the values in the PHC, the partitions can be generated once, and reused at every step in the MCMC method.

The full code for this partitioning process can be found in the appendix, either in the full repository or as Listing [A.2](#).

The partitions are used to calculate the probability, as shown in Equation [2.3](#). That probability is that of some candidate's and precinct's hypercube producing the observed electoral outcome. `prob_votes` uses this to calculate the log-probability. The more negative that value, the less likely the PHC generated the outcome of the election and is accepted in the Markov chain. This scoring method is more intuitive and perhaps more precise, but it is much more computationally expensive. The calculations of large products, exponents, and especially factorials makes this much slower than scoring by expectation with `expec_votes`.

## 2.4 ELECT

THE DISCRETE VOTER MODEL operates on data from elections, and the `elect` module prepares and packages the data for use.

This module implements the `Election` class, a data structure and accompanying methods that represents an election in a district. It has the following attributes:

1. `candidates`: the candidates
2. `dpp`: demographic distribution, per precinct
3. `dvp`: the demographic voting probabilities, per precinct
4. `mock`: whether the election is a mock election
5. `name`: the name of the election
6. `num_demo_groups`: the number of demographic groups in the district
7. `outcome`: the vote outcome, by demographic group, of the election, per precinct
8. `precincts`: the precinct IDs
9. `vote_totals`: the vote totals of each candidate, per precinct
10. `winner`: the winner of the election

An `Election` object for a mock election is initialized by providing the candidates, demographic distribution per precinct, and demographic voting probabilities (DVP) per precinct. Upon initialization, an election is simulated based on how each demographic votes, and other attributes, like the outcome, vote totals, and winner, are calculated.

That DVP is typically the desired result of the Discrete Voter Model, and by providing it to mock elections that then use it to generate election data, one can run the model on the object and attempt to recover the “truth” that generated the results of the election. This is covered in depth in Section 2.6.

An `Election` object for a real election is initialized by providing the candidates, demographic distribution per precinct, and vote totals per precinct. Upon initialization, the election is “decided” by calculating which of the candidates won.

The `elect` module also provides a function, `create_elections`, that takes demographic data per precinct and voting data per precinct, both as the ubiquitous `Pandas DataFrames`, to create an `Election` object. This allows DVM to be deployed quickly and easily to real election data. Section 2.6 employs it for the case studies of North Carolina and Chicago.

The Discrete Voter Model uses an `Election` object’s candidates, demographic distribution per precinct, and vote totals to create probabilistic hypercubes and run the Markov chain Monte Carlo method to find a distribution of PHCs that could have resulted in that `Election`’s outcomes. For mock elections, the accuracy of DVM can be compared to the original DVP that generated the `Election` object.

## 2.5 DVM

THE FINAL MODULE, and the one that runs the entire process, is `dvm`. It implements the Markov chain Monte Carlo method with Hamiltonian Monte Carlo (HMC) and Random Walk Metropolis

(RWM) kernels, with both probability-based and expectation-based state scoring.

The Discrete Voter Model's MCMC method used probabilistic programming Python package **TensorFlow Probability** (TFP) as a substrate. TFP is a relatively new probabilistic programming package maintained by Google and the TensorFlow team. It is built on their popular open source machine learning platform, TensorFlow, and offers many advantages over a pure Python implementation of MCMC.<sup>19</sup> The top three advantages to the Discrete Voter Model are:

1. Nearly all of the operations are vectorized, allowing for much more efficient computation on high-dimensional objects like probabilistic hypercubes.
2. It automatically scales to different platforms and architectures. Without any difference in the code, this allows the Discrete Voter Model to be able to be run on CPUs (central processing units) found in most computers, GPUs (graphics processing units) found in more powerful computers and sometimes optimized for machine learning, and TPUs (tensor processing units) developed by Google as AI accelerator application-specific integrated circuits (ASIC). As one increases the granularity of their PHC or the complexity of an election, they will be able to benefit from the increased computing power of modern machines.
3. TensorFlow, and thus TensorFlow Probability, has lazy evaluation built in, through the Graph structure. TensorFlow Graphs employ the dataflow programming paradigm to increase the efficiency of parallel computations.<sup>24</sup> This allows for:
  - (a) easy identification and execution of operations that can be computed in parallel.
  - (b) distributed computing to take advantage of several machines at once.
  - (c) saving the model as a language-agnostic Graph. The Graph can then be executed in

any of the many TensorFlow-compatible languages. For example, a graph generated in Python can be executed in C++.

(d) visualizing the Graph with built-in tools.

Specifically, TFP is used to create the transition kernels and sample from the Markov chain.

The final subroutine, `mcmc`, runs the Markov chain Monte Carlo method with the Metropolis-Hastings algorithm, employing either `expes_votes` or `prob_votes` to score candidates.

The pseudocode for the MCMC is as follows:

1. Choose a kernel: either Random Walk Metropolis or Hamiltonian Monte Carlo
2. Choose a scoring function,  $g$ : either `expes_votes` or `prob_votes`
3. Initialize some probabilistic with `phc`,  $X_i$  where  $X_{i=0}^t$  is the full sample. Score  $X_i$  it with the kernel of choice to get  $s_i$
4. Iterate  $t$  times. At each iteration,  $i$ :
  - (a) Generate a candidate hypercube,  $Y_{i+1}$ , with the kernel
  - (b) Calculate the score for that candidate,  $s_{i+1} = g(Y_{i+1})$
  - (c) Accept  $X_i$  with probability  $\alpha = \frac{s_{i+1}}{s_i}$ 
    - i. If accepted,  $X_{i+1} \leftarrow Y_{i+1}$
    - ii. If rejected,  $X_{i+1} \leftarrow X_i$
5. Burn a predetermined fraction of the observations at the beginning of the sample. Removing that early “burn-in” period has been shown by some studies to be effective in creating a better sample.<sup>20</sup>
6. Output the sample as the distribution of probability hypercubes.

The HMC kernel, `tfp.mcmc.HamiltonianMonteCarlo`, is initialized with an adaptive step size for the first 60% of the burn-in period. The adaptive step size is controlled by `tfp.mcmc.SimpleStepSizeAdaptation`, and is based on research by Andrieu and Thoms that shows that when the first few steps of an HMC-directed chain are allowed to adapt the step size towards the direction of the log probability of the state, the succeeding steps are more fruitful.<sup>8</sup>

The RWM kernel, `tfp.mcmc.RandomWalkMetropolis`, is initialized with a proposal function that chooses two cells of the probability hypercube at random, and adds a small value,  $\epsilon$ , to one while subtracting it from the other. That proposal function is reversible and allows for good mixing, while also being simple enough to be performant.

Both kernels trace the log-probability and log-acceptance rate of the chain during sampling.

`dvm` implements the lower level functions for running the MCMC, `rwm` and `hmc`. It also implements the higher level function `dvm_election` that applies the Discrete Voter Model to an `Election` object.

Finally, the module implements functions, `chain_mle` and `mean_phc`, for finding the Maximum Likelihood Estimate and the sample mean of distributions/samples produced by the model. These are quite useful when determining the central tendencies of the sample, and evaluating it easily.

## 2.6 EVALUATION

THERE ARE TWO AUXILIARY MODULES, `dvm_plot` and `dvm_eval`, that provide visualizations and analysis of the Discrete Voter Model's results and behavior.

`dvm_plot` provides a function to generate trace plots of the chains' results, as well as a bevy of functions to visualize the distribution of PHCs and the PHCs themselves. Those functions generated Figures 2.1, 2.2, and 2.3.

All of the methods for inferring racially polarized voting noted in Chapter 1, including ecological regression, homogeneous precincts, and King's EI, are necessarily imperfect – as is the Discrete Voter Model. The extent of that imperfection, however, can be evaluated and compared.

`dvm_eval` is crucial to evaluating the Discrete Voter Model's performance generally, as well as how PHC granularity, the number of MCMC iterations, and kernel and score function choices interact with differences in election sizes and types to produce different results. Those results can be compared to the mock `Election` object's `dvp` attribute which generated the election's data. The `dvm_evaluator` function does this comparison by running DVM and calculating the mean squared error (MSE), a common loss metric, between DVM's result and the `Election`'s DVP. The lower the MSE, the more accurate the model's predictions. `dvm_evaluator` also records the time taken to run the Discrete Voter Model.

In addition, for the  $2 \times 2$  case, `dvm_eval` has the function `kei_evaluator` that can similarly evaluate the performance of the canonical King's Ecological Inference method on mock elections. Those

results can then be compared to the Discrete Voter Model's.

A set of experiments were created and run to determine the viability of the Discrete Voter Model.

Those experiments' specifications and results are in [Results](#).



*Nulla facilisi. In vel sem. Morbi id urna in diam dignis-  
sim feugiat. Proin molestie tortor eu velit. Aliquam erat  
volutpat. Nullam ultrices, diam tempus vulputate egestas,  
eros pede varius leo.*

Quoteauthor Lastname

# 3

## Results

### EXPERIMENT SPECIFICATION

#### **2 × 2** Case

For the 2 demographic group case, 5 random elections were generated for each of the demographic distributions of the precinct in Table 3.1 for precincts of size 100. Each of these distributions is labeled

$d_i \forall i \in [1, 3]$  (without loss of generality) for reference in Table 3.5.

The true demographic voting patterns for candidates  $a$  and  $b$  in Table 3.2. Each of these voting patterns is labeled  $v_i \forall i \in [1, 2]$  (without loss of generality) for reference in Table 3.5.

**Table 3.1:** Demographic Distributions Tested in the  $2 \times 2$  Case

Label	Group 1 (%)	Group 2 (%)
$d_1$	50	50
$d_2$	25	75
$d_3$	10	90

**Table 3.2:** Demographic Voting Patterns Tested in the  $2 \times 2$  Case

Label	Group	Cand. $a$ (%)	Cand. $b$ (%)
$v_1$	1	50	50
	2	30	70
$v_2$	1	25	75
	2	20	80

### $3 \times 2$ Case

For the 3 demographic group cases, 5 random elections were generated for each of the demographic distributions of the precinct in Table 3.3 precincts of size 100. Each of these distributions is labeled  $d_i \forall i \in [4, 6]$  (without loss of generality) for reference in Table 3.5.

The true demographic voting patterns for candidates  $a$  and  $b$  in Table 3.4. Each of these voting patterns is labeled  $v_i \forall i \in [3, 4]$  (without loss of generality) for reference in Table 3.5.

## ACCURACY AND RUNTIME

Accuracy is reported as the mean squared error (MSE) between the models' outputted demographic voting distribution and the true demographic voting distribution, averaged over all trials. The lower

**Table 3.3:** Demographic Distributions Tested in the  $3 \times 2$  Case

Label	Group 1 (%)	Group 2 (%)	Group 3 (%)
$d_4$	50	50	0
$d_5$	25	25	50
$d_6$	33	33	34

**Table 3.4:** Demographic Voting Patterns Tested in the  $3 \times 2$  Case

Label	Group	Cand. $a$ (%)	Cand. $b$ (%)
$v_3$	1	20	80
	2	80	20
	3	40	60
$v_4$	1	30	70
	2	40	60
	3	50	50

the MSE, the more accurate the model. The runtime is in seconds. The lower the runtime, the faster the model.

Table 3.5 summarizes the results of running DVM and King’s EI on all of the experiments, for a total of 12 experiments replicated 5 times each, in Python 3.7.4 on a 2018 MacBook Pro with a 2.3 GHz Quad-Core Intel Core i5 processor, 16 GB of 2133 MHz LPDDR3 RAM, and a Intel Iris Plus Graphics 655 1536 MB graphics card. All numbers are given with 3 significant figures of precision. The DVM’s MCMC was configured with 200 steps, and the hypercubes had a granularity of 10.

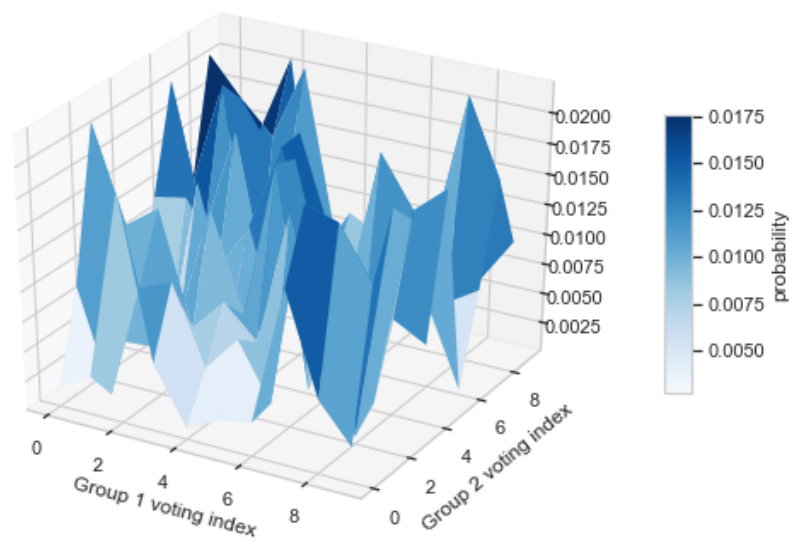
In addition to the presentation of the numeric evaluation, a presentation of select visualizations of the models produced by the Discrete Voter Model show its expressive power.

Figure 3.1 shows the distribution of probability over the 2D hypercube for a model in the  $2 \times 2$  case. The figure is 3-dimensional, with the third dimension being the probability of being in one of the cells of the hypercube. This is multinomial and much more expressive of possible outcomes.

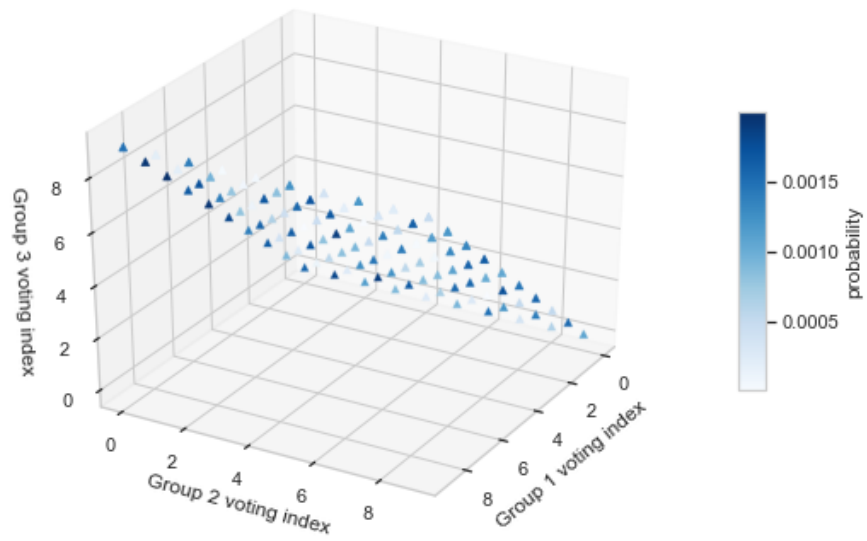
**Table 3.5:** Accuracy and Runtime of the Discrete Voter Model (DVM) and King's EI (KEI)

Label	DVM		KEI	
	MSE	Runtime (s)	MSE	Runtime (s)
$d_1 \times v_1$	0.0825	3.08	0.0253	13.9
$d_1 \times v_2$	0.0893	3.04	0.0589	14.0
$d_2 \times v_1$	0.153	4.05	0.0665	16.9
$d_2 \times v_2$	0.194	2.89	0.0604	11.5
$d_3 \times v_1$	0.107	2.74	0.0134	9.77
$d_3 \times v_2$	0.103	2.72	0.0620	9.55
$d_4 \times v_3$	0.0692	1500	-	-
$d_4 \times v_4$	0.0825	1470	-	-
$d_5 \times v_3$	0.236	1570	-	-
$d_5 \times v_4$	0.156	1530	-	-
$d_6 \times v_3$	0.149	1530	-	-
$d_6 \times v_4$	0.125	1500	-	-

Figure 3.2 shows the distribution of probability over the 2D hypercube for a model in the  $3 \times 2$  case. Since there are 3 dimensions for demographic groups, Figure 3.2 presents the fourth dimension, probability, as directed arrows across a cross section of the dimensions. The colors of these arrows give the corresponding probability.



**Figure 3.1:** Distribution of Probability over the 2D Hypercube for a Model in the  $2 \times 2$  Case



**Figure 3.2:** Distribution of Probability over the 3D Hypercube for a Model in the  $3 \times 2$  Case

*Nulla facilisi. In vel sem. Morbi id urna in diam dignis-  
sim feugiat. Proin molestie tortor eu velit. Aliquam erat  
volutpat. Nullam ultrices, diam tempus vulputate egestas,  
eros pede varius leo.*

Quoteauthor Lastname

# 4

## Discussion

The results presented in Table 3.5 show three primary trends:

1. the Discrete Voter Model performs nearly as well as, but a little worse than, King's Ecological Inference in the  $2 \times 2$  case (two demographic groups and two candidates)
2. the Discrete Voter Model is remarkably faster than King's Ecological Inference in the  $2 \times 2$  case (two demographic groups and two candidates)

3. unlike King's EI, the Discrete Voter Model can produce reliable results for the  $3 \times 2$  case (three demographic groups and two candidates)

The Discrete Voter Model was limited to 200 steps, with a hypercube granularity of 10 (that is, each dimension of the cube went from 0 to 9). Further testing corroborates the theory that, generally, with more iterations, Markov chain Monte Carlo converges to the true distribution. Hence, allowing DVM to run longer may increase the accuracy of the model, at the expense of runtime.

With that said, DVM was able to use those 200 steps to perform comparably to King's EI in the  $2 \times 2$  cases. This also corresponded to a lower runtime than King's EI in all cases, which is a testament to improvements in computing power and algorithms.

DVM was also able to do what King's EI cannot in any amount of time: run on the larger  $3 \times 2$  tables. Higher dimensional tables are also possible with DVM, with no changes to the statistical theory or implementation.

If King's EI is taken to be a standard method, as it has been by some courts, DVM's performance on these generated elections show that it is a good contender. Not only does it perform comparably, but:

- it is inherently extensible and can be tuned. Generally, at the expense of time:
  - one can increase the number of iterations to get closer to the true distribution
  - one can increase the size of the hypercube to get more granular estimates for the racial voting patterns
- it can produce visualizations of more complex voting pattern distributions

The Discrete Voter Model leverages the power of discretization and Markov chain Monte Carlo to



provide more expressive models within a more malleable and sound structure. Its flexibility has been shown on generated election data, and improvements continue to be made to its speed.

As Jim Greiner said: "Racial bloc voting, often a hotly contested issue in an "ordinary" redistricting dispute, may now become even more so, and as a result, using the best methodology is more critical."<sup>15</sup>

The Discrete Voter Model, as presented here and in future iterations, is a direct response to the need for more flexible and robust methods in the redistricting and voting rights space.

# 5

## Conclusion

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor

ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetur erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit, vel lacinia enim metus eu nunc.

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Vestibulum tortor quam, feugiat vitae, ultricies eget, tempor sit amet, ante. Donec eu libero sit amet quam egestas semper. Aenean ultricies mi vitae est. Mauris placerat eleifend leo. Quisque sit amet est et sapien ullamcorper pharetra. Vestibulum erat wisi, condimentum sed, commodo vitae, ornare sit amet, wisi. Aenean fermentum, elit eget tincidunt condimentum, eros ipsum rutrum orci, sagittis tempus lacus enim ac dui. Donec non enim in turpis pulvinar facilisis. Ut felis.

Cras sed ante. Phasellus in massa. Curabitur dolor eros, gravida et, hendrerit ac, cursus non, massa. Aliquam lorem. In hac habitasse platea dictumst. Cras eu mauris. Quisque lacus. Donec ipsum. Nullam vitae sem at nunc pharetra ultricies. Vivamus elit eros, ullamcorper a, adipiscing sit amet, porttitor ut, nibh. Maecenas adipiscing mollis massa. Nunc ut dui eget nulla venenatis aliquet. Sed

luctus posuere justo. Cras vehicula varius turpis. Vivamus eros metus, tristique sit amet, molestie dignissim, malesuada et, urna.

Cras dictum. Maecenas ut turpis. In vitae erat ac orci dignissim eleifend. Nunc quis justo. Sed vel ipsum in purus tincidunt pharetra. Sed pulvinar, felis id consectetur malesuada, enim nisl mattis elit, a facilisis tortor nibh quis leo. Sed augue lacus, pretium vitae, molestie eget, rhoncus quis, elit. Donec in augue. Fusce orci wisi, ornare id, mollis vel, lacinia vel, massa.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi commodo, ipsum sed pharetra gravida, orci magna rhoncus neque, id pulvinar odio lorem non turpis. Nullam sit amet enim. Suspendisse id velit vitae ligula volutpat condimentum. Aliquam erat volutpat. Sed quis velit. Nulla facilisi. Nulla libero. Vivamus pharetra posuere sapien. Nam consectetur. Sed aliquam, nunc eget euismod ullamcorper, lectus nunc ullamcorper orci, fermentum bibendum enim nibh eget ipsum. Donec porttitor ligula eu dolor. Maecenas vitae nulla consequat libero cursus venenatis. Nam magna enim, accumsan eu, blandit sed, blandit a, eros.

Quisque facilisis erat a dui. Nam malesuada ornare dolor. Cras gravida, diam sit amet rhoncus ornare, erat elit consectetur erat, id egestas pede nibh eget odio. Proin tincidunt, velit vel porta elementum, magna diam molestie sapien, non aliquet massa pede eu diam. Aliquam iaculis. Fusce et ipsum et nulla tristique facilisis. Donec eget sem sit amet ligula viverra gravida. Etiam vehicula urna vel turpis. Suspendisse sagittis ante a urna. Morbi a est quis orci consequat rutrum. Nullam egestas feugiat felis. Integer adipiscing semper ligula. Nunc molestie, nisl sit amet cursus convallis, sapien lectus pretium metus, vitae pretium enim wisi id lectus. Donec vestibulum. Etiam vel nibh. Nulla facilisi. Mauris pharetra. Donec augue. Fusce ultrices, neque id dignissim ultrices, tellus mauris dictum elit,

vel lacinia enim metus eu nunc.



# Appendix

The full set of code for this research and method can be found [here](#). To completely reproduce these experiments, run `methods.ipynb` in full. The following are notable code snippets mentioned in the body of the article.

**Listing A.1:** Python Implementation of King's Ecological Inference

```
import numpy as np
```

```

import pymc3 as pm

def ei_two_by_two(demo_pcts, candidate_pcts, precinct_populations, lambda
=0.5):
    """
    Run King's Ecological Inference method on
    a 2x2 example (2 demographic groups).

    group_demo_pcts (NumPy array): the percentage of people in the
    demographic group for each precinct

    group_voting_pcts (NumPy array): the percentage of people in the
    precinct who voted for a candidate

    precinct_populations (NumPy array): the populations of the
    precincts

    lambda (float):

    return: the probabilistic model
    """
    demo_counts = candidate_pcts * precinct_populations

    # Number of populations

    p = len(precinct_populations)

    with pm.Model() as model:
        c_1 = pm.Exponential('c_1', lambda)

```

```

d_1 = pm.Exponential('d_1', lambda)

c_2 = pm.Exponential('c_2', lambda)

d_2 = pm.Exponential('d_2', lambda)


b_1 = pm.Beta('b_1', alpha=c_1, beta=d_1, shape=p)

b_2 = pm.Beta('b_2', alpha=c_2, beta=d_2, shape=p)


theta = demo_pcts * b_1 + (1 - demo_pcts) * b_2

Tprime = pm.Binomial('Tprime', n=precint_populations , p=theta,
                      observed=demo_counts)

return model

```

#### Listing A.2: Integer Partitioning

```

from itertools import chain, permutations

def integer_partition(n, k, min_size=0):
    """
    Partition an integer.

    n (int): the integer to partition
    k (int): the number of elements in a partition
    min_size (int): the minimum size of an element
    in the partition

```



```

    return: a generator of partitions as tuples

    """

    if k < 1:

        return

    if k == 1:

        if n >= min_size:

            yield (n,)

        return

    for i in range(min_size, n // k + 1):

        for result in integer_partition(n - i, k - 1, i):

            yield (i,) + result

def permute_integer_partition(n, k, min_size=0):

    """

    Partition an integer, with all permutations

    n (int): the integer to partition

    k (int): the number of elements in a partition

    min_size (int): the minimum size of an element

    in the partition

    return: a generator of all permutations of partitions as tuples

```

```

"""

return chain.from_iterable(set(permutations(p)) for p in
                           integer_partition(n, k, min_size))

```

**Listing A.3:** Generate a Random Election

```

import numpy as np

def generate_random_election(candidates, demo, beta):
    """
    Generate a random election.

    candidates (string list): the candidates
    demo (dict): the demographics of the electorate
    beta (dict): the theoretical voting percentages of
    each demographic group, for each candidate

    return: a dictionary of candidates and the vote breakdowns by
    demographic group
    """
    # Set up the result dictionary
    num_groups = len(demo)
    result = {'a': (0, [0] * num_groups),
              'b': (0, [0] * num_groups),

```

```

        'c': (0, [0] * num_groups)}

    # Iterate through each demographic group
    for group_index, group in enumerate(demo):

        # Simulate each voter
        for voter in range(demo[group]):

            vote = np.random.choice(candidates, 1, beta[group])[0]

            prev_total, prev_breakdown = result[vote]

            prev_breakdown[group_index] += 1

            result[vote] = prev_total + 1, prev_breakdown

    return result

```

**Listing A.4:** Evaluate and Compare DVM and King's EI

```

import numpy as np

from operator import mul

import functools

import random

import pymc3 as pm

import time

from itertools import chain, permutations

def dvm_king_evaluator(election, demo, beta, label, n_iter=1, met_iter=200)

```

```

:

"""
Run and compare the results of the Discrete
Voter Model and King's EI on generated
election data.

election (dict): the random election to
evaluate on

demo (dict): the demographic dictionary of
a precinct

label (string): the label of the experiment

n_iter (int): the number of times to repeat
the experiment

met_iter (int): the number of iterations to use
for Metropolis-Hastings

return: a dictionary of the label and times and MSEs for
the Discrete Voter Model and King's EI
"""

results = {}

dvm_total_time = 0

dvm_total_mse = 0

```

```

king_total_time = 0

king_total_mse = 0

temp_pcts_array = [pcts[0] for group, pcts in beta.items()]

true_pcts = np.fromiter(temp_pcts_array, dtype=float)

for _ in range(n_iter):

    # Get the observed votes for candidate a

    candidate_a_obs = random_election_1_1['a'][0]

    # Run Metropolis-Hastings and time it

    dvm_total_time -= time.time()

    initial_grid = make_grid(len(demo), 10)

    met_result = metropolis_hastings(met_iter, initial_grid,

        candidate_a_obs, demo, scoring_type='prob')

    dvm_total_time += time.time()

    # Find the best grid and output the result

    best_grid = met_result['best_grid']

    best_cell = get_most_probable_cell(met_result['best_grid'])

    vote_pcts = get_vote_pcts(best_cell, 10, demo)

    # Find the MSE of the result

```

```

dvm_mse_array = np.fromiter(vote_pcts.values(), dtype=float)

dvm_total_mse += mse(dvm_mse_array, true_pcts)

# Run King's EI and time it, if at the right dimension

if len(demo) > 2:

    continue

king_demo = list(demo.values())[0] / 100

king_cand_vote = candidate_a_obs / 100

king_prec_pop = np.array([100])

king_total_time -= time.time()

king_model = ei_two_by_two(king_demo, king_cand_vote, king_prec_pop

    )

with king_model:

    king_trace = pm.sample()

king_total_time += time.time()

# Find the MSE of the result

king_mse_array = np.fromiter([king_trace.get_values('b_1').mean(),

                                king_trace.get_values('b_2').mean()],

                                dtype=float)

king_total_mse += mse(king_mse_array, true_pcts)

```

```
return {'label': label,  
        'dvm_time': dvm_total_time / n_iter,  
        'dvm_mse': dvm_total_mse / n_iter,  
        'king_time': king_total_time / n_iter,  
        'king_mse': king_total_mse / n_iter}
```

# References

- [Hol] *Holder v. Hall*, 512 U.S. 874, 904 n.13 (1994).
- [cle] *League of United Am. Citizens v. Clements*, 999 F.2d 831, 860 (5th Cir. 1993) (*en banc*).
- [LUL] *League of United Latin American Citizens v. Perry*, 126 S. Ct. 2594 (2006).
- [Lew] *Lewis v. Alamance County*, 99 F.3d 600, 604 n.3 (4th Cir. 1996).
- [tho] *Thornburg v. Gingles*. *Oyez*.
- [vra] Voting Rights Act (1965). *Our Documents - Voting Rights Act (1965)*.
- [7] (2015). Introduction to federal voting rights laws. *The United States Department of Justice*.
- [8] Andrieu, C. & Thoms, J. (2008). A tutorial on adaptive mcmc. *Statistics and Computing*, 18(4), 343–373.
- [9] Betancourt, M., Byrne, S., Livingstone, S., & Girolami, M. (2017). The geometric foundations of hamiltonian monte carlo. *Bernoulli*, 23(4A), 2257–2298.
- [10] Christian, R. & Casella, G. (2012). A short history of markov chain monte carlo: Subjective recollections from incomplete data. *arXiv.org*, 26(1).
- [11] Commons, W. (2014). File:traviscountydistricts.png — wikimedia commons, the free media repository. [Online; accessed 15-November-2019].
- [12] Deruelle, N. & Uzan, J.-P. (2018). Hamiltonian mechanics. In *Relativity in Modern Physics*. Oxford University Press.
- [13] Ferree, K. E. (2004). Iterative approaches to r x c ecological inference problems: Where they can go wrong and one quick fix. 12(2).
- [14] Gelman, A., Park, D., Ansolabehere, S., Price, P., & Minnite, L. (2001). Models, assumptions and model checking in ecological regressions. *Journal Of The Royal Statistical Society Series A-Statistics In Society*, 164, 101–118.



- [15] Greiner, D. J. (2007). Ecological inference in voting rights act disputes: Where are we now, and where do we want to be? *Jurimetrics*, 47(2), 115–167.
- [16] King, G. (1997). *A Solution to the Ecological Inference Problem: Reconstructing Individual Behavior from Aggregate Data*. Princeton, New Jersey: Princeton University Press.
- [17] King, G. (2013). *A Solution to the Ecological Inference Problem: Reconstructing Individual Behavior from Aggregate Data*. Princeton: Princeton University Press.
- [18] King, G., Rosen, O., & Tanner, M. A. (1999). Binomial-beta hierarchical models for ecological inference. *Sociological Methods & Research*, 28(1), 61–90.
- [19] Lao, J., Suter, C., Langmore, I., Chimisov, C., Saxena, A., Sountsov, P., Moore, D., Saurous, R., Hoffman, M., & Dillon, J. (2020). tfp.mcmc: Modern markov chain monte carlo tools built for modern hardware. *arXiv.org*.
- [20] Liang, F. F. (2010). *Advanced Markov chain Monte Carlo methods: learning from past samples*. Wiley series in computational statistics. Chichester, England: Wiley.
- [21] Neal, R. (2012). Mcmc using hamiltonian dynamics. *arXiv.org*.
- [22] Roberts, G. O. & Rosenthal, J. S. (2004). General state space markov chains and mcmc algorithms. *Probab. Surveys* 1, 1(20-71), 20–71.
- [23] Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*. Wiley series in probability and mathematical statistics. Applied probability and statistics. New York ;: Wiley.
- [24] Wongsuphasawat, K., Smilkov, D., Wexler, J., Wilson, J., Mane, D., Fritz, D., Krishnan, D., Viegas, F. B., & Wattenberg, M. (2018). Visualizing dataflow graphs of deep learning models in tensorflow. *IEEE Transactions on Visualization and Computer Graphics*, 24(1), 1–12.



**T**HIS THESIS WAS TYPESET using  $\text{\LaTeX}$ , originally developed by Leslie Lamport and based on Donald Knuth's  $\text{\TeX}$ . The body text is set in 11 point Egenolff-Berner Garamond, a revival of Claude Garamont's humanist typeface. The above illustration, *Science Experiment 02*, was created by Ben Schlitter and released under [CC BY-NC-ND 3.0](#). A template that can be used to format a PhD dissertation with this look & feel has been released under the permissive [AGPL](#) license, and can be found online at [github.com/suchow/Dissertate](https://github.com/suchow/Dissertate) or from its lead author, Jordan Suchow, at [suchow@post.harvard.edu](mailto:suchow@post.harvard.edu).