

ELECTIONS ANALYSIS & PREDICTIVE MODELING

ELECTION
PREDICTIVE MODELING

59%

28% 12%

Prediction Model of Electoral Results in Mexico:

A Case Study using PREP, INEGI and INE Data

Data Science Intensive Capstone Project

August 3rd, 2024

By Javier Jorge Pérez Ontiveros

Abstract

This project uses the PREP, INEGI and INE databases, this provides detailed data on previous elections, including voter demographics, polling station information, economic variables, etc.

Analyzing this data helps identify hidden patterns and trends, which can enhance campaign strategies, policymaking, and voter engagement efforts.

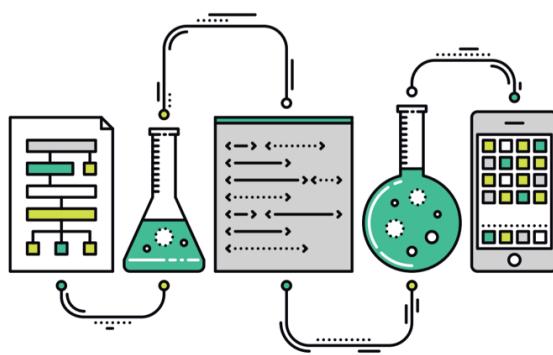
The goal is to understand why the winning party (MORENA) achieved a 59% acceptance rate, while the opposition (PAN-PRI-PRD) secured only 28% of the votes.

Objectives:

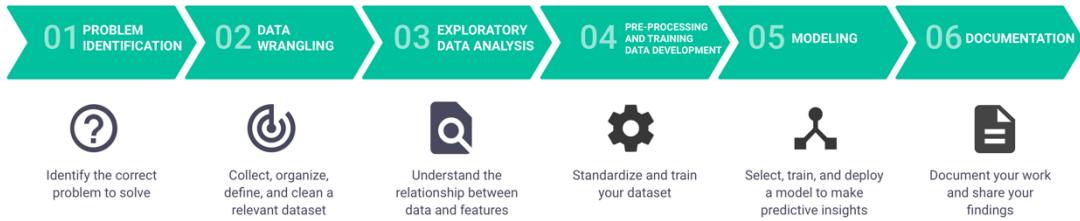
- 1) To generate a Model that predicts votes for MORENA
- 2) To provide actionable insights to support the Opposition Party in developing a strategy to improve their results in the 2030 elections.

Results:

- Geographic heatmap and interesting insights
- 2 Main Voters Profiles were discovered with unsupervised learning (k-Means Clustering).
- A final Machine Learning Model with 98% accuracy was developed.



Methodology followed



● Problem definition

Elections are a critical component of democratic societies, and understanding the factors that influence voting behavior is essential for all stakeholders involved. The PREP (INE) database provides comprehensive data on past elections, including voter demographics, polling station information, and real-time results. Analyzing this data can uncover patterns and trends that are not immediately apparent, helping to improve campaign strategies, policymaking, and voter engagement efforts.

Criteria for Success

- Identify key factors that influence voter turnout, while generating valuable insights to optimize Opposition's campaign strategies.
- Develop an accurate predictive model with an R² score greater than 70%.
- Deliver findings and actionable recommendations in a detailed final report and presentation.

Scope of Solution Space

- Data Analysis and Clustering: Use clustering techniques, such as K-Means, to segment voter groups based on demographics, voting history, and other relevant factors during the exploratory data analysis (EDA) phase.
- Data Quality and Validation: Ensure rigorous data cleansing and validation processes to maintain the accuracy and reliability of the models.
- Predictive Modeling: Develop a state-of-the-art predictive model to accurately forecast election results and voter turnout.

- Reporting and Visualization: Summarize key findings and recommendations in a final report and slide deck, including visualizations that highlight important trends and insights.

Data Sources

- PREP Database: The primary data source, providing detailed election data, including voter demographics, polling station information, and real-time results.
 - Census Data: Data from INEGI (Instituto Nacional de Estadística y Geografía) can provide detailed demographic information such as age, gender, education level, income, and employment status. This data can help identify demographic trends and correlations with voting patterns.
 - INE Database: Data from INE (Instituto Nacional Electoral) to show the budget spent per political party, the public events and marketing strategies.

Acquiring, Cleaning and Merging

<u>PREP</u> 2024 and 2018 election results exported from database to CSV file. 55,976,881 votes 170,944 polling stations 1,580 municipalities 32 states	<u>INEGI</u> Socioeconomic inputs for the 32 states in csv: Hospitals, Schooling Years, vehicles purchased, catholic believers, murders per year, population ages, Poverty ratios, Average Income, etc.	<u>INE</u> Performance and Behavior of Parties per state in csv: # of Campaign events, Marketing Budget, Federal Welfare,
--	---	---

Final Dataframe:

A 40 columns x 32 rows Data frame was created to predict the votes for MORENA Entity.

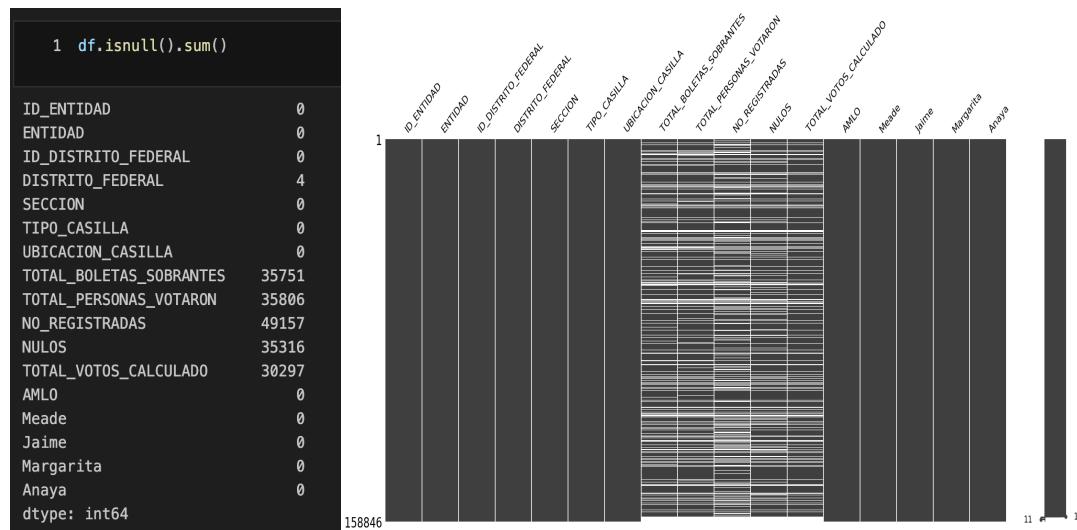
● Data Engineering / Wrangling

This step involves gathering the data, organizing it, and ensuring it is well-defined.

Careful attention to these tasks will yield significant benefits later.

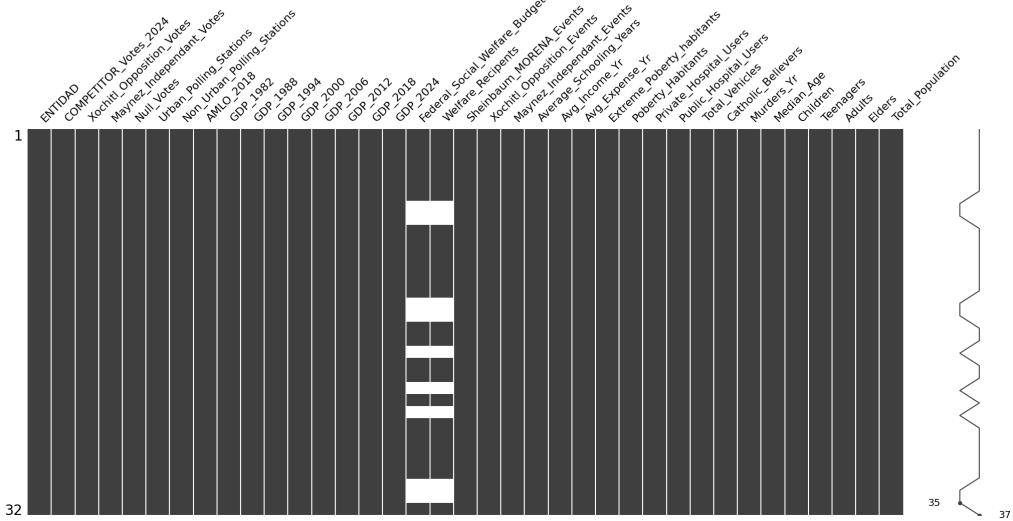
Some initial data cleaning might be performed at this stage, while ensuring that the data frame does not have Null values. Also, during this first stage, we explore the data to gain a better understanding of it.

From Prep Data source up to 49,157 null values were identified and classified as “Ceros”, as basically this values represented the documents of population with null votes.



Once that a single data frame was merged from the 3 data sources, it was identified that Federal Social Welfare and Welfare Recipients Features had null values and in the exact same rows.

I realized that the null values of the merged data frame, were caused due to incompatibility of formats in the foreign keys of the left join due to typographic differences ('apostrophes').



After all the cleansing and wrangling, the final data frame had 32 rows and 37 columns, all with numeric dtypes, except for the entities.

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32 entries, 0 to 31
Data columns (total 37 columns):
 #   Column           Non-Null Count Dtype  
--- 
 0   ENTIDAD          32 non-null    object  
 1   COMPETITOR_Votes_2024 32 non-null    int64  
 2   Xochitl_Opposition_Votes 32 non-null    int64  
 3   Maynez_Independant_Votes 32 non-null    int64  
 4   Null_Votes         32 non-null    int64  
 5   Urban_Polling_Stations 32 non-null    int64  
 6   Non_Urban_Polling_Stations 32 non-null    int64  
 7   AMLO_2018          32 non-null    int64  
 8   GDP_1982           32 non-null    int64  
 9   GDP_1988           32 non-null    int64  
 10  GDP_1994          32 non-null    int64  
 11  GDP_2000          32 non-null    int64  
 12  GDP_2006          32 non-null    int64  
 13  GDP_2012          32 non-null    int64  
 14  GDP_2018          32 non-null    int64  
 15  GDP_2024          32 non-null    int64  
 16  Sheinbaum_MORENA_Events 32 non-null    int64  
 17  Xochitl_Opposition_Events 32 non-null    int64  
 18  Maynez_Independant_Events 32 non-null    int64  
 19  Average_Schooling_Years 32 non-null    int64  
 20  Avg_Income_Yr       32 non-null    int64  
 21  Avg_Expense_Yr      32 non-null    int64  
 22  Extreme_Poverty_habitants 32 non-null    int64  
 23  Poverty_Habitants   32 non-null    int64  
 24  Private_Hospital_Users 32 non-null    int64  
 25  Public_Hospital_Users 32 non-null    int64  
 26  Total_Vehicles      32 non-null    int64  
 27  Catholic_Believers   32 non-null    int64  
 28  Murders_Yr          32 non-null    int64  
 29  Median_Age          32 non-null    int64  
 30  Children             32 non-null    int64  
 31  Teenagers            32 non-null    int64  
 32  Adults               32 non-null    int64  
 33  Elders               32 non-null    int64  
 34  Total_Population     32 non-null    int64  
 35  Federal_Social_Welfare_Budget_2024 32 non-null    int64  
 36  Welfare_Recipients   32 non-null    int64  
dtypes: int64(36), object(1)

```

● Exploratory Analysis (EDA) and Visualization

The Exploratory Data Analysis (EDA) phase is a crucial step in our Data Science project. In this phase, we delve into the dataset to uncover patterns, spot anomalies, test hypotheses, and check assumptions with the help of summary statistics and graphical representations.

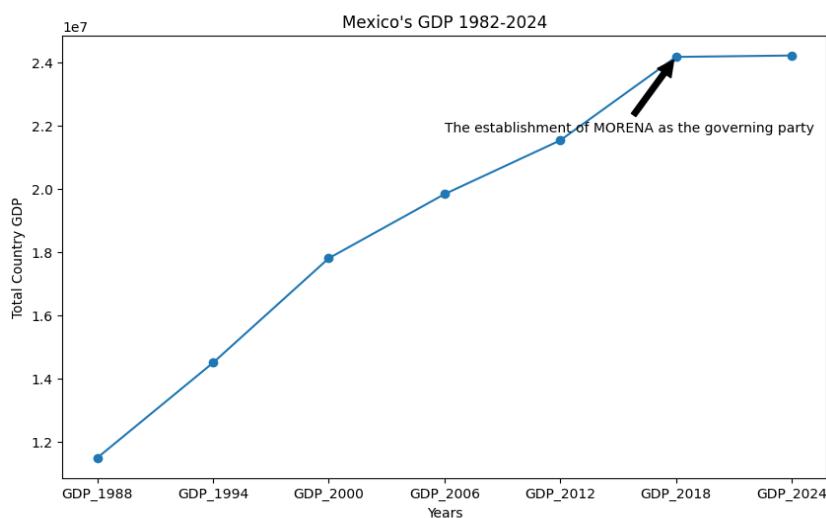
The primary goal of EDA is to gain a deeper understanding of the data, which will guide us in making informed decisions during the modeling phase. This involves identifying key features that influence voter decisions and understanding the underlying structure of the data. By doing so, we can develop strategies for the opposition (PAN-PRI-PRD) to enhance their chances in upcoming election cycles.

By thoroughly exploring the data, we aim to build a robust predictive model that can provide actionable insights for the opposition, enabling them to develop effective campaign strategies and allocate resources efficiently for future elections.

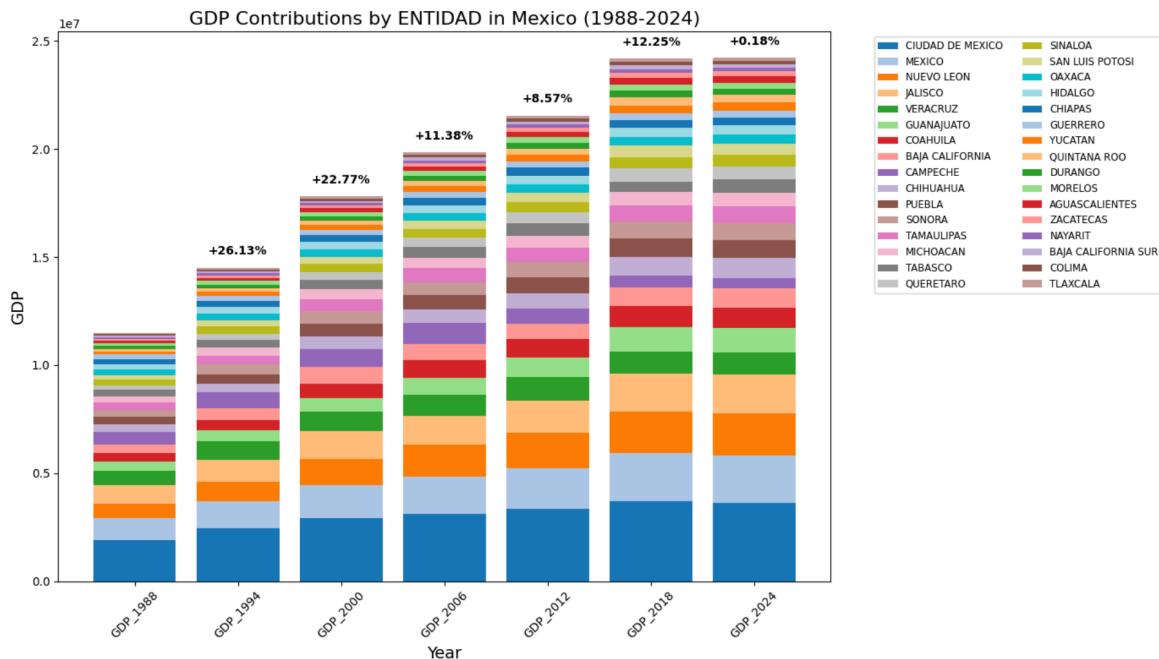
During EDA, I uncovered 4 interesting insights:

1) Economic Slowdown

The data shows an Economic Slowdown During the Period When MORENA Assumed Power



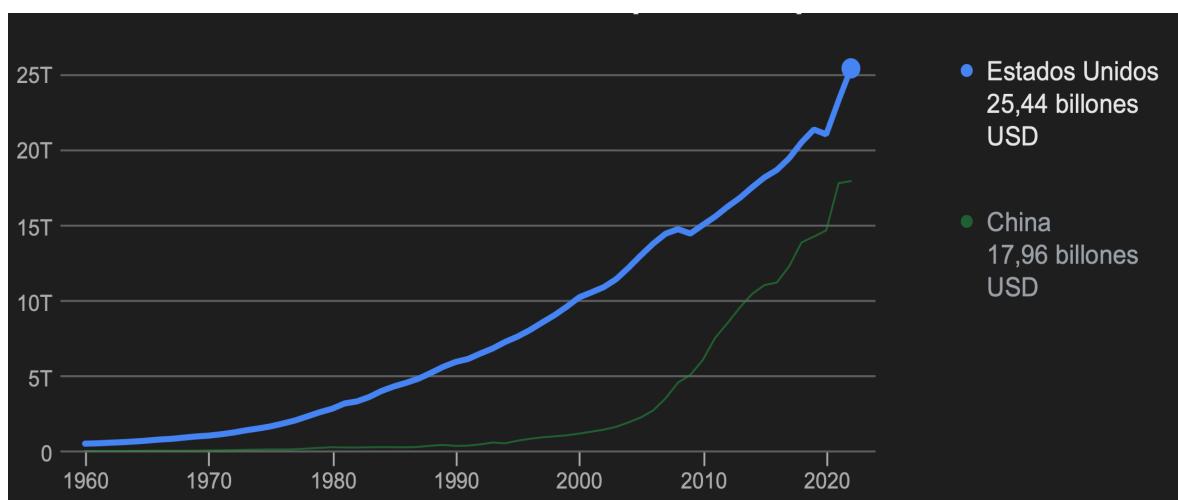
Same data can be split by Entity as follows:



As it can be observed, GDP Growth from 2018 to 2024 was only 0.18%, while the Opposition parties had a 10%-20% growth from 1988 to 2018.

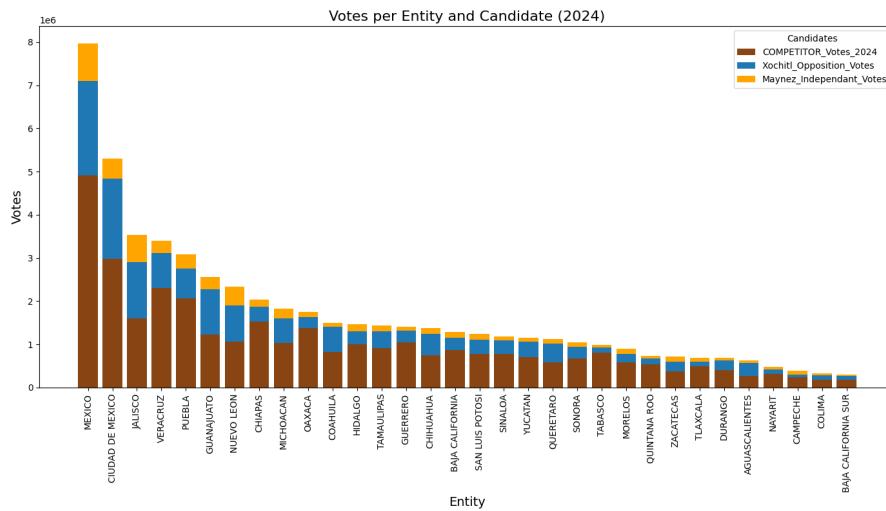
When I shared this insight with some political Experts, the first question that popped to their mind was: “Was this effect due to Covid 19?”.

For that reason, I compared the performance of the top 2 economic powers to review if those countries had a similar “slow down” or “stagnation”. But I observed that they didn’t have a similar effect from 2018-2024, so something else is causing the slowdown in México.



2) Effect of Independent Party “Movimiento Naranja”.

Contrary to what many expected, that the Independent Candidate Maynez would weaken the opposition, this effect can only be observed in states where the “Movimiento Naranja” party currently holds power, such as Nuevo León and Jalisco.”



3) Geographical Heatmap:

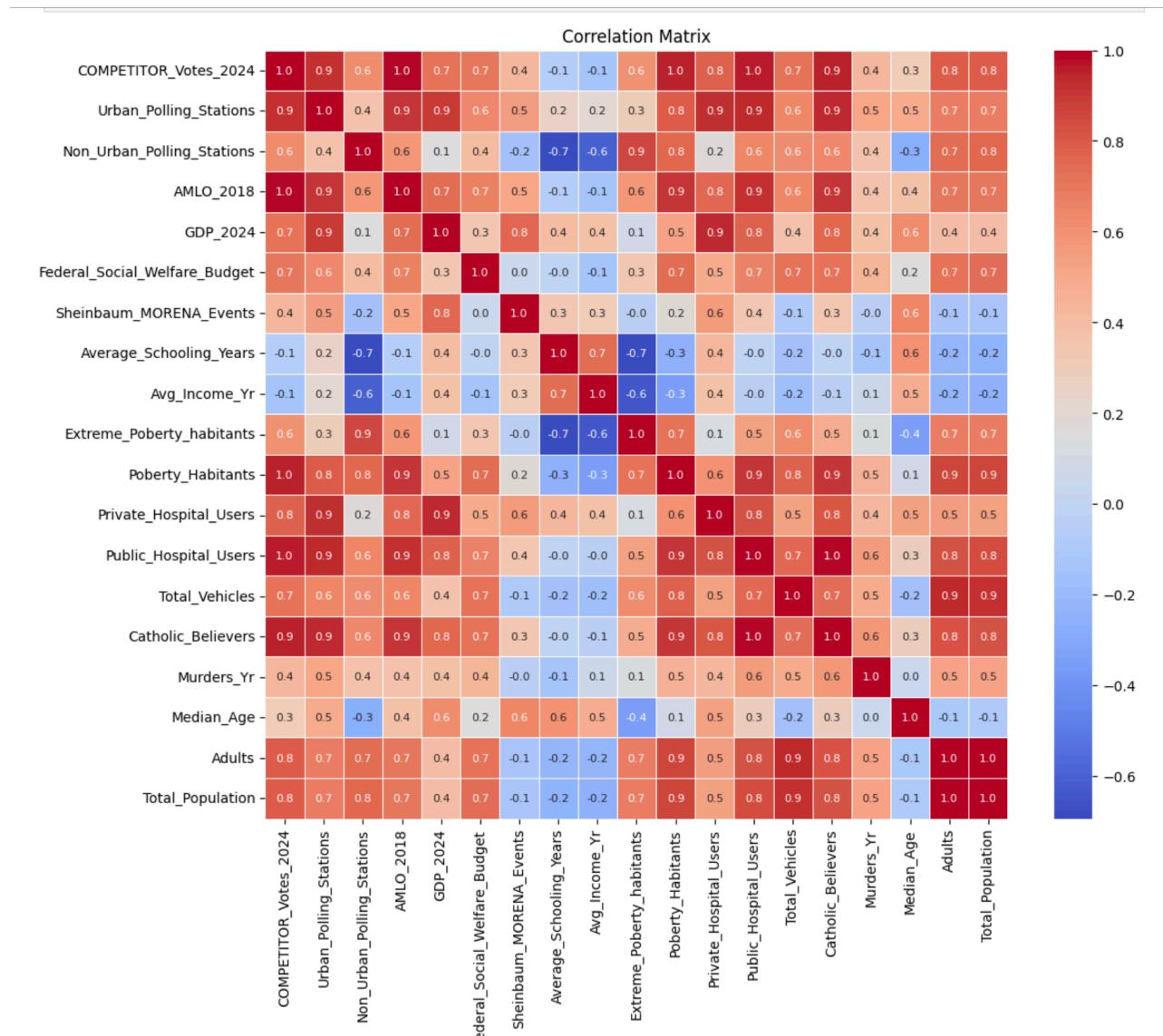
The strongest resistance to MORENA is observed in the Bajío region, known for its industrial and economic strength, while the strongest support for MORENA is in the less industrially and economically developed Southeast region. To have a chance of winning from 2024-2030, the opposition must focus on helping the less privileged population.



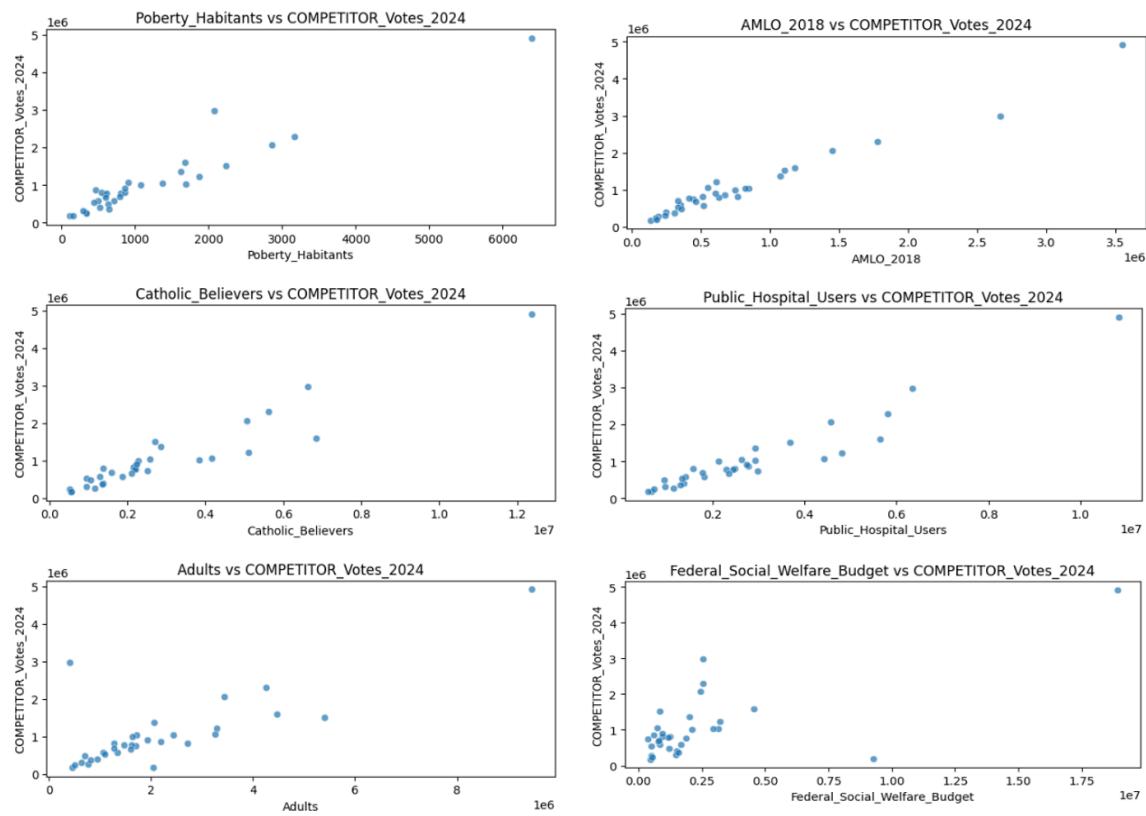
4) Features Correlated to Votes:

The following variables showed a positive correlation with MORENA votes:

- Preference for AMLO in 2018
- Federal Warefare budget
- Poberty Habitants
- Public Hospitals
- Catholic Believers
- Adults Amount



Their Linear correlations show that the votes for MORENA increase when these variables increase:

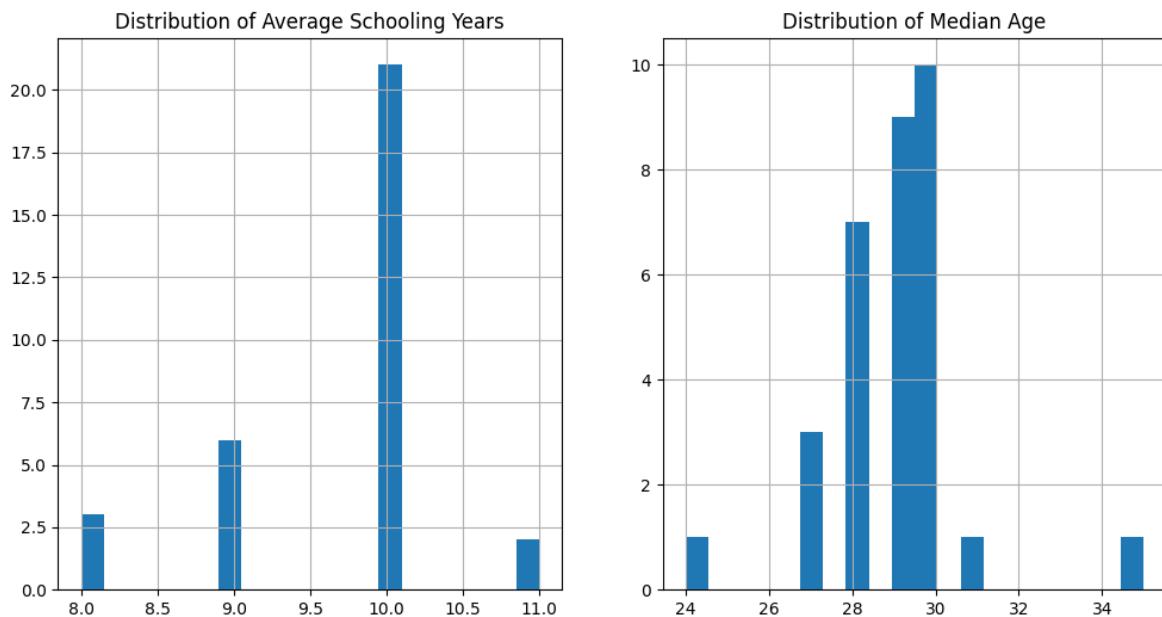


● Feature Engineering and Clusterization

2 features have a behavior of categorical variables (despite they are numerical):

- * Average_schooling days

- * Median Age



These 2 variables were encoded using “one-hot encoding” technique:

Schooling_Bins_Low	Schooling_Bins_Medium	Schooling_Bins_High	Schooling_Bins_Very High	Age_Bins_Young Adult	Age_Bins_Adult	Age_Bins_Middle Age	Age_Bins_Senior	Age_Bins_Elder
False	True	False	False	True	False	False	False	False
False	True	False	False	False	False	True	False	False
False	True	False	False	False	True	False	False	False
False	True	False	False	False	True	False	False	False
False	False	False	False	False	False	False	False	False
False	True	False	False	False	True	False	False	False
False	False	True	False	False	False	False	False	True
False	True	False	False	False	True	False	False	False
False	True	False	False	False	False	True	False	False
False	True	False	False	True	False	False	False	False

Ratios:

Some Features in my dataframe might seem to have very high levels or values depending on the population of citizens living in each Entity. In

order to get a more realistic feature, some ratios where calculated from those type of features that depend on the size of the population.:

```

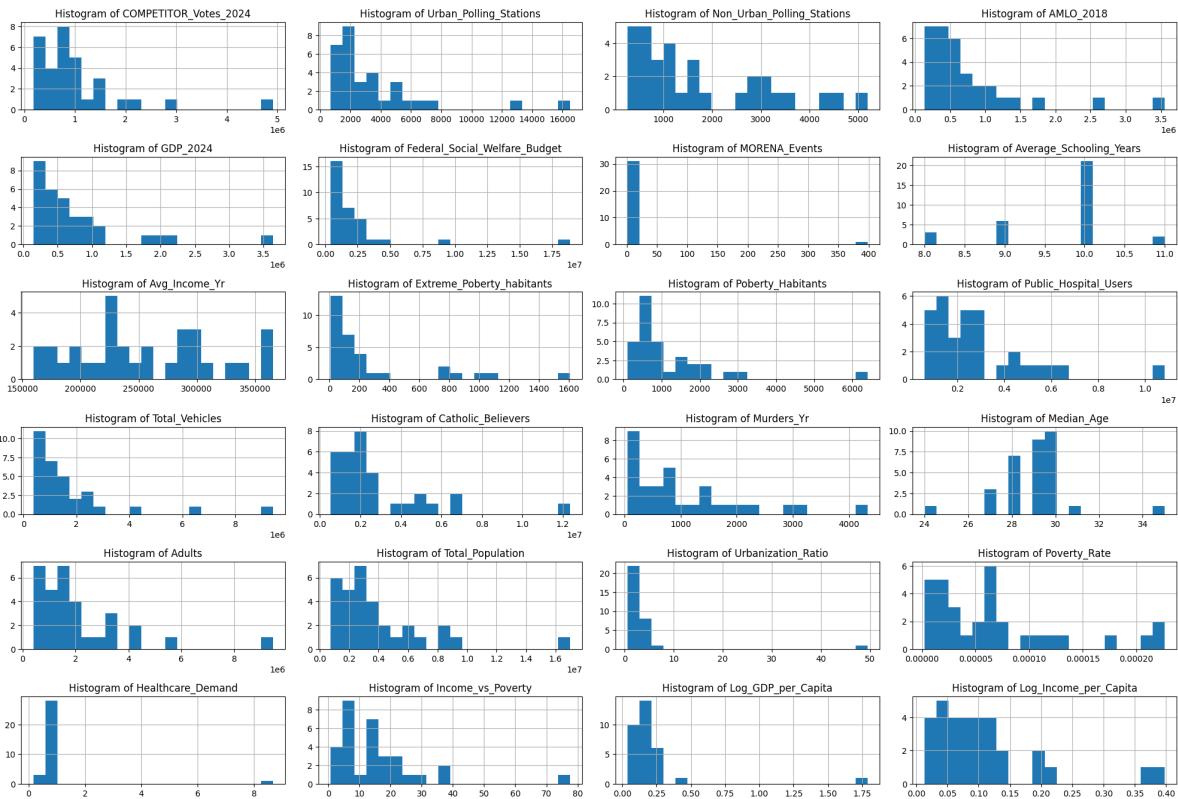
1 # Create new features with Ratios
2 df_encoded['Urbanization_Ratio'] = df_encoded['Urban_Polling_Stations'] / (df_encoded['Non_Urban_Polling_Stations'] + 1)
3 df_encoded['Poverty_Rate'] = df_encoded['Extreme_Poverty_habitants'] / df_encoded['Total_Population']
4 df_encoded['Healthcare_Demand'] = df_encoded['Public_Hospital_Users'] / df_encoded['Total_Population']
5 df_encoded['Income_vs_Poverty'] = df_encoded['Avg_Income_Yr'] * df_encoded['Poverty_Rate']
6 df_encoded['Vehicles_per_Capita'] = df_encoded['Total_Vehicles'] / df_encoded['Total_Population']

1 # Other proportions vs Total Population
2 df_encoded['GDP_per_Capita'] = df_encoded['GDP_2024'] / df_encoded['Total_Population']
3 df_encoded['Income_per_Capita'] = df_encoded['Avg_Income_Yr'] / df_encoded['Total_Population']
4 df_encoded['Vehicles_per_Capita'] = df_encoded['Total_Vehicles'] / df_encoded['Total_Population']

```

Standardization

As the data still seemed to be very skewed, normalization techniques were used in this study:



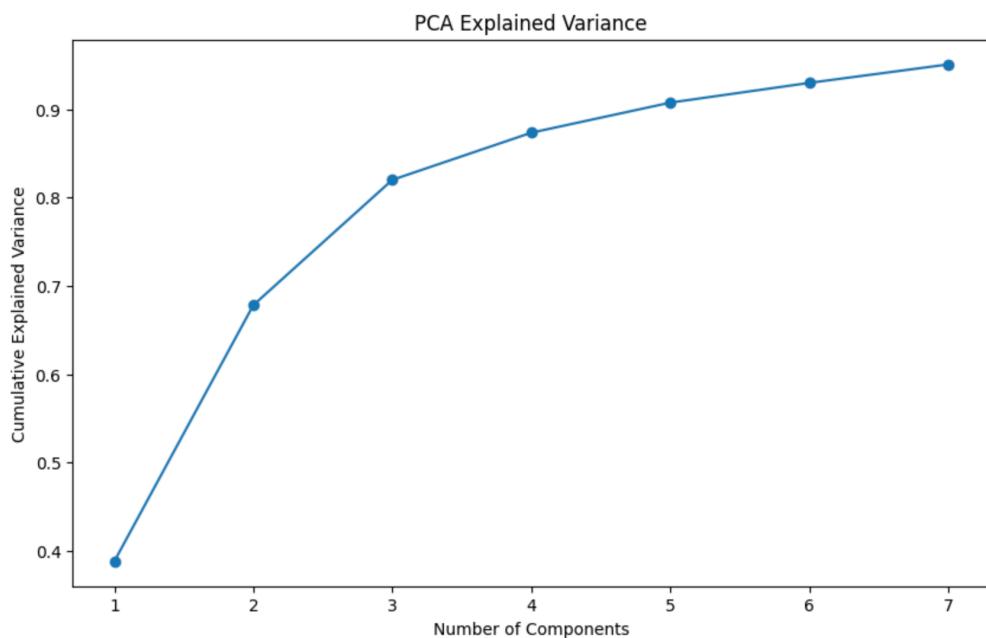
Median_Age	Adults	Total_Population	Urbanization_Ratio	Poverty_Rate	Healthcare_Demand	Income_vs_Poverty	Log_GDP_per_Capita	Log_Income_per_Capita
-1.16692	-0.760959	-0.778683	-0.179094	-0.818624	-0.107315	-0.684212	-0.054838	1.050425
0.58346	0.034264	-0.052355	0.344803	-0.898666	-0.158091	-0.752846	0.013539	-0.172075
0.00000	-0.931093	-0.973068	-0.106661	-0.991260	-0.088872	-0.888999	-0.054285	3.063130
0.00000	-0.904481	-0.932801	-0.273967	0.483175	-0.130636	0.495079	0.670421	1.305532
-2.91730	1.828805	1.634029	-0.381746	1.698576	-0.401373	0.844107	-0.608122	-0.995469

Principal Components Analysis (PCA)

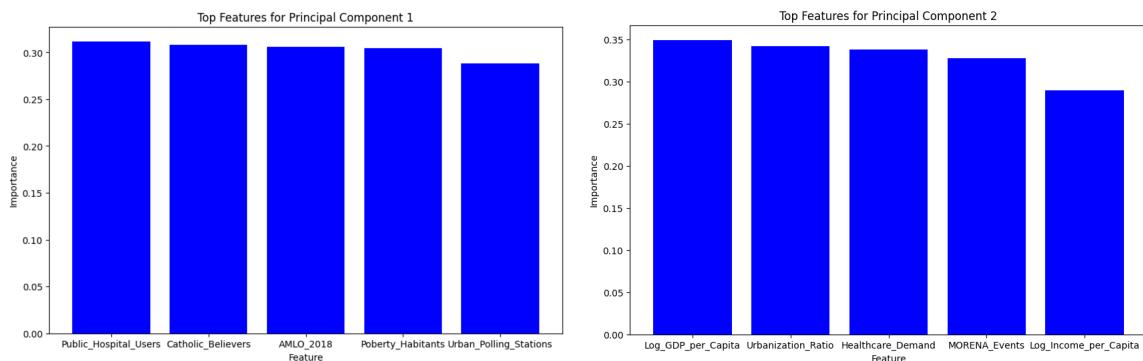
I conducted PCA to identify the most influential features in our dataset. The top features for the first principal component (PC1) included Public_Hospital_Users, Catholic_Believers, AMLO_2018, Poberty_Habitants, and Urban_Polling_Stations.

For the second principal component (PC2), the top features were Log_GDP_per_Capita, Urbanization_Ratio, Healthcare_Demand, Sheinbaum_MORENA_Events, and Log_Income_per_Capita.

PCA helped me to reduce dimensionality while retaining the most significant variance in the data.



- 2 Components explain 67% of the Total Variation
- 3 Components explain 87% of Total Variation

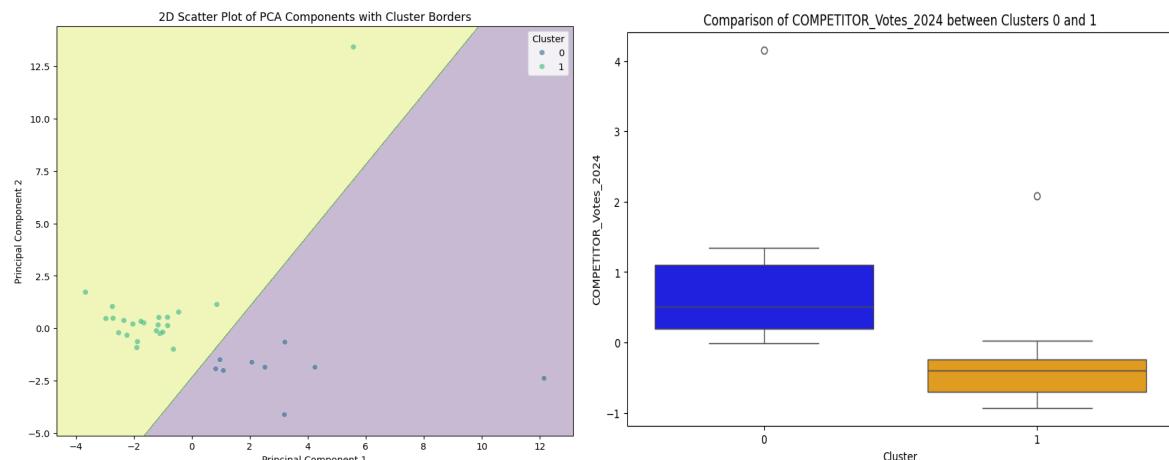


Clustering by using K-Means

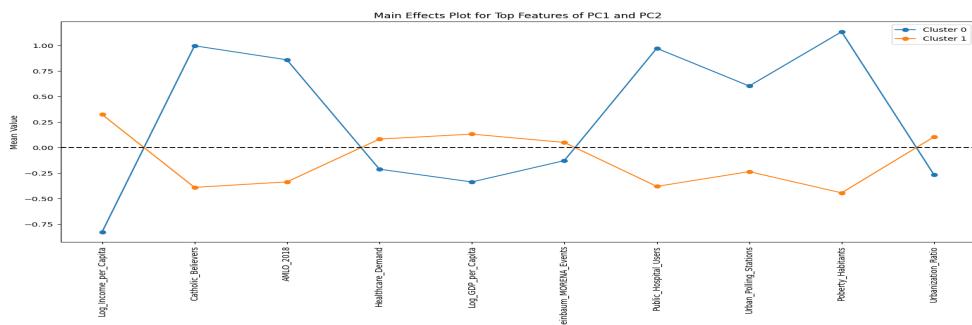
I applied K-means clustering on the PCA-transformed data to identify potential clusters within our dataset. 2 clusters were visualized, revealing distinct groups that indicate an underlying pattern in the data.

Orange Cluster #1 have these characteristics in general:

- LESS VOTES FOR MORENA
- Higher Income per Capita.
- Less chance of being catholic.
- Voted less for AMLO in 2018.
- Preference for Private Hospitals.
- Live in cities that have a lower Ratio of Poor Habitants



Detailed behavior of the features of both clusters:



The graph shows the normalized Average of different characteristics of both clusters.

- **Cluster #0 (BLUE)** has a **higher** preference for **MORENA**
- **Cluster #1 (ORANGE)** has a **lower** voting preference for **MORENA**.

● Predictive Modeling

The modeling phase is one of the most crucial stages in a data science project. This phase involves selecting and applying appropriate modeling techniques to our data to build a model that can make accurate predictions or provide valuable insights.

Choosing the right model is fundamental. Different problems require different types of models. For example:

- **Classification Models:** When the goal is to categorize data (e.g., spam/not spam).
- **Regression Models:** When the goal is to predict a continuous value (e.g., house price).
- **Clustering Models:** To find natural groups within the data (e.g., customer segmentation).

In this specific project the purpose is to predict continuous values and the objective is to test at least 3 different models to measure their performance.

Once the model is selected and trained, it's essential to evaluate its performance using appropriate metrics. This helps to understand whether the model is accurate enough and if it will generalize well to unseen data. Common metrics include accuracy, recall, F1-score for classification, and RMSE, MAE for regression.

In this study as our output variable is continuous, so we used R^2 and MSE.

As a result of this phase, 5 **Models where developed:**

- **Decision Tree**
- **Random Forest**
- **Linear Regression**
- **Ridge**
- **Lasso.**

The Lasso model had the best performance (98% accuracy in testing and training)

```
# Initialize models
models = {
    'Random Forest': RandomForestRegressor(random_state=42),
    'Lasso': Lasso(alpha=0.1),
    'Linear Regression': LinearRegression(),
    'Ridge': Ridge(alpha=1.0),
    'Decision Tree': DecisionTreeRegressor(random_state=42)
}

# Train and evaluate models
conclusions = []

for name, model in models.items():
    model.fit(X_train_scaled, y_train)
    y_train_pred = model.predict(X_train_scaled)
    y_test_pred = model.predict(X_test_scaled)

    mse_train = mean_squared_error(y_train, y_train_pred)
    r2_train = r2_score(y_train, y_train_pred)
    mse_test = mean_squared_error(y_test, y_test_pred)
    r2_test = r2_score(y_test, y_test_pred)

    conclusions.append({
        'Model': name,
        'Train MSE': mse_train,
        'Train R222 in descending order
conclusions_df = conclusions_df.sort_values(by='Test R2', ascending=False)
```

Metrics of the Model:

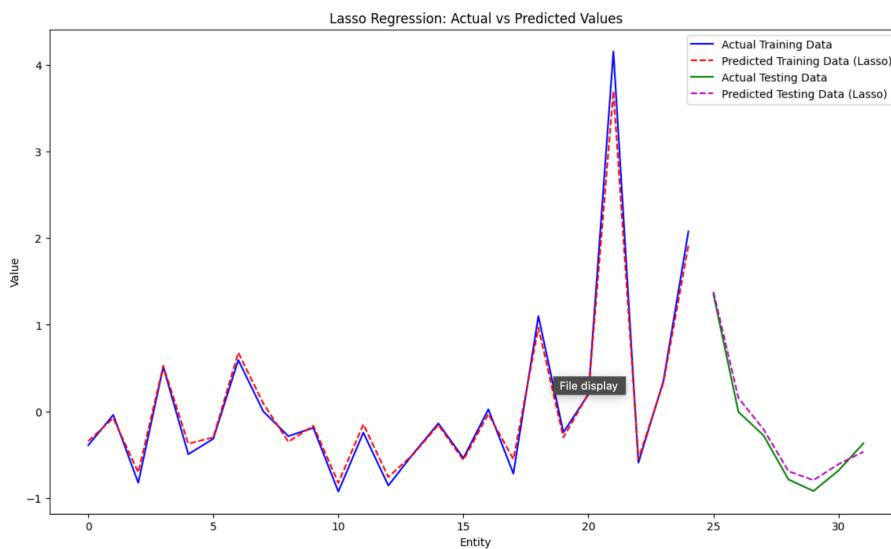
Model	Lasso
Train MSE	0.014998
Train R ²	0.986565
Test MSE	0.010197
Test R ²	0.97997
Model	Ridge
Train MSE	0.002201
Train R ²	0.998029
Test MSE	0.047715
Test R ²	0.906278
Model	Linear Regression
Train MSE	0.0
Train R ²	1.0
Test MSE	0.054846
Test R ²	0.892273
Model	Random Forest
Train MSE	0.100472
Train R ²	0.910003
Test MSE	0.067189
Test R ²	0.868028
Model	Decision Tree
Train MSE	0.0
Train R ²	1.0
Test MSE	0.127489
Test R ²	0.749587

Testing Predictions:

```
Lasso Coefficients:  
          Feature   Coefficient  
11      AMLO_2018    0.567453  
18      Poberty_Habitants  0.202205  
19      Public_Hospital_Users  0.209553
```

Lasso Equation:

```
Y = + 0.57 * AMLO_2018 + 0.20 * Poberty_Habitants + 0.21 * Public_Hospital_Users
```



Lasso Equation ML Model developed had a great performance that can be graphically seen in the above equation.

I suggest to choose the Lasso Model due to 3 main points:

- 1) The simplicity of its equation:

Lasso Equation: $Y = + 0.57 * \text{AMLO_2018} + 0.20 * \text{Poberty_Habitants} + 0.21 * \text{Public_Hospital_Users}$

- 2) The great accuracy it has in training and Testing Datasets (98% R^2)

- 3) With a low amount of features I can still predict the votes accurately.

Model Metrics

An aditinal trial was run in a CSV file with all the dataframe features to confirm the predictive power of the Model, showing similar results, MSE 0.007 and R2 of 99%.

ModelMetricsfile

AMLO_2018	Poverty_Habitants	Public_Hospital_Users	COMPETITOR_Votes_2024	Predicted	Difference
-0.7847173190976150	-0.7082027433568750	-0.7970364622246030	-0.8262555701128480	-0.7563070776241820	-0.06994849248866610
-0.1143552063371590	-0.5959998662049440	-0.0202714265554901	-0.1910465079640970	-0.18863944042982300	-0.002407067534274330
-0.8593764649741140	-0.8910518764933560	-1.0291730519087100	-0.9277501104378760	-0.8841813012347460	-0.04356880920313
-0.8116595205902790	-0.6990602867000510	-1.0032360720022700	-0.8581043968505120	-0.8131375591969460	-0.044966837653565900
0.4811805121176670	0.8751045231203730	0.4056895938703400	0.5134736491721350	0.5344886112439160	-0.021014962071780900
-0.4375800733621880	-0.4896149160164470	0.0709569178833855	-0.3174272341599740	-0.3324426722642260	0.0150154381042516
2.642901014090320	0.742123353847510	1.6801366939321400	2.077905693252200	2.0077069508341800	0.07019874241802350
-0.3359137617271650	-0.5320026696071760	-0.1614775186467470	-0.243844389546412	-0.33178165702173600	0.08793726747532420
-0.8022673464089400	-0.8553131822894070	-1.0581142342428600	-0.9209848087453530	-0.8505590131019790	-0.07042579564337440
-0.7069592395177220	-0.5444696559573910	-0.6955931623545990	-0.6844963307639760	-0.6579352618110460	-0.02656106895293010
-0.2070246582596980	0.5750657182918760	0.9429471643874960	0.1962412740405790	0.19502799297172200	0.0012132810688574300
0.1182099543517430	0.1619929038881000	-0.0926610771222894	-0.0016202055644616	0.0803194285624325	-0.0819396341268941
-0.0197990262678175	-0.0815288961527572	-0.331575231356973	-0.0416457380890304	-0.09722202278817170	0.05557628469914130
0.5850652843209090	0.4138260281624340	1.3388947538335000	0.5921765873621260	0.6974203160004410	-0.10524372863831500
3.862305164225730	4.335939933939930	3.8129722090114400	4.156590755587110	3.8694260942890600	0.28716466129805000
0.0855674840844155	0.4262930145126490	0.0496048731858473	-0.0076036580284104	0.14444909219967500	-0.15205275022808500
-0.3275471936153710	-0.3907101576380780	-0.6770813851807150	-0.4960974790652660	-0.40703102277632700	-0.08906645628893870
-0.7163624714464410	-0.738954643020738	-0.8857242922056600	-0.7876092668172880	-0.7421196386918080	-0.04548962812548040
-0.2855885711799470	-0.2253148053918980	0.756275346944859	0.0227608349200244	-0.049030623792529000	0.07179145871255340
0.4408957561919710	0.3706071421483570	0.0438952682863672	0.350638800508142	0.33465001579923200	0.015988784708910200
0.9569898530282180	1.3987179498293800	0.8322610811917750	1.0999570951993400	1.000002633242230	0.09995446195710380
-0.5643557648615520	-0.5685724962344720	-0.4787425493632280	-0.4881266924167510	-0.5359332205842570	0.047806528167506200
-0.5861782293170290	-0.6159470443652880	-0.7074536620406230	-0.5412452501572880	-0.6058762686122950	0.06463101845500760
-0.4771626623292210	-0.3084280477266620	-0.2504312662919310	-0.2871152934916490	-0.3862588929942940	0.09914359950264460
-0.1740172823288130	-0.4704988702794510	-0.1802742756829030	-0.2773180469854340	-0.23114722287672300	-0.0461708241087109
-0.4097214610553630	-0.4813035917829710	-0.2256461197084190	-0.3925845021858100	-0.37718763629691900	-0.015396865888891200
0.0110189156817218	-0.2618846320191940	-0.5957949301845360	-0.2461631248443250	-0.17121307980401000	-0.07495004504031510
-0.2090813992724770	-0.2643780292892370	-0.0407613833419057	-0.1399195743901760	-0.18061189394495900	0.04069231955478300
-0.5527658633883330	-0.4580318839292370	-0.8941550258903860	-0.5939991594235920	-0.5944554743541780	0.0004563149305858440
1.4132352743857700	1.6522133389504100	1.419281949632180	1.350763388271900	1.434035983612730	-0.0832725953408271
-0.5910768114066040	-0.3250506961936150	-0.5006049844161070	-0.3693070595836330	-0.5070509684678700	0.1377439088842370
-0.6238588857346110	-0.4455648975790220	-0.7271077414383660	-0.7202436786893890	-0.5974051700865890	-0.1228385086027990
Mean Squared Error	0.007706573204951880				
R^2 Score	0.9922934267950480				

Conclusion: Key Insights

- Target Population: MORENA focuses on impoverished populations, which are their major voters. However, these voters do not necessarily contribute to GDP growth or industrial/economic development, presenting an opportunity for the opposition to highlight long-term benefits.
- Campaign Efforts: MORENA held significantly more public events (617) compared to the opposition (223), despite spending 3x less on campaign budgets. The opposition should focus on engaging directly with impoverished populations rather than spending excessively on propaganda.
- Key Demographics: The opposition should prioritize engaging with populations that have a strong preference for AMLO, federal welfare recipients, impoverished communities, Catholic believers, and adults.
- Socioeconomic Clusters: There are two distinct clusters: one with a high socioeconomic profile (less likely to vote for MORENA) and another with a low socioeconomic profile (MORENA followers).
- Political Polarization: AMLO has polarized these clusters with rhetoric of “Fifis vs Chairos.” This polarization should be addressed as all Mexicans are equal, and unity should be emphasized.
- Future Projects: For future projects, incorporating demographic, economic, and social information per city could help create a more robust model that generalizes the data better.