

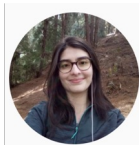


Prediction Model of Electoral Results in Mexico: A Case Study using PREP, INEGI and INE Data.

Data Science Intensive Capstone Project, August 3rd, 2024

By Javier Jorge Pérez Ontiveros

Mentor



Ale Berbesi

Executive Summary:

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.

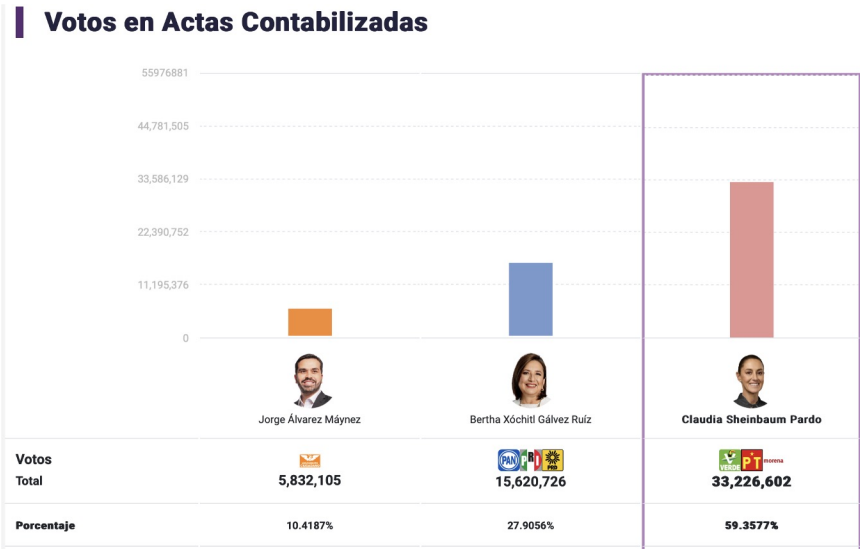
The goal is to understand why the winning party (**MORENA**) achieved a **59% acceptance rate**, while the opposition (PAN-PRI-PRD) secured only 28% of the votes.

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.

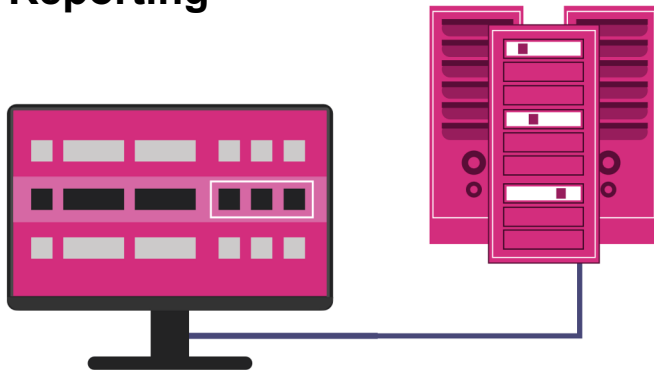
Results:

- 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

- Data Engineering
- Data Wrangling
- Exploratory Analysis (EDA) and Visualization
- Feature Engineering and Clusterization
- Predictive Modeling
- Reporting



Data Sources

Data Engineering and Data Wrangling: Acquiring, Cleaning and Merging

PREP

2024 and 2018 election results exported from database to CSV file.

55,976,881 votes
170,944 polling stations
1,580 municipalities
32 states

INEGI

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.

INE

Performance and Behavior of Parties per state in csv:

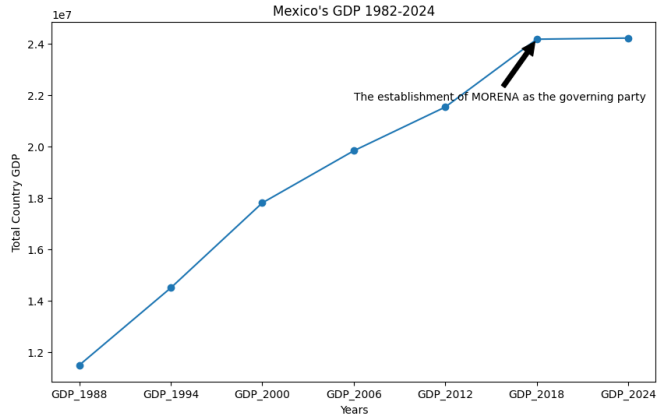
of Campaign events, Marketing Budget, Federal Welfare,

Final DataFrame

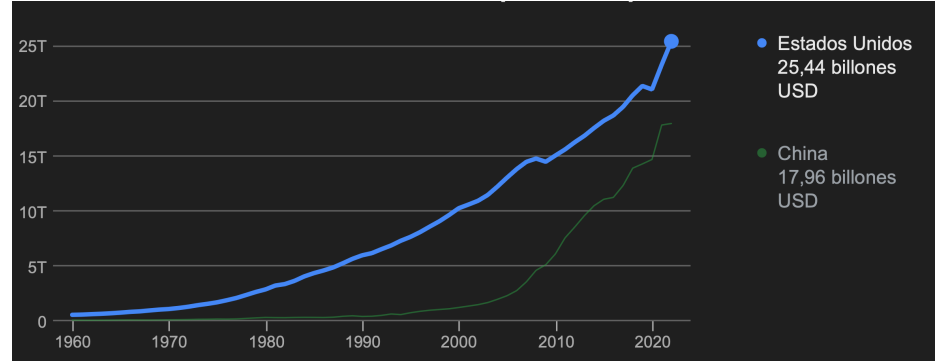
|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|

A 40 columns x 32 rows DataFrame was created to predict the votes for MORENA Entity
This DataFrame was further modified and expanded with the use of some ratios and feature engineering.

Insight #1: GDP has stagnated since 2018



Since AMLO president won elections in 2018, the economy in Mexico has slowed down.

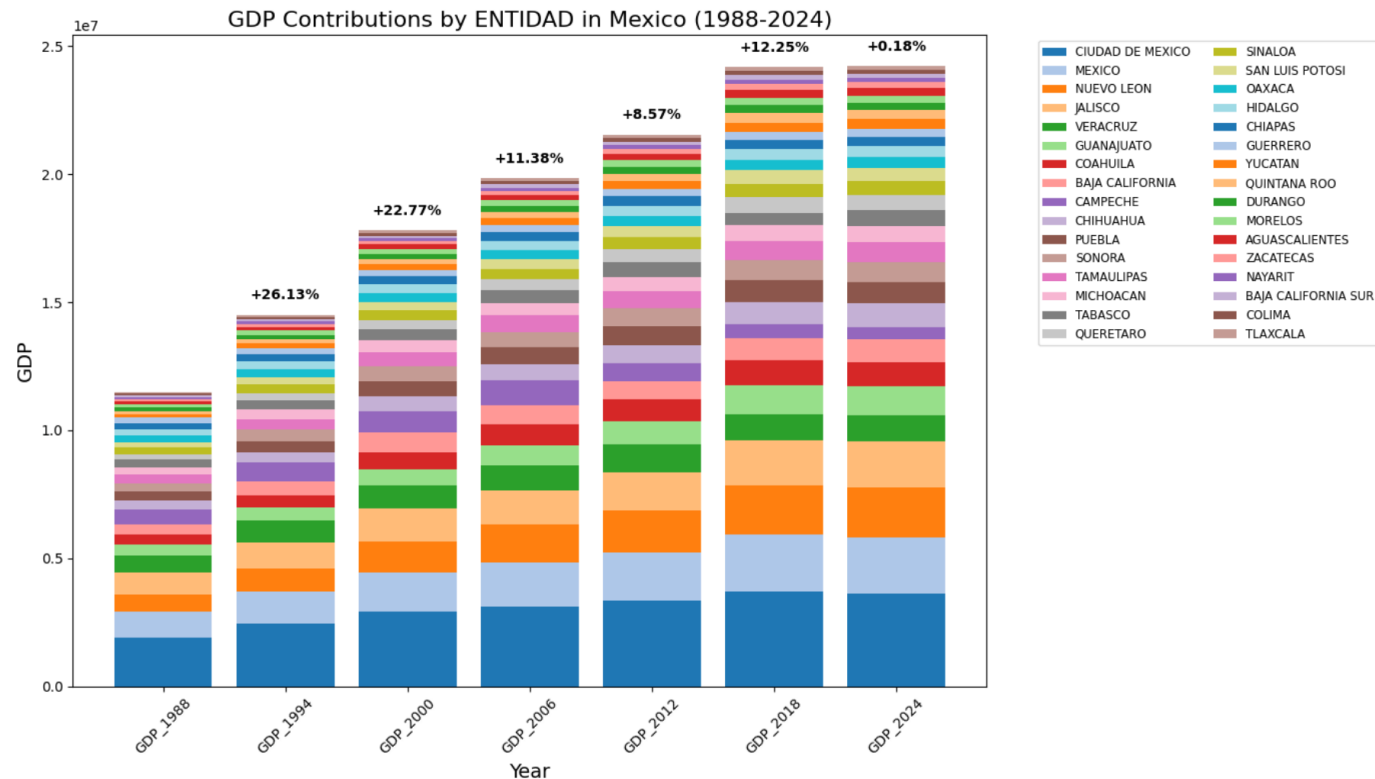


In comparisson, USA and China have a very positive trend (even after Covid)

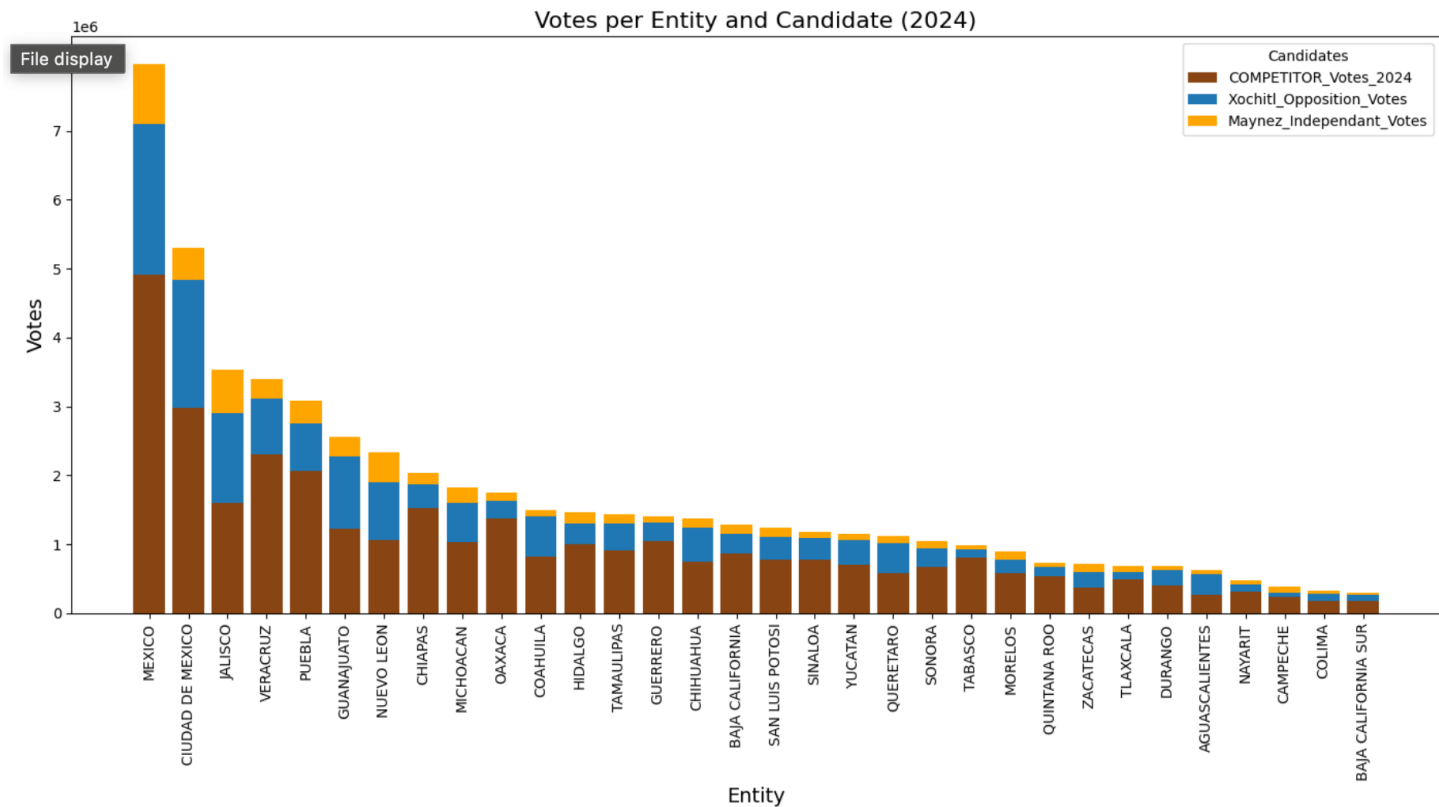
In the other hand, with Opposition parties, From 1988-2018, the GDP had a 10%-20% growth in average

This is a Key argument that needs to be “Sold” to the mexicans to gain more votes.

Crime Scene: Economy stagnation per State



Votes per Entity



Geographical Heatmap Distribution

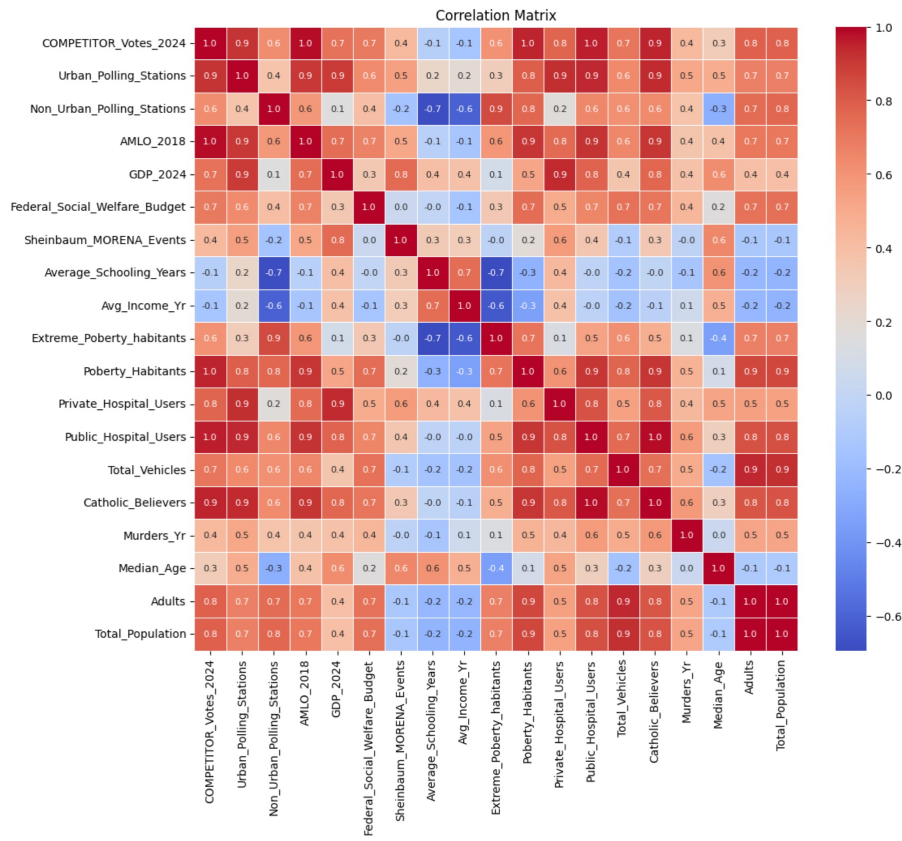


The map shows a significant trend where the highest resistance (blue color) is observed in the Bajío region, known for its industrial and economic strength. While the stronger MORENA states are in the South-East Region (Less Industrial and Economical development)

This is interesting as votes for MORENA, the competitor party, correlate highly with poverty ratios.

Opposition Party MUST focus on helping the less priviledged population from 2024-2030 if they want to have a chance of winning.

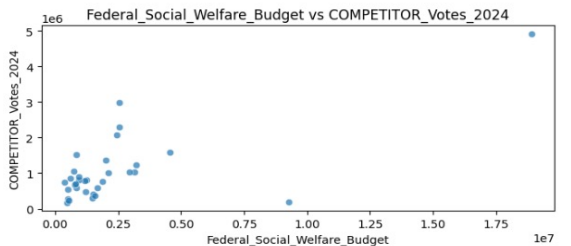
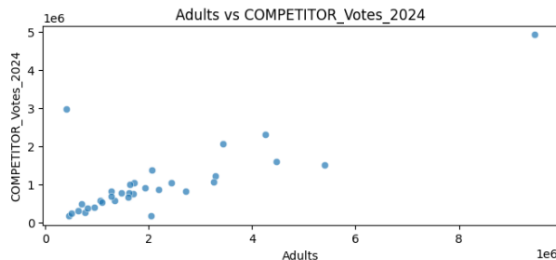
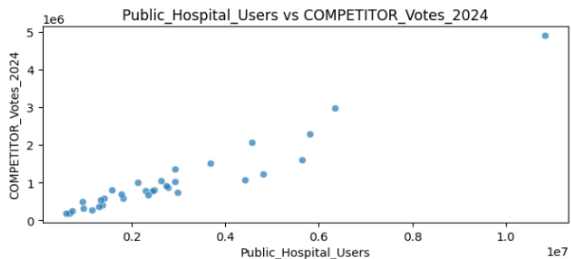
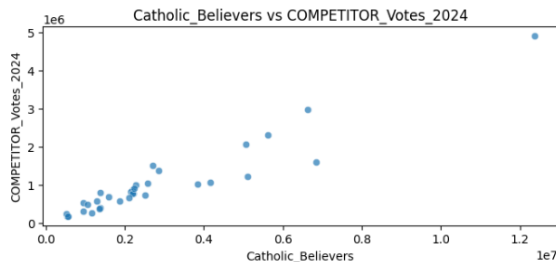
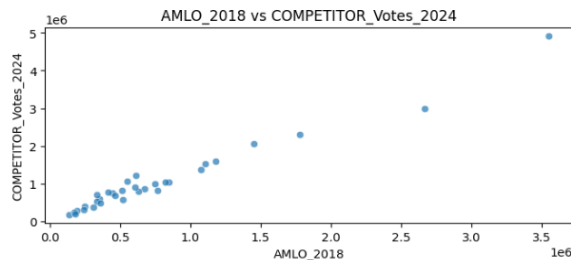
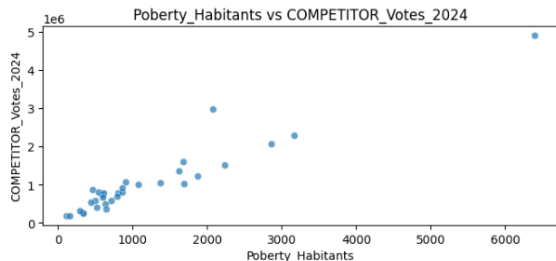
Features vs Response Heatmap



Since the begining of this study these variables showed a positive correlation with MORENA votes:

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

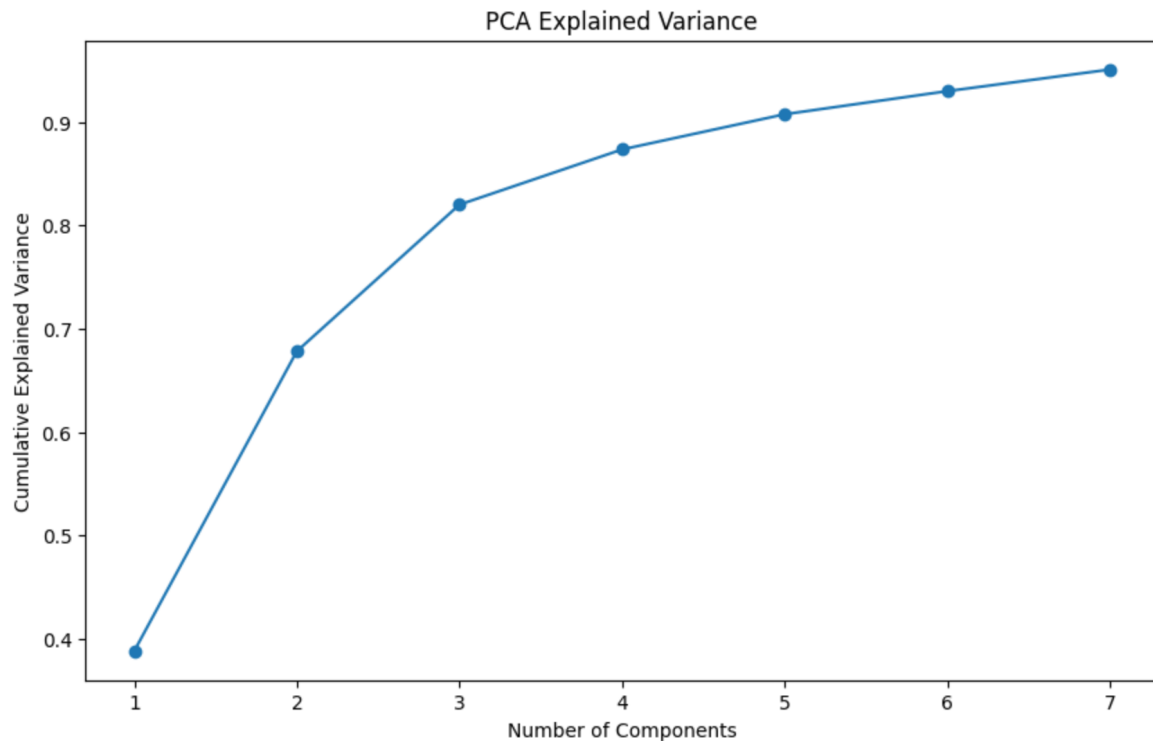
Linear Correlations



In General the votes for MORENA increase when these variables increase:

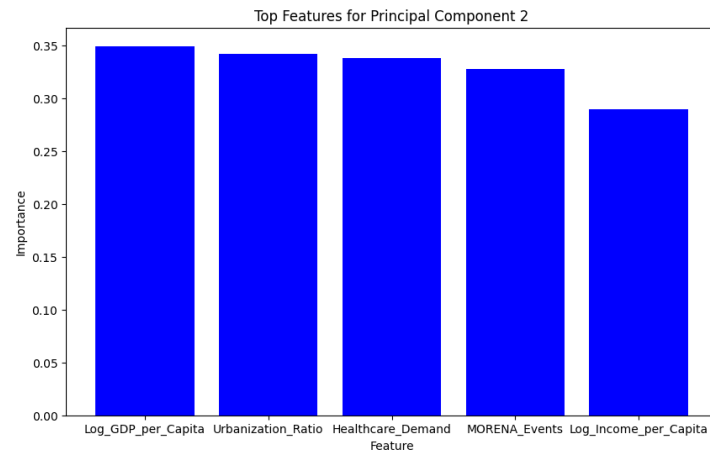
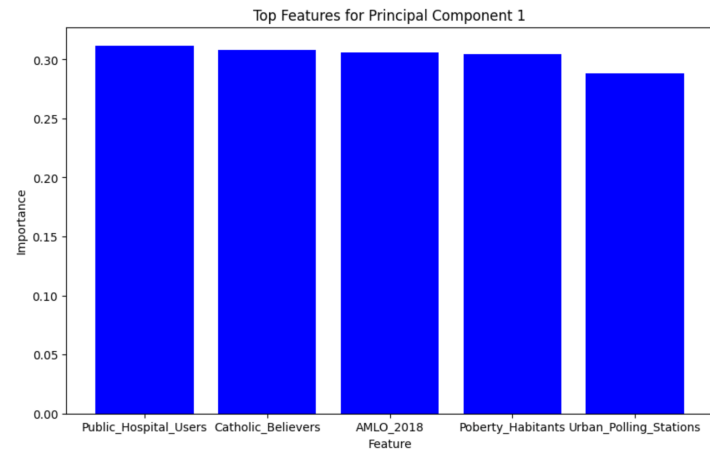
- Poberty Habitants
- Catholic Believers
- Adults 40-50 yrs old
- AMLO followers since 2018
- Public Hospital Users
- Federal Welfare Budget

Clustering (PCA)

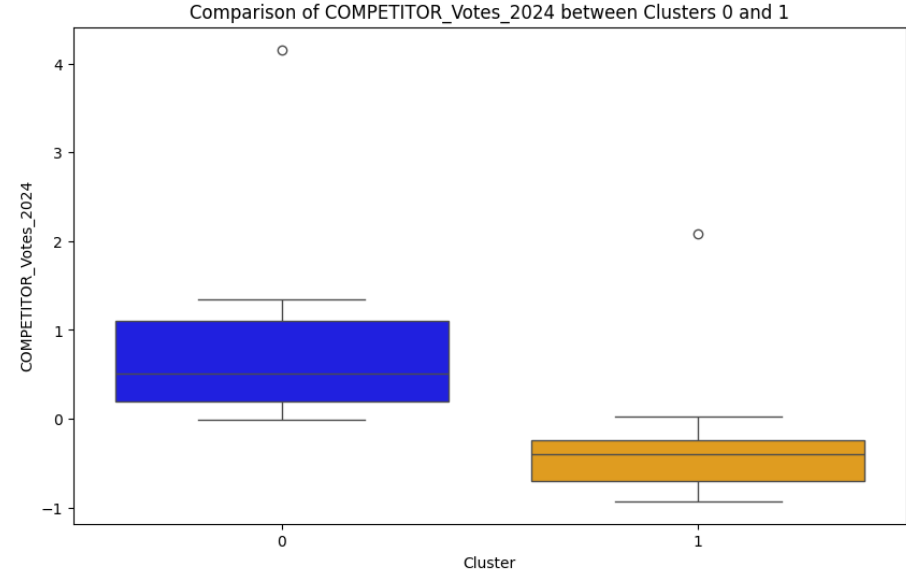
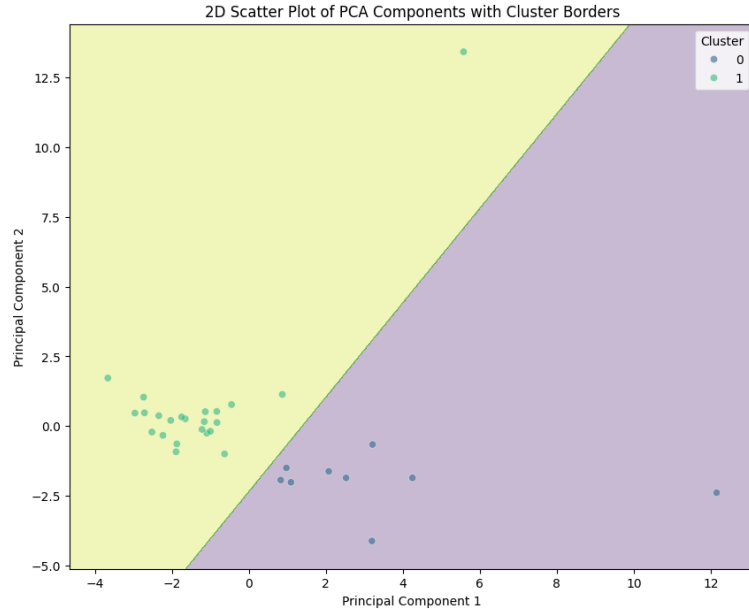


2 Componenten explain 67% of the Total Variation

3 Components explain 87% of Total Variation



Clustering

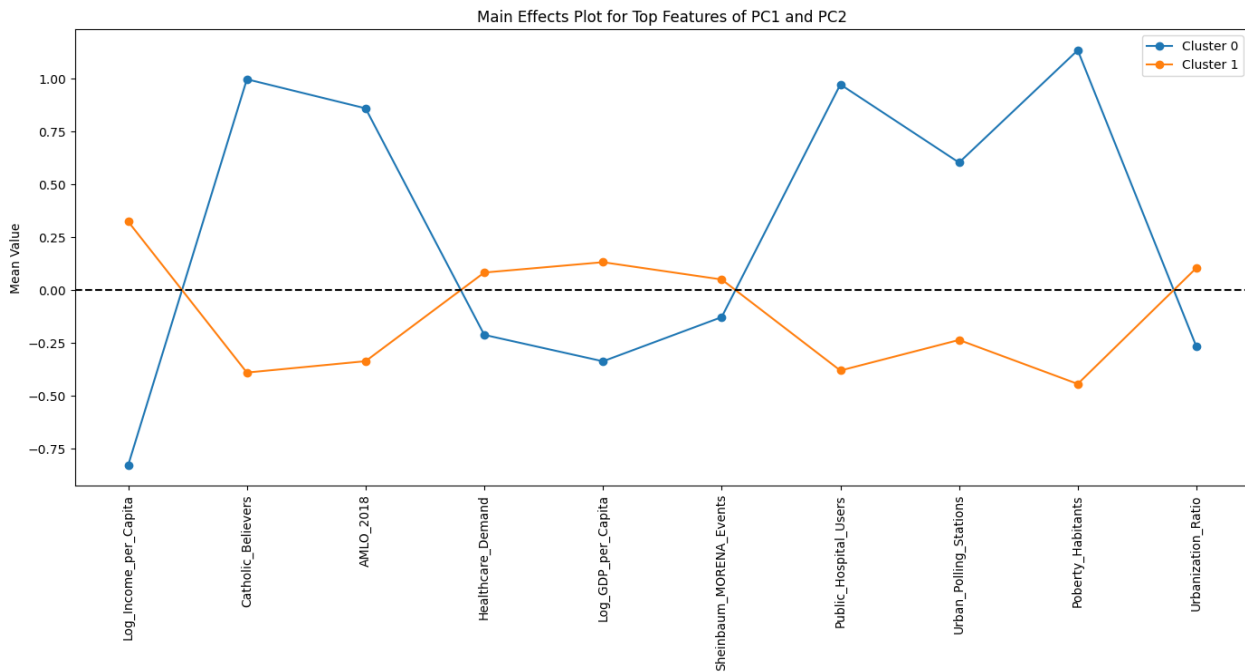


Cluster #1 (PCA 1) has a lower voting preference for MORENA.

While has this features as key characteristics: High GDP, High Urbanization Ratio, Higher Healthcare demand and Higher Income per Capita.

Also is important to mention that in this cluster is where MORENA performed more campaing events!

Behavior of the 2 clusters identified (Morena and Opposition)



Cluster #0 has a **higher** preference for **MORENA** while,
Cluster #1 (PCA 1) has a **lower** voting preference for **MORENA**.

Modeling

```
# 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 R²': r2_train,
        'Test MSE': mse_test,
        'Test R²': r2_test
    })

# Convert conclusions to DataFrame
conclusions_df = pd.DataFrame(results)

# Print the results:

# Sort the results by Test R² in descending order
conclusions_df = conclusions_df.sort_values(by='Test R²', ascending=False)
```

Model	Lasso
Train MSE	0.014998
Train R²	0.986565
Test MSE	0.010197
Test R²	0.97997
Model	Ridge
Train MSE	0.002201
Train R²	0.998029
Test MSE	0.047715
Test R²	0.906278
Model	Linear Regression
Train MSE	0.0
Train R²	1.0
Test MSE	0.054846
Test R²	0.892273
Model	Random Forest
Train MSE	0.100472
Train R²	0.910003
Test MSE	0.067189
Test R²	0.868028
Model	Decision Tree
Train MSE	0.0
Train R²	1.0
Test MSE	0.127489
Test R²	0.749587

5 Models were developed: Decision Tree, Random Forest, Linear Regression, Ridge and Lasso. The last model had the best performance (98% accuracy in testing and training)

Results



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 representation

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 Chairios.” 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.

Thank You!



<https://www.linkedin.com/in/jjpo/>



javierjorge77@gmail.com



<https://github.com/javierjorge77/Springboard/tree/main/Capstones/Capstone2>



(+521) 8711777903