

**Datasheet for climate-health mortality collected at Iganga -Mayuge HDSS in Uganda  
under Lacuna Project**

We present the climate and health datasheet created by INSPIRE Network at APHRC, Kenya. We followed the dataset datasheet framework developed by Gebru et al. (2021).

<b>Motivation</b>	
Who created the dataset?	The dataset was created by a team from Iganga-Mayuge Health and Demographic Surveillance Site (IMHDSS) in Uganda.
For what purpose was the data set created?	The dataset was created to understand the causes of deaths in Iganga Mayuge district in Uganda
Was there a specific task in mind?	Yes, we wanted to estimate the impact of climate change on mortality and use data science methods to predict the effect of climate change on health outcomes.
Who funded the creation of the dataset?	This work was conducted with the support of the Lacuna Fund
<b>Composition</b>	
What do the instances that comprise the dataset represent?	The instances contain information on deaths that occurred among residents within the Health and Demographic Surveillance System (HDSS) study area.
How many instances are there in total (of each type, if appropriate)?	The dataset comprises 5,054 observations detailing the causes of death in the Iganga-Mayuge district.
Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?	Yes, the dataset represents a sample of instances from a larger population.
What data does each instance consist of? “Raw” data or features?	The data consists of verbal autopsy data and climate data. Each data point is accompanied by attributes; gender, birthdate, deathdate, age, deathyear, birthyear, codmethod, ICD10causeofdeath, ICD10codes, General CoD, causeofdeath, climate record date, average

	temperature, maximum temperature. Minimum temperature, and precipitation.
Is there a label or target associated with each instance? If so, please provide a description.	No
Is any information missing from individual instances?	No
Are relationships between individual instances made explicit?	No
Are there recommended data splits (for example, training, development/validation, testing)?	We do not specify any data splits.
Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.	None
Is the dataset self-contained, or does it link to or otherwise rely on external resources?	The verbal autopsy data has been linked to climate data.
Does the dataset contain data that might be considered confidential?	No
Does the dataset contain data that, if viewed directly, might be of- offensive, insulting, threatening, or might otherwise cause anxiety?	No
<b>Collection Process</b>	
How was the data associated with each instance acquired?	The dataset was collected through Verbal Autopsy procedures, which involve interviewing surviving relatives or caregivers of the deceased to gather

	detailed information about the circumstances and events leading to the death.
What mechanisms or procedures were used to collect the data?	Verbal Autopsy procedures, which involve interviewing surviving relatives or caregivers of the deceased to gather detailed information about the circumstances and events leading to the death.
If the dataset is a sample from a larger set, what was the sampling strategy?	The dataset is not from a larger set.
Who was involved in the data collection process?	The Iganga Mayuge Data Team
Over what timeframe was the data collected?	January 2007 to December 2022
Were there any ethical review processes conducted (for example, by an institutional review board)?	The study was approved by the Makerere University School of Public Health Internal Review Board and the Uganda National Council of Science and Technology.

### **Preprocessing, Cleaning, and Labelling**

Was any pre-processing/cleaning/labelling of the data done	Yes.
Was the “raw” data saved in addition to the pre-processed/cleaned/ labelled data (for example, to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.	None
Is the software that was used to pre-process/clean/label the data available? If so, please provide a link or other access point.	We used R programming language for data preprocessing.

<b>Uses</b>	
Has the dataset been used for any tasks already? If so, please provide a description.	Yes, we have used the dataset.
Is there a repository that links to any or all papers or systems that use the dataset?	No
What (other) tasks could the dataset be used for?	The dataset can be utilized for building time-series models and conducting trend analysis.
Is there anything about the composition of the dataset or the way it was collected and pre-processed/cleaned/labelled that might impact future uses?	No
<b>Distribution</b>	
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.	Yes. The data provided in the Harvard dataverse represents a random 10% of the entire dataset. If you require access to the complete dataset, please follow this link for more information: <a href="https://microdataportal.aphrc.org/index.php/catalog/165">https://microdataportal.aphrc.org/index.php/catalog/165</a>
How will the dataset be distributed (for example, tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?	The dataset and its associated metadata are hosted on the Harvard Dataverse, an open-source data repository. However, the version available there represents only 10% of the full dataset. For access to the complete dataset, please follow this link for more information: <a href="https://microdataportal.aphrc.org/index.php/catalog/165">https://microdataportal.aphrc.org/index.php/catalog/165</a> The dataset has been assigned a Digital Object Identifier (DOI) for easy referencing: <a href="https://doi.org/10.7910/DVN/DHDNIC">https://doi.org/10.7910/DVN/DHDNIC</a>

When will be the dataset be distributed?	Upon approval, you will receive a link to download a compressed file containing the dataset in Stata format.
Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?	The dataset is licensed under the CC BY license
Have any third parties imposed IP-based or other restrictions on the data associated with the instances?	No
Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?	Yes. The data privacy and protection regulations in the host country.

#### **Maintenance**

Who will be supporting/hosting/maintaining the dataset?	The dataset will be maintained by the INSPIRE Network. The team will support, host, and maintain the dataset.
How can the owner/curator/manager of the dataset be contacted (for example, email address)?	The dataset owner can be contacted via email (INSPIRE Network).
Is there an erratum?	No
Will the dataset be updated (for example, to correct labelling errors, add new instances, delete instances)?	All updates to the dataset will be documented and communicated via the INSPIRE Network.
Will older versions of the dataset continue to be supported/hosted/ maintained? If so, please describe how.	Yes, previous versions will be stored on the APHRC microdata portal.

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?	Researchers interested in contributing to the dataset can reach out to INSPIRE Network for further discussions.
--	---

## References

Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Iii, H. D., & Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM*, 64(12), 86-92.