# Election Data: Analyzing the Effectiveness and Accuracy of Ipsos Forecasting Models

**Project Director:** Clifford Young

President of Ipsos Public Affairs in the U.S.

**Supervisor:** Prof. Sharon O'Halloran

**Team Members:** Matthew Kane McMahon

Juan Manuel Puyana

Adam Christopher Stoddard

Xiaozhi Wang

Na Wei

# Executive Summary

The Columbia Capstone Group prepared the following election forecasting project for Ipsos, a leading global market research firm. The objective of this capstone project is to review the capabilities and limitations of the original Ipsos election prediction model, expand the dataset by including new variables, and enhance the robustness of the estimates through rigorous analysis and diagnostics of the model's constraints. With the addition of a new data set and added variables, a new model was created to predict the outcome of an election.

The election model estimates the likelihood that the incumbent party wins the election. The forecasting model is a logistic regression estimated with maximum likelihood. In addition to the data provided by Ipsos, the team expanded the data set to include additional independent variables, such as government approval ratings, GDP growth rate, inflation rate, and employment growth. The group also included a war variable that measures whether the country in which the election takes place is involved in a military conflict. The scope of the assembled election data was also increased to include elections in 87 countries from 1980 until the end of 2015. This additional data provides for more robust analysis by the Ipsos prediction model.

Seven different models were created and tested using the foundation of the original Ipsos forecasting model. The Columbia Capstone Group was able to successfully enhance the predictive capabilities of the original Ipsos model by addition of the new variables and collection of a more expansive election data set. The predictive power of the original model increased from 74% accuracy to 76%.

The following report includes:

- A review of academic literature needed for an in-depth study of predicting elections;
- Methodology, data acquisition, and summarized codebook;
- The original model from Ipsos and its results using the updated data set;
- Estimates from new models using the additional variables; and
- Appendices, including an in-depth literature review, Stata code for each model, and a complete codebook.
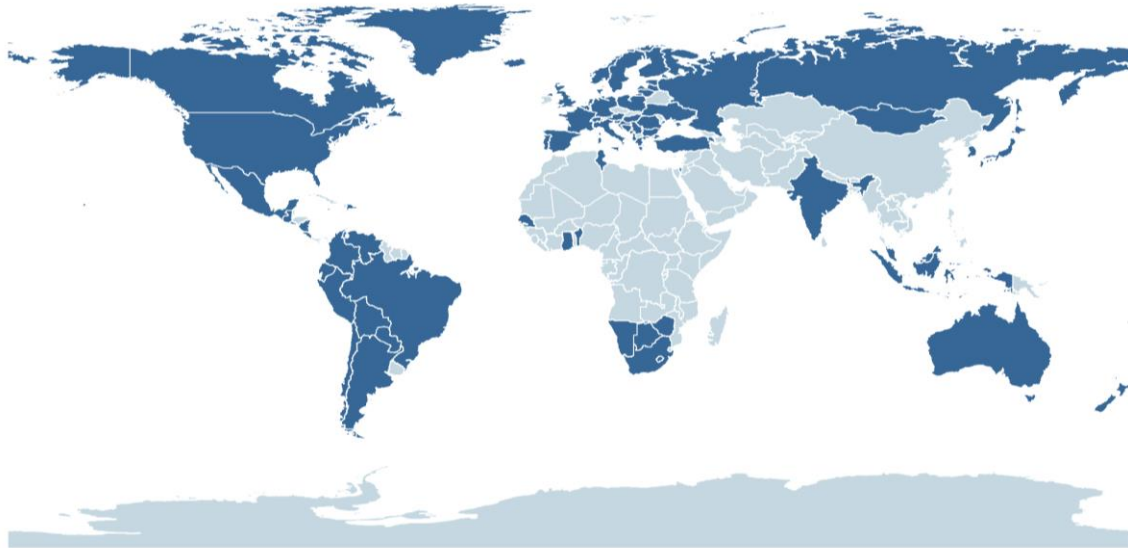
# Table of Contents

# I.    Introduction

This Columbia Capstone Group project is an election forecasting model originally designed by and updated for the global research firm Ipsos.  The project seeks to identify the likelihood of an incumbent electoral candidate's reelection based on contemporary public opinion, macroeconomic conditions, and whether the country is involved in a military dispute.

This election forecasting project uses election data collected both in the United States and in foreign countries to test and analyze the effectiveness and accuracy of the Ipsos forecasting model.  An analysis of different models was completed through a review of academic literature and best practice.  In places where no polling data exists or weak polling is available, aggregate data from across multiple levels was combined to determine vote share and probability of victory of the incumbent.



 Map 1: Country Coverage of the Election Dataset

We built upon the historic elections dataset previously developed by the Ipsos Public Affairs team to enhance its scope and accuracy.  There were previously 99 countries which were included in the original Ipsos dataset, with 540 observations, including state level elections.  The total number of countries was decreased based on the Freedom House Index and on election types.  Additional democratic countries were selected and included in the updated data set.  The total number of countries used are 87 countries with observations on 749 elections.   We also used a standardized thirty-five year time scale over which the data was collected, and a higher number of independent variables.  This project includes global elections in various contexts with a focus on free elections.  The goal of the project is to create an improved database accompanied by a more thorough and accurate statistical model.  Our hypothesis is that with greater scope of the data and a higher precision of the statistical model, we can enhance the accuracy of the election model and provide for a greater predictive capability of the Ipsos elections forecast.

# II.    Literature Review

The Columbia Capstone Group reviewed seventeen academic studies that covered election prediction models and methodology in order to understand the existing state of the field. Among them we chose the following five as the most salient in terms of variables and methodology, and we have drawn on them to create our new prediction model. It should be noted that almost all election prediction literature is based on the US electoral system.

## 1. Lewis-Beck/Tien Jobs Model Forecast

Michael S. Lewis-Beck and Charles Tien have written numerous peer-reviewed articles on predicting presidential elections through their Jobs Model Forecast. The key finding of their study is that there is a high correlation between economic performance during a president's term, his approval rating, and whether he gets re-elected or his successor is elected. These two primary variables are broken down into four measureable sub-variables. Economic performance is measured through change in GNP in the fourth quarter prior to the election year and through the percentage change of jobs growth over the first 3.5 years of the president's term. Approval rating is measured by the first Gallup poll in July of an election year, and incumbent or successor status is scored as 1 if the incumbent is running for reelection, 0 if the incumbent and successor have a tolerable relationship, and -1 if the incumbent party candidate and the president are not united. The authors apply ordinary least squares regression to these data to predict the vote share of the incumbent candidate or party, and correctly predict all U.S. elections after 1984 with the exception of the controversial 2000 election. This model has correctly picked winners for the past few elections, with the exception of the 2000 election, which forecast a Gore win (by popular vote, however, the model was correct). Each variable is significant.

## 2. Fair's Model

Ray C. Fair created one of the earliest presidential election vote-share models in 1978. He includes both economic and political variables. Economic variables include GDP growth rate per capita in first three quarters of election year, the inflation rate, and the number of quarters in the first 15 quarters in which the GDP growth is > 3.2%. Political variables include incumbency status, duration of time in office of the incumbent party, and whether the nation is in a military conflict. The author applies ordinary least squares regression to these variables in order to predict the Democratic Party vote share in the presidential election. The sample is comprised of U.S. elections from 1916-2006. The result shows that economic variables do affect vote-shares in presidential, on-term house elections, and midterm house elections

## 3. Lichtman "Keys" Model

Allan J. Lichtman developed an index forecasting model that evaluates key election issues. Through the application of pattern recognition methodology on American presidential elections from 1860 to 1980.  He uncovered thirteen key

indicators and a simple decision rule that accounted retrospectively for the popular vote winners of each of these contests, which include no polling data and consider a much wider range of performance indicators than economic concerns. The thirteen statements favor the re-election of the incumbent party including party mandate, contest, incumbency, third party, short-term economy, long-term economy, policy change, social unrest, scandal, foreign/military failure, foreign/military success, incumbent charisma, and challenger charisma. When five or fewer statements are false, the incumbent party wins. When six or more are false, the challenging party wins.

## 4. John E. Mueller's Model Indicating Presidential Popularity

John E. Mueller studies the indicators which affect presidential popularity in the 24 years' period from Truman to Johnson. The dependent variable is presidential popularity, the percentage of approval in the Gallup Poll question, "Do you approve or disapprove of the way (the incumbent) is handling his job as President?"

The key take away of this paper is that four major variables impact popularity: coalition of minorities, rally around the flag, economic slump, and war. The first means that each president will experience in each term a general decline of popularity. "Rally around the flag" means this decline will be interrupted from time to time with temporary upsurges associated with international crises and similar events. "Economic slump" means that the decline will be accelerated in direct relation to increases in unemployment rates over those prevailing when the President began his term, but that improvement in employment rates will not affect his popularity one way or the other. The "war" variable means that the president will experience an additional loss of popularity if a war is ongoing. The study is conducted in two phases, the first phase without the war variable, and the second with the Korean War and Vietnam War variable affecting the presidents in the time period. The fit of the resulting equation was very good: it explained 86 percent of the variance in presidential popularity. But the effect of each variable are very different.

## 5. Abramowitz: Forecasting the 2008 Presidential Election with the Time-for-Change Model

The argument presented in this article is that one of the largest factors in US Presidential elections is that there is a need to change parties every few elections. Electoral results can be determined by popularity of the incumbent president, the state of the economy and the length of time that the president's party has controlled the white house. Abramowitz found that there is a strong relationship that exists between approval and vote choice. The key findings of this article are that the leading indicators of presidential elections are the growth rate of the economy, particularly during the second quarter of the election year. Another is the incumbent president's approval rating at mid-year. The last is the length of time the incumbent's party has controlled the White House, which is called the Time-for-Change Factor. Abramowitz's Time-for-Change Model states that electoral results will be influenced by the natural change from one party to the other. The incumbent party tends to be blamed for whatever contentions and disagreements exist in the political landscape at that time, and that the recourse taken is generally the choosing of the other party in the following election. The Time-for-Change Model is ultimately highly accurate in predicting presidential elections in the United States. It predicted within two percentage points accuracy the voting results for the 2008 election of President Obama.

## *Hypothesis*

The original model focused on incumbent data, such as approval ratings or popularity of the incumbent president. In order to improve the original prediction model in global context, the literature suggests new fundamental variables that should impact electoral outcomes across country boundaries.  All of the five models above include macroeconomic data (GDP growth, employment growth, inflation etc.) can influence election prediction significantly. International crisis and military conflicts are also frequently mentioned in Fair's model, Mueller's model and Lichtman's model. In this report, we want to test whether the old model could be improved and work across different countries, time periods, and election types by adding these new variables. We hypothesize that GDP growth, employment growth, changes in the inflation rate, and whether a country is in a military conflict will influence the outcome of elections in any country and therefore should be included in the electoral model.

# III. Data and Methods

The original dataset Ipsos provided consists of 85 countries and 525 elections with a winner in total, with 266 national elections and 212 state elections. In terms of time frame, although the original dataset included several US national elections from time as early as the 1936, 87.6% of the total data is from after 2000. The original model consists of the following variables, which were filled in for each election year: country, election winner, party winner, name of candidates, vote of candidates, if incumbent or successor wins, and government approval rate.

Based on the models we listed in the literature review and our hypothesis, we selected the following variables to collect data on: GDP, inflation employment, and conflict data from the Correlates of War dataset. In order to subset the data further into relevant groups, the Freedom House's Freedom in the World indicators and types of elections were also included. We also selected these variables because they have higher data accessibility, as many organizations such as IMF, Oxford Economics, Freedom House etc. publicize their database on the internet. In addition to the new variables we include in the model, we also adjusted the scope of the sample. Particularly we removed all state elections from the sample as we believed the prediction method of state election and national election are different. Also, we added Germany, Greece and Sweden in our sample countries, and also adopted a larger time frame from 1980 to 2015 for all the 88 countries, resulting 727 elections in total in our database. Detailed description of the variables is listed in the codebook in Appendix 2
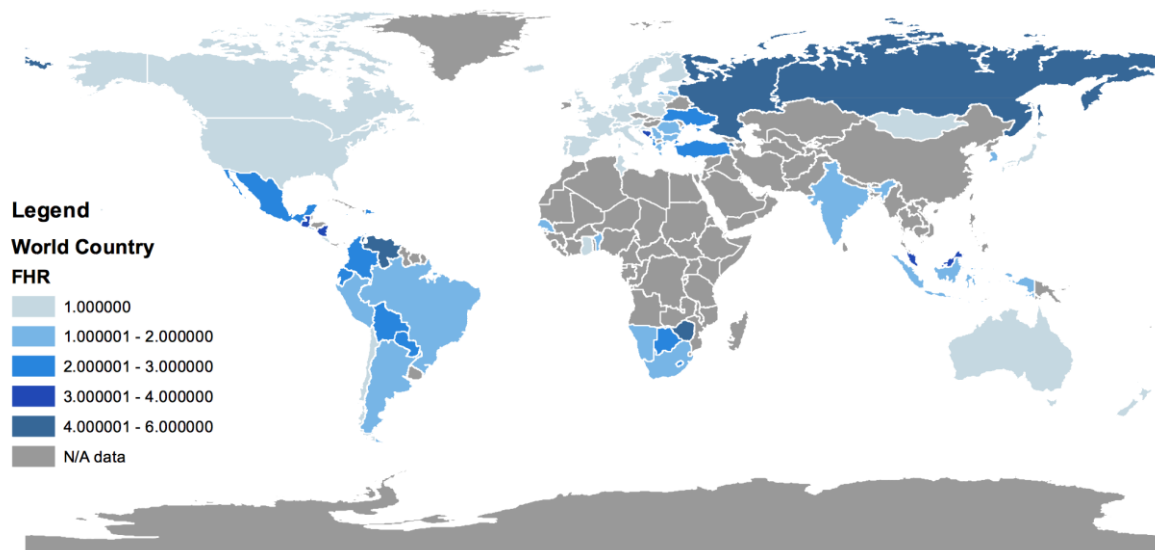
## Collection Method

### Additional data for original model

*Incumbency* - There is no single, comprehensive dataset for incumbency data. Therefore, we use data from each country's Wikipedia page to find the election result; if there was a discrepancy, we would check it with the country's election website to ensure the accuracy of the data. Detailed source of each polling data point is listed in the *Source* section in the codebook.

*Government Approval* - In terms of governmental approval data, we use data from various polling agencies and media outlets. Most of the data comes from open source, local media and polling firms. Some come from US based research centers and international polling agencies such as Gallup and Pew, or from multinational polling agencies such as Ipsos. Detailed source of each polling data point is listed in the *Source2* section in the codebook.

### New Variables

*FHR* - We use Freedom House's Freedom in the World (FIW) indicators to determine the level of democracy of a country. Specifically, we use the Political Rights ratings consisting of ratings on electoral process, political pluralism and participation, and functioning of the government, as we believe this is more relevant to outcome of elections than the Civil Liberties score. The scores range from 1 to 7, where a score 1 indicates countries which enjoy a wide range of political rights, including free and fair elections, and a score 7 indicates countries which have few or no political rights.



Map 2: Freedom House Political Rights Ratings of Test Countries in 2015
            Source: FIW in 2015

*GDP* - We used quarterly real GDP data in US Dollars of the each country between 1980 to 2015. We then calculated the quarterly percentage change. The model specifically looks at the year-on-year GDP growth of the quarter before the election. Most data comes from three sources: the IMF database; the Oxford Economics Global Economic Databank, and the national statistics bureau of each country. Most data are quarterly data, and a few countries are available for only annual data.

*CPI (Inflation)* - We used quarterly data of year-on-year CPI change of the election countries during 1980 to 2015. The model specifically looks at the year-on-year CPI growth of the quarter before the election. Most data comes from three sources: the IMF database; the Oxford Economics Global Economic Databank, and the national statistics bureau of each country. Most data are quarterly data, and a few countries are available for only annual data.

*Employment* - The employment growth data is found through the Thomson Reuters data acquisition program DataStream. DataStream is an information source which gathers economic data from individual government statistical surveys such as IMF database, Oxford Economics Global Economic Databank and the national

statistics bureau of each country. The data collected to build the employment growth was monthly employment numbers for each country.  Once that data was acquired, it was converted to quarterly data.  The employment data on a quarterly basis was used in the analysis of country trends prior to elections. The data was measured on a differential basis for the quarter prior to the election.  For countries in which monthly data was not available, annual data was used.

*War* - Conflict data comes from the Correlates of War Militarized Interstate Dispute dataset. The data contains conflicts ranging from threats to begin a war to beginning or entering a war. Previous literature did not use the lowest levels of conflicts in the dataset, so we chose to include only conflicts that resulted in sustained military operations with or against a foreign state. These data eliminate all conflicts coded from 0 (no conflict) through 3 (border violations); conflicts labeled 4 (blockade, occupation of territory, etc) and above remained. Missing years were filled in after researching conflicts, primarily in the years 2010-2015. The final variable is a dummy that indicates whether a country is in a conflict: 0 if no conflict, 1 if in conflict.

## Summary Statistics

| Statistic | Observations | Mean | St. Dev. | Minimum | Maximum |
|---|---|---|---|---|---|
| Current Government Wins | 727 | 51.90% | 0.5 | 0 | 1 |
| Government Approval | 189 | 42.20% | 0.17 | 3% | 87% |
| Incumbency Status | 664 | 53.00% | 0.5 | 0 | 1 |
| Freedom House Rating | 686 | 1.98 | 1.34 | 1 | 7 |
| Country in War | 728 | 10.60% | 0.31 | 0 | 1 |
| Change in Quarterly Employment | 296 | 0.01 | 0.03 | -0.09 | 0.3 |
| Change in Quarterly Inflation | 547 | 0 | 0.51 | -7.37 | 6.21 |
| Growth of Quarterly GDP | 358 | 0.02 | 0.09 | -0.14 | 1.18 |
| Election Type (0=Pres, 1=Parl) | 717 | 68.50% | 0.47 | 0 | 1 |

# IV.  Models

We started with a logistic model created by Ipsos, developed six additional logistic regressions models in order to evaluate the data, and tested those models on our newly created dataset. The first model was developed by Ipsos and provides a baseline for comparison. The second model runs robust standard errors and run the original explanatory variables in our new dataset. The third model tests the validity of our additional variables. The fourth model combines the original model and new variables. The fifth, sixth, and seventh models attempt to impute the government approval rate due to missing observations.

Table 1. Regression Table

| Incumbent party wins election | | | | | | | |
|---|---|---|---|---|---|---|---|
| Model | (1) | (2) | (3) | (4) | (5) | (6) | (7) |
| Gov Approval | 0.07*** | 6.90*** | | 7.72*** | | | |
| | (0.009) | (1.341) | | (1.443) | | | |
| Gov Approval + | | | | | 6.35*** | 7.00*** | 3.44*** |
| | | | | | (1.232) | (1.282) | (0.833) |
| Candidate is incumbent | 1.02*** | 1.29*** | 1.78*** | 1.39*** | 1.67*** | 1.72*** | 1.20*** |
| | (0.262) | (0.366) | (0.308) | (0.393) | (0.260) | (0.273) | (0.193) |
| Country is at war | | | -0.07 | 0.95 | | -0.25 | |
| | | | (0.443) | (0.637) | | (0.445) | |
| Inflation | | | -0.33 | -8.97 | | -1.21 | |
| | | | (0.956) | (6.173) | | (0.919) | |
| Additional Variables | No | No | Yes | No | No | No | No |
| Constant | -3.65*** | -3.54*** | -0.45 | -4.12*** | -3.28*** | -3.59*** | -1.98*** |
| | (0.479) | (0.607) | (0.511) | (0.715) | (0.535) | (0.565) | (0.370) |
| Observations | 383 | 173 | 233 | 164 | 321 | 307 | 497 |
| Pseudo R Squared | 0.203 | 0.232 | 0.147 | 0.284 | 0.189 | 0.212 | 0.0886 |
| Correctly Predicted | 74.41% | 75.14% | 71.24% | 76.22% | 71.96% | 73.62% | 64.99% |
| Prediction US 2016 | 52.21% | 49.54% | 45.81% | 43.32% | 48.90% | 49.11% | 44.41% |
| Prediction US 2016 in War | | | 44.08% | 66.51% | | 42.97% | |

Standard errors in parentheses. Additional Variables include Freedom House Rating, increase in employment the quarter before the election, GDP growth the quarter before the election and election type. Gov Approval + is the government approval with some missing values created through imputation methods. *** p<0.01, ** p<0.05, * p<0.1

## Model 1: Original model with original dataset

Ipsos created a dataset and model to predict the outcome of elections overtime and across countries prior to the beginning of the project. The Ipsos model is a logit regression and relied primarily on government approval rate and whether the election had an incumbent or successor running. The sample covered 85 countries from 1936 to 2015 for a total of 383 elections. This dataset includes state-level

data and years not included in our final model. Equation 1 of Table 1 shows the result of the regression with government approval and incumbency status being statistically significant. The model yielded a pseudo R-squared statistic of 20 percent. Table 2 shows the prediction classification of the model. On average, 74.4% of the elections were correctly predicted

Table 2. Prediction Classification of Model 1

| Observation Classification Model 1- Original model and data | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 184 | 61 | 245 |
| Negative | 37 | 101 | 138 |
| Total | 221 | 162 | 383 |
| Correctly classified | 74.41% | | |

## Model 2: Original model with new 1980-2015 dataset

Our first step in expanding the dataset was to add additional data points to Ipsos' original model and to find whether the model performed differently. In this case, we added additional countries, expanded the years up to 2015, and input more elections. In order to compare across subsequent models, we limited the years to 1980 to 2015 and eliminated state-level elections, reducing the number of observations to 173. This reduced number is mainly because of missing observations in the government approval variable. In models 5, 6 and 7 we use imputation methods to improve our regression.

Previous models in literature were based on national elections, and we want to give an international focus to the data gathered. Model 2 of Table 1 shows the result of this regression using the expanded data set. Government approval and incumbency status still remains statistically significant. Table 3 shows a higher correct classification rate, with 75.1% of the elections correctly predicted.

Table 3. Prediction Classification of Model 2

| Observation Classification Model 2 - Original model with new data | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 70 | 24 | 94 |
| Negative | 19 | 60 | 79 |
| Total | 89 | 84 | 173 |
| Correctly classified | 75.14% | | |

There are two specific measures of classification that improve our understanding of the precision of the model: sensitivity and specificity. Sensitivity is the probability of classifying a positive occurrence as positive. Sensitivity in our model is the probability of predicting that the election will go to the incumbent party given that the incumbent party did win (70/89=78.65%). Specificity is the probability of predicting that the incumbent party will lose given that it lost (60/84=71.43%). Even though our model does well predicting observations in general, it does a better job at predicting positive than negative results.

This classification is done assuming a positive result if the predicted probability is above 50%. A sensitivity and specificity analysis measures the specificity and sensitivity of the model for each probability cutoff. The following graph shows that the probability cutoff that make positive and negative predictions equally powerful is closer to 52%. This implies that the model should classify an incumbent party as winner only if the predicted probability of winning is above 52%.
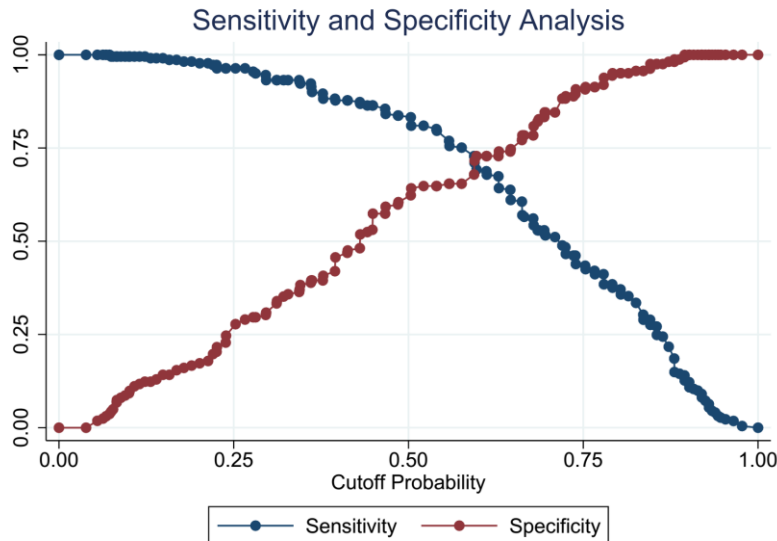


Figure 1. Sensitivity and Specificity Analysis of Model 2

The resulting classification is shown in Table 4. The positive predictions are now slightly less powerful, but we gained a great deal of power from the negative prediction, increasing the model's overall predictive power to 75.72%.

Table 4. Prediction Classification of Model 2 with Cutoff of 52%

| Observation Classification Model 2 - Cutoff = 52% | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 67 | 20 | 87 |
| Negative | 22 | 64 | 86 |
| Total | 89 | 84 | 173 |
| Correctly classified | 75.72% | | |

## Model 3: Additional Variables on full 1980-2015 data

We determined that the model would have the most validity by expanding the data to include additional countries and variables, and to reduce the years to 1980 onward. Additional countries increase the robustness of the data, and ensures that there is sufficient data to analyze in difficult situations, for instance in the case of new democracies.

A review of salient literature, described above, pointed to numerous potential additional variables. To decide which variables would be the best for our model, we looked at two criteria: differentiation between variables, so they would not

measure similar attributes of an election, and availability of data, so we could apply the variable across countries and time. After applying these criteria to the salient variables in the literature review, we determined that four additional variables would be valuable: GDP growth, employment growth, inflation growth, and military conflicts. Some descriptive variables were also necessary in order to categorize countries for further analysis, including government type and Freedom House ratings. Model 3 in Table 1 does not include government approval, demonstrating the validity of these additional variables. Only incumbency status is statistically significant at the 95% level, although the model holds global significance. The pseudo R squared decreased to 14% in this model, with 233 observations. Table 5 shows the prediction power of the model, with 71.2% of elections being predicted correctly.

Table 5. Prediction Classification of Model 3

| Observation Classification Model 3 - New data with new variables | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 100 | 34 | 134 |
| Negative | 33 | 66 | 99 |
| Total | 133 | 100 | 233 |
| Correctly classified | 71.24% | | |

## Model 4: 1980-2015 data and a new model

Our most complete model combines inflation, war status, incumbency status and government approval, the most salient variable. This model achieves our goal of boosting the classification rate up to 76.22%, increasing the pseudo R squared up to 28%, with the drawback of losing some observations. Table 6 shows the prediction classification of the model.

Table 6. Prediction Classification of Model 4

| Observation Classification Model 4 - New data and new model | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 66 | 22 | 88 |
| Negative | 17 | 59 | 76 |
| Total | 83 | 81 | 164 |
| Correctly classified | 76.22% | | |

Conducting a sensitivity and specificity analysis brings us to the conclusion that the best probability cutoff for the prediction is 52%, just as in Model 2. The new election classification is shown in Table 7 and the analysis can be seen in Figure 2.
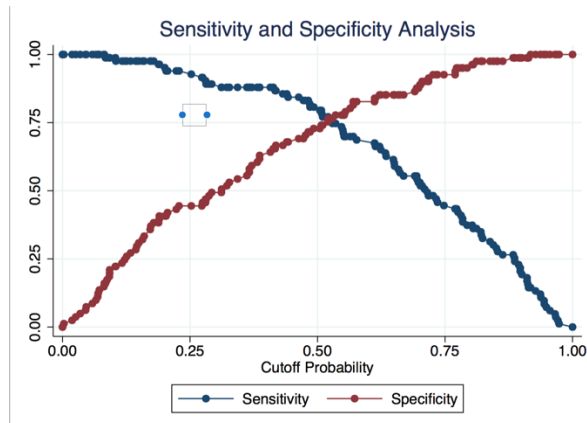
Figure 2. Sensitivity and Specificity Analysis of Model 4

Table 7. Prediction Classification of Model 4 with Cutoff of 52%

| Election classification of Model 4 with cutoff of 52% | | | |
|---|---|---|---|
| **Classified** | **TRUE** | **FALSE** | **Total** |
| **Positive** | 64 | 20 | 84 |
| **Negative** | 19 | 61 | 80 |
| **Total** | 83 | 81 | 164 |
| **Correctly classified** | 76.22% | | |

Table 7 shows a slight increase in the classification power of negative results, and a decrease in classification power of positive results. The prediction power of the model remains the same, but fits the data better.

## Model 5: Imputed government approval - Original Model - Method 1

We next attempt two imputation methods with government approval. Government approval only has information for 189 elections out of the 727. Given the missing data problems, if we find a viable method to impute data on the missing observations of the variables, we might be able to gain more information from other explanatory variables, like incumbency status.

The first method predicts government approval by regressing government approval against 3 country characteristics: war status, change in employment in the last quarter and the type of election. Table 8 shows the result of this regression. Even though war status is marginally significant, employment and election type are strong predictors. The model has global significance (checked through the F statistic), and the R squared is 15%.

Table 8. Government Approval Prediction

| Government Approval | Coeff | St. Error | p-value | [95% Conf. Interval] | |
|---|---|---|---|---|---|
| War Status | 0.056 | 0.034 | 0.106 | -0.012 | 0.123 |
| Increase in Employment | 1.125 | 0.529 | 0.035 | 0.079 | 2.171 |
| Election Type | -0.09 | 0.03 | 0.003 | -0.149 | -0.031 |
| Constant | 44.90% | 0.028 | 0 | 0.394 | 50.40% |

| | |
|---|---|
| Observations | 137 |
| F(3, 133) | 5.65 |
| Prob > F | 0.001 |
| R-squared | 0.1563 |

Based on this regression we predict government approval, and replace the values of the variable for which we have real data. Then we use this variable in the original model. The result can be seen in model 5 of Table 1. The main gain from this model is the increase in the number of observations, which increase to 321. The predictive power goes down slightly to 72% and the R-squared is now 18.9%. The prediction classification for the model can be seen in Table 9.

Table 9. Prediction Classification of Model 5

| Observation Classification Model 5 - New data and original model with input | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 138 | 52 | 190 |
| Negative | 38 | 93 | 131 |
| Total | 176 | 145 | 321 |
| Correctly classified | 71.96% | | |

If we assume that the classification power for the observations of Model 2 is kept untouched, the classification power for the new observations is 68.2%. This model still has substantial predictive power, while covering a much wider range of countries.

## Model 6: Imputed government approval - New Model - Method 1

Following the previous imputation method we also include the explanatory variables of Model 3, war status and inflation. Model 6 in Table 1 show the results. The pseudo R-squared is improved to 21% and the prediction power is improved to 73%. Table 10 shows the prediction classification for the model.

Table 10. Prediction Classification of Model 6

| Observation Classification Model 6 - New data and new model with input | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 134 | 48 | 182 |
| Negative | 33 | 92 | 125 |
| Total | 167 | 140 | 307 |
| Correctly classified | 73.62% | | |

# Model 7: Imputed government approval - Original Model - Method 2

In a final imputation model we calculate means for government approval for different groups. We use Freedom House Rating, war status, election type and the decade of the election (to somewhat control for the trend) to generate the means. Once the means are created, we replace the missing values of the government approval variable. Model 7 of Table 1 shows the result of the regression. Even though the number of observations climbed to 497, the prediction power and pseudo R-squared fell dramatically to 65% and 8.8%, respectively. Table 11 shows the prediction classification for the model.

Table 11. Prediction Classification of Model 7

| Observation Classification Model 7 - New data and original model with input 2 | | | |
|---|---|---|---|
| Classified | TRUE | FALSE | Total |
| Positive | 189 | 100 | 289 |
| Negative | 74 | 134 | 208 |
| Total | 263 | 234 | 497 |
| Correctly classified | 64.99% | | |

Assuming that the classification power of the original model is maintained for the unchanged observations, the percentage of correctly classified observations of the new data is 59%.

## Predictions

To illustrate the strength of our analysis we conducted a prediction for the United States presidential election for 2016, considering the latest economic variables and approval of the current government. The last two rows of Table 1 show the probability of the Democratic Party winning if the election were to occur tomorrow. The last row shows what the probability would be if the country was officially at war. The probability goes from 43.3% to 52.2%, and goes up to 66% if we were at war with Model 4.

# V.   Conclusions and Recommendations

The purpose of this project was to create a robust analysis of a previously analyzed problem but on a global scale.  The value added for this project is that it will increase the predictive power of the Ipsos forecasting model.  We evaluated and enhanced the original model by adding new variables and a more comprehensive data range based on the number of countries' collected and the time span over which those country elections took place.  The enhanced model is more thorough and accurate.  It also covers a much wider portion of the globe.  Additionally, it can help predict the result of the coming election in the United States with greater accuracy than the previous Ipsos model.

The conclusion drawn from this analysis is that it is possible to accurately predict the outcome of elections by considering a selection of factors that have an impact on that outcome.  The analysis finds that there is a higher likelihood for predicting the result of an election if one takes into account the incumbency status of the government party candidate, the current economic variables, the government approval rating and the current military commitment if that country is in a conflict.  This result supports our original hypothesis by confirming that by providing a more comprehensive set of data and a more precise model, we can enhance the accuracy of the election model and provide for greater predictive capability for Ipsos election forecasting.

The results of the analysis tell us that in the original Model 1 received from Ipsos, there is a 74% correct classification rate.  As the added variables and the 35 year time series were included to the expanded data set, the accuracy of Model 2 increased to 75%.  Using the most complete dataset in Model 4 and each of the new exogenous variables for the study, the rate of elections correctly classified jumped to 76%.  Our analysis reveals that our models are more precise than the original Ipsos model (Model 1).  Our model enhances the capacity to reveal False Negatives, while keeping True Positives almost unchanged.  On average, the total percentage of correct predictions increased.  This added precision in the model is due to the additional elections and new variables included, as well as the fact that the new model does not use state level data, it is only national level data that is included.  In addition, according to the specificity and sensitivity analysis in Model 2, we find that the model can accurately classify an incumbent party winner if the predicted probability of winning is at or greater than 52%.

## *Recommendations*

There is room to expand the project further, both in economic analysis and further data collection.  Of the many aspects that can be modified, these are the most important:

- Finding government approval rating data was the greatest challenge in assembling the data for this project, and more data is needed.  There are many countries and years where polling does not exist for government approval.  Additionally, the polling questions of government approval are not standardized metrics throughout the world, and every polling firm, if they exist at all in a country, often uses a slightly different method of

measuring government approval rather than an individual representative number.

- Add to the elections database as more elections take place in the countries listed. Recent and upcoming elections are very likely to include better polling and accurate macroeconomic indicators, and their addition to the dataset will greatly increase the predictive power of the model. It will also be possible to continue to add countries as more democracies develop and individual countries' Freedom House ratings approach a more democratic status.

- The model could also be modified to include other variables and change the way some are analyzed:

  - The macroeconomic data could be changed to look at different time periods, such as the entirety of an incumbent leader's term.
  - Conflict data could be looked at as a scale, rather than a zero to one binary variable.
  - Exogenous variables that affect the outcome of an election could be the effect of a scandal or an increase in sympathy for the prospective candidate, which would be influenced by Lichtman's difficult to quantify "Keys" model.
  - Measuring the corruption index in a country may have an effect on election results.

# Works Cited

Abramowitz, Alan I. "Forecasting the 2008 Presidential Election with the Time-for-Change Model." *PS: Political Science and Politics* 41.4 (2008): 691-95. Web.

Ajzen, Icek. "The Theory of Planned Behavior." *Organizational Behavior and Human Decision Processes* 50.2 (1991): 179-211. Web.

Congleton, Roger D. "Median Voter Model." *Encyclopedia of Public Choice* (2002): n. pag. Print.

Fair, Ray C. "Presidential and Congressional Vote-Share Equations." *American Journal of Political Science* 53.1 (2009): 55-72. Web.

"Freedom House." *|Freedom House Country Ratings.* N.p., n.d. Web. 04 May 2016.

"Global Economic Databank." *Global Economic Databank.* Oxford Economics, n.d. Web. 2 Apr. 2016.

Graefe, Andreas, and J. Scott Armstrong. "Predicting Elections from the Most Important Issue: A Test of the Take-the-best Heuristic." *ScholarlyCommons* (2010): n. pag. Print.

Healy, Andrew J., Neil Malhotrab, and Cecilia Hyunjung Mo. "Irrelevant Events Affect Voters' Evaluations of Government Performance." *Irrelevant Events Affect Voters' Evaluations of Government Performance.* PNAS, 20 July 2010. Web. 04 May 2016.

"IMF Data." *IMF Data.* International Monetary Fund, n.d. Web. 02 Apr. 2016.

Jackman, Simon. "Estimation and Inference Are Missing Data Problems: Unifying Social Science Statistics via Bayesian Simulation." *Political Analysis* 8.4 (2000): 307-32. Web.

Jackman, Simon. "Pooling the Polls over an Election Campaign." *Australian Journal of Political Science* 40.4 (2005): 499-517. Web.

Jennings, W., and C. Wlezien. "Distinguishing Between Most Important Problems and Issues?" *Public Opinion Quarterly* 75.3 (2011): 545-55. Web.

Lewis-Beck, MIchael, and Charles Tien. "The Job of President and the Jobs Model Forecast: Obama for '08?" *PS: Political Science & Politics* 41.4 (2008): 687-90. Print.

Lichtman, Allan J. *The Keys to the White House: Updated Forecast for 2008.* 2008. MS. American University, Washington.

Linzer, Drew A. "Dynamic Bayesian Forecasting of Presidential Elections in the States." *Journal of the American Statistical Association* 108.501 (2013): 124-34. Web.

Montgomery, J. M., F. M. Hollenbach, and M. D. Ward. "Improving Predictions Using Ensemble Bayesian Model Averaging." *Political Analysis* 20.3 (2012): 271-91. Web.

Mueller, John E. "Presidential Popularity from Truman to Johnson." *The American Political Science Review* 64.1 (1970): 18. Web.

Mueller, John E. "Public Expectations of War During the Cold War." *American Journal of Political Science* 23.2 (1979): 301. Web.

Palmer, Glenn, Vito D'Orazio, Michael Kenwick, and Matthew Lane. 2015. "The MID4 Data Set: Procedures, Coding Rules, and Description." Conflict Management and Peace Science. Forthcoming.

"Pollsters May Be Herding." *VOTAMATIC*. Votamatic, 05 Nov. 2012. Web. 04 May 2016.

Wlezien, Christopher, and Robert S. Erikson. "Temporal Horizons and Presidential Election Forecasts." *American Politics Research* 24 (1996): 492. Web.

# Appendix 1: Literature Review Summary Table

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|---|---|---|---|---|---|---|---|
| Abramowitz | Forecasting 2008 | Electoral results can be determined by: Popularity of the incumbent president, the state of the economy and the length of time that the president's party has controlled the white house. Strong Relation between approval and vote choice | Leading Indicators: Growth Rate of the Economy – during the second quarter of the election year. The incumbent president's approval rating at mid-year. The length of time the incumbent's party has controlled the White House (the Time-for-Change Factor) | Incumbent Party Vote | Incumbent Populartiy; Economic Factors, Growth Rate of Economy Approval Rating Length of Time in Office | Time for Change Model | The outcomes found were that the main factors in an election of an incumbent president were the popularity of the incumbent, the condition of the economy, and the lengh of time that the president's party has controlled the government. Outside factors that impact elections are: the rising effect of polarization and race. |
| Allan J. Lichtman | The Keys to the White House: Updated Forecast for 2008 | The Keys to the White House are a historically-based prediction system, which give specificity to the theory that that presidential election results turn primarily on the performance of the party controlling the White House and that politics as usual by the challenging candidate will have no impact on results. The Keys include no polling data and consider a much wider range of performance indicators than economic concerns. | The Keys to the White House show that a pragmatic American electorate chooses a president according to the performance of the party holding the White House as measured by the consequential events and episodes of a term — economic boom and bust, foreign policy successes and failures, social unrest, scandal, and policy innovation. | Winning of the incumbent party | The 13 Keys To The White House : Party Mandate, Contest, Incumbency, Third party, Short-term economy, Long-term economy, Policy change, Social unrest, Scandal, Foreign/military failure, Foreign/military success, Incumbent charismam, Challenger charisma | The Keys are statements that favor the re-election of the incumbent party. When five or fewer statements are false, the incumbent party wins. When six or more are false, the challenging party wins. | This system correctly picked the winner of all six presidential elections from 1984 to 2004, |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|------|-------|----------|-------------|--------------------|---------------------|--------|----------|
| Andrew J. Healya, Neil Malhotrab,1, and Cecilia Hyunjung Mob | Irrelevant events affect voters' evaluations of government performance | Information irrelevent to government performance affects voting behavior, especially before the election day. However, making people more aware of the reasons of the current state of mind will reduce such bias on making voting decisions. | Study1, the successes and failures of the local college football team before Election Day significantly influence the electoral prospects of the incumbent party<br><br>Study2, 2009 NCAA men's college bas-ketball tournament win experienced by respondents significantly increased approval of President Obama's job performance | Study 1, Outcome of Presidential, Gubernatorial, and Senate Elections (incumbent vote share)<br><br>Study 2, Approval of President Obama's job performance | 1, outcome of local football team (win, lose, tie) 10d before election<br><br>2, Outcome of 2009 NCAA men's college bas-ketball tournament | 1, We first performed simple difference of means tests, comparing the change in in-cumbent party vote share between counties in which the football teamwon to counties where the teamlost or tied. We further demonstrate robustness by using point spreads from the betting market, to con-struct an independent variable that isolates the surprise component of game outcomes.<br><br>2, Similar, without betting for surprise component | 1, For games 10 d before the election, a victory increases the incumbent party's vote shareby 1.13percentagepoints (P=0.05).<br><br>2, each additional adjusted win experienced by respondents significantly increased approval of President Obama's job performance, with the effect size being 2.3 percentage points (P= 0.04) |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|------|-------|----------|-------------|-------------------|---------------------|--------|----------|
| Congleton | Median Voter Model | A Median Voter Model is a theoretical model that predicts the policy stance of candidates based on the stances of the population. The model has powerful conclusions based on a very simple economic model. | The basic models has a general finding: whatever the median voter votes for, wins the election. The model has strong implications on policy adoption by political parties. Implications include that policy tends to be moderate, many people are usually displeased by elected policy, and policies are stable over time. Equilibrium in this kind of model is not pareto-efficient. The model runs into problems when the preferences of individuals are more than 1-dimensional, or when transitivity doesnt hold. | No model | | | |
| Drew | Pollsters May be Herfing | The herding behavior of pollsters (The possibility that polling firms, out of fear of being wrong, are looking at the results of other published surveys and weighting or adjusting their own results to match) could make the results of polls collectively biased. | Herding around the wrong value is potentially much worse than any one or two firms having an unusual house effect. But even if the variance of the polls is decreasing, they might still have the right average. | Absolute survey error | Time | Plot the absolute value of the state polls' error, over time. (The error is the difference between a pollís reported proportion supporting Obama, and my modelís estimate of the ìtrueî population proportion.)<br><br>Herding would be indicated by a decline | The underlying trend reveals that the average error in the polls started at 1.5% in early May, but is now down to 0.9%. |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|------|-------|----------|-------------|--------------------|----------------------|--------|----------|
| | | | | | in the average survey error towards zero ñ representing no difference from the consensus mean ñ over the course of the campaign. | | |
| Fair | Presidential congressional | Previous studies have found that economic variables affect voting behavior in presidential elections. This paper argues that economic variables also affect House elections. | Economic variables do affect vote-shares in presidential, on-term house elections, and midterm house elections | V = Democratic share (in presidential, on-term House, and midterm House votes) | I = incumbent president, DPER = Democratic incumbant president running again (-1 if Rep, 0 otherwise), DUR = 0 if either party in White House for one term, 1 if Democratic party in WH for 2 consecutive (-1 for Reps) and so on, WAR = US at war, G = GDP growth rate per capita in first three quarters of election year, P = absolute value of GDP growth rate, Z = number of quarter in the first 15 quarters in which the GDP growth is > 3.2%. G, P, and Z seperately include mid-term House vote variables | Three vote-share models are developed by the author, one each for presidential, on-term House elections, and midterm elections. Each equation estimates a voters utility function based on the dependent variables. | The economic variables affect all three models, and they are roughly the same. There is no evidence of presidential "coat-tail" effects; the election of a president and House representatives of the same party is by correlation, not causation. There is a positive serial correlation in the House vote share models from the previous election's vote-share, possibly due to an incumbant effect. Finally, the presidential vote share has a negative effect on the next midterm House vote share; the author notes that voters may prefer a balanced government system, and don't want one party to become too dominant. |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|---|---|---|---|---|---|---|---|
| Icek Ajzen | Theory of Planned Behavior | A review of various aspects of the theory of planned behavior and discussion about some unresolved issues. | Intention, perception of behavioral control, attitude toward the behavior, and subjective norm each reveals a different aspect of the behavior and each can serve as a point of attack in attempts to change it. The underlying foundation of beliefs provides the detailed descriptions needed to gain substantive information about a behavior's determinants.<br><br>Unsolved issues: the exact form of relations between behavioral beliefs and attitudes toward the behavior, between normative beliefs and subjective norms, and between control beliefs and perceptions of behavioral control is still uncertain. | 1 Behavior<br><br>2 Intention | 1 Intention, Perceived behavioral control<br><br>2 Attitudes toward the behavior, subjective norms with respect to the behavior, and perceived control | Empirical Findings | Attitudes toward the behavior, subjective norms with respect to the behavior, and perceived control over the behavior are usually found to predict behavioral intentions with a high degree of accuracy. In turn, these intentions, in combination with perceived behavioral control, can account for a considerable proportion of variance in behavior. |
| **Name** | **Title** | **Argument** | **Key Finding** | **Dependent Variable** | **Independent Variable** | **Method** | **Outcomes** |
| J. Scott Armstrong and Andreas Graefe | Predicting Elections from the Most Important Issue: A Test of the Take- | The big-issue (BI) voting model to predict the outcome of U.S. presidential elections. The model is based on information about how | The take-the-best heuristic generated accurate forecasts based on voters' | The actual two-party popular vote share received by the candidate of the incumbent party (V) | The incumbent party candidate's 3-day rolling average of voter support on the most important issue. (S) | **BI-H Model**: (1) Identify the issue seen as most important by voters, (2) calculate the two-party | **BI-H Model:**Over all 1,000 forecasts, BI-H correctly predicted the winner 88% of the times.**BI-M Model:**Over all 1,000 forecasts, BI-M correctly predicted the winner 97% of the times. |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | the-Best Heuristic | voters expect the candidates to deal with the issue seen as most important. The model relies solely on information from polls and uses a heuristic similar to Take-the-Best Heuristic (TTB) to determine the winner of the popular vote. The goal was to develop a model that can provide fast advice on which issues candidates should stress in their campaign. | perceptions on how the candidates will handle the single most important issue facing the country. | | | shares of voter support for the candidates on this issue and average them for the last three days, and (3) predict the candidate with the higher voter support to win the popular vote.**BI-M Model** :A simple linear regression to build the BI model (BI-M). V = 27.0 + 0.50 * S | Combining forecasts of BI-H and BI-M was expected to further increase accuracy |
| **Name** | **Title** | **Argument** | **Key Finding** | **Dependent Variable** | **Independent Variable** | **Method** | **Outcomes** |
| Jackman | Pooling the Polls | Pooling polls has the advantage of improving the precision of polls, but this might lead to issues due to bias from house effects. The author develops a statistical model to correct house effect. | House effects are significant and most likely dependent on method for interviewing. These effects can't be calculated without the election result. Recommendation to use this effects to calibrate real-life forecasts or only use phone interviews. | Variable that we are looking for: Vote share and house effects. | Variables used to find it: results from polls. | Bayesian Simulation Methods | Predictions are much more precise, house effects are significant and dependent on method of questioning, results can be used to improve real time forecasting. |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|---|---|---|---|---|---|---|---|
| Jennings and Wlezien | Distinguishing Between Most Important Problems and Most Important Issues: AAPOR 2011 | the MII and the MIP are highly correlated | Comparison between the "Most Important Problem" and the "Most Important Issue." Gallup has run surveys since the 1930's asking people what their most critical issues of importance are in determining who is elected. They determined a list of 18 different topics which were most critical. Independent Variable, percentage of respondents identifying the particular category as the most important problem or the most important issue. | Most Important Issue | Most Important Problem | this is a survey which measures the relative importance of the difference between asking individuals in a survey format about critical items dealing with politics by using the words "problem" or the word "issue." Traditionally, polling has taken place using "Problems." This article measure the results by using the word "Issues." | The results found in this article show that the MII and the MIP are highly correlated |
| John Mueller | Presidential Popularity | This investigation has applied multiple regression analysis to the behavior of the responses to the Gallup Poll's Presidential popularity question in the 24 year period from the Truman to the Johnson administration. | 1) Each President will experience in each term a general decline of popularity; 2) that this decline will be interrupted from time to time with temporary upsurges associated with international crises and similar events; 3) that the decline will be accelerated in direct relation to increases in unemployment rates over those prevailing when the | Presidential Popularity | 1, coalition of minorities 2, Rally round the flag 3, Economic Slump 4, War | multiple regression. (1)The first time without the "war" variable, adding dummy variable for each president, (2)and the second time with all four variables to see the effect of Korea and Vietnam war on the popularity. | The fit of the resulting equation was very good: it explained 86 percent of the variance in presidential popularity. But the effect of each variable are very different. |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | President began his term, but that improvement in unemployment rates will not affect his popularity one way or the other. 4) the president will experience an additional loss of popularity if a war is on. | | | | |
| **Name** | **Title** | **Argument** | **Key Finding** | **Dependent Variable** | **Independent Variable** | **Method** | **Outcomes** |
| John Mueller | Public Expectations | American people's anticipations of world war during the cold war period, specifically from after the WW2 till 1962. | 1) , Generally, American's expectation of world war is affected mostly by belligerent behaviors by the two blocs, and the "mood" of the American public was briefly sensitive to short-term crises, but most determined by longer term (but not long-term) forces; 2) Educational level also account for expectation of war, well educated people are more optimistic of war in the short term (a few years), but differences diminish greatly. | Expectation by American public of major war | 1) Time (short term and long term), 2)Incident of International conflicts (refractory behaviour and reconcilatory behavior by both western and soviet union ) 3) Educational level | 1) The study uses both the polling data and the Gamson-Modigliani variable to run a mutiple regression to see the major impact of major refractory and reconcilatory behaviors of the two sides in a certain period of time before a poll on the poll result of war expectiation. 2) Then they seperate the people in three education levels, college and higher, high school level, and grade school level, use dummy variables for each educational group, | (1)belligerent behavior: the total amount of belligerent behavior (Western plus Soviet) occurring in the two months before the poll was taken has the strongest relation towards war expectation; (2)the American popular expectations of war react favorably to evidence of consistent conciliatory behavior on the part of the Soviet bloc given at least a 2-month period before the poll is taken; (3) differences are great in shorter period (5 year less), more educated people are less likely to expect a war, however in longer term, difference diminishes or even reverse. |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|---|---|---|---|---|---|---|---|
| | | | | | and time period from half year to 50 years and life time's expectation of war | | |
| Lewis-Beck/Tien | The Job of President and the Jobs Model Forecast: Obama for '08? | There is a high correlation between economic performance during a president's term, his approval rating, and whether he gets re-elected or his successor is elected. The two variables are broken down into five sub-variables that are broadly applicable. | The Jobs Model predicts an Obama two-party popular vote forecast of 50.1% (actual: 52.9%) | Vote share | Presidential popularity, economic growth (GNP), elected president running or not, job growth, incumbant status | OLS | This model has correctly picked winners for the past few elections, with the exception of the 2000 election, which forecast a Gore win (by popular vote, however, the model was correct). Each variable is significant. |
| Linzer | Dynamic Bayesian Forecasting | Dynamic Bayesian forecasting is superior to structural models because it updates in real-time as new polls arrive. The author applies this method to individual states, giving it | The author shows how a sequence of state-level polls can be used to estimate both current voter preferences and forecasts of the election outcome for | Election outcome - based on electoral votes | Pre-election polls | Polls are input into a Bayesian model, which updates its forecast based on the added information. This information is added to structural | The model was tested on the 2008 election, and was able to predict an Obama victory two months prior to the election. The model estimated nearly 100% chance of victory, with electoral vote counts in a range of 338 to 387; the actual electoral vote count was 365. |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | superior prediction of electoral college votes. | each state on any given day, regardless of whether a poll was conducted on that day. | | | forecasts in order to create a combined model. | |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|---|---|---|---|---|---|---|---|
| Mongetomery, Hollenbeck, and Ward | Improving predictions | Ensemble Bayesian methods are the best forecasting models, as they combine multiple diverse models into a coherent whole model. | Presidential election forecasting is greatly improved by combining multiple models in Ensemble Bayesian Model Averaging (EBMA). | Model prediction accuracy. | Six distinct forecasting models developed by academics: Campbell, Abromowitz, Hibbs, Fair, Lewis-Beck/Tien, EWT2C2. All but Hibbs were simple regression models. | EBMA creates forecasts by creating weighted averages of component predictions, or component predictive probability distribution functions (PDFs). The weight assigned to each component forecast reflects two aspects of the components' past forecasts. First, ceteris paribus, the EBMA model will give greater weight to forecasts that were more accurate in the past. Second, ceteris paribus, it will assign a greater weight to models that made unique (but correct) predictions. | The EBMA model outperformed any one compenent model. In the presidential election portion of the paper, it outperformed by having more correct predictions and, therefore, lower MSE than any component model. |
| Simon Jackman | Estimation and Inference Are Missing Data Problems: Unifying | Bayesian models improve estimations by being able to calculate "auxiliary variables", like missing data. This advantage is | We can use bayesian models to track auxiliary variables. | Example about incumbacy estimation: Vote share | Previous winning party and party affiliation | OLS, Maximum Likelihood and Bayesian Methods | Republican incumbents receive a bigger advantage than democrats. Bayesian methods say that incumbency advantage is the same. |

| Name | Title | Argument | Key Finding | Dependent Variable | Independent Variable | Method | Outcomes |
|---|---|---|---|---|---|---|---|
| | Social Science Statistics via Bayesian Simulation | analyzed in several contexts and examples. | | | | | |
| Wlezien and Erikson | Temporal Horizons | Using Leading economic indicators substantially improves electoral forecast | The two variables used are Cumulative Per Capita Income and Presidential Approval. At the Election Time, The net economic effect is the most important issue in determining elections. Each quarter's income growth rate is weighted 1.25 times the one before. Income growth is more important toward the end of a President's term than in the beginning. Measures of the leading indicators may capture economic performance better than simple summary indicators taken after the fact. Using Leading economic indicators substantially improves electoral forecast. | Incumbent Party Vote | LEIG: Leading Economic Indicators Growth Income Growth (W : The Nation's Economic Welfare) Presidential Approval | Regression Analysis: | The measure of leading indicators and growth in income capture economic performance better than simple summary indicators. Using leading economic indicators substantially improves the electoral forecast |

# Appendix 2: Codebook

| Code | Definition of Code |
|---|---|
| Country | Country |
| ISO3 | ISO 3166-1 alpha-3 codes 3-letter country codes |
| Year | Year (1980-2015) |
| FHR | Freedom House Rating (Political Rights Scores) |
| Region | Region of the World |
| Location | Election Level (State, National) |
| Election Winner | Who won the election? |
| Party (winner) | The winning party of the election that year |
| Incumbent Party | The party in office the time of the election |
| Incumbent Leader | The leader in office the time of the election |
| First Round/Second Round | 1 = First Round of Voting, 2 = Second Round of Voting |
| name.govt.cand | Name of the candidate supported by the current government. Usually the name of the incumbent or the successor. |
| name.cand1 | Candidate 1 is the person with largest support who is NOT the government candidate |
| name.cand2 | Candidate 2 is the person with the second largest number of votes who is NOT the government candidate |
| name.cand3 | Candidate 3 is the person with the third largest number of votes who is NOT the government candidate |
| name.cand4 | Candidate 4 is the person with the 4th largest number of votes who is NOT the government candidate |
| vote.govt.cand | Percentage of votes for government supported candidate |
| vote.cand1 | Percentage of votes for candidate 1 |
| vote.cand2 | Percentage of votes for candidate 2 |
| vote.cand3 | Percentage of votes for candidate 3 |
| vote.cand4 | Percentage of votes for candidate 4 |
| date.election | What date did the election occur? |
| election.quarter | The quarter of the year that the election occur |
| govapproval | Government or candidate approval rating before election took place - For Example, in 1992 President Bush was the incumbent and the candidate, ideally, we would want President Bush's approval rating in October 1992. In the 2000 Presidential Election, we would want President Clinton's approval rating in October 2000, even though Al Gore was the candidate. |

| Excellent/Good | Alternative way to measure approval for countries that measure it based on "excellent or good" instead of "approval", most applicable for Brazil. |
|---|---|
| approval.date | When was this approval rating reported? It should be BEFORE the election. |
| Incumbant_or_successor_wins | The Incumbent Party OR Incumbent Candidate in office wins election 1 = yes |
| Sucessor (new candidate from the party in office running) | 1= yes, there is a successor for the government party, instead of an incumbent. In Presidential election in 2000, Al Gore was the democratic successor to the party currently in office. |
| Incumbent (Candidate in office running for re-election) | Is there an incumbent in the race? 1= yes, In 2008 Presidential elections, there was no incumbent in the race, but there was a successor. |
| Government in office on ballot? | Is the government in office currently on the ballot? 1 =yes |
| Government in Office Supporting a canddiate | Is the party or candidate currently in office known to support any of the candidates? 1 = yes, For Example, President Bill Clinton support Al Gore for the Presidency in 2000 election OR President Obama supported his own election to the presidency in 2012. Both would be 1. |
| New Democracy - First Elections after Non-democractic government | Is this the first time a democratic election was held in the country? 1 = yes |
| source | Where did you get the information to support approval ratings and vote counts? |
| source 2 | Where did you get the information to support approval ratings and vote counts? |
| Notes Section | Any pertinent information |
| etype | election type: presidential or parliamentary |
| GDP | quarterly GDP data of year-on-year growth, the model specifically looks at GDP growth a quarter before the election |
| Inflation | quarterly CPI data of year-on-year growth, the model specifically looks at CPI change a quarter before the election |
| Employment | quarterly employment data of year-on-year growth, the model specifically looks at employment change a quarter before the election |
| War | a dummy that indicates whether a country is in a conflict: 0 if no conflict, 1 if in conflict |

# Appendix 3: Do Files

## 1. Data organization and file merges

```
clear all
cd "C:\Users\puyan\Desktop\Data"

import excel Source\electiond.xlsx, sh(Data) first case(l)
drop if location!="National"
sort country year

*Elimination Duplicates
duplicates tag country year, g(tag)
drop if tag>0 & incumbentparty==""
drop if tag==2 & incumbentleader==""
duplicates tag country year, g(tag2)
drop if tag2==1 & fhr==.
drop tag tag2
destring incumbant_or_successor_wins, force replace

save data1.dta, replace

*Merge with annual employment
import excel Source\employmenta.xlsx, sh(data) first case(l) clear
destring employ, force replace
sort country year
bysort country: gen increase=employ-employ[_n-1]
bysort country: gen y_employ=increase*100/employ[_n-1]
label var y_employ "Increase in Annual Employment (%)"
drop if year==2016
keep iso3 year y_employ
save data2.dta, replace

use data1.dta, clear
merge 1:1 iso3 year using data2.dta
drop if _merge==2
drop _merge

save data1.dta, replace

*Merge with quarterly employment
import excel Source\employmentq.xlsx, sh(data) first case(l) clear
destring q1 q2 q3 q4, force replace
bysort country: replace q1=. if year==1980
sort country year
forvalues i=1/4{
      bysort country: gen increaseq`i'=q`i'-q`i'[_n-1]
      bysort country: gen q`i'_employ=increaseq`i'*100/q`i'[_n-1]
      label var q`i'_employ "Increase in Quarterly Employment (%, Q`i')"
}
drop if year==2016
keep iso3 year q*employ q4
rename q4 empq4
save data2.dta, replace

use data1.dta, clear
merge 1:1 iso3 year using data2.dta
drop if _merge==2
drop _merge

save data1.dta, replace

*Merge with quarterly inflation
import excel Source\inflationq.xlsx, sh(data) first case(l) clear
destring q1 q2 q3 q4, force replace
```

```
bysort country: replace q1=. if year==1980
sort country year
forvalues i=1/4{
        bysort country: gen q`i'_inflation=q`i'-q`i'[_n-1]
        label var q`i'_inflation "Increase in Quarterly Inflation (%,
Q`i')"
}
drop if year==2016
keep iso3 year q*inflation
save data2.dta, replace

use data1.dta, clear
merge 1:1 iso3 year using data2.dta
drop if _merge==2
drop _merge

save data1.dta, replace

*Merge with quarterly gdp
import excel Source\gdpq.xlsx, sh(data) first case(l) clear
destring q1 q2 q3 q4, force replace
bysort country: replace q1=. if year==1980
sort country year
forvalues i=1/4{
        rename q`i' q`i'_gdp
        label var q`i'_gdp "Quarterly GDP Growth (%, Q`i')"
}
drop if year==2016
keep iso3 year q*gdp
save data2.dta, replace

use data1.dta, clear
merge 1:1 iso3 year using data2.dta
drop if _merge==2
drop _merge

save data1.dta, replace

*Merge with annual freedom house rating
import excel Source\fhra.xlsx, sh(data) first case(l) clear
destring rating, force replace
sort country year
rename rating fhra
label var fhra "Freedom House Rating (1-7)"
drop if year==2016
keep iso3 year fhra
save data2.dta, replace

use data1.dta, clear
merge 1:1 iso3 year using data2.dta
drop if _merge==2
drop _merge
drop fhr
rename fhra fhr

save data1.dta, replace

*Merge with war data
import excel Source\midb.xlsx, sh(data) first case(l) clear
destring war, force replace
sort country year
label var war "Country is at war"
drop if year==2016
keep iso3 year war
save data2.dta, replace

use data1.dta, clear
merge 1:1 iso3 year using data2.dta
replace war=0 if _merge==1 & war==.
drop if _merge==2
drop _merge
```

```
replace war=1 if war>1 & war!=.

save data1.dta, replace

*Create decade variable
drop if year<1980
gen decade=1 if year<=1989
replace decade=2 if year>=1990 & year<=1999
replace decade=3 if year>=2000 & year<=2009
replace decade=4 if year>=2010 & year<2019

*Merge PennTables employment
merge 1:1 iso3 year using Source\pwt81.dta, keepus(emp)
drop if _merge==2

replace empq4=empq4/1000000
gen emp1=emp if emp!=.
replace emp1=empq4 if emp==.
drop empq4 emp
rename emp1 emp
bysort country: gen emp_a=(emp[_n-1]-emp[_n-2])*100/emp[_n-2]
drop emp
rename emp_a emp

drop if electionwinner==""

*Create appropriate variable with quarter
destring electionquarter, force replace

gen employment_quarter=.
forvalues i=1/4{
        replace employment_quarter=q`i'_employ if electionquarter==`i'
}

gen inflation_quarter=.
forvalues i=1/4{
        replace inflation_quarter=q`i'_inflation if electionquarter==`i'
}
replace inflation_quarter=. if inflation_quarter>1000 |
inflation_quarter<-1000

gen gdp_quarter=.
forvalues i=1/4{
        replace gdp_quarter=q`i'_gdp if electionquarter==`i'
}

gen result_gov=1 if partywinner== incumbentparty
replace result_gov=0 if result_gov!=1
label var result_gov "Incumbent party wins"

gen candidate=1 if incumbent==1
replace candidate=1 if namegovtcand==incumbentleader & incumbent==.
replace candidate=0 if sucessor==1 & candidate!=1
replace candidate=0 if namegovtcand=="" & candidate!=1
label var candidate "Candidate is incumbent"

replace etype="" if etype=="#N/A"
encode etype, gen(etype2)
drop etype
rename etype2 etype
recode etype (2=0)

recode govapproval (0.44=44)

*Change in units
replace govapproval=govapproval/100
replace employment_quarter=employment_quarter/100
replace inflation_quarter=inflation_quarter/100
replace gdp_quarter=gdp_quarter/100
```

```
*Create US observation
local new = _N + 1
set obs `new'
replace year=2016 in `new'
replace country="United States" in `new'
replace iso3="USA" in `new'
replace govapproval=0.51 in `new'
replace candidate=0 in `new'
replace fhr=1 in `new'
replace war=1 in `new'
replace employment_quarter=0.0211 in `new'
replace inflation_quarter=0.01 in `new'
replace gdp_quarter=0.014 in `new'
replace etype=0 in `new'


save data1.dta, replace
erase data2.dta

*Summary statistics
tabstat result_gov govapproval candidate fhr war employment_quarter
inflation_quarter gdp_quarter etype, s(n mean sd min max) format(%9.3f
```

## 2. Regressions, predictions and tables

```
clear all
cd "C:\Users\puyan\Desktop\Data"

***************************************
* Model 1 - Original Model and Data *
***************************************
use dataoriginal.dta, clear

local new = _N + 1
set obs `new'
replace year=2016 in `new'
replace country="United States" in `new'
replace govapproval=51 in `new'
replace candidate=0 in `new'

logit result_gov govapproval candidate
predict prob1
outreg2 using Results\Table.xls, se bdec(2) excel replace e(r2_p) label
estat clas
sum prob1 if country=="United States" & year==2016

***************************************
* Model 2 - Original Model, New Data *
***************************************

use data1.dta, clear

logit result_gov govapproval candidate, robust
outreg2 using Results\Table.xls, se bdec(2) excel e(r2_p) label
predict probability1
estat clas
lsens, graphr(c(white)) title("Sensitivity and Specificity Analysis")
xtitle("Cutoff Probability") ytitle("")
gr export Results\SS_M2.pdf , replace
window manage close graph

*Generating predictions
gen prediction1=1 if probability1>=0.5 & result_gov==0 & probability1!=.
replace prediction1=0 if probability1<0.5 & result_gov==1 &
probability1!=.
tab prediction1

label var probability1 "Predicted P Model 2"
label var prediction1 "Predicted Result Model 2
```

```
        sum probability1 if country=="United States" & year==2016


        ***********************************
        * Model 3 - New Model, New Data *
        ***********************************

        logit result_gov candidate fhr war employment_quarter inflation_quarter
        gdp_quarter etype, robust
        outreg2 using Results\Table.xls, se bdec(2) excel e(r2_p) label
        predict probability2
        estat clas

        label var probability2 "Predicted P Model 3"

        sum probability2 if country=="United States" & year==2016


        ****************************************
        * Model 4 - New Model 2, New Data *
        ****************************************

        logit result_gov govapproval candidate war inflation_quarter, robust
        outreg2 using Results\Table.xls, se bdec(2) excel e(r2_p) label
        predict probability3
        estat clas
        lsens, graphr(c(white)) title("Sensitivity and Specificity Analysis")
        xtitle("Cutoff Probability") ytitle("")
        gr export Results\SS_M4.pdf , replace
        window manage close graph

        *Generating predictions
        gen prediction3=1 if probability3>=0.5 & result_gov==0 & probability3!=.
        replace prediction3=0 if probability3<0.5 & result_gov==1 &
        probability3!=.
        tab prediction3

        label var probability3 "Predicted P Model 4"
        label var prediction3 "Predicted Result Model 4"

        sum probability3 if country=="United States" & year==2016

        ****************************************
        * Determinants of government approval *
        ****************************************

        reg govapproval war employment_quarter etype, robust
        predict p_govapproval
        gen p_govapproval2=govapproval if govapproval!=.
        replace p_govapproval2=p_govapproval if govapproval==.

        *Model 5 - Predicted Government Approval and Original Model
        logit result_gov p_govapproval2 candidate, robust
        outreg2 using Results\Table.xls, se bdec(2) excel e(r2_p) label
        predict probability4
        estat clas
        lsens, graphr(c(white)) title("Sensitivity and Specificity Analysis")
        xtitle("Cutoff Probability") ytitle("")
        gr export Results\SS_M5.pdf , replace
        window manage close graph

        gen prediction4=1 if probability4>=0.5 & result_gov==0 & probability4!=.
        replace prediction4=0 if probability4<0.5 & result_gov==1 &
        probability4!=.
        tab prediction4

        label var probability4 "Predicted P Model 5"
        label var prediction4 "Predicted Result Model 5"

        sum probability4 if country=="United States" & year==2016
```

```
*Model 6 - Predicted GA and Model 4
logit result_gov p_govapproval2 candidate war inflation_quarter, robust
outreg2 using Results\Table.xls, se bdec(2) excel e(r2_p) label
predict probability5
estat clas

sum probability5 if country=="United States" & year==2016

***************************
* Input through averages *
***************************
bysort war fhr etype decade: egen a_gov=mean(govapproval)
gen a_govapproval=govapproval if govapproval!=.
replace a_govapproval=a_gov if govapproval==.

logit result_gov a_govapproval candidate, robust
outreg2 using Results\Table.xls, se bdec(2) excel e(r2_p) label
sortvar(govapproval a_govapproval p_govapproval2 candidate)
predict probability6
estat clas

sum probability6 if country=="United States" & year==2016
```