



UNIVERSIDAD DE CHILE
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS
DEPARTAMENTO DE INGENIERÍA INDUSTRIAL

**OUTLIER DETECTION FOR MULTICANDIDATE ELECTIONS WITH
DEMOGRAPHIC GROUPS**

TESIS PARA OPTAR AL GRADO DE MAGÍSTER EN GESTIÓN DE OPERACIONES

MEMORIA PARA OPTAR AL TÍTULO DE INGENIERO CIVIL INDUSTRIAL

PABLO UBILLA PAVEZ

PROFESOR GUÍA:

Charles Thraves Cortés-Monroy

MIEMBROS DE LA COMISIÓN:

Denis Sauré Valenzuela

Fernando Ordóñez Pizarro

Esta tesis ha sido parcialmente financiado por: FONDO PUENTE DAI - FCFM 2023

Powered@NLHPC: Esta tesis fue parcialmente apoyada por la infraestructura de
supercómputo del NLHPC (ECM-02)

SANTIAGO DE CHILE

2023

DETECCIÓN DE ANOMALÍAS EN ELECCIONES CON MÚLTIPLES CANDIDATOS Y GRUPOS DEMOGRÁFICOS

La inferencia ecológica es una técnica para estimar el comportamiento individual utilizando datos agregados. Un caso particular se encuentra en las elecciones políticas, donde en cada mesa electoral conocemos los votos de los candidatos y el número de votantes de distintos grupos demográficos (como la edad, el sexo y la nacionalidad). En este trabajo, aplicamos el algoritmo EM para estimar las probabilidades de voto de los grupos demográficos para cada candidato en un distrito determinado. Desafortunadamente, el *E-step* escala exponencialmente en el número de candidatos.

En este estudio proponemos cuatro métodos polinomiales alternativos para estimar las probabilidades del paso *E-step*: (1) simulación de escenarios utilizando un método de *hit-and-run*, (2) aproximación utilizando una distribución normal multivariada con integración de Monte Carlo o (3) una distribución normal multivariada utilizando su FDP, y (4) aproximación mediante una única multinomial. Mostramos a partir de experimentos numéricos que el método de aproximación multinomial es el más rápido, ejecutándose en menos de una centésima de segundo. Además, el error absoluto promedio de la probabilidad estimada con este método es muy similar al obtenido al realizar el algoritmo EM con la probabilidad exacta.

Implementamos los métodos propuestos en la primera vuelta de la elección presidencial de Chile de 2021. Presentamos una metodología que utiliza programación entera mixta para agregar grupos y estimar intervalos de confianza en las probabilidades estimadas mediante el uso de *bootstrapping*, de modo que el modelo pueda identificar correctamente las probabilidades de voto. Observamos que los distritos con más mesas electorales se benefician de conjuntos de grupos menos agregados en contraste con los distritos con menos mesas electorales. Finalmente, calculamos p-valores utilizando simulación con la aproximación multinomial, obteniendo 27 urnas electorales con un p-valor menor o igual a 10^{-8} .

OUTLIER DETECTION FOR MULTICANDIDATE ELECTIONS WITH DEMOGRAPHIC GROUPS

Ecological inference is a technique to estimate individual behaviour by using aggregated data. A particular case is found in political elections where in each ballot box we know the candidates' votes and the number of voters for different demographic group (like age, sex and nationality). In this work, we apply the EM-algorithm to estimate the voting probabilities of demographic groups for each candidate at a particular district. Unfortunately, the E-Step scales exponentially in the number of candidates.

We propose four alternative polynomial methods to estimate the E-Step probabilities: (1) sample scenarios using hit-and-run, (2) approximate using a multivariate normal with Monte Carlo integration or (3) a multivariate normal using its PDF, and (4) approximate by a single multinomial. We show from numerical computations that the multinomial approximation method is the fastest, running in less than a hundredth of a second. In addition, the mean absolute error of the estimated probability with this method is very similar to the one obtained when performing the EM-algorithm with the exact probability.

We run the proposed methods in the first round Chilean Presidential Election of 2021. We present a methodology that uses mixed integer programming to aggregate groups, and estimate confidence intervals on the estimated probabilities by using bootstrapping, so that the model can correctly identify the voting probabilities. We observe that districts with more ballot boxes benefit from less aggregated group sets in contrast to districts with less ballot boxes. Finally, we compute p-values using simulation with the multinomial approximation, obtaining 27 ballot-boxes with a p-value lower or equal than 10^{-8} .

Pa' mi weli Ruchi

Agradecimientos

Primero que todo agradecer a mi maire y mi paire, que me han apoyado siempre. También a mis hermanos Jesu y Carlos, y al resto de mi hermosa familia.

A todas las lindas amistades que he tenido a lo largo de mi vida. Mis amigos del colegio, los del Plan Común, los de Industrias, las tierras de Quilpué, los de Leiden y tantos otros que les guardo bonitos recuerdos. Al Nico que lo quiero mucho.

A la Dorkita.

Al Luca y al Pepi, que ojalá algún día se lleven bien.

A José Miguel Alvarado por motivarme a estudiar Ingeniería Industrial.

A Charles Thraves por guiarme en este proceso de tesis. A Denis Sauré y Fernando Ordóñez por participar de la comisión.

No tengo palabras para agradecerle a Linda y toda su ayuda en la gestión de titulación. A Fernanda por su buena disposición y su motivación con el magíster.

Cómo no agradecerle a *lofi hip hop radio - beats to relax/study to* que me mantuvo concentrado incontables horas.

Quedaron un poco desordenados los agradecimientos y seguro me faltó alguien, se hace lo que se puede. Saludos a quienes lean esta tesis, ojalá les guste.

TABLE OF CONTENT

1. Introduction	1
2. Model	3
2.1. E-Step	4
2.2. M-Step	5
2.3. EM-algorithm	5
2.4. Log-likelihood	6
3. Computing the Conditional Probability	7
3.1. Exact Method (EXACT)	7
3.2. Hit-and-Run Method (H&R)	10
3.3. Multivariate-Normal Approximation Method Using the CDF (MVN-CDF) .	12
3.4. Multivariate-Normal Approximation Method Using the PDF (MVN-PDF) .	13
3.5. Multinomial Approximation Method (MULT)	14
4. Computing the p-value	16
4.1. Exact Computation	16
4.2. Simulation Using the Multinomial Approximation	17
5. Numerical Results	18
5.1. Numerical Instances	18
5.2. EM-Algorithm Probability Estimation	20
5.2.1. Time Results	20
5.2.2. Error Results	22
5.3. p-value Estimation	23
5.3.1. Probability Gap	23

5.3.2. Method Gap	25
6. Case Study: a Real Election	28
6.1. Historical Data	28
6.2. Group Aggregation	30
6.2.1. Log-Likelihood	33
6.2.2. Bootstrapping: Estimated Probability Standard Deviation	33
6.2.3. Adding More Groups: Trade-off	34
6.2.4. Best Group Aggregation	34
6.3. Results	37
6.3.1. Group Aggregation Sets	37
6.3.2. Probabilities	39
6.3.3. p-values and Outliers	40
6.4. General Recommendations	44
7. Conclusions	45
7.1. Future Work	46
Bibliography	48
ANNEXES	50
Annex A. Hit and Run Method	50
A.1. Starting Point for the Hit-and-Run Algorithm	50
A.2. Choosing the Step-Size	50
A.3. Another Formulation for the Discrete Polytope	52
Annex B. EM-algorithm Probability Estimation with 200 Voters per Ballot-Box	53
B.1. Time results	53
B.2. Error results	55
Annex C. Detailed Results for Polling Places with Low p-value Ballot-Boxes	56

Table Index

5.1.	Mean running time over 20 scenarios in seconds for the EM-algorithm for varying instances, with fixed values of $I_b = 100$ and $B = 50$	21
5.2.	Comparison of mean simulation time and mean EM running time in seconds over 20 scenarios for varying instances, with fixed values of $I_b = 100$ and $B = 50$	22
5.3.	Mean absolute prediction error over 20 scenarios for varying instances, with fixed values of $I_b = 100$ and $B = 50$	23
6.1.	p-value count in the 2021-GCE.	42
B.1.	Mean running time over 20 scenarios in seconds for the EM-algorithm for varying instances, with fixed values of $I_b = 200$ and $B = 50$	53
B.2.	Comparison of mean simulation time and mean EM running time in seconds over 20 scenarios for varying instances, with fixed values of $I_b = 200$ and $B = 50$	54
B.3.	Mean absolute prediction error over 20 scenarios for varying instances, with fixed values of $I_b = 200$ and $B = 50$	55

Figure Index

5.1.	Mean absolute error over 20 scenarios for varying values of mix λ in different instances. $I_b = 100, B = 50$	20
5.2.	Correlation between p-values mean absolute error and estimated probability mean absolute error.	24
5.3.	Comparison between p-values computed with real probability and p-values computed estimated probability.	25
5.4.	Effect of covariance determinant increase in the mean absolute error in the p-values.	26
5.5.	Comparison between p-values computed with the exact method and p-values simulated with the multinomial approximation.	27
6.1.	District count for ranges of length 25 for number of ballot-boxes.	29
6.2.	Proportions for different age ranges in polling places from the district of <i>Puente Alto</i>	30
6.3.	Proportions for different group aggregations in the district of <i>Río Negro</i>	31
6.4.	Estimated probabilities for different group aggregations.	32
6.5.	Standard deviation of estimated probability vs log-likelihood over districts of different ballot-boxes number.	35
6.6.	Fraction of districts where the age-range is included as a single group.	38
6.7.	Statistics for aggregated group sets	39
6.8.	Probability estimation for the districts of <i>Calama, Maipú</i> and <i>Villarrica</i>	40
6.9.	Results in the polling place of: <i>Colegio Filipense</i>	41
6.10.	Estimated p-value accross all ballot-boxes.	43
A.1.	Mean correlation for various step-sizes over different instances. The graph is log-scaled.	51

C.1.	Results in the polling place of: <i>Liceo de Hombres de Antofagasta Mario Bahamonde Silva</i>	56
C.2.	Results in the polling place of: <i>Escuela Mariano Latorre</i>	57
C.3.	Results in the polling place of: <i>Colegio Alborada del Mar</i>	58
C.4.	Results in the polling place of: <i>Colegio D-417 Guardiamarina Guillermo Zañartu Irigoyen</i>	59
C.5.	Results in the polling place of: <i>Escuela Presidente Carlos Ibáñez del C.</i>	60
C.6.	Results in the polling place of: <i>Colegio Clementinos</i>	60
C.7.	Results in the polling place of: <i>Escuela Gerónimo Lagos Lisboa</i>	61
C.8.	Results in the polling place of: <i>Liceo Juan Pablo II de Las Condes Local: 1</i> . .	62
C.9.	Results in the polling place of: <i>Escuela Básica Algarrobal</i>	63
C.10.	Results in the polling place of: <i>Colegio Particular Ozanam Local: 2</i>	64
C.11.	Results in the polling place of: <i>Liceo Javiera Carrera Local: 2</i>	64
C.12.	Results in the polling place of: <i>Colegio Diferencial Madre Tierra</i>	65
C.13.	Results in the polling place of: <i>Colegio Los Alerces Local: 2</i>	65
C.14.	Results in the polling place of: <i>The Mayflower School Local: 2</i>	66
C.15.	Results in the polling place of: <i>The Newland School Local: 2</i>	67
C.16.	Results in the polling place of: <i>Colegio Excelsior Local: 1</i>	68
C.17.	Results in the polling place of: <i>Colegio Excelsior Local: 2</i>	68
C.18.	Results in the polling place of: <i>Liceo Darío Salas</i>	69
C.19.	Results in the polling place of: <i>Escuela Monseñor Carlos Oviedo</i>	70

Chapter 1

Introduction

Election processes are a key part of the political systems in most countries. They are used to determine presidents, senators, mayors; to approve new laws or to make new constitutions. It is necessary that transparency and impartiality are upheld to ensure the seamless execution of these procedures and to ensure outcomes that most faithfully approximate the true will of the electorate.

Elections that are apparently clean may not be absent of unintended errors or fraudulent acts in low scale (Leemann & Bochsler, 2014). These cases may not be detected since they are not necessarily systematic or produced on a large scale. It is still an important task to detect such acts for various reasons: in close elections they could alter the final outcome, these actions could damage the trust in the democratic institutions, and prompt detection can improve the voting process for future elections (Fortin-Rittberger, Harfst, & Dingler, 2017).

In this study we will propose methods to determine how likely the results of an election are, thus, assessing its transparency. We will construct our models using the Chilean election system as a reference, however, these models could be applied to different scenarios with slight modifications.

In most Chilean election processes, voters are assigned to specific ballot-boxes in the area they live in. Ballot-boxes are located in what are call *mesas* (tables in Spanish), where a group of randomly selected citizens are in charge of assisting the voters, registering election attendance and counting the votes at the end of the election. These randomly selected citizens are called *vocales de mesa* (ballot-box committee), and their labour is fundamental in order

to have an honest procedure.

The demographic composition of voters between ballot-boxes may greatly differ. For example, there could be a ballot-box where mostly young people vote and another where is mostly old people. These differences between ballot-box composition could make the election results in each ballot-box completely different. If we chose to omit this difference between ballot-boxes we would make wrong estimations of how likely the results are in each. This is why our model that should take into account the possibility of having different demographic group composition for each ballot-box.

The number of candidates is also an important factor, and it may scale in some elections. In Chile (and the rest of the world) there are many election processes that have multiple-candidates, for example the election of councilors, senators and deputies. In a two-candidate election it is easier to compute probabilities associated with the voting distribution. However, when the number of candidates increases, it also increases the number of different combination outcomes. When we also consider demographic groups the number of different combination outcomes also increases.

Considering that we will have two sets of aggregate data: ballot-box voter composition and ballot-box outcome, we will be working with ecological inference, where one would like to estimate individual behavior by using aggregate data. There has been some applications of ecological inference to elections in other contexts like voting migration (Antweiler, 2007) and political studies (King, Tanner, & Rosen, 2004). Ecological inference also arises in different fields such as sociology (Duncan & Davis, 1953; Goodman, 1953, 1959; Glynn & Wakefield, 2010), geography (Cleave, Brown, & Payne, 1995; Withers, 2001; Anselin & Cho, 2002), epidemiology (Morgenstern, 1995; Jackson, Best, & Richardson, 2006; Wakefield, 2008). Most models in the literature focus on the 2×2 case, where for both classes of aggregate data there are only 2 categories. The $R \times C$ case has gained more attention due to its wide applications to different fields, being a more general case. Some of the studies have approached the $R \times C$ case with parametric approaches, where prior distributions are assumed.

The complexity of the general case is its main challenge, meaning that approximate solutions could be a promising approach to this problem. In this study we will focus on the general case with an arbitrary number of demographic groups and an arbitrary number of candidates, estimating a non-parametric distribution.

Chapter 2

Model

Consider an election process with candidates from a set \mathcal{C} , where $|\mathcal{C}| = C$. Voters belong to a set \mathcal{I} , where $|\mathcal{I}| = I$. Each voter belongs to a group from a set \mathcal{G} and is assigned to a designated ballot-box from a set \mathcal{B} , where $|\mathcal{G}| = G$ and $|\mathcal{B}| = B$. In this context, \mathcal{I}_g represents the subset of individuals belonging to group g , while \mathcal{I}_b represents the subset of individuals assigned to ballot-box b , where $|\mathcal{I}_b| = I_b$ and $|\mathcal{I}_g| = I_g$.

The groups delineated by \mathcal{G} could encompass demographic attributes of the voters, such as age brackets and gender. The amount of voters from group $g \in \mathcal{G}$ that vote in ballot-box $b \in \mathcal{B}$ is known and is denoted by w_{bg} . The amount of votes to candidate c in each ballot-box b is also known and is denoted by x_{bc} .

We define the probability that a voter from group $g \in \mathcal{G}$ votes for candidate $c \in \mathcal{C}$ as p_{gc} . Notably, this probability remains latent due to our access solely to the aggregate election outcome for each candidate. Our aim is to obtain an estimator \hat{p}_{gc} that maximizes the likelihood of the election result given the known data.

Let us consider \mathcal{X} to be the observed data: $\mathcal{X} = \{X_{bc} = x_{bc}\}_{b \in \mathcal{B}, c \in \mathcal{C}}$. We may also define the unobserved data as $\mathcal{Y} = \{\mathcal{Y}_{bic}\}_{b \in \mathcal{B}, i \in \mathcal{I}_b, c \in \mathcal{C}}$, where \mathcal{Y}_{bic} is a binary random variable that indicates if the i -th person from ballot-box b did vote for candidate c . We will also define $Z_{bgc} = \sum_{i \in \mathcal{I}_g \cup \mathcal{I}_b} Y_{bic}$, this is the total votes that candidate c received from group g in ballot-box b . We may also use \mathbf{X}_b , \mathbf{Y}_b and \mathbf{Z}_b to refer to the respective multidimensional random variables of ballot-box b .

2.1. E-Step

Let $\mathbf{p} = \{p_{gc}\}_{g \in \mathcal{G}, c \in \mathcal{C}}$. We may write the likelihood function and derive the log-likelihood function as follows:

$$\begin{aligned} l(\mathbf{p}; \mathcal{X}, \mathcal{Y}) &= \sum_{b \in \mathcal{B}, i \in \mathcal{I}_b, c \in \mathcal{C}} Y_{bic} \ln(p_{g(i)c}) \\ &= \sum_{b \in \mathcal{B}, g \in \mathcal{G}, c \in \mathcal{C}} Z_{bgc} \ln(p_{gc}) \end{aligned}$$

We would like to know the distribution of Y_{bic} conditional to the result in ballot-box b . We will define q_{bgc} as the probability that someone from ballot-box b and group g did vote for candidate c , conditional to the result in that ballot-box. We will use the notation $Y_{bi(g)c}$ to indicate any voter i in ballot-box b that belongs to group g . This probability may be calculated with Bayes' Theorem.

$$\begin{aligned} q_{bgc} &:= P(Y_{bi(g)c} = 1 | X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC}; \mathbf{p}) \\ &= \frac{P(X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC} | Y_{bi(g)c} = 1; \mathbf{p}) P(Y_{bi(g)c} = 1; \mathbf{p})}{P(X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC}; \mathbf{p})} \quad (2.1) \\ &= \frac{P(X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC} | Y_{bi(g)c} = 1; \mathbf{p}) \cdot p_{gc}}{P(X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC}; \mathbf{p})} \end{aligned}$$

For some purposes in following methods we could also write the denominator using the law of total probabilities:

$$q_{bgc} = \frac{P(X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC} | Y_{bi(g)c} = 1; \mathbf{p}) \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} P(X_{b1} = x_{b1}, \dots, X_{bC} = x_{bC} | Y_{bi(g)c} = 1; \mathbf{p}) \cdot p_{gc'}}$$

Once we calculate q_{bgc} , we can do the E-step, considering $Q(\mathbf{p}; \mathbf{p}^{(old)}) := \mathbb{E} \left[l(\mathbf{p}; \mathcal{X}, \mathcal{Y}) | \mathcal{X}, \mathbf{p}^{(old)} \right]$.

$$\begin{aligned} Q(\mathbf{p}; \mathbf{p}^{(old)}) &= \mathbb{E} \left[\sum_{b \in \mathcal{B}, i \in \mathcal{I}_b, c \in \mathcal{C}} Y_{bic} \ln(p_{g(i)c}) | \mathcal{X}, \mathbf{p}^{(old)} \right] \\ &= \sum_{b \in \mathcal{B}, g \in \mathcal{G}, c \in \mathcal{C}} q_{bgc} \ln(p_{gc}) \\ &= \sum_{b \in \mathcal{B}, g \in \mathcal{G}} w_{bg} \sum_{c \in \mathcal{C}} q_{bgc} \ln(p_{gc}) \end{aligned}$$

In this case $\mathbf{p}^{(old)}$ will be a prior probability we assume for the distribution.

2.2. M-Step

In the M-Step we maximize $Q(\mathbf{p}; \mathbf{p}^{(old)})$ on \mathbf{p} . We have to add a constraint so that the probabilities add up to 1 and are non-negative. We can see that $\mathbf{p}^{(old)}$ (the probability we consider as a prior) is expressed using \mathbf{q} .

$$\begin{aligned} \max_{\mathbf{p}} \quad & \sum_{b \in \mathcal{B}, g \in \mathcal{G}} w_{bg} \sum_{c \in \mathcal{C}} q_{bgc} \ln(p_{gc}) \\ \text{subject to:} \quad & \sum_{c \in \mathcal{C}} p_{gc} = 1 \quad \forall g \in \mathcal{G} \\ & p_{gc} \geq 0 \quad \forall g \in \mathcal{G}, \forall c \in \mathcal{C} \end{aligned}$$

To solve this optimization we can use the Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\lambda}) = \sum_{b \in \mathcal{B}, g \in \mathcal{G}} w_{bg} \sum_{c \in \mathcal{C}} q_{bgc} \ln(p_{gc}) - \sum_{g \in \mathcal{G}} \lambda_g \sum_{c \in \mathcal{C}} (p_{gc} - 1)$$

We state the first order conditions:

$$\frac{\partial \mathcal{L}}{\partial p_{gc}} = 0 \Leftrightarrow \frac{1}{p_{gc}} \sum_{b \in \mathcal{B}} w_{bg} q_{bgc} = \lambda_g \quad (2.2)$$

$$\Leftrightarrow p_{gc} = \frac{1}{\lambda_g} \sum_{b \in \mathcal{B}} w_{bg} q_{bgc}$$

$$\Leftrightarrow \sum_{c \in \mathcal{C}} \frac{1}{\lambda_g} \sum_{b \in \mathcal{B}} w_{bg} q_{bgc} = 1$$

$$\Leftrightarrow \lambda_g = \sum_{b \in \mathcal{B}, c \in \mathcal{C}} w_{bg} q_{bgc}$$

$$\Leftrightarrow \lambda_g = \sum_{b \in \mathcal{B}} w_{bg} \quad (2.3)$$

Replacing 2.3 in 2.2 we obtain:

$$p_{gc} = \frac{\sum_{b \in \mathcal{B}} w_{bg} q_{bgc}}{\sum_{b \in \mathcal{B}} w_{bg}} \quad \forall g \in \mathcal{G}, \forall c \in \mathcal{C} \quad (2.4)$$

2.3. EM-algorithm

The EM-algorithm consists of an iterative alternation between the E-step and the M-step. The formulation for this model goes as follows:

Algorithm 1 EM-algorithm

Input: $\mathcal{I}, \mathcal{G}, \mathcal{C}, \mathcal{B}, \mathbf{x}, \mathbf{w}, \epsilon$ **Output:** $\hat{\mathbf{p}}$ **Initialize:** $\mathbf{p}_{gc}^{(old)} = \frac{\sum_{b \in \mathcal{B}} w_{bg} \cdot x_{bc}}{I_g}$

converge = False

while not converge **do** **E-step:** Compute q_{bgc} using $\mathbf{p}_{gc}^{(old)}$ as described in 2.1. **M-step:** Compute $\hat{\mathbf{p}}$ as described in 2.4 **if** $\max_{g \in \mathcal{G}, c \in \mathcal{C}} \{\hat{\mathbf{p}}_{gc} - \mathbf{p}_{gc}^{(old)}\} \leq \epsilon$ **then**

converge = True

 $\mathbf{p}^{(old)} = \hat{\mathbf{p}}$

2.4. Log-likelihood

We can work the following expression for the log-likelihood.

$$l(\mathbf{p}; \mathcal{X}) = \sum_{b \in \mathcal{B}} \ln P(\mathbf{X}_b = \mathbf{x}_b | \mathbf{p}) \tag{2.5}$$

$$\begin{aligned} &= \sum_{b \in \mathcal{B}} \ln \left[\sum_{i \in \mathcal{I}_b, c \in \mathcal{C}} P(\mathbf{X}_b = \mathbf{x}_b, Y_{bg(i)c} = 1 | \mathbf{p}) \right] \\ &= \sum_{b \in \mathcal{B}} \ln \left[\sum_{i \in \mathcal{I}_b, c \in \mathcal{C}} P(Y_{bg(i)c} = 1 | \mathbf{X}_b = \mathbf{x}_b; \mathbf{p}^{(old)}) \cdot \frac{P(\mathbf{X}_b = \mathbf{x}_b, Y_{bg(i)c} = 1 | \mathbf{p})}{P(Y_{bg(i)c} = 1 | \mathbf{X}_b = \mathbf{x}_b; \mathbf{p}^{(old)})} \right] \\ &= \ln \mathbb{E} \left[\left(\frac{P(\mathcal{X}, \mathcal{Y} | \mathbf{p})}{P(\mathcal{Y} | \mathcal{X}, \mathbf{p}^{(old)})} \right) | \mathcal{X}, \mathbf{p}^{(old)} \right] \\ &\geq \mathbb{E} \left[\ln \left(\frac{P(\mathcal{X}, \mathcal{Y} | \mathbf{p})}{P(\mathcal{Y} | \mathcal{X}, \mathbf{p}^{(old)})} \right) | \mathcal{X}, \mathbf{p}^{(old)} \right] \tag{Jansen's Inequality} \\ &= \mathbb{E} \left[l(\mathbf{p}; \mathcal{X}, \mathcal{Y}) | \mathcal{X}, \mathbf{p}^{(old)} \right] - \mathbb{E} \left[\ln P(\mathcal{Y} | \mathcal{X}; \mathbf{p}^{(old)}) | \mathcal{X}, \mathbf{p}^{(old)} \right] \\ &= Q(\mathbf{p}; \mathbf{p}^{(old)}) - \mathbb{E} \left[\ln P(\mathcal{Y} | \mathcal{X}; \mathbf{p}^{(old)}) | \mathcal{X}, \mathbf{p}^{(old)} \right] \tag{2.6} \\ &= \sum_{b \in \mathcal{B}, g \in \mathcal{G}} w_{bg} \sum_{c \in \mathcal{C}} q_{bgc} \ln(p_{gc}) - \sum_{b \in \mathcal{B}, g \in \mathcal{G}} w_{bg} \sum_{c \in \mathcal{C}} q_{bgc} \ln(q_{bgc}) \\ &= \sum_{b \in \mathcal{B}, g \in \mathcal{G}} w_{bg} \sum_{c \in \mathcal{C}} q_{bgc} \ln \left(\frac{p_{gc}}{q_{bgc}} \right) \end{aligned}$$

We can see that 2.6 works as a lower bound for the log-likelihood. As in each iteration we are maximizing over \mathbf{p} , we only focus in $Q(\mathbf{p}; \mathbf{p}^{(old)})$, thus, increasing the lower bound for each iteration. This process is repeated until we get close enough to convergence ($\mathbf{p} = \mathbf{p}^{(old)}$), where 2.5 inequality becomes an equality.

Chapter 3

Computing the Conditional Probability

In order to perform the EM-algorithm we will need to compute the conditional probability q_{bgc} described in 2.1. In this chapter we will present an algorithm to compute this probability exactly and approximate methods that reduce its complexity.

3.1. Exact Method (EXACT)

A naive approach to calculate q_{bgc} would be to just consider each possible scenario that satisfies the observed data and compute using the law of total probabilities.

$$\begin{aligned} P(Y_{bi(g)c} = 1 | \mathbf{X}_b = \mathbf{x}_b) &= \sum_{z_b \in \Omega_b} P(Y_{bic} | \mathbf{Z}_b = z_b) P(\mathbf{Z}_b = z_b) \\ &= \sum_{z_b \in \Omega_b} \frac{z_{bgc}}{I_b} P(\mathbf{Z}_b = z_b) \end{aligned}$$

Where Ω_b is the set of all possible outcomes of \mathbf{Z}_b given the election result in that ballot-box. This can be represented by the following restrictions.

$$\begin{aligned}
\sum_{g \in \mathcal{G}} z_{bgc} &= x_{bc} & \forall c \in \mathcal{C} \\
\sum_{c \in \mathcal{C}} z_{bgc} &= w_{bg} & \forall g \in \mathcal{G} \\
z_{bgc} &\geq 0 & \forall c \in \mathcal{C} \forall g \in \mathcal{G} \\
z_{bgc} &\in \mathbb{Z} & \forall c \in \mathcal{C} \forall g \in \mathcal{G}
\end{aligned} \tag{3.1}$$

This method presents two complications: first, the occurrence of redundant calculations due to overlapping probabilities, and second, the lack of a straightforward means to acquire all feasible outcomes from Ω_b . To address these challenges, we introduce a recursive strategy. This methodology centers on tracking the aggregate election result up to a fixed group $f \in \mathcal{G}$.

For this purpose, we introduce the notation $T_{bf}(\mathbf{k})$, which denotes the probability of observing the aggregate result $\mathbf{k} = (k_1, \dots, k_C)$ within ballot-box b up to the inclusion of group f . By using this concept and applying the law of total probabilities, we arrive at the following derivation:

$$\begin{aligned}
T_{bf}(\mathbf{k}) &= \sum_{\mathbf{h} \in \mathcal{H}_{bf} : \mathbf{h} \leq \mathbf{k}} T_{b,f-1}(\mathbf{k} - \mathbf{h}) \cdot \binom{w_{bf}}{h_1, \dots, h_C} \prod_{c \in \mathcal{C}} p_{fc}^{h_c} \\
T_{b0}(\mathbf{k}) &= \mathbf{1}\{\mathbf{k} = \mathbf{0}\}
\end{aligned} \tag{3.2}$$

Where \mathcal{H}_{bf} is the set of all possible voting outcomes for group f . We will also define \mathcal{K}_{bf} as the set of all possible voting aggregated outcomes of what has happened until group f . This last set determines all the values of \mathbf{k} that we will need to compute with this recursive approach, starting from $f = 0$ until $f = G$. The formulation for these sets is as follows:

$$\begin{aligned}
\mathcal{H}_{bf} &= \{\mathbf{h} \in \mathbb{Z}_+^C : \sum_{c \in \mathcal{C}} h_c = w_{bf} \quad \forall c \in \mathcal{C}\} \\
\mathcal{K}_{bf} &= \{\mathbf{k} \in \mathbb{Z}_+^C : \sum_{c \in \mathcal{C}} k_c = \sum_{f' \leq f} w_{bf'}; k_c \leq x_{bc} \quad \forall c \in \mathcal{C}\}
\end{aligned} \tag{3.3}$$

We can see that calculating $P(\mathbf{X}_b = \mathbf{x}_b)$ is equivalent to calculating $T_{bG}(\mathbf{x}_b)$.

Furthermore, it will be necessary to compute the conditional probabilities for the case where an individual from group g has already voted for candidate c . Employing analogous

notation, we introduce $U_{bfgc}(\mathbf{k})$ to represent the probability of observing the partial result $\mathbf{k} = (k_1, \dots, k_C)$ within ballot-box b up to the inclusion of group f , conditional to $Y_{bi(g)c} = 1$. This can be expressed as follows:

$$U_{bfgc}(\mathbf{k}) = \sum_{\mathbf{h} \in \mathcal{H}_{bf} : \mathbf{h} \leq \mathbf{k}} U_{b,f-1,gc}(\mathbf{k} - \mathbf{h}) \cdot \binom{w_{bf}}{h_1, \dots, h_C} \cdot \left(\prod_{c' \in \mathcal{C}} p_{fc'}^{h_{c'}} \right) \cdot \left(\frac{h_c}{p_{fc} w_{bf}} \right)^{\mathbb{1}_{\{f=g\}}}$$

$$U_{b0gc}(\mathbf{k}) = \mathbf{1}\{\mathbf{k} = 0\}$$

To determine the probability of the complete conditional outcome, it suffices to compute $U_{bGgc}(\mathbf{x}_b)$ for each group and candidate. By incorporating both discussed components, we arrive at an expression for q_{bgc} .

$$q_{bgc} = \frac{U_{bGgc}(\mathbf{x}_b) \cdot p_{gc}}{T_{bG}(\mathbf{x}_b)}$$

We can compute the conditional probability for the last candidate (or an alternative candidate) using the complement probability: $U_{bGgC}(\mathbf{x}_b) = 1 - \sum_{c < C} U_{bGgc}(\mathbf{x}_b)$.

To calculate \mathbf{q}_b , we propose the following algorithm employing dynamic programming.

Algorithm 2 Exact \mathbf{q}_b computation with dynamic programming.

Initialize $\mathcal{H}_{bf}, \mathcal{K}_{bf}$ as in 3.3
 $T_{bf}(\mathbf{k}) = 0$; $U_{bfgc}(\mathbf{k}) = 0$; $\forall \mathbf{k} \in \mathcal{K}_{bf}$
for $f = 1$ to G **do**
 for $\mathbf{k} \in \mathcal{K}_{bf}$ **do**
 for $\mathbf{h} \in \mathcal{H}_{bf}$ **do**
 if $h_c \leq k_c \forall c \in \mathcal{C}$ **then**
 $a = \binom{w_{bf}}{h_1, \dots, h_C} \prod_{c \in \mathcal{C}} p_{fc}^{h_c}$
 $T_{bf}(\mathbf{k}) = T_{bf}(\mathbf{k}) + T_{bf}(\mathbf{k} - \mathbf{h}) \cdot a$
 for $c \in \mathcal{C}$ **do**
 for $g = 1$ to $f - 1$ **do**
 $U_{bfgc}(\mathbf{k}) = U_{bfgc}(\mathbf{k}) + U_{b,f-1,gc}(\mathbf{k} - \mathbf{h}) \cdot a$
 $U_{bggc}(\mathbf{k}) = U_{bggc}(\mathbf{k}) + U_{b,g-1,gc}(\mathbf{k} - \mathbf{h}) \cdot a \cdot \left(\frac{h_c}{p_{fc} w_{bf}} \right)$
 Output $q_{bgc} = \frac{U_{bGgc}(\mathbf{x}_b) \cdot p_{gc}}{T_{bG}(\mathbf{x}_b)}$

We can see that for each group $f \in \mathcal{G}$ we are doing $C \cdot |\mathcal{K}_{bf}| \cdot |\mathcal{H}_{bf}|$ operations, thinking of a worst case scenario this could be $C \cdot (I_B)^{2C}$ operations. We conclude that the complexity of the algorithm is $O(B \cdot G \cdot C \cdot (I_B)^{2C})$.

3.2. Hit-and-Run Method (H&R)

Given the complexity of the EXACT method to compute q_{bgc} we may use an approach using simulation so that we do not need to consider each possible outcome of the election.

We shall further examine Ω_b , the set of all possible outcomes for $\mathbf{Z}_b = \mathbf{z}_b$, as previously defined in 3.1. The depiction of an outcome within ballot-box b can be illustrated using matrix representation.

$$\begin{array}{rcc}
 & & \text{Candidate} \\
 & & 1 \quad \cdots \quad C \\
 \text{Group} \quad 1 & \boxed{\begin{array}{ccc} z_{b11} & \cdots & z_{b1C} \end{array}} & w_{b1} \\
 \quad \vdots & \boxed{\begin{array}{ccc} \vdots & \ddots & \vdots \end{array}} & \vdots \\
 \quad G & \boxed{\begin{array}{ccc} z_{bG1} & \cdots & z_{bGC} \end{array}} & w_{bG} \\
 & x_{b1} \quad \cdots \quad x_{bC} &
 \end{array}$$

Where the sum of row g is set to w_{bg} and the sum of column c is set to x_{bc} .

Our objective is to generate sample values for \mathbf{z}_b while adhering to the conditions outlined in the matrix. One approach to achieve this is by iteratively updating the values within the matrix using a stochastic process. This process should ensure that we end up in a point that is independent of the initial position.

We may note that shifting two components of \mathbf{z}_b that are neither in the same row nor column results in the alteration of two other components of \mathbf{z}_b . To illustrate, consider the scenario where we move z_{bgc} and $z_{bg'c'}$ downward by one unit; to maintain feasibility within the polytope, it necessitates moving $z_{bgc'}$ and $z_{bg'c}$ upward by one unit. Additionally, when opting to reduce a component of \mathbf{z}_b , we should verify that this modification does not infringe upon the non-negativity constraints.

We propose the following algorithm to sample values for \mathbf{z}_b .

Algorithm 3 Hit and Run

Input: $\mathbf{w}_b, \mathbf{x}_b, M, S$ Initialize $\mathbf{z}_b = \alpha(\mathbf{w}_b, \mathbf{x}_b)$, $\mathcal{S}_b = \phi$ **for** $s = 1$ to S **do** **for** $m = 1$ to M **do** Sample without replacement $\{g, g'\} \leftarrow \mathcal{G}$ and $\{c, c'\} \leftarrow \mathcal{C}$ **if** $z_{bgc} > 0$ and $z_{bg'c'} > 0$ **then** $z_{bgc} \leftarrow z_{bgc} - 1$ $z_{bg'c'} \leftarrow z_{bg'c'} - 1$ $z_{bgc'} \leftarrow z_{bgc'} + 1$ $z_{bg'c} \leftarrow z_{bg'c} + 1$ $\mathcal{S}_b = \mathcal{S}_b \cup \{\mathbf{z}_b\}$ **Output:** \mathcal{S}_b

Where S is the number of samples, M is the step size and α is any function that gives a starting point for \mathbf{z}_b (see Annex A.1). In Annex A.2 we discuss the election of M .

The complexity of the Hit and Run algorithm is $O(M \cdot S)$, as the size of the instance does not affect the number of operations (we always change 4 components when moving).

We could also consider methods of sampling that use polytopes defined with inequalities, where there are hit and run approaches that guarantee a uniform distribution (Metz & Zabinsky, 2012). The formulation for this approach is described in Annex A.3. It should be noted that this algorithm was not chosen as the number of iterations to get uncorrelated points was noticeably higher in this context.

As we sample we obtain a subset $\mathcal{S}_b \subseteq \Omega_b$. It is now possible to condition the event $Y_{bi(g)c} = 1$ to both \mathcal{S}_b and \mathbf{X}_b . We may note that $P(\mathbf{X}_b = \mathbf{x}_b | \mathcal{S}_b) = 1$, since we built \mathcal{S}_b so that all its element satisfy the given outcome of the ballot-box. Considering this method we would get the following approximation.

$$\begin{aligned} q_{bgc} &= P(Y_{bi(g)c} = 1 | \mathbf{X}_b = \mathbf{x}_b, \mathcal{S}_b; \mathbf{p}) \\ &= P(Y_{bi(g)c} = 1 | \mathcal{S}_b; \mathbf{p}) \\ &= \sum_{\mathbf{z}_b \in \mathcal{S}_b} P(Y_{bi(g)c} = 1 | \mathbf{Z}_b = \mathbf{z}_b) P(\mathbf{Z}_b = \mathbf{z}_b) \\ &= \sum_{\mathbf{z}_b \in \mathcal{S}_b} \frac{z_{bgc}}{w_{bg}} \cdot \prod_{g' \in \mathcal{G}} \binom{w_{bg'}}{z_{bg'1}, \dots, z_{bg'C}} \prod_{c' \in \mathcal{C}} p_{g'c'}^{z_{bg'c'}} \end{aligned}$$

The complexity of this calculation is: $O(S \cdot B \cdot G \cdot C)$. It should be noted that the term

$P(\mathbf{Z}_b = \mathbf{z}_b)$ can be precomputed for each ballot-box and does not need to be calculated in each combination.

3.3. Multivariate-Normal Approximation Method Using the CDF (MVN-CDF)

In this section we will consider a method based on the Normal-Approximation proposed by Zhengzhi Lin (Lin, Wang, & Hong, 2022). Lin's application is done for a more general case, where each component of a matrix has a distinct categorical probability, described as the Poisson Multinomial Distribution. We adapted this implementation to the case of having groups.

Since the last component of \mathbf{X}_b is determined by the rest, we will consider \mathbf{X}_b^* as the reduced version where we omit the last component: X_{bc} .

We will consider the following asymptotic distribution:

$$\mathbf{X}_b^* \sim \mathcal{N}(\boldsymbol{\mu}_b, \boldsymbol{\Sigma}_b)$$

Consider \mathbf{p}^* as the reduced version of \mathbf{p} . In order to obtain the mean of the distribution we just consider that the expected voting outcome for one person of group g is the vector \mathbf{p}_g^* (we are not considering the last candidate).

$$\begin{aligned} \boldsymbol{\mu}_b &= \mathbf{p}^{*T} \mathbf{w}_b \\ \boldsymbol{\Sigma}_b &= \text{diag}(\boldsymbol{\mu}_b) - \mathbf{p}^{*T} \text{diag}(\mathbf{w}_b) \mathbf{p}^* \end{aligned} \tag{3.4}$$

We can also get the asymptotic distribution to the conditional random variable. Using the same notation for reduced variables we get:

$$\begin{aligned} \mathbf{X}_b^* | Y_{bj(g)c} = 1 &\sim \mathcal{N}(\boldsymbol{\mu}_b^g, \boldsymbol{\Sigma}_b^g) + \mathbf{e}_c \\ \boldsymbol{\mu}_b^g &= \boldsymbol{\mu}_b - \mathbf{p}_g^{*T} \\ \boldsymbol{\Sigma}_b^g &= \boldsymbol{\Sigma}_b - \text{diag}(\mathbf{p}_g^{*T}) - \mathbf{p}_g^{*T} \mathbf{p}_g^* \end{aligned} \tag{3.5}$$

We can use a continuous approximation to compute the desired probability. Let us consider F_{bg} to be the cumulative density function of the multivariate normal approximation just

described (note that it does not depend on c , since it is only added as a constant). Consider the following hypercube:

$$\mathcal{A}_{bc} := [x_{b1} - 0.5, x_{b1} + 0.5] \times \dots \times [x_{bc} - 1.5, x_{bc} - 0.5] \times \dots \times [x_{bC-1} - 0.5, x_{bC-1} + 0.5]$$

Then we have the following approximation.

$$\begin{aligned} P(\mathbf{X}_b = \mathbf{x}_b | Y_{bi(g)c} = 1) &\approx P_{NA}(\mathbf{x}_b - 0.5 \leq \mathbf{X}_b \leq \mathbf{x}_b + 0.5 | Y_{bi(g)c} = 1) \\ &=: F_{bg}(\mathcal{A}_{bc}) \end{aligned}$$

Using this method we will need to compute the integral $F_{bg}(\mathcal{A}_{bc})$ for all combinations of b, g, c . This integral can be computed using Monte Carlo integration methods such as the one proposed by (Genz, 1992). Using this approximation we would end up with the following result for the conditional probabilities.

$$q_{bgc} = \frac{F_{bg}(\mathcal{A}_{bc}) \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} F_{bg}(\mathcal{A}_{bc'}) \cdot p_{gc'}}$$

The complexity of this algorithm is as described in Monte Carlo integration, where we consider an error parameter ϵ , obtaining a complexity of $O(B \cdot G \cdot C^2 \cdot \epsilon^{-\frac{1}{2}})$ for running each integral. We also get a complexity of $O(B \cdot G \cdot C^3)$ for doing the Cholesky decomposition for each group. Then the complexity of the method is $O(B \cdot G \cdot C^2 \cdot [\epsilon^{-\frac{1}{2}} + C])$. It should be noted that the number of iterations needed to compute the integral may depend on the size of the instance. We are considering a scale of C^2 for the integrals as for each combination of C the integration method considers a simulation of $C - 1$ random uniform variables.

3.4. Multivariate-Normal Approximation Method Using the PDF (MVN-PDF)

We may note that in order to calculate $P(\mathbf{X}_b = \mathbf{x}_b | Y_{bi(g)c} = 1)$, we do not necessarily need to calculate $P(Y_{bi(g)c} = 1 | \mathbf{X}_b = \mathbf{x}_b)$, but we may calculate something proportional to it. What we are really interested in is the ratio between $P(Y_{bi(g)c} = 1 | \mathbf{X}_b = \mathbf{x}_b)$ for the different values of c for a fixed g .

Using this idea we may choose to calculate something that is easier to compute than the actual probability. Based on the approximation in 3.4, we may use the probability density function instead of the cumulative density function. This can be done as we are only interested in the ratios between values to calculate the probability, and we do not need this quantity to represent a probability on itself. We could also interpret it as taking the limit where we make the hypercube as small as possible.

Using the same notation as the previous method, let us consider f_{bg} as the probability density function for the random variable described in 3.5. Using this notion we can get this simple approximation.

$$\begin{aligned}
q_{bgc} &= \frac{f_{bg}(x_{b1}, \dots, x_{bc} - 1, \dots, x_{bC}) \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} f_{bg}(x_{b1}, \dots, x_{bc'} - 1, \dots, x_{bC}) \cdot p_{gc'}} \\
&= \frac{\sqrt{(2\pi)^{C-1} |\Sigma_b^g|} \cdot e^{-\frac{1}{2}(\mathbf{x}_b - \mathbf{e}_c - \boldsymbol{\mu}_b^g)^\top (\Sigma_b^g)^{-1} (\mathbf{x}_b - \mathbf{e}_c - \boldsymbol{\mu}_b^g)} \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} \sqrt{(2\pi)^{C-1} |\Sigma_b^g|} \cdot e^{-\frac{1}{2}(\mathbf{x}_b - \mathbf{e}_{c'} - \boldsymbol{\mu}_b^g)^\top (\Sigma_b^g)^{-1} (\mathbf{x}_b - \mathbf{e}_{c'} - \boldsymbol{\mu}_b^g)} \cdot p_{gc'}} \\
&= \frac{e^{-\frac{1}{2}(\mathbf{x}_b - \mathbf{e}_c - \boldsymbol{\mu}_b^g)^\top (\Sigma_b^g)^{-1} (\mathbf{x}_b - \mathbf{e}_c - \boldsymbol{\mu}_b^g)} \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} e^{-\frac{1}{2}(\mathbf{x}_b - \mathbf{e}_{c'} - \boldsymbol{\mu}_b^g)^\top (\Sigma_b^g)^{-1} (\mathbf{x}_b - \mathbf{e}_{c'} - \boldsymbol{\mu}_b^g)} \cdot p_{gc'}}
\end{aligned}$$

Computing the inverse covariance matrix for each combination of g and c has a complexity of $O(B \cdot G \cdot C^3)$, as inverting a matrix has a cubic order in its dimension. Then, computing q_{bgc} considering the inverse covariance matrix as a precomputed parameter has a complexity of $O(B \cdot G \cdot C)$. We see that the complexity of the entire computation is $O(B \cdot G \cdot C^3)$.

3.5. Multinomial Approximation Method (MULT)

Considering that \mathbf{X}_b is a sum of multinomial random variables, we may consider an approximation so that \mathbf{X}_b distributes multinomial. To get this approximation we observe that each group g has a probability of voting for candidate c , so we can weight the contribution of each group with the number of voters in that group w_{bg} . We will consider r_{bc} to be the weighted probability for candidate c in ballot-box b .

$$\begin{aligned}
r_{bc} &= \frac{\sum_{g \in \mathcal{G}} w_{bg} p_{bg}}{\sum_{g \in \mathcal{G}} w_{bg}} = \frac{\sum_{g \in \mathcal{G}} w_{bg} p_{bg}}{I_b} \\
\mathbf{X}_b &\sim \text{Multinomial}(I_b, \mathbf{r}_b)
\end{aligned} \tag{3.6}$$

We can also get approximations for the conditional probabilities. If we consider that someone from group g already voted for candidate c , then its weight in the probability is decreased by one unit. Considering r_{bgc} to be the weighted probability for candidate c in ballot-box b , when someone from group g already voted for c :

$$r_{bgc} = \frac{\sum_{g' \in \mathcal{G}} w_{bg'} p_{g'c} - p_{gc}}{\sum_{g' \in \mathcal{G}} w_{bg'} - 1} = \frac{\sum_{g' \in \mathcal{G}} w_{bg'} p_{g'c} - p_{gc}}{I_b - 1}$$

$$\mathbf{X}_b | Y_{bi(g)c=1} \sim \text{Multinomial}(I_b - 1, \mathbf{r}_b^{(g)}) + \mathbf{e}_c$$

Using these approximations we can work out an expression for q_{bgc} .

$$\begin{aligned} q_{bgc} &= \frac{P(\mathbf{X}_b = \mathbf{x}_b | Y_{bj(g)c} = 1) \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} P(\mathbf{X}_b = \mathbf{x}_b | Y_{bj(g)c'} = 1) \cdot p_{gc'}} \\ &= \frac{\binom{I_b - 1}{x_{b1}, \dots, x_{bc} - 1, \dots, x_{bC}} \cdot \prod_{d \in \mathcal{C}} r_{bgd}^{x_{bd}} \cdot r_{bgc}^{-1} \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} \binom{I_b - 1}{x_{b1}, \dots, x_{bc'} - 1, \dots, x_{bC}} \cdot \prod_{d \in \mathcal{C}} r_{bgd}^{x_{bd}} \cdot r_{bgc'}^{-1} \cdot p_{gc'}} \\ &= \frac{\frac{(I_b - 1)!}{x_{b1}! \dots (x_{bc} - 1)! \dots x_{bC}!} \cdot \prod_{d \in \mathcal{C}} r_{bgd}^{x_{bd}} \cdot r_{bgc}^{-1} \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} \frac{(I_b - 1)!}{x_{b1}! \dots (x_{bc'} - 1)! \dots x_{bC}!} \cdot \prod_{d \in \mathcal{C}} r_{bgd}^{x_{bd}} \cdot r_{bgc'}^{-1} \cdot p_{gc'}} \\ &= \frac{x_{bc} \cdot \prod_{d \in \mathcal{C}} r_{bgd}^{x_{bd}} \cdot r_{bgc}^{-1} \cdot p_{gc}}{\sum_{c' \in \mathcal{C}} x_{bc'} \cdot \prod_{d \in \mathcal{C}} r_{bgd}^{x_{bd}} \cdot r_{bgc'}^{-1} \cdot p_{gc'}} \\ &= \frac{\frac{x_{bc} \cdot p_{gc}}{r_{bgc}}}{\sum_{c' \in \mathcal{C}} \frac{x_{bc'} \cdot p_{gc'}}{r_{bgc'}}} \\ &= \frac{(x_{bc} \cdot p_{gc}) / r_{bgc}}{\sum_{c' \in \mathcal{C}} (x_{bc'} \cdot p_{gc'}) / r_{bgc'}} \\ &= \frac{x_{bc} \cdot p_{gc} \cdot r_{bgc}^{-1}}{\sum_{c' \in \mathcal{C}} x_{bc'} \cdot p_{gc'} \cdot r_{bgc'}^{-1}} \end{aligned}$$

This complexity of this algorithm is straightforward as it relies on basic operations, scaling as: $O(B \cdot G \cdot C)$.

Chapter 4

Computing the p-value

4.1. Exact Computation

Once we have our estimate for \mathbf{p} it is in our interest to determine how likely the result of each ballot-box is. Given the result \mathbf{x}_b , we can determine $P(\mathbf{X}_b = \mathbf{x}_b)$ using algorithm 2. Since we want to know how likely is this result we would also need to compute the probability for all other possible outcomes $\mathbf{x}_b \in R_{\mathbf{X}_b}$, which can be done using the same set of recursive equations.

We will consider the p-value of ballot-box b as the probability of obtaining something as or less probable than \mathbf{x}_b . It can be computed as follows:

$$\begin{aligned} \text{p-val}(\mathbf{x}_b) &= \sum_{\mathbf{x}'_b \in R_{\mathbf{X}_b} | P(\mathbf{x}'_b) \leq P(\mathbf{x}_b)} P(\mathbf{x}'_b) \\ &= \sum_{\mathbf{x}'_b \in R_{\mathbf{X}_b}} P(\mathbf{x}'_b) \cdot \mathbb{1}_{\{P(\mathbf{x}'_b) \leq P(\mathbf{x}_b)\}} \\ &= \sum_{\mathbf{x}'_b \in R_{\mathbf{X}_b}} T_{bG}(\mathbf{x}'_b) \cdot \mathbb{1}_{\{T_{bG}(\mathbf{x}'_b) \leq T_{bG}(\mathbf{x}_b)\}} \end{aligned} \tag{4.1}$$

Where $T_{bG}(\mathbf{x}'_b)$ is as described in 3.2.

4.2. Simulation Using the Multinomial Approximation

We could approximate the p-value using simulation and the approximation proposed in 3.6. We may simulate H identical, independent and identically distributed random variables following a Multinomial(I_b, \mathbf{r}_b), obtaining: $\mathbf{x}_b^{(1)}, \dots, \mathbf{x}_b^{(H)}$. Using these simulations we may calculate the p-value by counting how many of this outcomes were less probable than \mathbf{x}_b :

$$\text{p-val}(\mathbf{x}_b) = \frac{1}{H} \sum_{h=1}^H \mathbb{1}_{\{P(\mathbf{x}_b^{(h)}) \leq P(\mathbf{x}_b)\}} \quad (4.2)$$

Where we we can compute the comparison as follows.

$$P(\mathbf{x}_b^{(h)}) \leq P(\mathbf{x}_b)$$

$$\binom{I_b}{x_{b1}^{(h)}, \dots, x_{bC}^{(h)}} \prod_{c \in \mathcal{C}} (r_{bc})^{x_{bc}^{(h)}} \leq \binom{I_b}{x_{b1}, \dots, x_{bC}} \prod_{c \in \mathcal{C}} (r_{bc})^{x_{bc}}$$

Chapter 5

Numerical Results

In this chapter we will compare the 5 methods studied for the EM-algorithm in Chapter 3: EXACT, H&R, MVN-CDF, MVN-PDF and MULT. We will also compare both methods for computing the p-values as described in Chapter 4.

We will describe the methodology for generating instances that represent the model described in Chapter 2. We will use instances with varying value of $G \in \{2, 3, 4\}$ and $C \in \{2, 3, 5, 10\}$, keeping fixed values of $B = 50$ and $I_b = 100$. We assume that all ballot-boxes will have the same amount of voters. For each instance we will run 20 scenarios.

5.1. Numerical Instances

In Algorithm 4 we describe the methodology to generate a scenario from an instance with parameters (I, G, C, B) . We also include a mixing parameter λ that will be discussed further. A scenario is completely described by the observed information \mathbf{x} and \mathbf{w} . We also output the real probability \mathbf{p} , used only to test estimations.

Algorithm 4 Election Outcome Generation

Input: I, G, C, B, λ

Output: $\mathbf{w} \in \mathbb{R}^{B \times G}, \mathbf{x} \in \mathbb{R}^{B \times C}, \mathbf{p} \in \mathbb{R}^{G \times C}$

Initialize: $\omega^0 : \omega_i^0 = \lceil \frac{i \cdot G}{I} \rceil, \omega = \omega^0, N = \frac{I}{B}, \mathcal{I} = \{1, \dots, I\}$

Sample a vector $\mathbf{v} \in \mathbb{N}^{\lceil \lambda \cdot I \rceil}$ where each component is drawn without replacement from \mathcal{I} .

Generate $\mathbf{v}' = \text{sort}(\mathbf{v})$

for $i = 1$ to $\lfloor \lambda \cdot I \rfloor$ **do**

$$\omega_{v_i} = \omega_{v'_i}^0$$

$$w_{bg} = \sum_{i=(b-1) \cdot N + 1}^{b \cdot N} \mathbb{1}_{\{\omega_i = g\}}$$

Simulate \mathbf{p}_g from $\text{Dirichlet}(\mathbf{1}_C)$

Simulate z_{bgc} from $\text{Multinomial}(w_{bg}, \mathbf{p}_g)$

$$x_{bc} = \sum_{g \in \mathcal{G}} z_{bgc}$$

Note that we choose to simulate \mathbf{p}_g from a Dirichlet distribution with parameter 1 on all its components. This is done so that we obtain different real probabilities for each scenario and better represent the set of different elections that may occur.

We include the λ factor as a mixing parameter that controls how groups are sorted through the ballot-boxes. For the EM-algorithm to get to a good estimator of the desired probabilities we would like the groups to have a heterogeneous distribution through the ballot-boxes. The intuition behind lies in how changes in the number of voters from a group in different ballot-boxes are affecting the outcome of votes.

We will test how the distribution of groups affects the mean absolute error (MAE) of the estimated probabilities. We calculate the MAE metric considering the absolute error of all components of the probability.

$$\frac{\sum_{g \in \mathcal{G}; c \in \mathcal{C}} |p_{gc} - \hat{p}_{gc}|}{G \cdot C} \quad (5.1)$$

Figure 5.1 shows how the value of λ affects the estimation of the probability \mathbf{p} using the EM-algorithm in different instances.

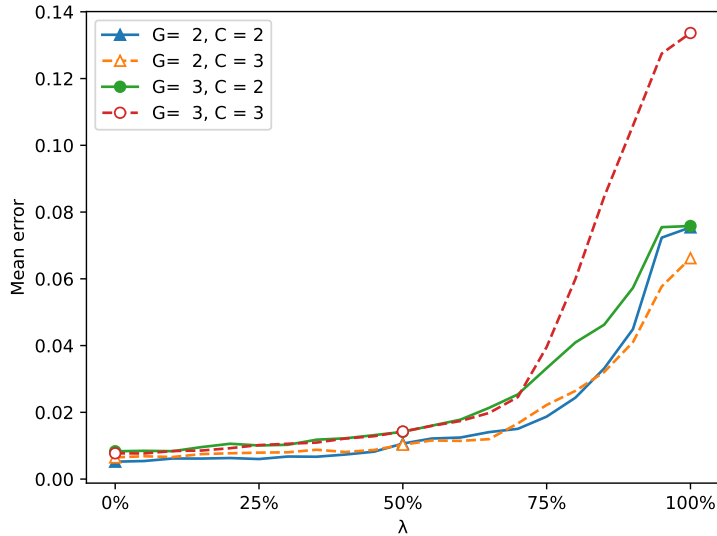


Figure 5.1: Mean absolute error over 20 scenarios for varying values of mix λ in different instances. $I_b = 100$, $B = 50$

We can observe that as the mixing parameter λ increases we get worse estimated probabilities. When we have small values of λ , ballot-boxes are more heterogeneous, thus, it is more apparent how the different composition in groups is affecting the outcome in votes.

According to these results we will consider that a mixing parameter of $\lambda = 50\%$ is sufficiently good for obtaining an estimator. We do not choose a lower value of λ since it may not be representative of what is really happening in the group assignment process of real Chilean elections.

5.2. EM-Algorithm Probability Estimation

In Annex B we present the results for the same instances considering $I_b = 200$.

5.2.1. Time Results

In Table 5.1 we present the average times that it takes the EM-algorithm to run. We set a time limit of 1 hour to each run, so we do not show results for those that did not converge. It should be noted that the number of iteration it takes the algorithm to converge varies across different methods, so the time presented takes into account how fast it converges and how fast the conditional probabilities are computed for each iteration.

Table 5.1: Mean running time over 20 scenarios in seconds for the EM-algorithm for varying instances, with fixed values of $I_b = 100$ and $B = 50$.

Instance		Method					
C	G	EXACT	H&R $S = 10^3$	H&R $S = 10^2$	MVN CDF	MVN PDF	MULT
2	2	14.510	719.939	72.335	0.851	0.140	0.003
	3	34.235	720.107	73.224	1.990	0.208	0.004
	4	51.563	723.778	72.621	3.421	0.233	0.003
3	2	549.902	709.502	71.909	3.604	0.217	0.003
	3	2,115.500	708.062	68.265	11.373	0.265	0.004
	4	2,808.304	676.996	67.471	15.927	0.320	0.005
5	2	-	683.890	67.924	12.315	0.173	0.003
	3	-	637.631	60.941	41.609	0.279	0.005
	4	-	603.475	59.644	47.176	0.281	0.004
10	2	-	594.104	60.332	61.702	0.181	0.003
	3	-	553.448	54.908	210.639	0.221	0.004
	4	-	514.156	50.178	355.949	0.308	0.004

The MULT method is the fastest and it is not highly impacted by C and G , as it is polynomial. We have a similar behavior for the MVN-PDF method. The EXACT method shows a clear exponential increase through C that makes it impossible to compute in the time limit for $C = 5$ and $C = 10$. Both H&R methods are not highly impacted by C and G , as the main time limitation is in the simulation, which will be discussed later on. The MVN-CDF method is also noticeably impacted by C but it was computable in the time limit.

We also show separate time results between simulation and EM-algorithm for the H&R methods in Table 5.2. It is clear that the bottle-neck of this algorithm comes in the simulation, as we will later show how the amount of samples affects the accuracy of the method. We can also observe that bigger instances tend to have lower simulation times due to the fact that the polytope is more restricted and it is more probable to hit the non-negativity constraints (thus, making less operations).

Table 5.2: Comparison of mean simulation time and mean EM running time in seconds over 20 scenarios for varying instances, with fixed values of $I_b = 100$ and $B = 50$.

Instance		Method			
C	G	H&R ($S = 10^2$)		H&R ($S = 10^3$)	
		Sim-time	EM-time	Sim-time	EM-time
2	2	71.869	0.466	719.469	0.471
	3	71.324	1.900	713.681	6.427
	4	70.213	2.408	702.951	20.827
3	2	70.307	1.603	704.637	4.865
	3	66.246	2.019	683.981	24.081
	4	64.956	2.515	651.587	25.409
5	2	66.042	1.882	666.138	17.752
	3	59.140	1.801	614.540	23.091
	4	57.597	2.047	581.373	22.102
10	2	58.708	1.624	578.243	15.861
	3	53.296	1.612	534.916	18.532
	4	48.476	1.701	492.142	22.014

5.2.2. Error Results

In Table 5.3 we show the MAE metric as described in 5.1. These results show that both MVN methods and the MULT method have similar MAE when compared to the EXACT method. Regarding the H&R method the magnitude of the error depends on the amount of points simulated. It is clear that bigger instances should require a bigger S , making the H&R method more demanding.

Table 5.3: Mean absolute prediction error over 20 scenarios for varying instances, with fixed values of $I_b = 100$ and $B = 50$.

Instance		Method					
C	G	EXACT	H&R $S = 10^3$	H&R $S = 10^2$	MVN CDF	MVN PDF	MULT
2	2	0.011	0.011	0.011	0.011	0.011	0.011
	3	0.014	0.014	0.016	0.014	0.014	0.015
	4	0.017	0.018	0.025	0.017	0.017	0.018
3	2	0.010	0.010	0.011	0.010	0.010	0.010
	3	0.014	0.015	0.025	0.014	0.014	0.015
	4	0.017	0.020	0.034	0.017	0.017	0.017
5	2	-	0.009	0.012	0.008	0.009	0.009
	3	-	0.016	0.025	0.013	0.013	0.013
	4	-	0.022	0.032	0.014	0.014	0.014
10	2	-	0.007	0.012	0.006	0.007	0.006
	3	-	0.015	0.020	0.009	0.009	0.009
	4	-	0.020	0.024	0.012	0.012	0.012

As we obtained similar MAE metrics for the MVN methods and the MULT method, we would recommend using the latter as it was shown to be the fastest.

5.3. p-value Estimation

In this section we will use a subset of the computed instances, considering $C = \{2, 3\}$ and $G = \{2, 3, 4\}$. We choose these ones as we need the EXACT method in order to do comparisons, where it is only reliable to run it until $C = 3$.

5.3.1. Probability Gap

It is in our interest to observe how the estimated probability $\hat{\mathbf{p}}$ changes the p-values obtained for each ballot-box compared to the original probability \mathbf{p} . In this case we will be comparing using the estimated $\hat{\mathbf{p}}$ that comes from the exact EM-method, we will call it $\hat{\mathbf{p}}_{(\text{EXACT})}$.

In Figure 5.2 we are showing the mean p-value absolute error and the mean absolute error between \mathbf{p} and $\hat{\mathbf{p}}_{(\text{EXACT})}$.

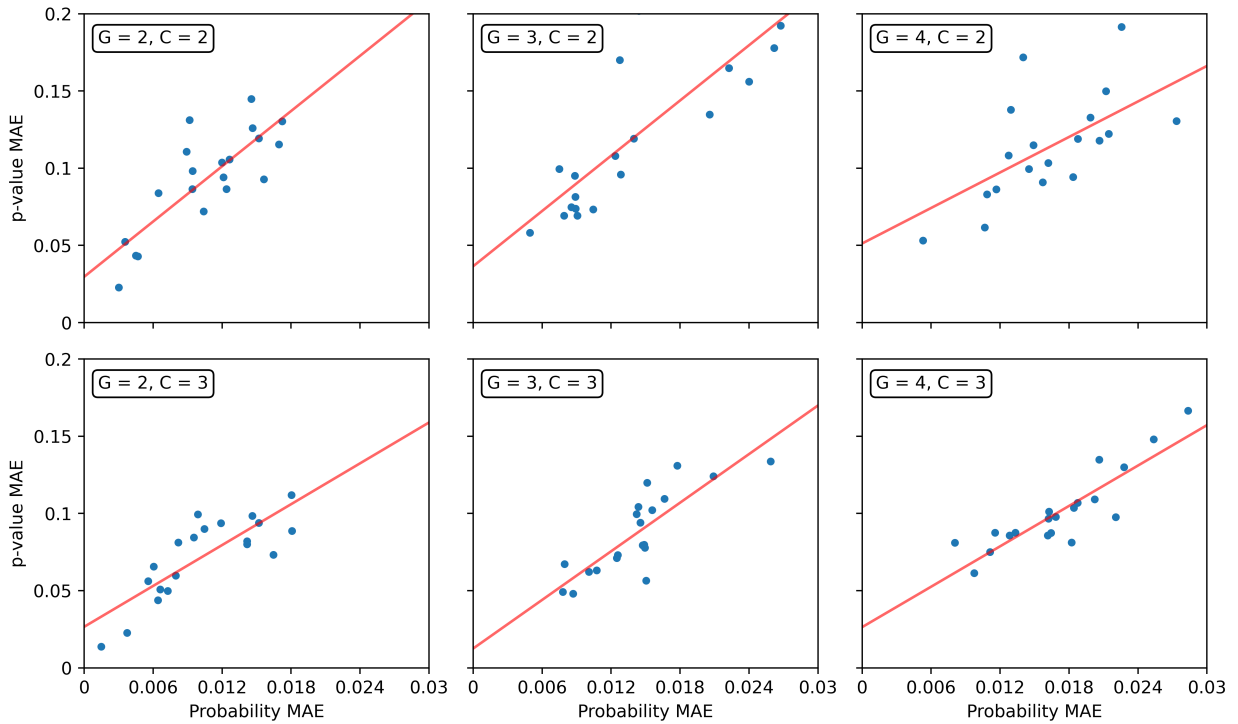


Figure 5.2: Correlation between p-values mean absolute error and estimated probability mean absolute error.

There is a clear correlation between both errors, which is the expected behavior.

We show a different approach in Figure 5.3 where we compare the p-values obtained using both probabilities \mathbf{p} and $\hat{\mathbf{p}}_{(\text{EXACT})}$. The color grading is now showing the magnitude of the MAE between both probabilities. We observe that the points with a smaller mean absolute error tend to be close to the middle, where both p-values are the same. We also observe that the p-values estimated with $\hat{\mathbf{p}}_{(\text{EXACT})}$ are overestimated and underestimated in similar amounts.

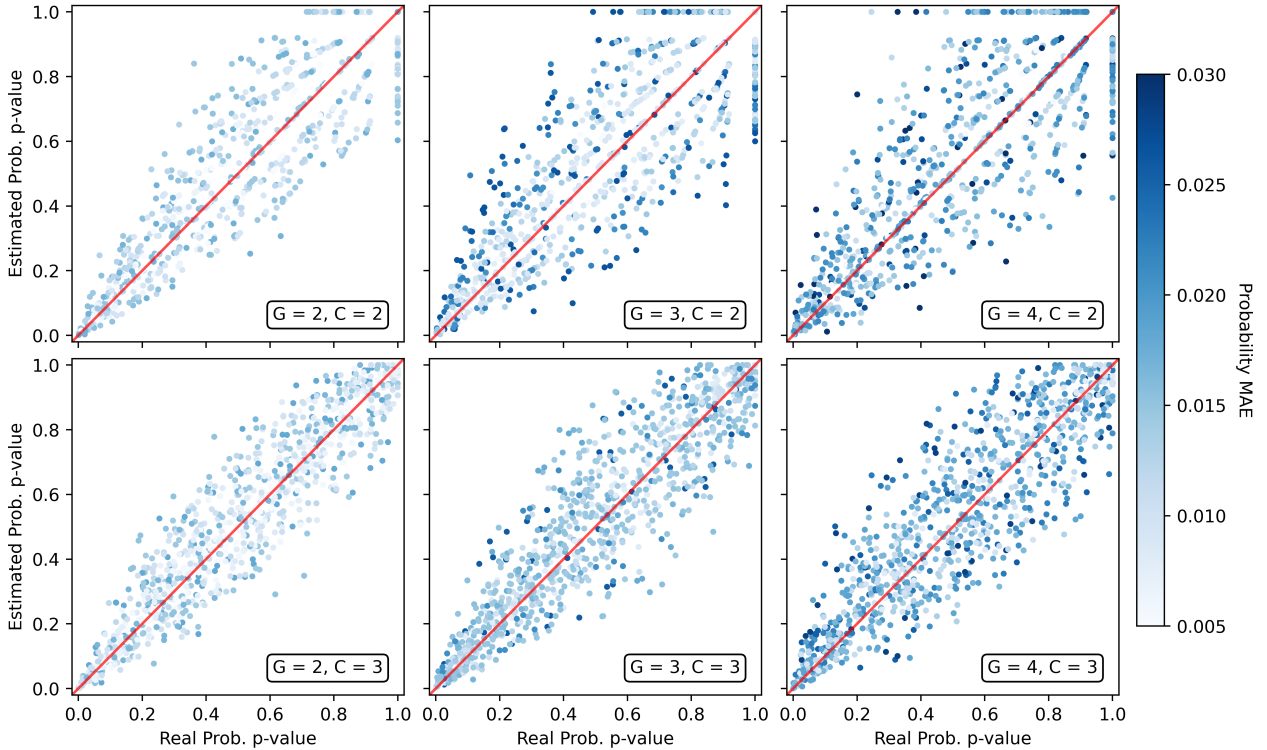


Figure 5.3: Comparison between p-values computed with real probability and p-values computed estimated probability.

5.3.2. Method Gap

We are also interested in studying how the simulation method using the multinomial approximation changes the computed p-values. For this comparison we will use the same probability \mathbf{p} and we will only change the method used to compute the p-value in each ballot-box.

It is intuitive to think that as we do the multinomial approximation of the distribution we are altering the variance of the original distribution. To test this hypothesis we will calculate the covariance matrix for both distributions.

The covariance matrix of the original distribution was already explained in Section 3.3 and displayed in 3.4. The covariance matrix of the multinomial approximation can be easily obtained using the parameters of the approximation as follows.

$$\Sigma_b^{(\text{Mult})} = I_b \cdot (\text{diag}(\mathbf{r}_b) - \mathbf{r}_b \mathbf{r}_b^\top)$$

To study how the variance changes in a ballot-box when using the multinomial approximation

we consider the ratio:

$$\frac{|\Sigma_b^{(Mult)}|}{|\Sigma_b|}$$

In Figure 5.4 we show how this ratio affects the absolute error for different p-values. We also show in color grading the magnitude of the original p-value.

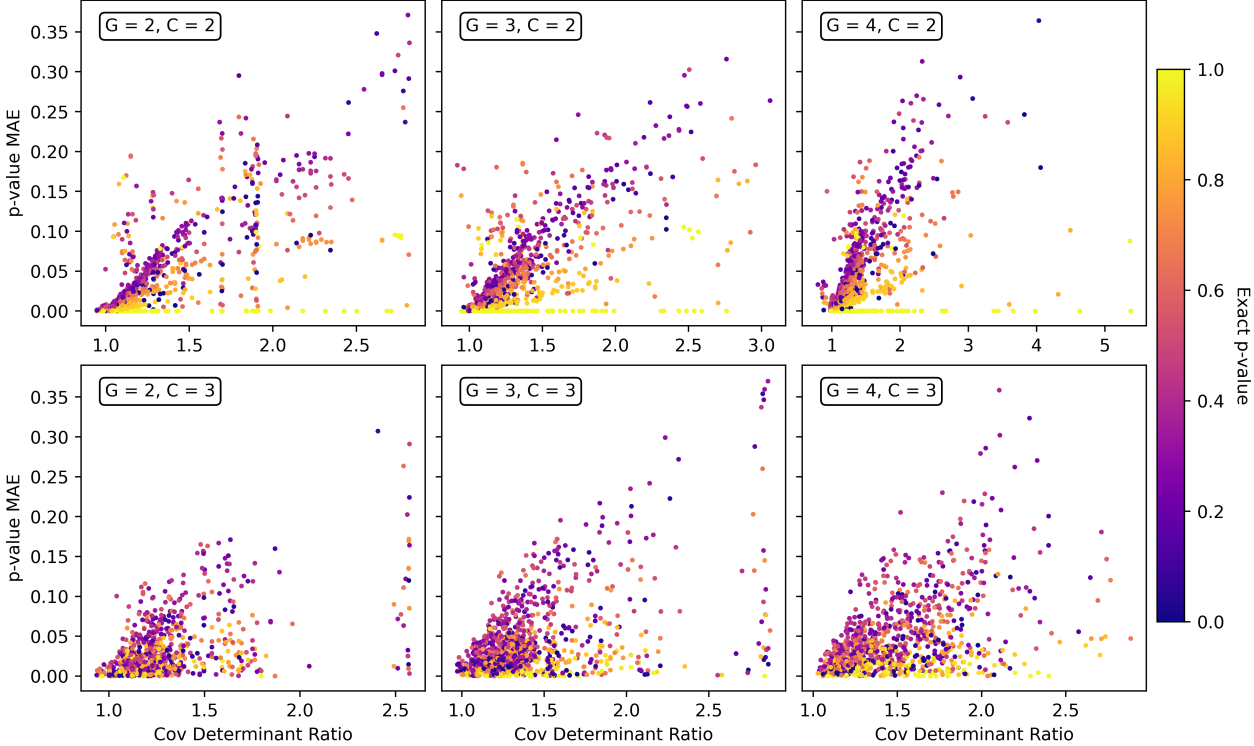


Figure 5.4: Effect of covariance determinant increase in the mean absolute error in the p-values.

We can observe that as the ratio between determinants increases we obtain higher mean absolute errors. We can also see that the error is more notable for p-values that are around the middle. This means that the main change in probability occurs to outcomes that are in the middle between the most probable and least probable outcome. The increase in variance we get through the multinomial approximation should increase the probability of this outcomes that are in the middle of the distribution.

In Figure 5.5 we show a different approach where we compare the p-value obtained with exact computation and the p-value obtained with simulation through the multinomial approximation.

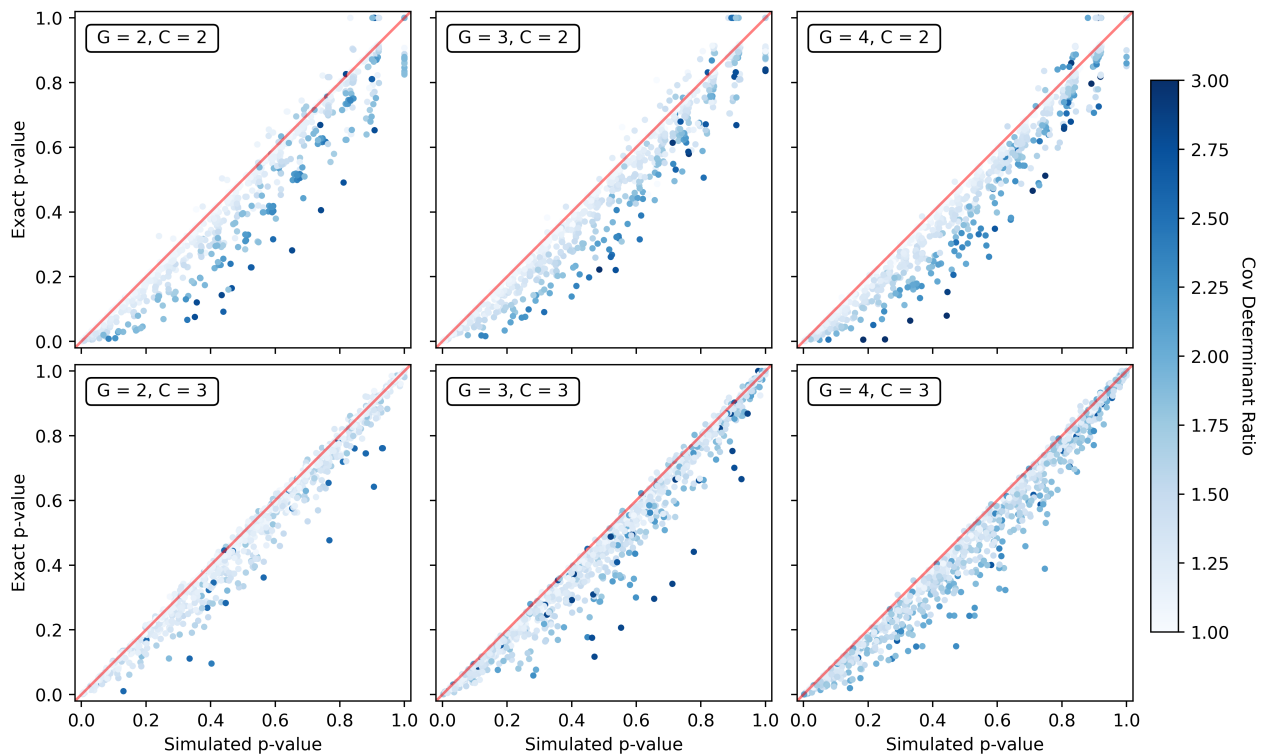


Figure 5.5: Comparison between p-values computed with the exact method and p-values simulated with the multinomial approximation.

It is clear that there is a bias that tends to increase the magnitude of the p-values computed with the approximation. We show in color grading the magnitude of the covariance determinant ratio. We can observe that those points with the biggest change in variance are the ones where the bias is more noticeable.

Chapter 6

Case Study: a Real Election

We will apply the EM-algorithm to the first round of the 2021 Chilean Presidential Election (2021-CPE) using the MULT method. This election adheres to the two-round system (TRS), which necessitates a subsequent round of voting featuring the two leading candidates if no single candidate secures an outright majority exceeding 50% during the first round.

6.1. Historical Data

The 2021-CPE election featured 7 candidates: Gabriel Boric, José Atonio Kast, Yasna Provoste, Sebastián Sichel, Eduardo Artés, Marco Enríquez-Ominami and Franco Parisi. Besides these 7 options, votes could be classified as blank votes (the ballot has no marking in it) or null votes (the ballot was filled out incorrectly). We will treat blank and null votes as a single category since different ballot-box committees may have applied varying criteria when labeling them.

Circunscripciones electorales are the territorial divisions used to assign voters to ballot-boxes. We will refer to them as districts. Across each district, there are distinct polling places, each of them with a fixed number of ballot-boxes. We will assume that two voters belonging to the same group within a district are considered comparable in terms of probability within that specific district. It is noteworthy that, in this election, the allocation of voters to ballot boxes did not consider geographical distance, a factor that was considered in subsequent elections starting from 2022.

In Figure 6.1 we show how the number of ballot-boxes varies across districts. We can observe that most districts have less than 50 ballot-boxes, which could make the estimation

of probabilities harder, as we have less information. We would expect that those districts with over a 100 ballot-boxes give us enough information for computing the estimating probabilities.

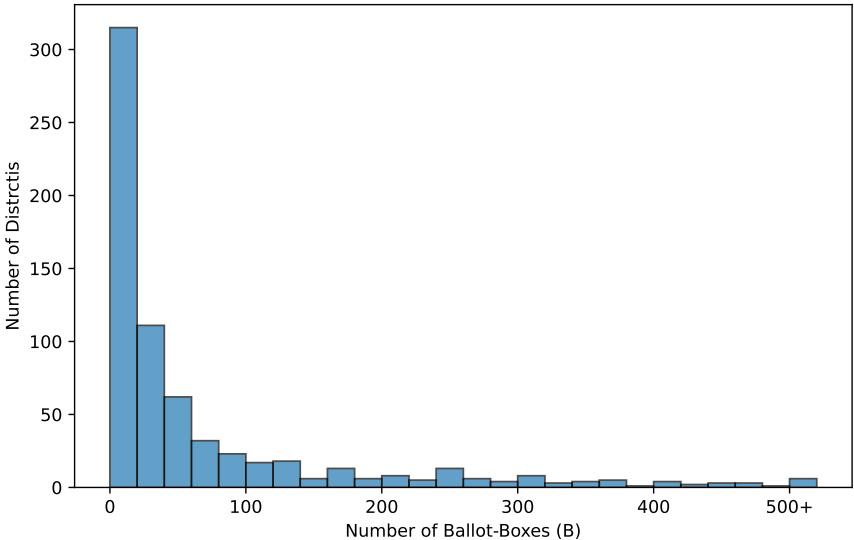


Figure 6.1: District count for ranges of length 25 for number of ballot-boxes.

The official website of the Electoral Service of Chile (SERVEL) offers detailed data of the election results (SERVEL, 2023). The detailed dataset for the votes emitted in Chile registers 7,080,276 votes emitted across 681 districts with a total of 46,639 ballot-boxes. We should note that this number may vary slightly after processing the data due to incomplete information in some ballot-boxes.

Additionally, SERVEL offers information about the aggregated demographic data of the voters in each ballot-box. There are three categorical features: sex, age range and nationality. There are 8 possible options for age range: 18-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80+. For this application, we will only consider age ranges as possible groups, as there is evidence to suggest a correlation between age and political ideology (Leigh, 2005; Peterson, Smith, & Hibbing, 2020). We will not consider nationality or gender as it would further increase the number of possible groups.

In Figure 6.2 we show an example that illustrates how the group age distribution in different polling places at the district of *Puente Alto* would look. This shows at first glance that there is heterogeneity in both: districts and polling places.

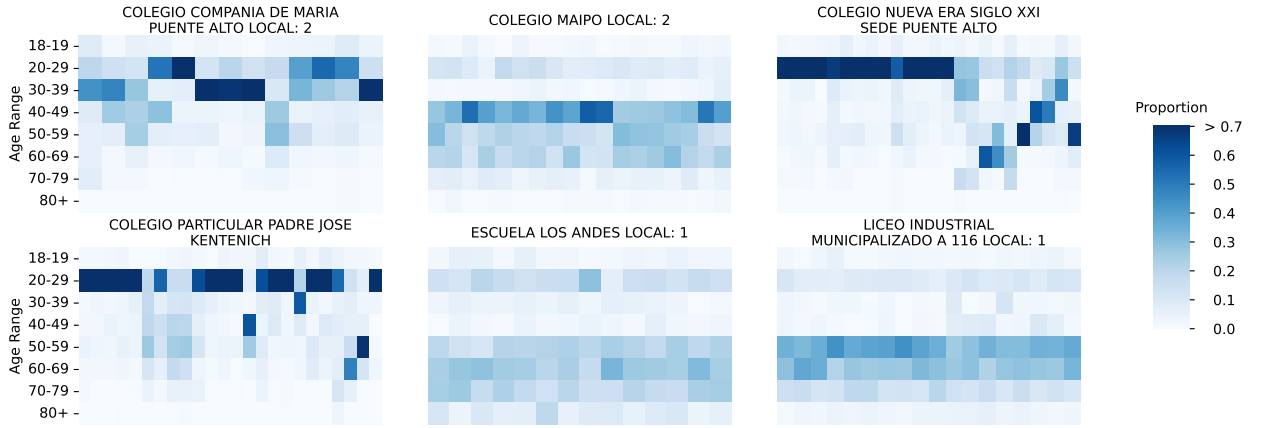


Figure 6.2: Proportions for different age ranges in polling places from the district of *Puente Alto*.

6.2. Group Aggregation

In order to run the EM-algorithm in the 2021-CPE we need to choose which groups to use. This is important because as we saw in Figure 5.1 the distribution of groups through ballot-boxes has an effect on the probability estimation.

The first group set that we could consider would be to use the full data available from age ranges, giving us the following:

$$\mathcal{G} = \{18-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80^+\}$$

However, we could aggregate some age ranges into one group aggregation a . This could be useful if the probabilities of that age range are difficult to identify given its distribution through the ballot-boxes. We may consider a group aggregation a as a set of groups from \mathcal{G} where we aggregate all voters from that group. Then, a group set \mathcal{A} is a set formed of group aggregations from \mathcal{G} .

For example, we could aggregate age ranges into 2 categories: voters up to 39 years old, and voter from 40 years old. This could be expressed as a set $\mathcal{A} = \{18-39, 40^+\}$. In Figure 6.3 we show how the distribution varies using this aggregation in the district of *Río Negro*.

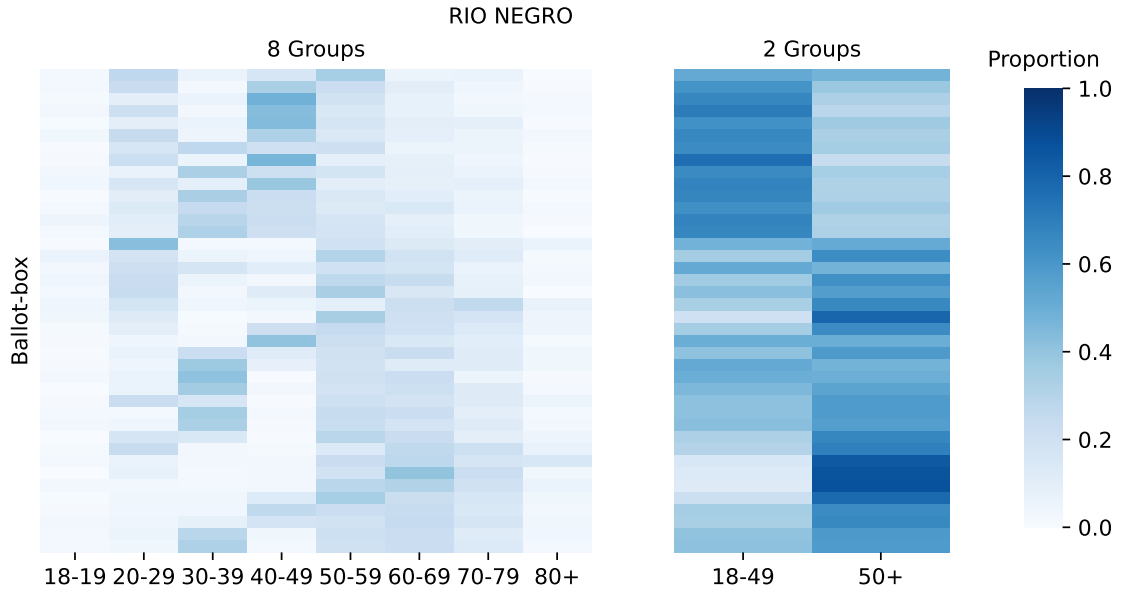


Figure 6.3: Proportions for different group aggregations in the district of *Río Negro*.

We can notice that when adding groups we may obtain distributions that are completely different. For example, we can observe that the extreme age ranges of 18-19 and 80⁺ have few voters per ballot-box and show little variance across the district. We can note that when we aggregate into 2 groups we obtain very identifiable categories, which we could expect to result in better estimations. It is not clear which aggregation is best, as we would like the estimated probabilities to better represent the reality of the district but we would not want them to be far off from the real latent probabilities.

To show an example of how the group aggregation affects the model, we will consider the original set \mathcal{G} and another set $\mathcal{A} = \{18-29, 30-59, 60^+\}$. We could run the EM-algorithm using both sets \mathcal{G} and \mathcal{A} to compare the estimated probabilities. In Figure 6.4 we show results for two different districts that noticeably vary in the number of ballot-boxes.

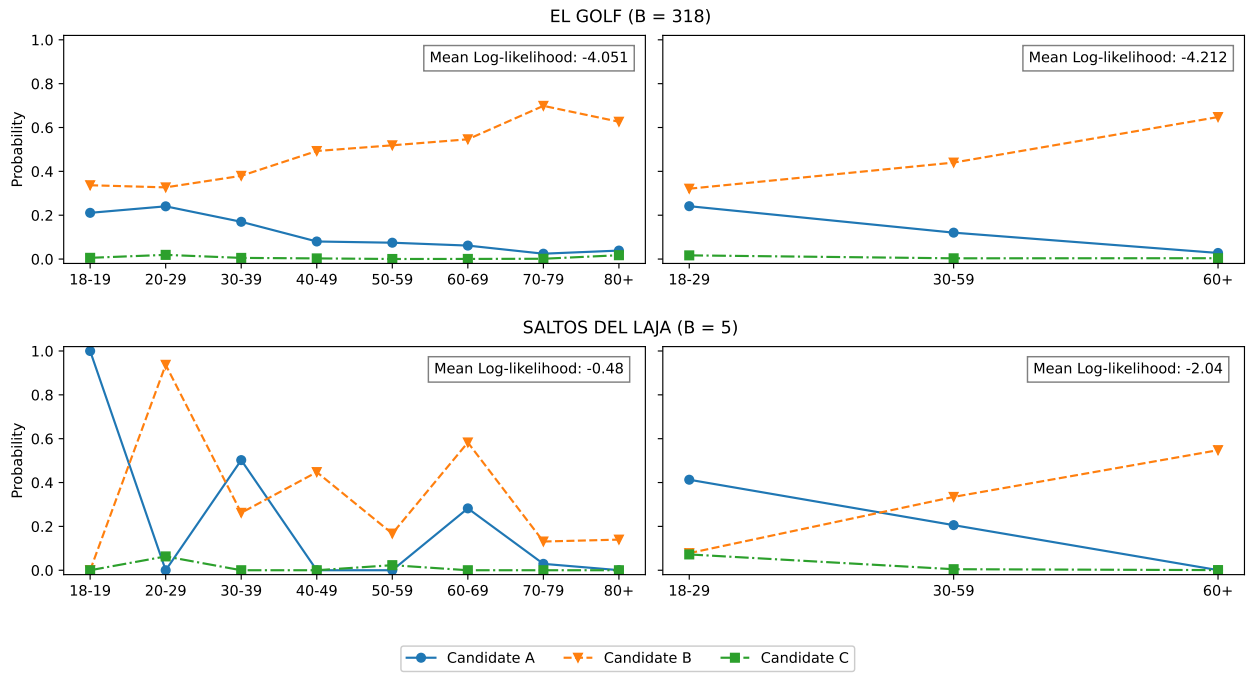


Figure 6.4: Estimated probabilities for different group aggregations.

The bottom of Figure 6.4 shows the estimated probabilities for the district *Salto del Laja*. We observe from the bottom left panel that the estimated voting probabilities have big variations among different age ranges, specially for candidates *A* and *B*. For example, the estimated voting probability for candidate *A* is 1 for voters in the 18-19 age range, whereas the estimate goes to 0 for voters in the 20-29 age range. In the bottom right panel we observe the estimated voting probabilities when considering 3 aggregated age groups. In this case, the estimated probabilities among the aggregated age groups show a more regular pattern. The district of *Salto del Laja* has only 5 ballot-boxes; as a result, it overfits the estimated probabilities when using eight demographic groups. Indeed, the log-likelihood of the bottom left case is considerably higher than the one of the bottom right panel. Note that the overfitting phenomenon, when using 8 age groups, does not occur in the district of *El Golf* (top left panel) since this district has a considerable amount of ballot boxes ($B = 318$). We observe that in the district of *El Golf* we get a similar probability in both the top left and top right panels, meaning that we should probably stay with the 8 age groups as they increase the log-likelihood.

These results show once again that it is not clear which group aggregation is better, as it may depend in the group distribution and number of ballot-boxes in the district. Indeed, there

are several other possible group aggregations we could have done.. In order to approach this issue we will introduce a metric of reliability which could be obtain through bootstrapping. Besides, we could use the log-likelihood to measure how the probability is fitting to the distribution of the district.

6.2.1. Log-Likelihood

As we ran the EM-algorithm in a district, we obtain an estimated probability matrix $\hat{\mathbf{p}}$. We are not able to use the MAE metric anymore as we do not know the real probability. However, we are able to calculate the log-likelihood of the estimation using equation 2.6. This log-likelihood will not necessarily have an interpretation of how good the estimator is, but it could be useful to study how the estimation changes as we add more groups.

6.2.2. Bootstrapping: Estimated Probability Standard Deviation

As we ran the EM-algorithm in a district, we obtain an estimated probability matrix $\hat{\mathbf{p}}$. It is our objective to determine how reliable this estimation is. We could consider a bootstrapping approach to see how this estimated probability changes as we change the subset of ballot-boxes that are considered for running the EM-algorithm.

We can formulate a very simple algorithm as follows.

Algorithm 5 Bootstrapping for estimated probabilities

Input: $\mathcal{I}, \mathcal{G}, \mathcal{C}, \mathcal{B}, \mathbf{x}, \mathbf{w}$
Output: $\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(S)} \in [0, 1]^{G \times C}$
for $s = 1$ to S **do**
 Sample: a_1, a_2, \dots, a_B from $\mathcal{U}\{1, B\}$ (with replacement)
 Create matrices: $\mathbf{x}' \in \mathbb{Z}_+^{B \times C}, \mathbf{w}' \in \mathbb{Z}_+^{B \times G}$
 Assign: $\mathbf{x}'_b = \mathbf{x}_{a_b}, \mathbf{w}'_b = \mathbf{w}_{a_b} \quad \forall b \in \mathcal{B}$
 Estimate $\mathbf{p}^{(s)}$ using \mathbf{x}' and \mathbf{w}' with Algorithm 1.

As we get a collection of estimated probabilities $\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(S)}$ we may use it to estimate a metric of standard deviation. This could be done for the probabilities estimated for an arbitrary group, as we are interested on observing if we should include that group or not.

The standard deviation of the estimated probability p_{gc} can be expressed as follows:

$$\text{std}(p_{gc}) := \sqrt{\frac{1}{S-1} \cdot \sum_{s=1}^S (p_{gc}^{(s)} - \mu_{gc})^2}$$

$$\text{std}(\mathbf{p}) := \frac{\sum_{g \in \mathcal{G}, c \in \mathcal{C}} \text{std}(\mathbf{p}_g)}{G \cdot C}$$

Where $\mu_{gc} := \frac{1}{S} \sum_{s=1}^S p_{gc}^{(s)}$.

This metric will give us an insight on how reliable the EM-algorithm is for estimating the probabilities associated for each combination of group g and candidate c .

6.2.3. Adding More Groups: Trade-off

We will study how considering more groups affects the estimation through the log-likelihood and the standard deviation. We will consider the following sets of aggregated groups.

- 1 Group: $\mathcal{A} = \{18^+\}$
- 2 Groups: $\mathcal{A} = \{18-49, 50^+\}$
- 4 Groups: $\mathcal{A} = \{18-29, 30-49, 50-69, 70^+\}$
- 8 Groups: $\mathcal{A} = \{18-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80^+\}$

We will select districts with varying number of ballot-boxes through Chile. For each selected district we are going to run the EM-algorithm with the 4 different group aggregations. For each group aggregation we will compute both the log-likelihood and the standard deviation, obtaining the result displayed in Figure 6.5.

6.2.4. Best Group Aggregation

Given the results in Figure 6.5, we would like to choose the best group combination for each district. We know that we would like to maximize the log-likelihood but without sacrificing so much standard deviation in the estimator. In this particular application we will consider a threshold (γ) for the standard deviation each group voting probability must satisfy given a group set. We will choose this threshold to be $\gamma = 0.05$.

We would like to try each possible group set for each district, however, this is not possible as the number of different group combinations scale exponentially. Consider we have G

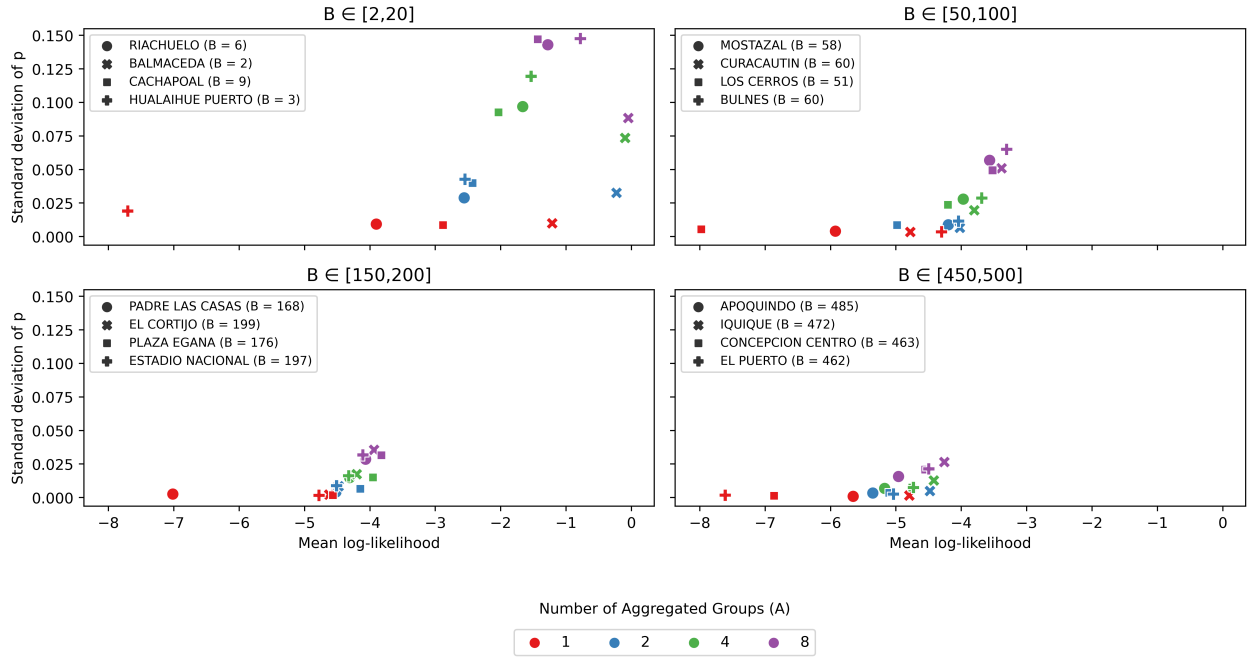


Figure 6.5: Standard deviation of estimated probability vs log-likelihood over districts of different ballot-boxes number.

possible age ranges and we can only merge neighbor age ranges, the total number of group sets of size A we can make are:

$$\binom{G-1}{A-1}$$

As trying each feasible group set is not a reliable option. We will propose an iterative approach to obtain a good group set for a given district. Firstly, we will determine a way to choose a group set given its size: A .

We know that we would like groups to be heterogeneous accross ballot-boxes in order to make good estimators, so we will choose the group set that maximizes the sum of the standard deviation of groups distribution in that set. This would be done through mixed-integer programming, where we set the number of groups as a constraint.

$$\begin{aligned}
 & \max_{\mathbf{a}} \sum_{m \in \mathcal{M}} a_m \cdot \sigma_m \\
 & \text{subject to: } \sum_{m \in \mathcal{M} | g \in m} a_m = 1 \quad \forall g \in \mathcal{G} \\
 & \sum_{m \in \mathcal{M}} a_m = A \\
 & a_m \in \{0, 1\} \quad \forall m \in \mathcal{M}
 \end{aligned} \tag{6.1}$$

The decision variable \mathbf{a} represents the following:

$$a_m = \begin{cases} 1 & \text{Group combination } m \text{ is included in the set} \\ 0 & \text{Otherwise} \end{cases}$$

We are considering σ_m to be the standard deviation of group combination m . Then, \mathcal{M} is the set of all possible group combinations given the full set \mathcal{G} . For the particular structure of this application we can define it as follows:

$$\mathcal{M} = \{m \in \mathcal{P}(\mathcal{G}) \mid \forall g_1, g_2 \in m, \{g \in \mathcal{G} \mid g_1 \leq g \leq g_2\} \subseteq m\} \quad (6.2)$$

We are considering the power set of \mathcal{G} but only taking those subsets that satisfy that each element has at least one neighbor in the set. Let us note that we are considering $\mathcal{G} = \{1, 2, \dots, G\}$ so that each index corresponds to an age range in ascending order.

Let us take into account that the cardinality: $|\mathcal{M}| = \sum_{k=1}^G = \frac{G \cdot (G+1)}{2}$, scales as $O(G^2)$. For the application in the 2021-CPE we have $|\mathcal{M}| = 36$ which is a reasonable instance of the problem.

Now that we have a way of choosing a group set given a size A , we will propose an iterative approach to get the final set for each district. The idea is that we would like to have as much groups as possible, as it is shown to improve the log-likelihood. This way we will start with the biggest set possible \mathcal{G} (as it is the original set) and observe how the standard deviation of the estimation is using Algorithm 5. If the standard deviation is below the threshold then we choose that set, else we try with a smaller set of size $G - 1$. As we have different choices for a set of this size we solve the optimization problem described in 6.1. If this new solution does not satisfy the threshold we keep decreasing the set size. If no group set satisfies the threshold we just use the group set of size 1 where we merge all age ranges into one.

Algorithm 6 Choosing a group set

Input: $\mathcal{I}, \mathcal{G}, \mathcal{C}, \mathcal{B}, \mathbf{x}, \mathbf{w}, \gamma$ **Output:** \mathcal{A} **for** $A = G$ to 1 **do** **Solve** the optimization problem described in 6.1, input \mathcal{G}, A $\mathcal{A} = \{m \in \mathcal{M} | a_m = 1\}$ $w'_{ba} = \sum_{g \in a} w_{bg}, \forall a \in \mathcal{A}$ **Apply** bootstrapping using Algorithm 5, input $\mathcal{I}, \mathcal{A}, \mathcal{C}, \mathcal{B}, \mathbf{x}, \mathbf{w}'$ **if** $\text{std}(p_{ac}) \leq \gamma, \forall a \in \mathcal{A}, \forall c \in \mathcal{C}$ **then** **Stop**

6.3. Results

We will evaluate the discussed models and methodology for each district of the 2021-PGE. We will first choose aggregated group sets according to algorithm 6, then we will perform the EM-algorithm according to 1 and finally we will compute the p-values for each ballot-box using the multinomial approximation as described in 4.2 using a precision up to 10^9 samples. We should take into account that these p-values may have biases according to the discussion in Section 5.3.

6.3.1. Group Aggregation Sets

As we run the methodology established in 6 we obtained a group set \mathcal{A} for each district in the election. We can study some statistics related to the chosen sets.

In Figure 6.6 we observe the fraction of districts where the original age range is selected as a single group. We can notice that the age ranges of 70-79 and 80⁺ are almost never selected as a single group, as their probabilities must be hard to identify. A similar phenomenon occurs with the younger age ranges of 18-19 and 20-29, but they do get selected as a single groups in some more instances. We observe that those age ranges in the middle: 30-39, 40-49 and 50-59; are the age ranges that tend to be selected as a single group, as they probably are well represented through the district and the estimated probability does not radically change when doing bootstrapping.

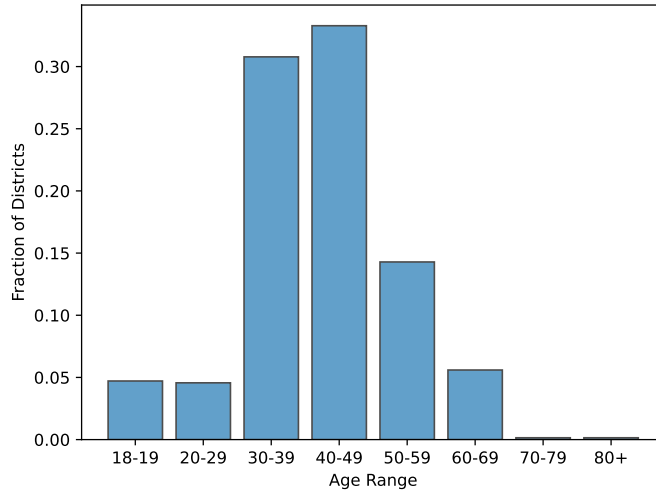


Figure 6.6: Fraction of districts where the age-range is included as a single group.

We can also study the size of the chosen group sets. In Figure 6.7 we display some statistics of the sets over the districts. In Figure 6.7.a we can observe how the number of ballot-boxes in a district affects the number of aggregated groups selected. We can observe that when there are more ballot-boxes the methodology tends to result in bigger sets, meaning that it is feasible to identify the probability of more age-ranges. In Figure 6.7.b we observe the proportion of districts that had each aggregated group size. We can observe that most districts get a set with a single group, this is most likely due to many districts containing few ballot-boxes, meaning that a lot of groups would be hard to identify. We can see that in general there are less districts as we increase the aggregated group set size. This indicates that there are few districts where the distribution of groups among the ballot-box is reliable enough to identify multiple age-range probabilities. We should also note that no district satisfied the threshold when trying the original group set where every age range is treated individually.

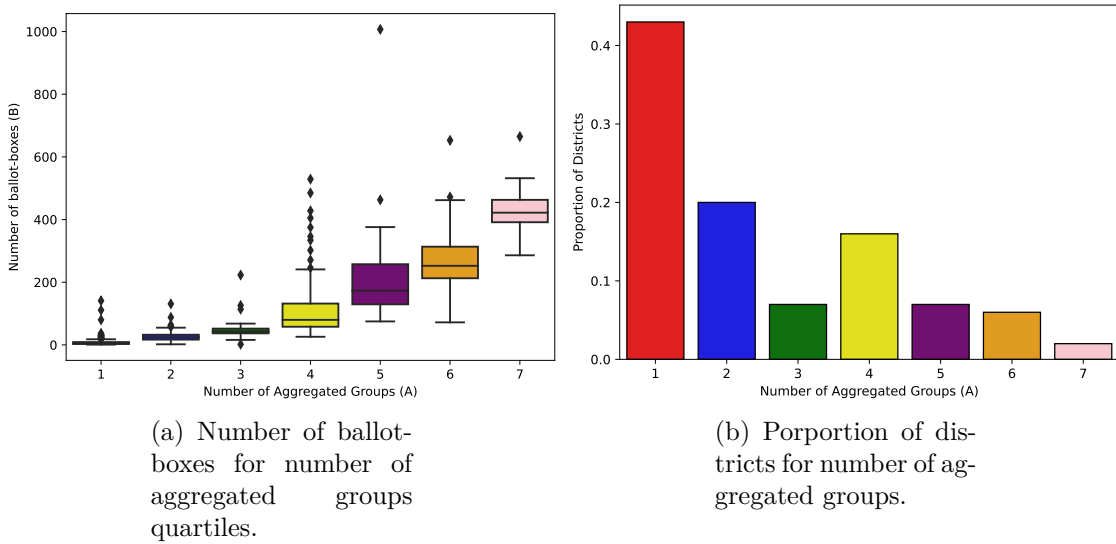


Figure 6.7: Statistics for aggregated group sets

6.3.2. Probabilities

As we run the EM-algorithm in each district we will obtain a probability matrix estimation \hat{p} . We would like this estimator to be a good representation of how voters actually vote in the given district.

In Figure 6.8 we observe a graded representation of how the estimated probabilities would look for the districts of *Calama*, *Maipú* and *Villarrica*. We can first notice that this districts vary in the number of selected age ranges. We can note that *Maipú* is the district with the higher number of ballot-boxes in 6.8, however we should note it was in the higher end of the distribution for districts with 6 age ranges according to Figure 6.7.a. Regarding the probabilities, we can observe decently regular patterns accross different candidates. We can see that *Candidate 1* and *Candidate 8* have low probabilities for all age groups accross all age ranges. Regarding *Candidate 2* we can observe higher voting probabilities in the district of *Calama* where there is quadratic correlation with age, in constrast with the other districts of *Maipu* and *Villarrica* where *Candidate 2* is less voted and it shows a decreasing correlation with age. In regards to *Candidate 3* we see a common decreasing correlation between probabilities and age in all districts, where the voting probabilities are higher in the district of *Maipú*. An opposite pattern is find in *Candidate 4* where there is an increase in the probabilities with age, with the higher probabilities found on the district of *Villarrica*.

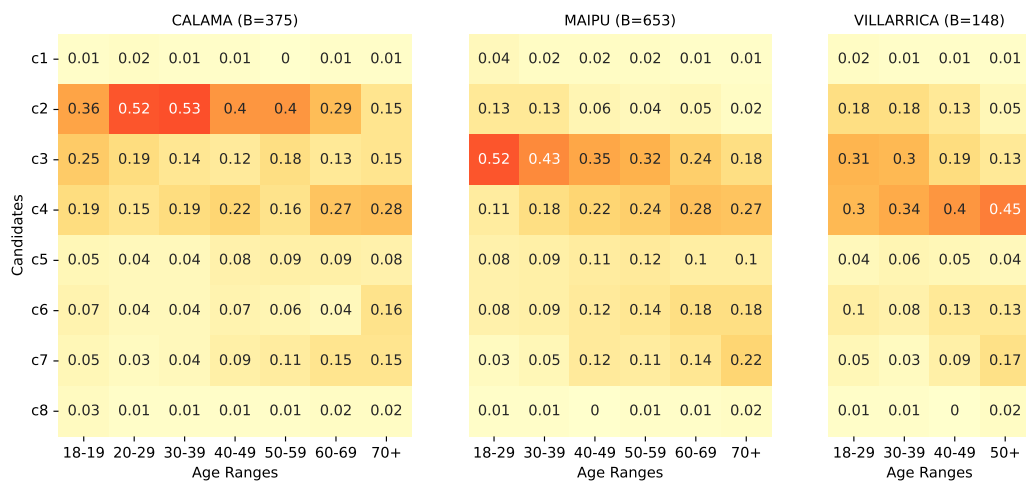


Figure 6.8: Probability estimation for the districts of *Calama*, *Maipú* and *Villarrica*

6.3.3. p-values and Outliers

We present some general results for the p-values obtained. We should note that obtaining a low p-value is an indicator that a ballot-box is an outlier but it does not states what happened in that ballot-box. Some reasons a ballot-box could show a low p-value are the following:

1. There were errors in the data, this may include: vote count, group count or ballot-box number mislabeling.
2. There was relevant information about the voter assignation of that ballot-box that was not considered in the model. This could happen if a particular set of voters (particularly different from the rest of the district) where assigned to a specific ballot-box or to a specific polling place due to some arbitrary criteria.
3. The vote count was altered and does not reflect the actual results of that ballot-box.

To visualize how a ballot-box with a low p-value looks, we present the scheme displayed in Figure 6.9. We show statistics for all ballot-boxes in the polling place in order to identify patterns related to the low p-values. The matrices displayed in Figure 6.9 represent the following from left to right: the vote count for each candidate, the vote difference to the expected value (rounded for visualization), the number of voters for each selected age range and the p-values. The difference to the expected value for a candidate c in ballot-box b , considering a estimated matrix probability \hat{p} and an aggregated group set \mathcal{A} can be calculated

as follows:

$$x_{bc} - \mathbb{E}(X_{bc}) = x_{bc} - \sum_{a \in \mathcal{A}} w_{ba} \cdot \hat{p}_{ac}$$

In Figure 6.9 we can observe that all ballot-boxes, besides those labeled 142 and 144, show normal p-values of at least 1 in a hundred. The ballot-box labeled 144 shows a p-value of less than 10^{-9} (less than 1 in a billion), meaning that the computed value with simulation was 0. In this ballot-box we can see an important deviation from the expected value in the votes for candidates 3, 6 and 8. We can also observe that this ballot-box was primarily composed by voters in the 18-29 range.

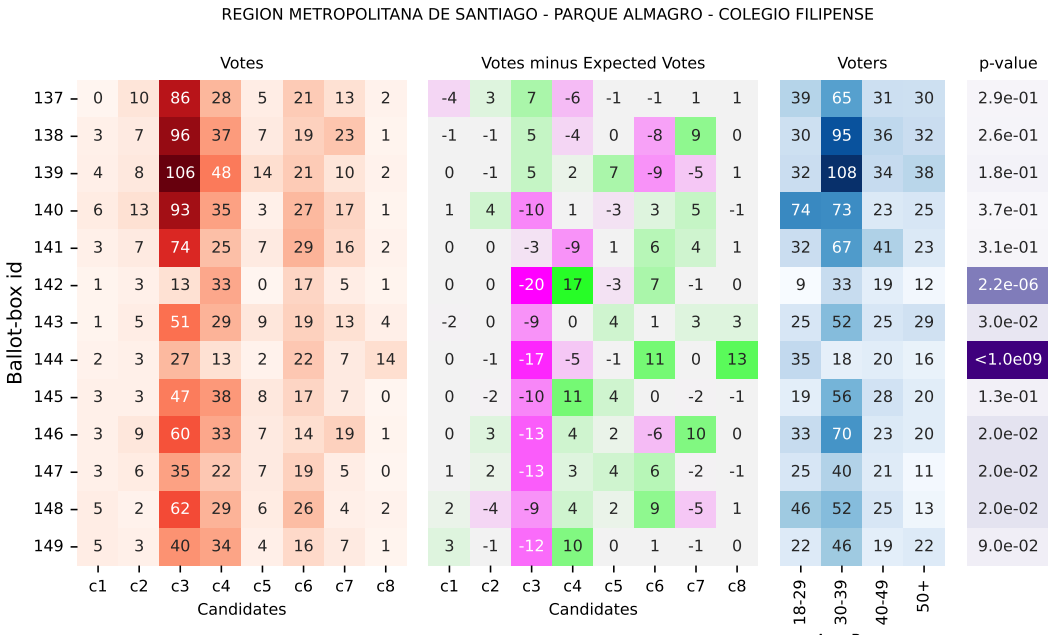


Figure 6.9: Results in the polling place of: *Colegio Filipense*

In Annex C we show a detailed result for the polling places containing the lowest encountered p-values in the range $\leq 10^{-8}$.

In Table 6.1 we present the p-value count for different ranges and in Figure 6.10 we present a scatter plot for the p-value and number of votes of each ballot-box.

We should note that by construction, p-values should approximate to a Uniform(0, 1) distribution. We observe a big concentration of p-values in the interval $[0.1, 1]$, at around 80%. We can note that the percentages do not match the expected quantity of each range. This could be attributed to different reasons: there is not enough information to get to good

estimators in some polling places, and the present outliers could alter both the probability distribution (specially in small polling places) and the interval count. We can observe that the count until 10^{-2} matches relatively well the expected results, however, from 10^{-3} we observe that the count does not decrease as it should.

We also observe that there are 15 ballot-boxes that have a p-value of less or equal than 10^{-9} . For those values that are strictly less than 10^{-9} it could be possible to further simulate samples in order to get an estimation to the order of magnitude. To further estimate, we would need to increase the sample size by 10 every time, making it computationally expensive. Even though we are not able to exactly compute these p-values, the chosen sample size is enough to determine those ballot-boxes as outliers.

Table 6.1: p-value count in the 2021-GCE.

By Intervals			Cumulative		
Range	Total	%	Range	Total	%
$(10^{-1}, 10^0]$	36584	78.50	$\leq 10^0$	46605	100.00
$(10^{-2}, 10^{-1}]$	7471	16.03	$\leq 10^{-1}$	10021	21.50
$(10^{-3}, 10^{-2}]$	2055	4.41	$\leq 10^{-2}$	2550	5.47
$(10^{-4}, 10^{-3}]$	259	0.56	$\leq 10^{-3}$	495	1.06
$(10^{-5}, 10^{-4}]$	111	0.24	$\leq 10^{-4}$	236	0.51
$(10^{-6}, 10^{-5}]$	54	0.12	$\leq 10^{-5}$	125	0.27
$(10^{-7}, 10^{-6}]$	28	0.06	$\leq 10^{-6}$	71	0.15
$(10^{-8}, 10^{-7}]$	16	0.03	$\leq 10^{-7}$	43	0.09
$(10^{-9}, 10^{-8}]$	12	0.03	$\leq 10^{-8}$	27	0.06
$(10^{-\infty}, 10^{-9}]$	15	0.03	$\leq 10^{-9}$	15	0.03

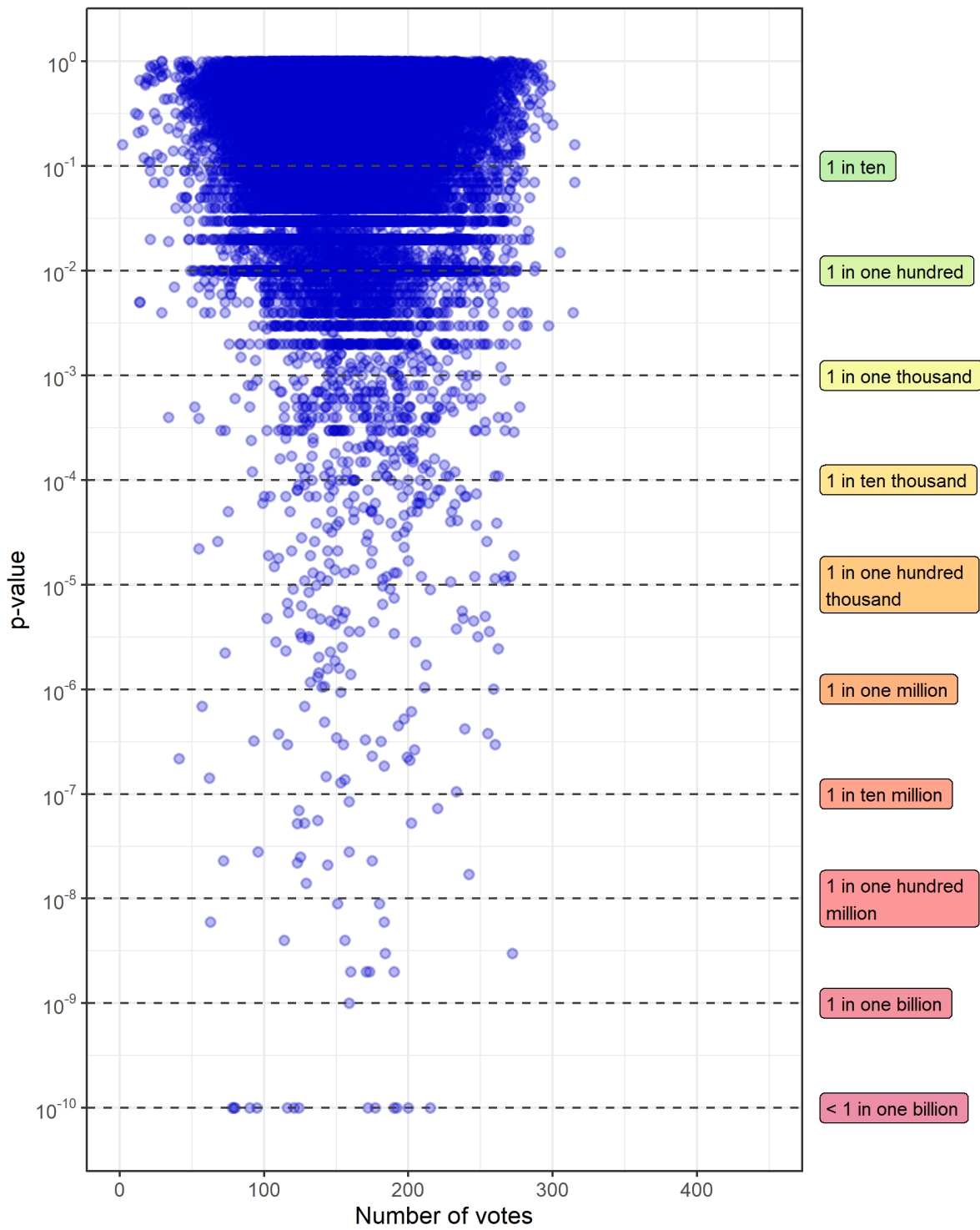


Figure 6.10: Estimated p-value accross all ballot-boxes.

6.4. General Recommendations

Identifying and addressing outlier ballot-boxes in future elections in Chile could help improve the overall integrity and transparency of the electoral process. Here are some actions that SERVEL could take in the short and medium term with the information the methods provide:

1. Hold the publication of official results for that ballot-box .
2. Verify that protocols were followed in that ballot-box and in the polling place containing that ballot-box.
3. Verify that the results in that ballot-box matches the official record.
4. Change the ballot-box committee for future elections.
5. Assign watchmen from the main parties to guard the counting of that ballot-box for future elections.

If these methods were to be put into practice, it would be beneficial to incorporate the voter assignment method directly into the model. This would assist in improving the identification of the territorial units for conducting the probability estimation and would also consider any arbitrary decisions that could have done in the allocation of voters.

Chapter 7

Conclusions

Through this study we analyzed different methods for computing the EM-algorithm in order to do ecological inference in the context of Chilean elections. It was shown that the EXACT method fails to compute medium and large instances, meaning it would be advisable to utilize approximated methods from $C = 3$ onwards. The amount of groups G does not affect the EXACT method exponentially, as the number of different combinations does not increase monotonically with G , so it could be possible to compute instances with a lot of demographic groups if the number of candidates is small. The MULT method was shown to be the fastest so it would be recommended to use it in most instances, specially if C is high. The MVN-PDF method is not as fast, but it is still polynomial and not highly impacted by the size of the instance, so it would also be a advisable to use it. The MVN-CDF gets similar estimations compared to the MULT and MVN-PDF methods, but is also highly impacted by the size of the instance so it would not be a good choice. The H&R method needs a lot of samples in order to have a good estimation in big instances meaning it is not reliable. It should be noted that these methods could be used for other applications where the distribution corresponds to a sum of multinomial random variable.

We observed that the methodology proposed for the EM-algorithm works best in heterogeneous ballot-boxes, so a previous assessment of the group distribution should be considered. This can be analyzed through the standard deviation and correlation of different groups in a ballot-box. It should be noted that in the case of having different amounts of voters per ballot-box, this analysis should be done through the proportion of voters from the group instead of the quantity, so that ballot-boxes are comparable.

In regards to the p-values, we showed that there is an important gap in the calculation when using the estimated probabilities from the EM-algorithm. This means that even though the MAE of the estimated probability obtained could be small, the error when computing the p-value could be amplified, as we are considering a sum of probabilities from a slightly different distribution. We also showed that there is a gap when using the multinomial approximation, specially when the variance of the original distribution is amplified. This shows that p-values should be carefully computed as there are errors attributed to both the estimated probability and the method used to compute the p-values.

We finally showed a methodology for aggregating groups in order to get reliable estimations. We used the 2021-PGE as a case study to test the proposed methods in this work. We showed that as there were more ballot-boxes in a district it was convenient to consider bigger sets (less aggregated). Using these group sets we computed estimated probabilities for each district and p-values for each ballot-box. We obtain 27 ballot-boxes with a p-value lower or equal than 10^{-8} . We discussed possible explanations to such low p-values: mislabeling, unmatched data, manipulation, among others. We stressed the fact that these models do not give an explanation to the low p-values, but they can work as an indicator to better improve the election process.

7.1. Future Work

This study leaves some open questions that would be interesting to address in the future.

- Consider different simulation approaches for the H&R method. As we have shown, there is an increase in the MAE of the estimation as we consider less samples, meaning it would be necessary to find a faster sampling method. It would be interesting to study how important the correlation between obtained points affects the estimation, so that we could allow a smaller step-size in some situations. It may even be possible to study a deterministic approach to come up with a subset of Ω_b that is representative of the aggregate outcome. It could also be interesting to study sampling approaches that take into account the original probability distribution and how it affects the final estimation.
- Consider different methods to approximate the p-values. As we have shown, the multinomial approximation has a bias estimating the p-values when there are big changes in

the original variance. We could study how simulations coming from the multivariate normal approximation compute the p-value, as it maintains the original covariance matrix. However, it should be noted that this method is considerably slower, and as we are interested in studying the presence of outliers, we do need a big sample size.

- Address special cases for elections in the Chilean context. The methods studied assume that voters from the same group are comparable within a district. This assumption may not be true for recent elections where the ballot-box assignation is based on distance. This factor could be taken into account in the probability distribution. It could also be interesting to study a model that uses different districts for the same ecological inference process. This would be particularly useful for small districts with a small number of ballot-boxes, where the estimation is less reliable. A general method could take into account both: aggregate demographic information from voters within a district (as shown in the application of this study) and aggregate demographic information from districts (such as average salary, unemployment rates and geographic location).

Bibliography

- Anselin, L., & Cho, W. K. T. (2002). Spatial effects and ecological inference. *Political analysis*, 10(3), 276–297.
- Antweiler, W. (2007). Estimating voter migration in canada using generalized maximum entropy. *Electoral Studies*, 26(4), 756–771.
- Cleave, N., Brown, P., & Payne, C. (1995). Evaluation of methods for ecological inference. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 158(1), 55–72.
- Duncan, O., & Davis, B. (1953). An alternative to ecological inference. *American Social Review*, 18, 665–666.
- Fortin-Rittberger, J., Harfst, P., & Dingler, S. C. (2017). The costs of electoral fraud: establishing the link between electoral integrity, winning an election, and satisfaction with democracy. *Journal of Elections, Public Opinion and Parties*, 27(3), 350–368.
- Genz, A. (1992). Numerical computation of multivariate normal probabilities. *Journal of computational and graphical statistics*, 1(2), 141–149.
- Glynn, A. N., & Wakefield, J. (2010). Ecological inference in the social sciences. *Statistical methodology*, 7(3), 307–322.
- Goodman, L. A. (1953). Ecological regressions and behavior of individuals. *American sociological review*, 18(6), 663.
- Goodman, L. A. (1959). Some alternatives to ecological correlation. *American Journal of Sociology*, 64(6), 610–625.
- Jackson, C., Best, N., & Richardson, S. (2006). Improving ecological inference using individual-level data. *Statistics in medicine*, 25(12), 2136–2159.
- King, G., Tanner, M. A., & Rosen, O. (2004). *Ecological inference: New methodological*

strategies. Cambridge University Press.

- Leemann, L., & Bochsler, D. (2014). A systematic approach to study electoral fraud. *Electoral studies*, *35*, 33–47.
- Leigh, A. (2005). Economic voting and electoral behavior: How do individual, local, and national factors affect the partisan choice? *Economics & Politics*, *17*(2), 265–296.
- Lin, Z., Wang, Y., & Hong, Y. (2022). The poisson multinomial distribution and its applications in voting theory, ecological inference, and machine learning. *arXiv preprint arXiv:2201.04237*.
- Mete, H. O., & Zabinsky, Z. B. (2012). Pattern hit-and-run for sampling efficiently on polytopes. *Operations Research Letters*, *40*(1), 6–11.
- Morgenstern, H. (1995). Ecologic studies in epidemiology: concepts, principles, and methods. *Annual review of public health*, *16*(1), 61–81.
- Peterson, J. C., Smith, K. B., & Hibbing, J. R. (2020). Do people really become more conservative as they age? *The Journal of Politics*, *82*(2), 600–611.
- SERVEL. (2023). *Resultados electorales históricos*. <https://www.servel.cl/centro-de-datos/resultados-electorales-historicos>. (Accesed: 2023-08-15)
- Wakefield, J. (2008). Ecologic studies revisited. *Annu. Rev. Public Health*, *29*, 75–90.
- Withers, S. D. (2001). Quantitative methods: advancement in ecological inference. *Progress in Human Geography*, *25*(1), 87–96.

Annex A

Hit and Run Method

A.1. Starting Point for the Hit-and-Run Algorithm

Algorithm 7 Starting Point for Hit and Run

Input: $w_b \in \mathbb{Z}_+^G, x_b \in \mathbb{Z}_+^C$
Initialize $z_b = \mathbf{0} \in \mathbb{Z}_+^{G \times C}$
for $g = 1$ to G **do**
 for $c = 1$ to C **do**
 $z_{bgc} = \min\{w_{bg}, x_{bc}\}$
 $w_{bg} = w_{bg} - z_{bgc}$
 $x_{bc} = x_{bc} - z_{bgc}$

A.2. Choosing the Step-Size

For algorithm 3 to work properly we would like to choose a step-size M so that the sampled points are independent from each other (in other words, we would like the starting point of each iteration not to determine the outcome).

For this study we will do an empiric evaluation on the effect of the step-size. We will evaluate 4 instances with varying values of G and C . For each instance we will simulate 20 scenarios containing 50 ballot-boxes with 100 voters each.

We will apply the algorithm using a step-size of $M = 100$ and $S = 10000$ for each ballot-box. Let us consider $\vec{\mathcal{S}}(b)$ as the ordered array containing all the points sampled in the algorithm (even if repeated). We will evaluate how correlated sampled points are using the following.

$$\rho_{m.M}(b) = \text{corr}(\vec{\mathcal{S}}_{m \rightarrow S}(b), \vec{\mathcal{S}}_{1 \rightarrow S-m}(b))$$

What we are doing is looking how one point correlates to the point obtained m iterations after, this allows us to evaluate the correlation for step-sizes in multiples of M .

Using this procedure we obtain the following average result for varying step-sizes over the different instances.

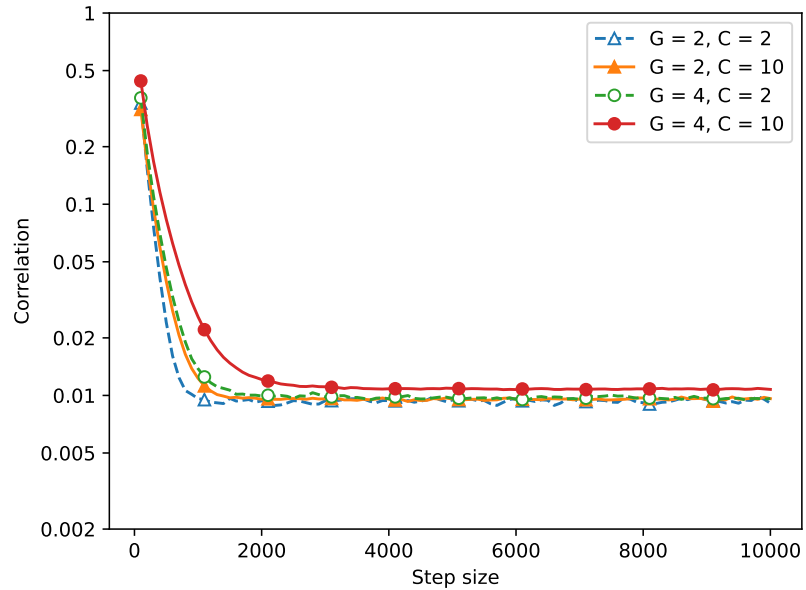


Figure A.1: Mean correlation for various step-sizes over different instances. The graph is log-scaled.

We observe a notable decrease for all instances when augmenting the step-size from 100 to 1000. Then we observe small decreases for some instances until they all hit a plateau around a step-size of 3000. This last value will be used as the step-size when running this method through this study

A.3. Another Formulation for the Discrete Polytope

$$\begin{aligned}
 \sum_{g=1}^{G-1} z_{bgc} &\leq x_{bc} && \forall c \in \{1, \dots, C-1\} \\
 \sum_{c=1}^{C-1} z_{bgc} &\leq w_{bg} && \forall g \in \{1, \dots, G-1\} \\
 \sum_{g=1}^{G-1} \sum_{c=1}^{C-1} z_{bgc} &\geq I_b - x_{bC} - w_{bG} \\
 z_{bgc} &\geq 0 && \forall c \in \{1, \dots, C-1\} \forall g \in \{1, \dots, G-1\} \\
 z_{bgc} &\in \mathbb{Z} && \forall c \in \{1, \dots, C-1\} \forall g \in \{1, \dots, G-1\}
 \end{aligned}$$

Annex B

EM-algorithm Probability Estimation with 200 Voters per Ballot-Box

B.1. Time results

Table B.1: Mean running time over 20 scenarios in seconds for the EM-algorithm for varying instances, with fixed values of $I_b = 200$ and $B = 50$.

Instance		Method					
C	G	EXACT	H&R $S = 10^3$	H&R $S = 10^2$	MVN CDF	MVN PDF	MULT
2	2	59.634	740.965	74.230	1.385	0.176	0.003
	3	100.076	736.357	75.210	2.469	0.242	0.004
	4	188.472	698.575	75.254	3.620	0.250	0.005
3	2	-	745.531	72.961	6.076	0.220	0.003
	3	-	718.089	70.551	12.326	0.249	0.004
	4	-	727.985	72.246	17.310	0.318	0.005
5	2	-	711.375	71.098	19.955	0.208	0.003
	3	-	576.660	67.952	42.570	0.265	0.003
	4	-	676.125	65.974	58.629	0.303	0.005
10	2	-	670.993	66.676	111.144	0.233	0.002
	3	-	534.038	61.340	204.632	0.266	0.003
	4	-	596.844	58.734	327.127	0.334	0.004

Table B.2: Comparison of mean simulation time and mean EM running time in seconds over 20 scenarios for varying instances, with fixed values of $I_b = 200$ and $B = 50$.

Instance		Method			
C	G	H&R ($S = 10^2$)		H&R ($S = 10^3$)	
		Sim-time	EM-time	Sim-time	EM-time
2	2	73.409	0.820	740.046	0.919
	3	72.915	2.295	722.202	14.155
	4	72.732	2.522	686.560	12.015
3	2	71.105	1.855	733.173	12.359
	3	68.670	1.881	699.799	18.290
	4	69.916	2.331	703.638	24.347
5	2	69.471	1.627	693.218	18.157
	3	65.756	2.196	567.549	9.111
	4	64.044	1.929	652.590	23.535
10	2	65.161	1.515	655.398	15.595
	3	59.866	1.474	525.634	8.404
	4	57.256	1.479	578.101	18.743

B.2. Error results

Table B.3: Mean absolute prediction error over 20 scenarios for varying instances, with fixed values of $I_b = 200$ and $B = 50$.

Instance		Method					
C	G	EXACT	H&R $S = 10^3$	H&R $S = 10^2$	MVN CDF	MVN PDF	MULT
2	2	0.009	0.009	0.009	0.009	0.009	0.009
	3	0.011	0.011	0.013	0.011	0.011	0.011
	4	0.012	0.013	0.024	0.012	0.012	0.013
3	2	-	0.009	0.010	0.009	0.009	0.009
	3	-	0.013	0.027	0.011	0.011	0.012
	4	-	0.022	0.040	0.013	0.013	0.013
5	2	-	0.008	0.013	0.007	0.007	0.007
	3	-	0.016	0.028	0.009	0.009	0.009
	4	-	0.025	0.036	0.011	0.011	0.011
10	2	-	0.008	0.013	0.005	0.005	0.005
	3	-	0.017	0.021	0.007	0.007	0.007
	4	-	0.022	0.026	0.008	0.008	0.008

Annex C

Detailed Results for Polling Places with Low p-value Ballot-Boxes

REGION DE ANTOFAGASTA - ANTOFAGASTA SUR - LICEO DE HOMBRES DE ANTOFAGASTA MARIO BAHAMONDE SILVA

Ballot-box id		Votes								Votes minus Expected Votes								Voters				p-value
		c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
100V	-1	40	49	42	7	15	14	2	-1	2	8	-2	-4	-2	-1	0	19	53	9	88	8.2e-01	
101V	-3	43	30	31	12	16	8	1	1	9	-7	-5	3	2	-3	0	18	59	10	58	3.0e-01	
81V	-1	39	38	35	7	12	9	0	-1	7	2	0	-2	-1	-3	-2	17	53	4	67	8.8e-01	
82V	-2	42	28	37	12	9	13	1	0	7	-10	2	3	-4	2	0	18	65	3	58	3.2e-01	
83V	-8	15	40	45	34	1	15	1	6	-23	-2	6	25	-14	3	-1	16	78	7	58	1.0e-09	
84V	-2	34	43	36	10	15	15	2	0	-3	3	-3	0	0	3	0	17	65	11	64	9.9e-01	
85V	-3	31	32	38	9	17	11	2	1	0	-3	1	0	3	-2	0	12	46	12	73	9.0e-01	
86V	-1	39	44	42	13	12	16	3	-1	2	3	-2	2	-5	0	1	16	50	5	99	8.8e-01	
87V	-0	37	41	27	5	8	15	2	-2	8	9	-8	-4	-6	3	0	13	40	6	78	8.0e-02	
88V	-4	42	37	20	11	5	12	8	2	10	2	-15	2	-8	1	6	17	54	4	64	1.2e-05	
89V	-3	42	36	36	20	18	11	2	1	5	-6	-7	9	1	-3	0	13	61	15	79	1.2e-01	
90V	-1	27	33	34	14	12	13	3	-1	-3	0	-1	5	-2	1	1	14	41	9	74	6.8e-01	
91V	-1	31	40	35	15	13	15	3	-1	-3	2	-4	5	-2	2	1	13	55	6	79	5.7e-01	
92V	-2	25	33	52	11	12	14	1	0	-8	-4	13	1	-3	1	0	6	61	10	74	3.0e-01	
93V	-1	34	39	21	15	17	17	0	-1	1	3	-15	6	3	5	-2	21	51	7	65	3.0e-02	
94V	-4	49	36	32	9	6	13	0	2	17	0	-7	-1	-9	0	-2	7	56	4	83	4.0e-03	
95V	-3	39	31	36	13	18	9	0	1	4	-7	-1	4	4	-3	-2	17	59	4	69	3.9e-01	
96V	-2	48	46	27	5	7	6	0	0	13	8	-7	-3	-6	-4	-1	15	75	3	48	4.0e-02	
97V	-0	28	39	24	12	15	14	0	-2	-3	5	-9	4	3	4	-2	22	49	4	57	1.7e-01	
98V	-1	36	37	32	10	18	10	2	-1	1	-1	-4	1	4	-1	1	21	59	7	59	9.3e-01	
99V	-0	25	39	47	11	12	18	1	-2	-9	1	8	1	-3	5	-1	15	58	3	77	2.1e-01	

Figure C.1: Results in the polling place of: *Liceo de Hombres de Antofagasta Mario Bahamonde Silva*

REGION DE LA ARAUCANIA - VILLARRICA - ESCUELA MARIANO LATORRE

Ballot-box id	Votes								Votes minus Expected Votes								Voters				p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
37V	2	19	24	78	8	20	19	3	0	-2	-12	10	0	0	3	1	8	51	55	59	2.8e-01
38V	2	22	40	70	8	22	15	0	0	1	2	0	-1	2	-2	-2	18	52	43	66	9.7e-01
39V	1	21	32	69	10	19	11	3	-1	1	-3	4	2	0	-4	1	14	46	54	52	7.8e-01
40V	0	27	37	67	14	15	11	4	-2	7	1	-2	6	-6	-6	2	16	35	64	60	5.0e-02
41V	1	19	28	60	9	18	14	1	-1	-1	-5	3	2	1	1	0	15	45	53	37	9.5e-01
42V	2	24	44	71	1	15	15	1	0	4	8	3	-7	-5	-2	-1	23	38	48	64	7.0e-02
43V	0	17	28	80	13	14	11	2	-2	-3	-7	16	5	-5	-4	0	15	43	58	49	8.0e-02
44V	2	16	32	76	11	13	15	2	0	-4	-3	11	3	-6	-1	0	23	37	46	61	4.3e-01
45V	1	26	30	15	4	17	15	6	0	13	7	-30	-2	4	3	5	14	37	53	70	4.0e-09
46V	2	15	37	77	7	21	18	1	0	-6	-1	8	-2	1	1	-1	16	49	51	62	8.8e-01
47V	2	25	34	67	7	18	17	0	0	3	-4	2	-1	0	2	-2	25	53	37	55	9.1e-01
48V	1	20	32	80	9	12	18	5	-1	-2	-7	12	0	-7	2	3	20	64	31	63	1.1e-01
49V	1	16	49	69	7	21	20	2	-1	-7	8	-2	-2	1	3	0	23	64	29	69	6.9e-01
50V	4	13	57	83	8	22	18	1	2	-12	13	3	-2	-1	-2	-1	20	61	50	75	7.0e-02
51V	4	20	39	84	4	15	16	0	2	-2	0	13	-5	-5	-1	-2	21	54	40	67	1.0e-01
52V	1	22	31	85	11	28	18	0	-1	-1	-10	9	2	6	-2	-3	35	41	38	82	3.4e-01

Figure C.2: Results in the polling place of: *Escuela Mariano Latorre*

REGION DE VALPARAISO - CONCON - COLEGIO ALBORADA DEL MAR

Ballot-box id	Votes								Votes minus Expected Votes								Voters				p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
63 - 0	12	46	112	5	40	17	2	-2	-8	-1	23	-6	-5	-1	0	52	43	53	85	1.0e-02	
64 - 2	18	40	121	4	60	19	2	0	-5	-14	21	-9	10	-3	0	52	46	52	116	7.0e-03	
65 - 5	11	40	90	8	44	8	1	2	-9	-11	15	2	5	-3	-1	121	26	19	41	4.0e-02	
66 - 3	13	47	84	7	52	17	5	1	-7	-4	-1	-2	9	1	3	89	31	37	70	2.0e-01	
67 - 7	19	49	73	4	54	12	3	4	-2	-8	-6	-2	12	1	1	155	20	17	28	1.2e-01	
68 - 3	13	42	108	3	38	17	1	0	-7	-11	26	-5	-5	3	-1	116	20	34	54	1.2e-02	
69 - 4	11	75	110	9	57	16	1	1	-14	17	0	-2	2	-3	-1	59	88	57	79	2.0e-02	
70 - 2	15	58	105	6	57	25	4	0	-9	2	1	-6	5	5	2	55	69	53	95	1.3e-01	
71 - 1	16	50	105	10	53	12	2	-1	-6	-1	9	0	5	-6	0	63	58	62	66	6.8e-01	
72 - 2	16	54	97	5	53	16	3	0	-6	3	2	-5	6	-1	1	60	64	51	72	4.8e-01	
73 - 1	5	30	142	15	63	14	2	-1	-15	-20	43	-4	13	-15	-1	14	1	51	206	3.0e-09	
74 - 3	17	62	109	5	61	9	3	0	-12	-4	5	1	9	0	1	120	138	5	6	1.2e-01	
75 - 1	21	52	110	11	49	14	0	-2	-2	-8	14	2	-1	-1	-2	126	41	45	46	4.8e-01	
76 - 0	18	39	77	5	46	17	0	-2	-1	-5	2	-3	8	3	-2	50	59	13	80	4.5e-01	
77 - 1	10	40	121	12	66	20	2	-1	-14	-16	16	1	13	1	0	61	72	64	75	4.0e-03	
78 - 1	9	47	131	6	65	13	1	-2	-15	-9	26	-5	12	-6	-1	61	80	55	77	2.9e-04	
79 - 2	15	57	115	8	64	14	4	0	-9	0	7	-4	10	-6	2	63	72	61	81	9.0e-02	
80 - 2	9	53	131	11	53	19	4	0	-16	-5	22	-1	-1	-1	2	64	73	63	81	4.0e-03	

Figure C.3: Results in the polling place of: *Colegio Alborada del Mar*

REGION DE VALPARAISO - EL BELLOTO - COLEGIO D-417 GUARDIAMARINA GUILLERMO ZANARTU IRIGOYEN

Ballot-box id		Votes								Votes minus Expected Votes								Voters				p-value
		c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
10	- 2	21	26	33	17	23	13	1	0	3	-6	-9	4	6	2	0	10	4	78	44	2.7e-01	
11	- 1	20	31	40	17	13	8	0	-2	0	1	3	3	-3	-1	-1	15	3	101	11	9.1e-01	
12	- 5	25	39	52	15	25	16	2	1	-3	-2	0	-4	3	4	1	15	7	142	15	7.2e-01	
13	- 3	23	43	51	16	17	13	1	-1	-2	6	2	-2	-4	1	0	13	2	138	14	9.3e-01	
14	- 3	30	42	35	18	16	8	1	-1	6	9	-10	1	-3	-2	0	12	0	133	8	3.6e-01	
15	- 3	27	34	50	18	22	7	0	-1	3	-4	4	1	2	-4	-1	23	3	120	15	9.0e-01	
16	- 4	35	23	37	19	16	10	1	1	13	-9	-6	3	-2	0	0	10	6	118	11	1.3e-01	
17	- 5	27	39	48	18	22	9	1	1	0	-2	1	0	2	-2	0	20	19	118	12	9.9e-01	
18	- 5	16	35	50	19	18	8	2	2	-7	0	5	3	-1	-3	1	16	3	113	21	4.7e-01	
19	- 8	13	39	45	2	23	22	4	5	-14	-3	4	-13	6	12	3	14	57	66	19	4.0e-09	
1M	- 5	34	59	36	16	16	13	2	2	14	7	-19	2	-4	-2	0	37	16	21	107	5.0e-03	
20	- 3	28	43	38	11	10	7	1	0	-1	-2	7	0	-3	-1	0	14	104	13	10	9.3e-01	
21	- 6	24	46	31	14	16	12	2	3	-7	-2	-3	2	2	4	1	20	99	18	14	2.3e-01	
22	- 1	32	39	44	13	17	5	0	-2	2	-8	9	1	3	-4	-1	11	108	8	24	3.6e-01	
2M	- 2	17	47	68	16	18	20	1	-1	-2	-4	8	1	-4	3	-1	26	16	23	124	9.0e-01	
3M	- 1	16	48	50	16	25	17	8	-1	-2	-2	-7	1	4	1	6	28	11	22	120	5.0e-02	
4M	- 1	23	66	44	17	29	24	3	-2	-4	8	-16	0	6	7	1	30	44	31	102	3.0e-02	
5M	- 2	13	35	39	9	8	16	2	0	1	-1	1	0	-6	5	0	27	7	7	83	5.6e-01	
6	- 3	15	21	28	13	6	11	3	2	6	-6	-5	5	-6	2	2	11	7	10	72	1.0e-02	
7	- 2	13	30	29	13	11	6	0	1	3	3	-5	4	-1	-4	-1	9	6	16	73	3.1e-01	
8	- 1	21	23	35	16	10	7	0	-1	9	-5	-2	6	-3	-3	-1	8	4	37	64	9.0e-02	
9	- 1	14	29	32	20	20	9	1	-1	1	-3	-9	9	5	-2	0	12	2	35	77	1.4e-01	

Figure C.4: Results in the polling place of: *Colegio D-417 Guardiamarina Guillermo Zanartu Irigoyen*

REGION DEL MAULE - LINARES - ESCUELA PRESIDENTE CARLOS IBANEZ DEL C.

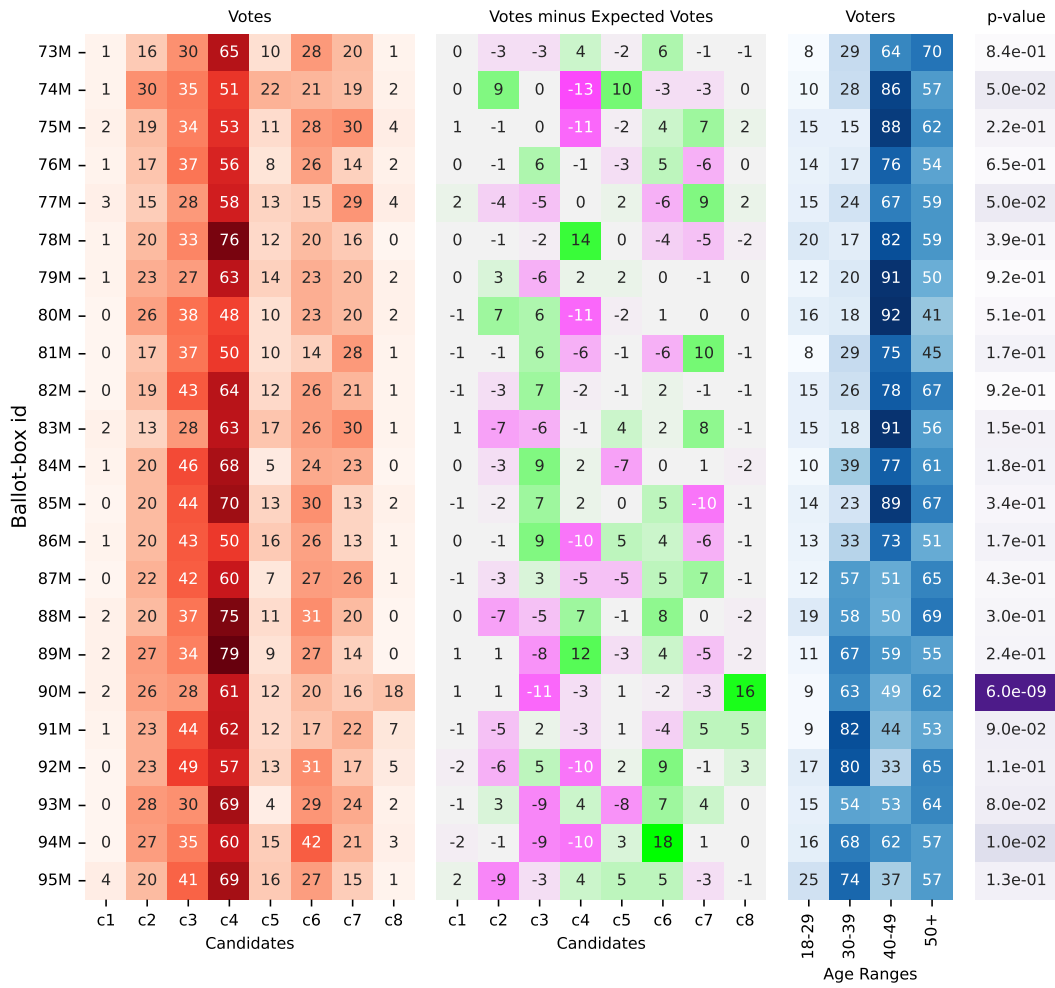


Figure C.5: Results in the polling place of: *Escuela Presidente Carlos Ibáñez del C.*

REGION DEL MAULE - SAN CLEMENTE - COLEGIO CLEMENTINOS

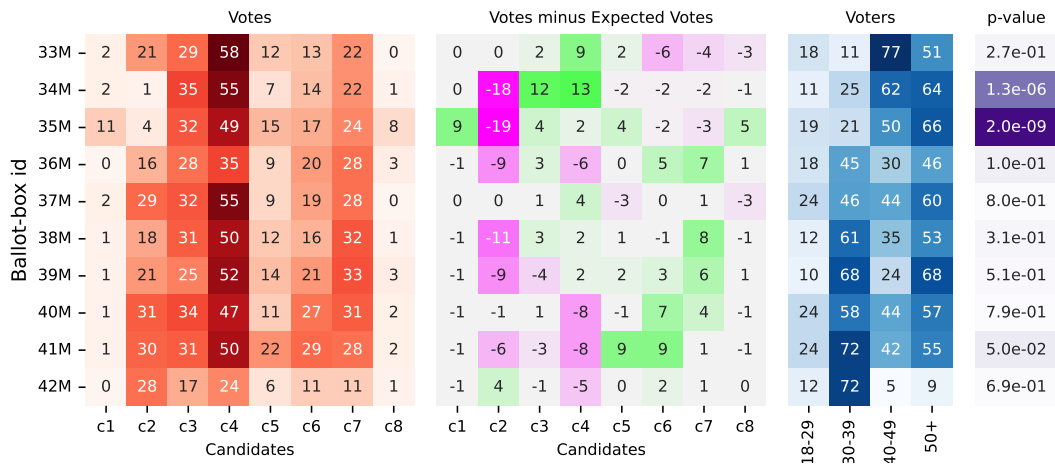


Figure C.6: Results in the polling place of: *Colegio Clementinos*

REGION DEL MAULE - SAN JAVIER - ESCUELA GERONIMO LAGOS LISBOA

Ballot-box id		Votes								Votes minus Expected Votes								Voters					p-value
		c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-19	20-29	30-39	40-49	50+	
43M	0	20	24	64	8	28	19	3	-2	-4	-6	9	-2	8	-4	1	1	17	61	36	51	1.8e-01	
43V	0	26	26	48	8	12	23	4	-2	5	-1	1	0	-7	2	2	3	35	12	39	58	3.3e-01	
44M	2	25	31	63	9	34	30	2	0	-2	-3	-2	-3	9	1	0	3	22	57	40	74	6.8e-01	
44V	1	27	29	47	13	21	19	1	-1	4	0	-4	4	1	-3	-1	2	33	23	43	57	7.9e-01	
45M	1	25	36	55	8	22	28	1	0	-1	4	-4	-3	1	4	-1	4	20	67	34	51	9.5e-01	
45V	1	31	18	55	4	17	25	1	-1	8	-10	5	-5	-1	5	-1	2	29	44	30	46	4.0e-02	
46M	0	21	24	58	11	22	30	1	-2	-2	-5	3	1	1	5	-1	3	20	45	30	69	7.8e-01	
46V	2	23	34	67	12	13	15	4	0	-2	3	11	2	-8	-8	2	3	24	59	35	49	8.0e-02	
47M	0	20	37	48	10	32	32	0	-1	-6	5	-12	-1	10	7	-2	6	16	69	30	58	5.0e-02	
47V	1	24	32	57	5	22	30	1	-1	-1	1	0	-5	1	6	-1	2	23	57	35	55	6.9e-01	
48M	1	22	39	61	12	25	31	0	-1	-6	5	-3	0	2	5	-2	3	14	83	39	52	7.0e-01	
48V	2	22	33	55	11	13	24	0	1	0	5	3	1	-8	0	-2	2	24	40	27	67	3.9e-01	
49M	3	27	37	50	11	19	13	0	1	-2	3	-1	2	2	-4	-1	3	62	49	16	30	7.8e-01	
49V	1	18	36	67	6	21	25	3	-1	-7	5	8	-5	-1	0	1	1	19	69	27	61	3.1e-01	
50V	3	11	30	73	36	20	16	1	1	-18	-5	9	24	-2	-9	0	4	16	96	31	43	2.0e-09	
51V	3	27	32	39	8	12	14	2	1	2	3	-3	1	-4	-1	1	3	72	10	13	39	8.1e-01	
52	2	34	24	35	5	11	4	0	0	7	-6	1	1	0	-2	-1	2	100	5	3	4	7.6e-01	
53	0	31	25	37	2	14	9	1	-2	3	-6	2	-2	3	2	0	1	105	2	2	8	4.2e-01	
54	4	18	33	38	5	17	11	1	2	-7	4	-1	-1	4	-1	0	1	78	13	10	26	4.2e-01	
55	1	28	26	37	7	13	12	2	0	4	-2	-4	2	0	-1	1	37	51	15	5	18	8.5e-01	
56	1	23	34	54	5	21	19	1	0	-4	3	1	-3	4	-1	0	53	25	35	16	29	9.5e-01	
57	2	33	36	55	8	14	19	1	0	1	0	0	1	-3	1	0	58	48	29	18	14	1.0e+00	
58	4	32	26	57	13	18	20	0	2	5	-7	2	3	-2	-1	-2	2	41	44	39	44	2.6e-01	

Figure C.7: Results in the polling place of: *Escuela Gerónimo Lagos Lisboa*

REGION METROPOLITANA DE SANTIAGO - APOQUINDO - LICEO JUAN PABLO II DE LAS CONDES LOCAL: 1

Ballot-box id	Votes								Votes minus Expected Votes								Voters				p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
111V	0	0	41	95	3	57	15	1	-1	-3	11	4	0	-11	0	0	52	48	48	64	2.4e-01
112V	0	1	25	107	5	13	64	0	-1	-2	-4	12	2	-54	48	-1	43	40	55	77	<1.0e09
113V	2	1	31	78	4	62	23	0	1	-2	2	-8	1	-1	9	-2	61	31	41	68	1.4e-01
114V	4	2	32	92	3	68	7	3	3	-1	2	2	0	1	-8	1	79	10	54	68	4.0e-02
115V	1	3	31	100	1	75	17	0	0	1	-1	1	-3	4	0	-2	78	6	67	77	8.4e-01
116V	2	1	32	81	4	90	20	2	1	-2	-1	-18	0	17	3	0	90	8	67	68	1.8e-01
117V	0	2	33	90	0	88	16	0	-1	-1	0	-7	-4	15	0	-2	94	9	66	60	2.3e-01
118V	1	5	32	114	2	55	11	1	0	2	2	17	-2	-13	-6	0	58	15	57	91	1.1e-01
119V	1	6	41	95	7	58	19	2	0	3	9	-3	3	-14	2	0	77	18	57	76	5.0e-02
120V	1	1	32	93	7	62	12	0	0	-1	2	4	4	-4	-3	-2	81	8	61	58	3.3e-01
121V	0	5	32	105	4	80	15	2	-1	2	-4	3	0	3	-3	0	102	9	53	79	8.9e-01
122V	0	1	20	109	3	66	13	0	-1	-2	-11	19	-1	0	-2	-2	82	12	41	77	1.6e-01
123V	1	9	37	86	2	54	9	1	0	7	9	0	-2	-7	-6	-1	68	7	44	80	6.0e-03
124V	2	4	39	101	0	62	15	1	1	1	8	3	-4	-6	-2	-1	68	7	52	97	9.0e-02
125V	3	4	31	105	1	71	17	3	2	1	-2	3	-3	-2	0	1	74	19	59	83	2.6e-01
126V	0	1	25	122	2	64	16	1	-1	-2	-7	22	-2	-8	-1	-1	75	10	55	90	2.6e-01
127V	2	1	33	106	1	74	19	2	1	-2	2	0	-3	2	0	0	62	5	57	114	5.5e-01
128V	0	1	27	111	2	74	28	0	-1	-2	-7	5	-2	-1	10	-2	75	10	63	94	2.0e-01
129V	0	2	30	108	5	79	17	3	-1	-1	-3	1	1	4	-2	1	72	9	62	101	9.2e-01

Figure C.8: Results in the polling place of: *Liceo Juan Pablo II de Las Condes Local: 1*

REGION METROPOLITANA DE SANTIAGO - CHICUREO - ESCUELA BASICA ALGARROBAL

Ballot-box id	Votes								Votes minus Expected Votes								Voters			p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-49	50+	
10 - 0	12	37	74	4	38	6	2	0	9	21	-15	2	-17	-1	1	75	78	20	2.0e-09	
11 - 0	8	26	78	3	75	6	3	-1	5	9	-26	1	12	-2	2	107	53	39	1.6e-04	
12 - 0	5	25	126	1	77	11	0	-1	2	5	-4	-2	3	-2	-1	56	114	75	7.5e-01	
1M - 0	2	5	165	3	88	9	0	0	0	-11	9	0	12	-9	-1	14	60	198	1.5e-02	
1V - 0	0	12	168	2	78	11	0	0	-2	-3	11	-1	3	-7	-1	17	41	213	4.8e-01	
2M - 0	4	18	148	2	67	15	1	0	1	2	4	-1	-5	-1	0	15	77	163	9.0e-01	
2V - 0	4	18	128	5	60	22	2	0	1	3	-6	2	-8	7	1	20	76	151	1.0e-01	
3M - 0	7	21	136	8	68	19	1	0	4	3	-8	5	-7	3	0	22	101	138	2.0e-02	
3V - 0	6	22	156	1	61	19	1	0	3	4	8	-2	-15	2	0	18	96	152	1.9e-01	
4M - 1	2	17	152	7	75	15	2	1	-1	-2	3	4	-4	-2	1	19	116	136	3.0e-01	
4V - 0	1	10	149	2	65	18	1	0	-2	-7	13	-1	-6	3	0	17	101	128	3.8e-01	
5M - 0	3	13	141	2	107	14	0	0	-1	-7	-11	-2	25	-3	-1	14	140	126	7.0e-02	
5V - 0	4	19	147	3	84	15	3	0	1	0	-5	0	4	-2	2	22	106	147	6.8e-01	
6M - 2	2	10	140	6	88	20	0	2	-2	-10	-4	2	9	4	-1	29	135	104	2.0e-02	
6V - 1	5	22	117	2	64	18	1	1	2	4	-6	-1	-5	5	0	20	130	80	3.1e-01	
7M - 0	4	32	96	4	56	9	2	0	1	16	-12	1	-4	-3	1	11	126	66	1.0e-02	
8 - 0	14	13	76	10	34	8	1	0	11	-1	-3	7	-14	0	0	10	126	19	1.4e-07	
9 - 0	9	15	73	7	23	6	2	0	6	2	6	5	-19	-1	1	14	118	3	2.6e-05	

Figure C.9: Results in the polling place of: *Escuela Básica Algarrobal*

REGION METROPOLITANA DE SANTIAGO - EL CENTRO - COLEGIO PARTICULAR OZANAM LOCAL: 2

Ballot-box id	Votes								Votes minus Expected Votes								Voters				p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
365	5	10	132	38	6	25	17	0	-2	-3	28	-13	-6	-2	0	-2	48	102	39	44	4.0e-02
366	9	15	124	27	13	27	16	5	1	2	17	-23	1	0	-1	3	68	91	38	38	9.0e-03
367	0	2	11	31	2	12	3	1	-2	-2	-19	19	-1	5	-1	1	10	40	8	5	1.4e-07
368	0	3	10	34	1	11	3	1	-2	-1	-20	21	-2	4	-1	1	15	36	6	6	6.0e-09
369	0	5	7	25	1	14	4	1	-2	2	-18	12	-2	7	0	1	11	27	7	12	6.9e-07
370	-1	2	10	40	3	20	3	0	-1	-2	-26	23	-1	11	-3	-1	12	39	18	10	<1.0e09
371	-1	2	12	27	5	23	2	0	-1	-2	-21	12	1	15	-3	-1	13	33	17	9	2.3e-08
372	-2	1	7	5	0	6	3	0	1	-1	-5	1	-1	3	2	0	12	7	4	1	7.0e-02
373	0	2	6	9	0	21	3	0	-1	0	-11	-1	-2	16	-1	0	2	19	10	10	2.2e-07
374	-2	13	147	34	12	25	19	1	-6	0	37	-22	-1	-6	-1	-1	63	86	46	57	3.0e-04
375	-13	13	130	38	12	16	24	2	5	0	21	-16	-1	-14	5	0	64	86	47	51	5.0e-03
376	-12	3	124	51	9	23	14	3	5	-10	19	-2	-4	-5	-4	1	60	86	42	51	1.1e-02
377	-8	12	123	43	14	26	15	1	0	-1	15	-9	1	-2	-3	-1	63	92	43	45	6.8e-01
378	-5	3	136	45	8	22	9	0	-2	-9	37	-6	-4	-5	-9	-2	50	85	41	52	1.1e-04
379	-8	12	137	40	12	19	14	2	0	-1	29	-13	-1	-10	-4	0	69	87	40	48	3.0e-02
380	-7	9	148	45	15	29	15	0	-1	-6	30	-14	1	-3	-5	-2	70	98	44	56	1.6e-02

Figure C.10: Results in the polling place of: *Colegio Particular Ozanam Local: 2*

REGION METROPOLITANA DE SANTIAGO - EL CENTRO - LICEO JAVIERA CARRERA LOCAL: 2

Ballot-box id	Votes								Votes minus Expected Votes								Voters				p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
18M	-3	4	31	22	5	20	12	2	1	1	1	-9	0	5	0	1	12	6	10	71	4.3e-01
19M	-1	2	26	31	3	23	11	1	-1	-1	-5	1	-2	8	0	0	15	11	5	67	4.8e-01
20M	-1	2	27	16	4	10	11	1	-1	0	3	-5	0	0	3	0	10	10	6	46	7.3e-01
21M	-3	1	21	26	3	19	14	0	1	-2	-6	0	-2	6	4	-1	12	5	10	60	2.5e-01
22M	-1	0	31	34	2	11	13	0	-1	-3	3	6	-3	-3	2	-1	8	10	9	65	3.6e-01
23M	0	1	22	21	5	20	9	2	-2	-2	-5	-2	1	8	1	1	15	10	6	49	5.0e-02
24M	-2	1	44	25	8	12	9	0	0	-2	12	-6	3	-3	-3	-1	13	11	5	72	2.3e-01
25M	-2	1	28	27	4	16	10	0	0	-2	1	0	-1	3	0	-1	8	8	8	64	9.5e-01
26M	-1	3	29	27	6	16	7	1	-1	0	0	0	1	3	-3	0	12	7	11	60	9.3e-01
27M	-2	1	26	34	4	18	12	1	0	-2	-4	4	-1	3	0	0	14	7	5	72	9.3e-01
28M	-1	4	22	25	8	10	9	0	-1	2	-1	0	4	-2	-1	-1	6	7	5	61	4.3e-01
29M	-5	3	17	29	5	23	8	0	3	0	-11	1	1	9	-2	-1	10	10	6	64	2.0e-02
30M	-1	4	32	15	3	14	11	2	-1	1	5	-9	-1	2	2	1	13	6	9	54	2.1e-01
31M	0	2	35	32	3	12	11	95	-4	-4	-23	-27	-7	-17	-11	93	11	9	6	69	<1.0e09
32M	0	4	31	32	5	14	13	0	-2	1	-2	3	0	-1	2	-1	13	15	7	64	8.4e-01
33M	-1	4	29	28	5	12	9	0	-1	1	2	1	0	-1	-1	-1	10	9	5	64	9.5e-01
34M	-2	4	27	30	3	16	15	0	0	1	-5	1	-2	2	4	-1	9	16	9	63	7.6e-01

Figure C.11: Results in the polling place of: *Liceo Javiera Carrera Local: 2*

REGION METROPOLITANA DE SANTIAGO - LO BARNECHEA - COLEGIO DIFERENCIAL MADRE TIERRA

Ballot-box id		Votes								Votes minus Expected Votes								Voters						p-value
		c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-19	20-29	30-39	40-49	50-59	60+	
82	-1	19	31	65	14	65	5	0	0	16	17	-39	12	-7	2	-1	8	72	101	5	9	5	<1.0e-09	
83	-0	4	12	106	1	107	11	1	-1	-2	-11	-10	-3	24	3	0	4	182	11	21	11	13	2.0e-02	
84	-1	18	29	74	11	49	10	0	0	13	12	-18	8	-19	5	-1	10	138	21	7	9	7	<1.0e-09	
85	-1	15	41	61	7	45	7	0	0	10	25	-23	5	-18	2	-1	6	146	6	10	7	2	<1.0e-09	
86	-0	0	7	143	0	102	5	2	-1	-7	-18	22	-4	11	-3	0	7	215	7	10	9	11	1.0e-06	
87	-0	1	21	134	1	86	8	0	-1	-6	-3	17	-3	-2	0	-2	6	202	10	8	13	12	7.0e-02	
88	-2	17	35	52	9	50	7	0	1	12	19	-29	7	-11	2	-1	3	149	3	9	6	2	<1.0e-09	
89	-1	17	27	81	5	64	6	2	0	12	8	-15	2	-9	1	1	4	166	19	9	4	2	1.5e-04	
90	-0	0	15	126	1	108	5	0	-1	-7	-10	7	-3	19	-3	-2	5	215	5	11	9	10	2.0e-03	
91	-0	2	12	143	1	111	4	1	-1	-6	-14	15	-3	14	-4	-1	8	237	3	11	10	5	2.0e-03	
92	-0	14	33	56	6	56	3	2	-1	9	17	-24	4	-4	-2	1	6	145	3	5	9	2	3.3e-07	
93	-1	10	36	80	6	55	5	0	0	5	18	-11	3	-13	-1	-1	5	160	4	9	9	6	5.0e-04	
94	-0	1	9	131	0	114	7	0	-1	-6	-16	8	-4	21	0	-2	8	224	13	6	4	7	1.1e-04	
95	-2	8	32	68	7	40	5	0	1	4	17	-11	4	-13	-1	-1	9	89	11	14	18	21	5.0e-05	

Figure C.12: Results in the polling place of: *Colegio Diferencial Madre Tierra*

REGION METROPOLITANA DE SANTIAGO - LO BARNECHEA - COLEGIO LOS ALERCES LOCAL: 2

Ballot-box id		Votes								Votes minus Expected Votes								Voters						p-value
		c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-19	20-29	30-39	40-49	50-59	60+	
17M	-2	5	25	78	8	48	13	2	1	0	6	-11	1	2	0	1	3	21	5	31	41	80	4.0e-01	
18M	-1	8	29	76	16	52	11	2	0	3	10	-23	10	1	-2	1	4	18	4	44	54	71	4.0e-04	
19M	-2	6	30	92	7	43	14	0	1	1	11	-6	0	-7	1	-1	2	16	7	44	48	77	1.8e-01	
20M	-3	6	15	81	7	44	14	2	2	1	-4	-3	0	0	3	1	6	28	2	34	27	75	3.2e-01	
21M	-0	5	9	83	10	56	7	1	-1	0	-11	1	3	13	-5	0	8	14	13	20	28	88	4.0e-02	
22M	-0	7	19	81	11	68	9	1	-1	2	-1	-15	4	14	-3	0	2	21	44	16	35	78	1.3e-01	
23M	-1	8	30	90	7	49	9	3	0	3	11	-7	1	-7	-3	2	6	21	39	11	50	69	8.0e-02	
24M	-2	13	34	66	12	44	10	3	1	9	18	-27	7	-10	0	2	5	19	51	13	41	55	3.0e-09	
25M	-0	1	23	97	11	50	13	0	-1	-4	4	2	4	-5	1	-1	10	25	37	8	43	72	3.0e-01	
26M	-1	4	25	103	11	52	11	1	0	-1	5	-1	5	-8	0	0	4	26	56	16	36	70	5.3e-01	
27M	-1	13	27	79	8	39	4	0	0	9	12	-7	3	-11	-5	-1	9	25	36	13	39	49	2.6e-05	
28M	-2	11	23	90	5	49	12	1	1	7	5	-6	0	-9	2	0	5	29	51	9	42	57	7.0e-02	
29M	-0	7	29	87	10	57	8	1	-1	2	10	-11	4	-2	-2	0	2	35	55	7	39	62	1.0e-01	
30M	-1	12	25	65	8	50	10	3	0	8	9	-22	3	-1	1	2	6	20	52	7	36	53	5.2e-05	
31M	-3	5	26	79	10	54	14	0	2	1	10	-17	5	-4	4	-1	5	28	48	9	55	46	2.0e-03	
32M	-3	12	28	79	5	58	5	0	2	9	13	-19	1	-1	-4	-1	8	28	55	16	48	35	7.5e-06	
33M	-1	8	30	93	4	63	7	1	0	4	14	-18	-1	4	-3	0	4	25	37	65	27	49	2.0e-02	

Figure C.13: Results in the polling place of: *Colegio Los Alerces Local: 2*

REGION METROPOLITANA DE SANTIAGO - LO BARNECHEA - THE MAYFLOWER SCHOOL LOCAL: 2

Ballot-box id	Votes								Votes minus Expected Votes								Voters						p-value	
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-19	20-29	30-39	40-49	50-59	60+		
175	-	0	3	14	141	1	90	16	0	-1	-2	-3	-4	-3	8	6	-1	8	45	62	80	43	28	4.8e-01
176	-	0	2	13	161	3	75	12	1	-1	-3	-5	16	-1	-7	1	0	7	41	61	74	49	35	5.3e-01
177	-	0	1	13	139	2	103	7	3	-1	-4	-5	-6	-2	20	-4	2	6	39	68	68	56	31	3.0e-02
178	-	0	3	11	156	0	85	9	1	-1	-2	-7	12	-4	4	-2	0	6	43	59	82	43	32	2.1e-01
179	-	0	3	7	146	3	102	12	1	-1	-2	-10	-4	-1	17	1	0	9	43	65	82	49	27	8.0e-02
180	-	1	2	14	146	3	97	9	1	0	-3	-4	-2	-1	12	-2	0	8	42	67	73	49	34	5.3e-01
70M	-	0	3	15	140	1	99	6	0	-1	-2	-7	3	-5	19	-6	-1	11	27	83	35	44	64	6.0e-02
71M	-	1	2	9	161	1	85	7	0	0	-3	-11	21	-4	3	-5	-1	5	27	87	40	52	55	1.0e-02
72M	-	0	0	11	153	4	100	4	0	-1	-5	-9	10	-1	12	-5	-1	10	37	118	27	37	42	2.0e-02
73M	-	0	1	14	141	6	105	8	2	-1	-3	-3	-8	3	11	0	1	7	39	133	28	49	21	1.5e-01
74M	-	1	1	14	135	2	99	10	0	0	-4	-4	-4	-2	14	1	-1	13	27	112	27	43	40	3.9e-01
75M	-	1	2	11	130	4	112	6	2	0	-2	-7	-13	0	24	-3	1	10	32	118	27	50	30	5.0e-02
76M	-	0	2	16	143	4	89	7	1	-1	-2	-2	3	0	5	-3	0	10	29	95	34	59	36	8.5e-01
77M	-	1	1	8	170	8	74	4	1	0	-4	-10	27	4	-10	-7	0	10	26	92	46	55	38	9.0e-04
78M	-	0	2	17	148	1	88	11	0	-1	-2	-1	4	-3	3	1	-1	8	32	93	48	52	34	7.9e-01
79M	-	0	1	18	133	0	83	11	1	-1	-3	2	-1	-3	4	2	0	7	24	99	48	35	34	2.8e-01
80M	-	1	5	28	119	5	75	13	1	0	0	9	-10	1	-3	3	0	14	49	66	33	46	39	4.3e-01
81M	-	1	12	33	59	8	46	9	3	0	8	18	-24	5	-12	3	2	6	114	12	14	16	9	2.0e-09

Figure C.14: Results in the polling place of: *The Mayflower School Local:*
2

REGION METROPOLITANA DE SANTIAGO - LO BARNECHEA - THE NEWLAND SCHOOL LOCAL: 2

Ballot-box id	Votes								Votes minus Expected Votes								Voters						p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-19	20-29	30-39	40-49	50-59	60+	
24V - 3	6	29	54	11	37	19	0	2	3	14	-25	6	-10	11	-1	7	14	55	4	24	55	2.8e-08	
25V - 0	3	15	81	7	55	10	0	-1	0	2	-8	4	1	3	-1	5	22	64	22	25	33	4.0e-01	
26V - 2	9	25	72	9	37	13	0	1	5	9	-11	4	-11	4	-1	3	17	46	7	36	59	1.0e-02	
27V - 0	6	36	82	7	34	11	1	-1	2	19	-6	1	-16	1	0	6	18	46	13	28	67	1.1e-03	
28V - 1	14	25	78	6	39	16	1	0	11	12	-16	3	-17	7	0	9	10	54	10	66	31	9.0e-09	
29V - 3	4	27	94	9	67	11	0	3	0	12	-22	6	0	2	-1	11	25	62	46	39	32	2.9e-04	
30V - 1	8	31	101	7	60	18	1	0	4	14	-21	2	-6	7	0	6	18	54	66	30	53	1.4e-03	
31V - 1	4	19	110	7	66	13	1	0	-1	-1	-2	0	3	1	0	3	27	43	38	39	71	9.9e-01	
32V - 1	6	31	104	13	57	12	0	0	1	10	-10	6	-6	0	-1	4	22	57	35	27	79	9.0e-02	
33V - 3	6	27	109	10	41	12	1	2	2	10	-2	5	-19	2	0	7	22	44	56	27	52	1.1e-03	
34V - 0	9	24	107	6	69	17	2	-1	5	7	-20	2	0	6	1	9	30	46	73	33	43	3.0e-02	
35V - 2	8	29	97	6	62	11	1	1	4	12	-19	2	-1	1	0	7	28	38	66	31	46	1.0e-02	
36V - 0	2	32	100	9	59	10	0	-1	-3	14	-12	4	-1	0	-1	2	25	46	61	22	56	4.0e-02	
37V - 0	5	26	102	8	57	7	5	-1	0	7	-6	2	-3	-3	4	6	36	43	44	20	61	6.0e-02	
38V - 1	16	21	108	10	67	10	1	0	11	2	-17	4	1	-1	0	4	22	48	78	17	65	1.8e-03	
39V - 1	7	22	116	7	48	15	1	0	2	5	0	2	-14	5	0	7	23	46	64	23	54	1.1e-01	
40V - 2	6	17	121	4	71	12	2	1	0	-5	0	-3	7	-1	1	8	26	29	59	33	80	5.5e-01	
41V - 2	2	24	123	5	55	10	0	1	-2	7	3	0	-7	-1	-1	5	23	34	85	22	52	3.7e-01	
42V - 0	10	28	129	9	60	12	1	-1	5	9	-6	3	-10	0	0	1	32	33	101	25	57	6.0e-02	
43V - 0	6	23	141	5	57	9	1	-1	1	6	7	0	-10	-3	0	4	28	25	113	22	48	5.6e-01	
44V - 0	4	8	170	0	63	9	1	-1	-1	-10	27	-5	-7	-3	0	5	24	20	134	23	49	1.0e-02	
45V - 1	3	8	144	2	65	6	0	1	-1	-8	17	-2	-1	-5	-1	3	32	21	99	38	36	1.1e-01	

Figure C.15: Results in the polling place of: *The Newland School Local: 2*

REGION METROPOLITANA DE SANTIAGO - PARQUE ALMAGRO - COLEGIO EXCELSIOR LOCAL: 1

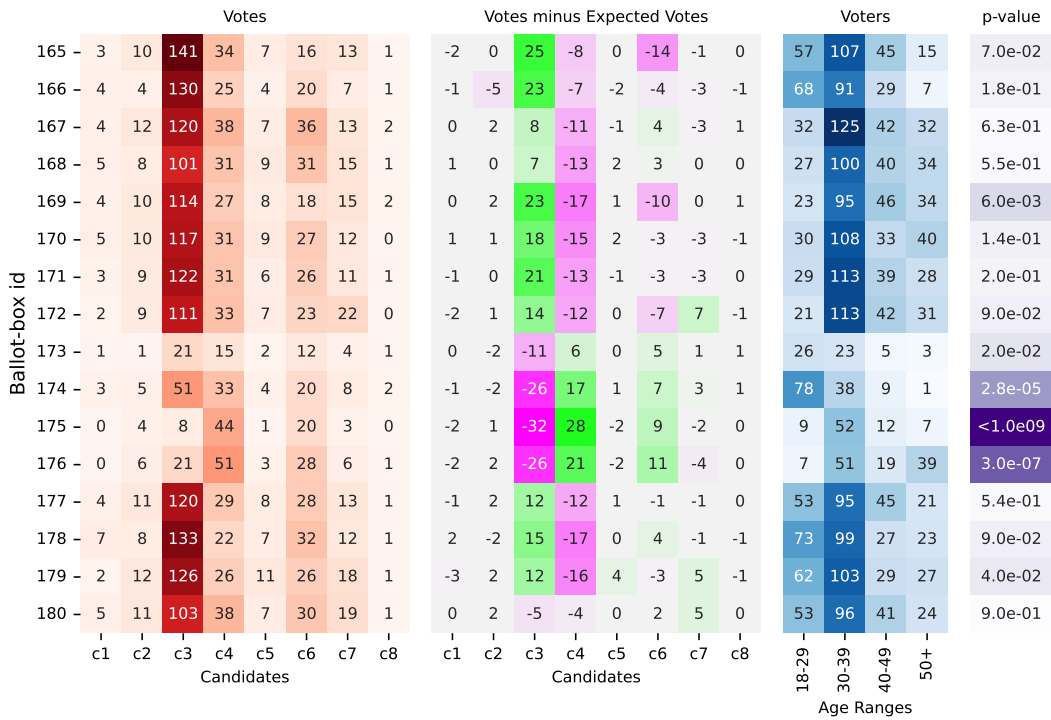


Figure C.16: Results in the polling place of: *Colegio Excelsior Local: 1*

REGION METROPOLITANA DE SANTIAGO - PARQUE ALMAGRO - COLEGIO EXCELSIOR LOCAL: 2

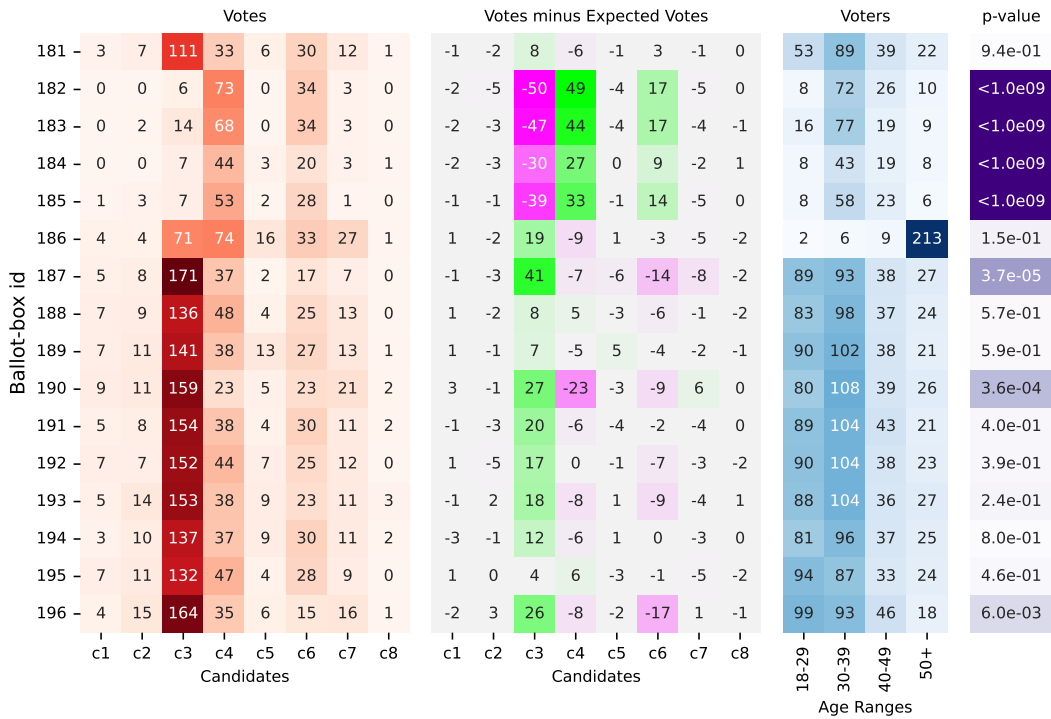


Figure C.17: Results in the polling place of: *Colegio Excelsior Local: 2*

REGION METROPOLITANA DE SANTIAGO - PARQUE ALMAGRO - LICEO DARIO SALAS

Ballot-box id	Votes								Votes minus Expected Votes								Voters				p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50+	
39M	2	19	71	40	4	28	11	2	-1	12	-4	-4	-3	3	-4	1	19	83	10	65	9.0e-03
40M	5	6	48	27	4	25	21	0	3	1	-7	-8	-2	5	9	-1	13	57	6	59	5.0e-02
41M	3	7	67	40	9	20	3	2	0	1	3	3	3	-2	-9	1	14	73	9	55	7.0e-02
42M	2	12	66	27	9	14	9	12	-1	6	0	-9	3	-7	-3	11	16	76	11	48	9.0e-09
43M	1	6	67	35	5	20	15	0	-2	0	5	-3	-1	-1	3	-1	14	69	7	59	9.5e-01
44M	4	11	66	36	15	22	15	0	1	4	-5	-6	8	-2	1	-1	16	80	12	61	4.0e-02
45M	3	4	70	39	9	22	13	1	0	-2	6	-3	2	-1	-2	0	18	63	9	72	9.4e-01
46M	4	8	72	40	6	22	15	1	1	1	2	-1	-1	-1	0	-1	42	43	8	75	9.8e-01
47M	5	6	64	33	12	25	13	0	1	-1	-8	-1	6	5	0	-2	65	22	12	59	2.5e-01
48M	7	13	66	37	6	21	9	2	3	6	-10	3	0	1	-3	0	72	22	8	59	1.4e-01
49M	5	15	72	35	9	22	13	2	1	7	-13	1	2	2	0	0	87	23	5	58	2.0e-01
50M	1	5	47	31	12	19	13	0	-2	0	-5	-1	6	2	1	-1	34	23	9	62	3.4e-01
51M	3	5	64	52	4	28	15	1	0	-1	5	2	-5	2	-3	0	12	38	36	86	7.7e-01
52M	2	7	57	47	11	33	21	2	-1	1	5	-11	1	6	-1	0	15	20	10	135	6.6e-01
53M	1	6	67	48	9	21	18	0	-2	0	8	-1	1	-4	0	-2	16	31	43	80	8.5e-01
54M	2	6	74	41	10	28	16	2	-1	-1	0	-4	3	2	0	1	23	63	31	61	9.5e-01

Figure C.18: Results in the polling place of: *Liceo Darío Salas*

REGION METROPOLITANA DE SANTIAGO - PUDAHUEL - ESCUELA MONSEÑOR CARLOS OVIEDO

Ballot-box id	Votes								Votes minus Expected Votes								Voters					p-value
	c1	c2	c3	c4	c5	c6	c7	c8	c1	c2	c3	c4	c5	c6	c7	c8	18-29	30-39	40-49	50-59	60+	
10M	3	8	38	22	17	18	16	1	0	0	5	-1	-2	1	-2	-1	17	4	1	30	71	9.6e-01
11M	2	9	28	20	17	13	14	2	-1	1	-1	0	1	0	0	0	20	3	4	23	55	1.0e+00
12M	0	14	47	23	17	13	16	2	-3	4	11	-3	-3	-4	-2	0	18	5	13	33	63	1.5e-01
13M	3	3	32	19	24	18	16	2	0	-5	1	-4	6	3	-1	0	16	3	7	29	62	3.8e-01
14M	5	9	27	32	19	17	17	2	2	0	-7	7	-1	0	-1	0	13	5	6	36	69	6.6e-01
15M	1	9	30	19	17	10	17	3	-2	1	1	-1	1	-4	3	1	17	4	6	27	52	7.5e-01
16M	3	13	32	28	21	8	24	1	0	3	-4	3	1	-8	6	-1	18	7	5	41	59	2.0e-01
17M	5	12	33	32	20	21	18	3	1	2	-6	4	-3	3	-2	1	20	2	10	48	63	6.9e-01
18M	2	10	40	28	19	19	21	1	-1	-1	1	1	-3	2	2	-1	21	6	15	38	61	9.9e-01
19M	3	13	29	24	12	19	17	1	0	4	-4	1	-6	5	1	-1	15	5	13	41	44	4.6e-01
1M	5	5	38	24	16	16	15	2	1	-7	1	3	-1	2	1	0	41	3	2	18	57	5.2e-01
20M	3	4	30	21	24	14	13	2	0	-4	-1	0	7	0	-2	0	17	4	6	34	49	4.6e-01
21M	4	12	40	26	10	18	13	1	1	2	4	3	-9	3	-3	-1	28	5	1	40	49	3.3e-01
22M	1	12	32	27	28	13	16	1	-2	4	-3	0	7	-3	-2	-1	8	5	4	66	47	5.0e-01
23M	5	5	24	22	16	14	20	1	2	-2	-5	1	-1	1	5	-1	11	3	0	48	45	5.3e-01
24M	7	9	35	25	21	21	11	3	4	0	-1	-1	0	4	-7	1	15	3	8	52	55	1.9e-01
2M	3	8	26	16	16	22	14	3	0	-1	-4	-3	0	8	-1	1	25	3	3	10	66	5.3e-01
3M	4	10	37	26	17	17	11	1	1	1	3	3	-2	1	-6	-1	24	3	5	16	75	7.9e-01
4M	4	6	31	20	17	17	20	3	1	-4	-2	-2	-1	2	5	1	26	4	4	25	60	6.5e-01
5M	6	10	38	20	21	21	17	1	2	0	1	-5	1	3	-1	-1	29	2	1	25	77	7.6e-01
6M	20	6	27	28	2	23	17	1	17	-2	-4	4	-17	5	-2	-1	13	3	1	22	85	<1.0e09
7M	4	12	46	20	18	14	12	2	0	1	9	-3	-1	-2	-4	0	32	6	3	22	66	7.2e-01
8M	5	6	32	25	20	22	21	3	1	-4	-4	0	0	4	2	1	25	1	1	22	85	7.1e-01
9M	3	8	26	25	22	19	10	5	0	0	-4	2	4	3	-8	3	15	0	5	21	77	1.6e-01

Figure C.19: Results in the polling place of: *Escuela Monseñor Carlos Oviedo*