

Datasheet for ‘US Voter File Dataset’

Hailey Jang

April 1, 2024

The ‘US Voter File Dataset’ includes detailed information about American voters, featuring their demographic profiles, voting histories, and socio-economic data. Sourced from a private entity, this comprehensive dataset awaits application in specific projects. It’s ripe for exploration in political science research and analysis of voter behavior, offering insights into election predictions and voter demographics. Users must handle this dataset with care to respect privacy concerns and legal standards concerning voter data.

Motivation 1. *For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.* - The dataset was Crafted to enhance the 2020 US Cooperative Election Study, this dataset enriches demographic accuracy at an individual level by incorporating voter file information from a private source. This enhancement supports in-depth voter behavior and preference studies, addressing a notable gap in understanding diverse electorate traits. 2. *Who created the dataset (for example, which team, research group) and on behalf of which entity (for example, company, institution, organization)?* - TA dedicated team in political science and data analytics compiled this dataset. Due to confidentiality and legal agreements, the team and their representing body, including the data-providing company, remain anonymous. 3. *Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.* - The dataset’s development was likely supported through academic and private sector grants focused on political research. Exact details on funders and grants are withheld for privacy and legal reasons. 4. *Any other comments?* - It’s crucial to acknowledge the dataset’s creation and use within strict legal and ethical boundaries, emphasizing transparent reporting on its limitations and ethical considerations for responsible data management.

Interpretation

1. *What do the instances that comprise the dataset represent?*

- The dataset captures individual US voters, including demographics, voting patterns, affiliations, locations, and socio-economic factors.

2. *How many instances are there in total?*
 - The dataset’s volume is tied to the breadth of voter files received, with the exact count pending.
3. *Does the dataset contain all possible instances or is it a sample of instances from a larger set?*
 - Representing a subset of the broader voter base, this dataset aims to reflect the larger US voter landscape through a strategic sampling approach.
4. *What data does each instance consist of?*
 - Comprising raw voter file data, each entry includes detailed demographic and socio-economic information, along with voting histories.
5. *Is there a label or target associated with each instance?*
 - While it may include labels or targets, such as voter behaviors or demographics, these specifics are yet to be finalized.
6. *Is any information missing from individual instances?*
 - Some entries might lack complete information, attributable to gaps in voter records or data availability constraints.
7. *Are relationships between individual instances made explicit?*
 - The dataset focuses on individual attributes without explicitly linking voter relationships or networks.
8. *Are there recommended data splits?*
 - Suggestions for training, validation, and testing splits may be provided, tailored to analysis goals and methodological rigor.
9. *Are there any errors, sources of noise, or redundancies in the dataset?*
 - Inherent voter record inaccuracies or data collection flaws may exist, with quality assurance efforts in place to mitigate these.
10. *Is the dataset self-contained or does it rely on external resources?*
 - While primarily self-contained, the dataset could leverage external data for enrichment or validation purposes.
11. *Does the dataset contain data that might be considered confidential?*
 - Given its inclusion of personal voter information, the dataset navigates strict privacy and ethical guidelines.

12. *Does the dataset contain data that might be considered offensive, threatening, or might otherwise cause anxiety?*
 - Political data’s delicate nature necessitates careful, ethical handling, with no expectation of offensive material.
13. *Does the dataset identify any sub-populations?*
 - It potentially delineates sub-groups based on demographic criteria, though sampling and representativeness validation vary.
14. *Is it possible to identify individuals, either directly or indirectly from the dataset?*
 - There’s a risk of indirect identification, prompting de-identification efforts to protect privacy.
15. *Does the dataset contain sensitive data?*
 - Containing sensitive political and demographic data, it underscores the need for strict data protection measures.

Collection process

1. *How was the data associated with each instance acquired?*
 - Data was gathered through surveys conducted by the providing company, focusing on voter demographics and histories.
2. *What mechanisms or procedures were used to collect the data?*
 - Utilizing a survey platform, data collection involved direct input from participants, overseen by human curators for quality assurance.
3. *If the dataset is a sample from a larger set, what was the sampling strategy?*
 - Aimed at representing the US electorate broadly, the sampling employed probabilistic techniques for demographic and geographic diversity.
4. *Who was involved in the data collection process?*
 - Comprising trained surveyors and analysts, the team’s efforts were compensated, though details remain confidential.
5. *Over what timeframe was the data collected?*
 - Data gathering spanned months before the 2020 election, aligning with its relevance to electoral studies.
6. *Were any ethical review processes conducted?*

- Internal ethical reviews were likely conducted to align with privacy standards, though specifics are not disclosed.
7. *Did you collect the data from the individuals in question directly?*
 - Individuals directly provided their data via surveys, ensuring firsthand accuracy and relevance.
 8. *Were the individuals in question notified about the data collection?*
 - Survey participants were informed about the data collection aims and privacy measures at the outset.
 9. *Did the individuals in question consent to the collection and use of their data?*
 - Voluntary survey completion indicated consent, with detailed statements provided, though not specifically depicted here.
 10. *If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent?*
 - Mechanisms for consent revocation were presumably in place, allowing participants to opt-out or remove their data upon request.
 11. *Has an analysis of the potential impact of the dataset and its use on data subjects been conducted?*
 - An internal impact analysis on subjects might have been performed to identify and mitigate privacy risks, details of which are undisclosed.
 12. *Any other comments?*
 - The dataset’s creation emphasized ethical data use and adherence to privacy laws, with a call for greater process transparency.

Preprocessing/cleaning/labeling

1. *Was any preprocessing/cleaning/labeling of the data done?*
 - Data underwent rigorous cleaning and standardization, including deduplication and missing value handling, to ensure integrity for analysis.
2. *Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data?*
 - Both raw and processed data versions are maintained, offering flexibility for future research applications.
3. *Is the software that was used to preprocess/clean/label the data available?*
 - Data processing utilized R and RStudio, with potential access to methodologies through open-source platforms for reproducibility.

4. *Any other comments?*

- Emphasizing data accuracy and usability, the processing phase is critical for supporting robust research findings.

Uses

1. *Has the dataset been used for any tasks already?*

- To date, the dataset remains untapped for specific research or analysis tasks, preserved for internal purposes.

2. *Is there a repository that links to any or all papers or systems that use the dataset?*

- No public repository exists due to its proprietary status, limiting external access and collaboration opportunities.

3. *What (other) tasks could the dataset be used for?*

- The dataset's rich content offers vast possibilities for political science studies, including electoral prediction and voter trend analysis.

4. *Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?*

- Future utilization must navigate the dataset's proprietary nature and sampling concerns to ensure analytical validity.

5. *Are there tasks for which the dataset should not be used?*

- Its sensitive nature and ownership terms dictate cautious use, particularly avoiding privacy intrusions or legal infractions.

6. *Any other comments?*

- Responsible use, guided by transparency, ethical standards, and legal compliance, is paramount for leveraging the dataset's value in political research.