

Benford's Law: Fraud Analysis of 2016 U.S. Presidential Election Data

Introduction

This report presents an analysis of the 2016 U.S. Presidential Election results using **Benford's Law**, a mathematical principle often applied to detect irregularities in large datasets, including financial statements and elections. Benford's Law provides a probabilistic framework for detecting anomalies by examining the frequency distribution of the first significant digits of numbers in naturally occurring data.

The primary objective of this report is to investigate whether the 2016 election results, recorded at the county level across the United States, adhere to Benford's Law. Deviations from the expected distribution may signal potential issues, such as vote manipulation or reporting errors, requiring further scrutiny.

Dataset Overview

The dataset contains 3,112 rows representing U.S. counties and 14 columns covering vote totals for the presidential elections of 2008, 2012, and 2016. Below is a summary of the key features included:

Key Features and Descriptions

- `fips_code`: A unique identifier for each county (used to identify and track counties geographically).
- `county`: Name of the county.
- Election Data Columns:
 - For each election year (2008, 2012, 2016), the following data points are available:
 - `total_{YY}`: Total votes cast in that year.
 - `dem_{YY}`: Votes received by the Democratic candidate.
 - `gop_{YY}`: Votes received by the Republican (GOP) candidate.
 - `oth_{YY}`: Votes received by other or independent candidates.

Rationale for Using Benford's Law

Benford's Law states that in naturally occurring datasets, the leading digits are distributed in a non-uniform pattern, with smaller digits (e.g., 1, 2, 3) appearing more frequently as the first significant digit than larger ones (e.g., 8, 9). Specifically, the probability $P(d)P(d)P(d)$ of a digit `ddd` (where `ddd` ranges from 1 to 9) appearing as the first digit is calculated as:

$$P(d) = \log_{10} \left(\frac{1}{1+d} \right) \quad P(d) = \log_{10} \left(1 + \frac{1}{d} \right)$$

In a fraud-free environment, vote counts should approximately follow this distribution, given the assumption that vote totals arise from natural processes. Deviations from Benford's distribution can indicate anomalies that may require further investigation, such as vote tampering or reporting errors.

Scope of the Analysis

The 2016 election is our primary focus, as we aim to assess whether any anomalies exist in the county-level vote data. The following vote categories are used for analysis:

- `total_2016`: Total votes cast in 2016.
- `dem_2016`: Votes for the Democratic candidate in 2016.
- `gop_2016`: Votes for the Republican candidate in 2016.
- `oth_2016`: Votes for other or independent candidates in 2016.

These features are suitable for Benford's Law analysis as they represent large, naturally varying numbers where human manipulation or irregularities might leave detectable traces.

Methodology

1. Data Preparation

- We extract the first significant digit from the relevant vote columns in the 2016 dataset.
- Any negative or zero values are excluded from the analysis, as Benford's Law applies only to positive, naturally occurring data.

2. Expected vs. Observed Distributions

- **Expected Distribution:** Calculated using Benford's formula for digits 1 through 9.
- **Observed Distribution:** The actual frequencies of the first digits from the 2016 vote data are computed.

3. Chi-Square Test for Goodness-of-Fit

- A **chi-square goodness-of-fit test** is performed to determine whether the observed distribution significantly deviates from the expected Benford distribution.
- **Hypotheses:**
 - Null Hypothesis (H_0): The observed data conforms to Benford's Law.
 - Alternative Hypothesis (H_1): The observed data does not conform to Benford's Law.
- **Significance Level (α):** 0.05
 - A **p-value < 0.05** suggests a statistically significant deviation, indicating potential anomalies.

Results and Interpretation

1. **Chi-Square Statistic and P-Value**
 - The chi-square statistic measures the difference between observed and expected frequencies.
 - The p-value indicates the probability of observing the given distribution if the data were consistent with Benford's Law.
2. **Found Results**
 - **Chi-Square Statistic:** 9.80
 - **P-Value:** 0.2793
3. **Interpretation of the Results**

P-Value Interpretation: Since the p-value (0.2793) is greater than 0.05, we fail to reject the null hypothesis that the data follows Benford's Law distribution. This suggests no significant anomalies or irregularities in the 2016 election results. The data aligns closely with the expected distribution, providing no evidence of fraud or manipulation.
4. **Visualization**
 - A bar chart displays the expected Benford distribution.
 - A line chart overlays the empirical distribution from the 2016 election data, allowing for a visual comparison.
 - Interpretation of Charts: Any significant visual deviation between the two distributions might indicate anomalies.

Conclusion

Based on the Benford's Law analysis, the 2016 U.S. Presidential election data shows no significant deviations from the expected distribution, as indicated by the chi-square test ($p = 0.2793$). This suggests that the election results align with typical patterns, reinforcing the integrity of the data. However, it is important to emphasize that Benford's Law is a probabilistic tool—while it can flag irregularities, deviations alone are not definitive proof of fraud or manipulation.

Recommendations and Next Steps

1. Further Investigations (If Needed):

- **Manual Audits:** If future analyses reveal anomalies, conducting manual audits in affected counties can help verify the reported vote counts.
- **Complementary Statistical Tests:** Additional techniques such as Z-tests, cluster analysis, or regression models could offer deeper insights and validate the findings.
- **Historical Comparison:** Comparing election data from other years (e.g., 2008, 2012) may help identify trends and determine if deviations are consistent over time or unique to 2016.

2. Caution in Interpretation:

- **Benford's Law as an Indicator:** While it is a powerful tool to detect irregularities, Benford's Law should not be used in isolation to make conclusive judgments. Results should always be considered alongside other quantitative and qualitative evidence, such as audits, surveys, and election processes.

Limitations of Benford's Law

- Not Definitive: A deviation from Benford's distribution does not confirm fraud; it only flags areas for further investigation.
- Data Size and Structure: Benford's Law performs best with large datasets covering diverse data ranges. In cases with small sample sizes or constrained values (e.g., capped vote counts), results may be less reliable.

Final Thoughts

This report provides a statistical assessment of the 2016 U.S. Presidential election data using Benford's Law. The findings suggest no evidence of manipulation in the analyzed data, supporting the election's overall integrity. However, statistical results should always be interpreted with caution—they are one part of a broader assessment. To draw well-founded conclusions, the insights from this analysis must be combined with other investigative methods, audits, and context-specific knowledge. The report serves as a starting point for further investigations, offering transparency and accountability in the analysis of electoral data.