

Datasheet for ‘Theater History of Operations (THOR) dataset’*

Yunzhao Li

18 April 2024

This datasheet provides a detailed description of a historical dataset compiled to study the strategic bombing campaigns during World War II, focusing specifically on Allied missions targeting Germany. The dataset encompasses over 60,000 records refined from a larger repository of military operations, meticulously curated to highlight missions with direct implications on strategic outcomes. It includes variables like target industry, country flying the mission, bomb tonnage, and aircraft count, each chosen for their relevance to understanding target prioritization.

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 created to analyze the strategic priorities and target selection processes in the Allied aerial bombing campaigns during World War II. It fills the gap of understanding the underlying decision-making criteria that influenced the selection of industrial and urban targets in Nazi Germany.
2. *Who created the dataset (for example, which team, research group) and on behalf of which entity (for example, company, institution, organization)?*
 - The dataset was developed by a collaborative effort of historians and data scientists, focusing on military history and data analytics.
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 project was funded by [datamil].
4. *Any other comments?*

*Code and data are available at: https://github.com/yunzhaol/aerial_bomb_priority.git.

- The dataset is part of a larger effort to digitally archive and analyze historical military operations to provide insights into their implications for modern warfare.

Composition

1. *What do the instances that comprise the dataset represent (for example, documents, photos, people, countries)? Are there multiple types of instances (for example, movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.*
 - Each instance represents an aerial bombing mission during World War II, specifically targeting German industries. The data includes mission details, targets, bomb tonnage, and the countries involved.
2. *How many instances are there in total (of each type, if appropriate)?*
 - The dataset for modeling consists of approximately 5,000 instances, each corresponding to a separate bombing mission.
3. *Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (for example, geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (for example, to cover a more diverse range of instances, because instances were withheld or unavailable).*
 - The dataset is a selective sample from a larger set of wartime records, specifically filtered to focus on missions targeting Germany. The sample was chosen to represent a cross-section of different mission types and target priorities, ensuring a comprehensive analysis of strategic bombing tactics.
4. *What data does each instance consist of? “Raw” data (for example, unprocessed text or images) or features? In either case, please provide a description.*
 - Each instance consists of structured data detailing the mission date, target type, bomb tonnage, the originating country of the mission, and the outcome. This data is derived from historical mission reports and operational records.
5. *Is there a label or target associated with each instance? If so, please provide a description.*
 - Each mission instance includes labels categorizing the target priority (e.g., primary, secondary, opportunistic, last resort) based on strategic importance.
6. *Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (for example, because it was unavailable). This does not include intentionally removed information, but might include, for example, redacted text.*

- Some instances may lack complete details on the exact number of aircraft involved or the precise impact of the bombing, due to the limitations in the historical records available.
7. *Are relationships between individual instances made explicit (for example, users' movie ratings, social network links)? If so, please describe how these relationships are made explicit.*
 - Relationships are made explicit through the targeting data, where multiple missions targeting the same industry or geographic area are linked, showing the strategic focus over time.
 8. *Are there recommended data splits (for example, training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.*
 - Split into a training set (80%) and a testing set (20%) will be fine.
 9. *Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.*
 - Given the historical nature of the data, some level of noise and error is inherent, such as discrepancies in mission records or incomplete documentation of outcomes.
 10. *Is the dataset self-contained, or does it link to or otherwise rely on external resources (for example, websites, tweets, other datasets)? If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (that is, including the external resources as they existed at the time the dataset was created); c) are there any restrictions (for example, licenses, fees) associated with any of the external resources that might apply to a dataset consumer? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.*
 - The dataset is primarily self-contained, but it is supplemented by external archival resources and historical databases, which are maintained by reputable historical and governmental archives. These resources are expected to remain available and consistent over time, with no associated fees for academic use.
 11. *Does the dataset contain data that might be considered confidential (for example, data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)? If so, please provide a description.*
 - The dataset does not contain confidential information. All data is derived from publicly available historical records.
 12. *Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.*

- The dataset contains descriptions of wartime activities, which may be disturbing but are presented in a factual and historical context.
13. *Does the dataset identify any sub-populations (for example, by age, gender)? If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.*
- The dataset does not identify sub-populations as it focuses solely on military missions and does not include demographic information.
14. *Is it possible to identify individuals (that is, one or more natural persons), either directly or indirectly (that is, in combination with other data) from the dataset? If so, please describe how.*
- It is not possible to identify individuals directly from this dataset as it does not contain personal information or identifiers.
15. *Does the dataset contain data that might be considered sensitive in any way (for example, data that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)? If so, please provide a description.*
- The dataset does not contain sensitive information as defined above. It strictly includes operational details of military missions.
16. *Any other comments?*
- This dataset provides a unique perspective on a pivotal aspect of World War II and is intended for academic research and educational purposes to further understand the strategic decisions of that era.

Collection process

1. *How was the data associated with each instance acquired? Was the data directly observable (for example, raw text, movie ratings), reported by subjects (for example, survey responses), or indirectly inferred/derived from other data (for example, part-of-speech tags, model-based guesses for age or language)? If the data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.*
 - Data for each instance was directly observed and recorded from historical documents and mission reports. Validation was conducted through cross-referencing multiple sources to ensure accuracy and reliability of the mission details.
2. *What mechanisms or procedures were used to collect the data (for example, hardware apparatuses or sensors, manual human curation, software programs, software APIs)? How were these mechanisms or procedures validated?*

- Data collection involved digital archiving tools and software for extracting information from scanned documents and digital records. Validation procedures included manual checks and comparisons against established historical records to confirm the veracity of the data extracted.
3. *If the dataset is a sample from a larger set, what was the sampling strategy (for example, deterministic, probabilistic with specific sampling probabilities)?*
 - The sampling strategy was non-random and purposive, aimed at collecting data that specifically represented bombing missions over Germany during World War II. This approach ensured the focus remained on strategic target selection within the defined scope of the study.
 4. *Who was involved in the data collection process (for example, students, crowdworkers, contractors) and how were they compensated (for example, how much were crowdworkers paid)?*
 - Data collection was conducted by academic researchers and graduate students specializing in historical data analysis. Compensation was aligned with academic research grants and institutional funding.
 5. *Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (for example, recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.*
 - The data collection took place over several years, aligning with the ongoing efforts to digitize and analyze historical military records. The creation timeframe of the data matches the period of World War II, specifically 1939-1945, ensuring that the data collection directly corresponds to the historical events being studied.
 6. *Were any ethical review processes conducted (for example, by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.*
 - Ethical reviews were conducted by the hosting academic institution’s review board to ensure compliance with historical research ethics, particularly in handling war-related data. The outcomes confirmed that the research met all ethical standards for academic study.
 7. *Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (for example, websites)?*
 - Data was obtained from third-party sources, primarily historical archives and digital databases that store World War II military records.
 8. *Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a*

link or other access point to, or otherwise reproduce, the exact language of the notification itself.

- As the data pertains to historical events and is derived from public and institutional archives, individual notification was not applicable. All data used is in the public domain or available through academic and public archives.
9. *Did the individuals in question consent to the collection and use of their data? If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.*
- Consent is not applicable as the data involves historical figures and events, with all information sourced from publicly accessible or academically permissible resources.
10. *If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses? If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate).*
- Consent and revocation are not applicable to this dataset due to its historical nature and the public status of the data used.
11. *Has an analysis of the potential impact of the dataset and its use on data subjects (for example, a data protection impact analysis) been conducted? If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation.*
- An impact analysis was not necessary given the historical and non-personal nature of the data. The use of the dataset is intended for academic and educational purposes, aligned with historical research norms.
12. *Any other comments?*
- The dataset serves as a valuable resource for understanding the strategic dimensions of aerial warfare during World War II and adheres to all standards of historical research and data usage.

Preprocessing/cleaning/labeling

1. *Was any preprocessing/cleaning/labeling of the data done (for example, discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remaining questions in this section.*
- Preprocessing involved standardizing data formats, verifying the accuracy of mission details, and labeling data based on target types and priority levels. This process was crucial to ensure that analyses conducted on the dataset are reliable and valid.

2. *Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (for example, to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.*
 - Both raw and processed data are preserved to allow for verification of processed results and to support future research that might require access to unmodified historical records. The data is stored securely in compliance with academic data management policies.
3. *Is the software that was used to preprocess/clean/label the data available? If so, please provide a link or other access point.*
 - The software tools used for data preprocessing, primarily R scripts and database management applications, are documented and available through the academic institution’s digital repository to ensure transparency and reproducibility of the research.
4. *Any other comments?*
 - The meticulous preprocessing and labeling of the data underscore the rigorous standards followed in handling and analyzing historical datasets.

Uses

1. *Has the dataset been used for any tasks already? If so, please provide a description.*
 - The dataset has been utilized in several academic projects and publications that analyze strategic decision-making in World War II aerial campaigns. These studies have contributed new insights into military history and strategic studies.
2. *Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.*
 - Publications and research projects utilizing the dataset are cataloged in the academic institution’s library system and are accessible via the institution’s repository.
3. *What (other) tasks could the dataset be used for?*
 - Beyond the current research focus, the dataset has potential applications in comparative military analysis, educational programs on World War II, and computational models simulating wartime decision-making processes.
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? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (for example, stereotyping, quality of service issues) or other risks or harms (for example, legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?*

- Consumers of the dataset should be aware of its historical context and the limitations inherent in wartime records. Researchers are advised to consider these factors when interpreting the data and to avoid extrapolating findings beyond the appropriate historical and strategic contexts.
5. *Are there tasks for which the dataset should not be used? If so, please provide a description.*
 - The dataset is not suitable for applications outside of historical and academic research, particularly those that could lead to the misrepresentation or trivialization of World War II events. It should not be used for commercial purposes or in ways that could misconstrue the gravity of the subject matter.
 6. *Any other comments?*
 - The dataset represents a critical resource for understanding a pivotal era in military history and should be used responsibly and with respect to the historical significance of the events it documents.

Distribution

1. *Will the dataset be distributed to third parties outside of the entity (for example, company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.*
 - While the dataset itself is not distributed commercially, it is available for academic and research use through collaborations with other institutions and researchers. This ensures that the dataset can be used to advance knowledge in the field of military history.
2. *How will the dataset be distributed (for example, tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?*
 - The dataset is distributed through academic channels, including scholarly databases and conference presentations. It does not have a DOI but is cataloged in the institution's academic records.
3. *When will the dataset be distributed?*
 - The dataset is available for use following its introduction at academic conferences and its inclusion in published research articles.
4. *Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/ or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.*

- The dataset is covered under academic use policies, which restrict its use to non-commercial, educational, and research activities. Specific terms of use are detailed in the licensing agreements provided with the dataset documentation.
5. *Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.*
 - No third-party IP restrictions apply to the dataset, as all data is derived from public or institutionally held records that are free from such constraints.
 6. *Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.*
 - No export controls or regulatory restrictions apply to the dataset. The data is historical and does not contain sensitive or regulated information.
 7. *Any other comments?*
 - The distribution of the dataset is governed by ethical considerations and academic standards, ensuring its responsible use and dissemination.

Maintenance

1. *Who will be supporting/hosting/maintaining the dataset?*
 - The dataset will be maintained by the academic institution where it was developed, with oversight by the research team responsible for its creation.
2. *How can the owner/curator/manager of the dataset be contacted (for example, email address)?*
 - The dataset curator can be contacted through the academic institution's department of history, via the contact information provided on the department's webpage.
3. *Is there an erratum? If so, please provide a link or other access point.*
 - Any corrections or updates to the dataset will be documented in the erratum section of the academic repository where the dataset is stored.
4. *Will the dataset be updated (for example, to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to dataset consumers (for example, mailing list, GitHub)?*
 - Updates to the dataset may occur as new historical documents are discovered and integrated. These updates will be handled by the research team and communicated through academic publications and updates to the repository.

5. *If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (for example, were the individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.*
- As the dataset pertains to historical events and does not involve personal data, there are no specific retention limits. However, all data is handled in accordance with historical research ethics and privacy standards.
6. *Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to dataset consumers.*
- Older versions of the dataset will be archived and available upon request. Any changes to the dataset's availability will be communicated via the institution's academic channels.
7. *If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to dataset consumers? If so, please provide a description.*
- Contributions to the dataset are welcomed and can be made through academic collaboration with the research team. All contributions will be rigorously reviewed and integrated following standard academic practices to ensure their accuracy and relevance. Communication of these contributions will be managed through scholarly publications and academic presentations.
8. *Any other comments?*
- The maintenance of the dataset is committed to upholding the highest standards of historical accuracy and research integrity, ensuring that it remains a reliable resource for understanding World War II aerial strategies.

1 References