

# 2019 StackOverflow Developer’s Survey Datasheet\*

Happy Developer Happy Company

Ray Wen

28 April 2022

## Abstract

This datasheet contains information regarding the dataset used by the report ‘Happy Developer Happy Company’.

Extract of the questions from Gebru et al. (2021)

## 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 purpose of this dataset is to gain insight into the Developer’s community. Specifically, this datasheet was created to investigate the factors that influence a developer’s satisfaction with their current occupation. It contains several factual information regarding the developer’s current job status as well as their subjective opinion on their preferred job setting
2. *Who created the dataset (for example, which team, research group) and on behalf of which entity (for example, company, institution, organization)?*
  - This dataset was created by Ray Wen while studying at the University of Toronto
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.*
  - No funding was received for the creation of this dataset
4. *Any other comments?*
  - None

## 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 row consists of a single survey response done by a user of the StackOverflow Community
2. *How many instances are there in total (of each type, if appropriate)?*
  - There are a total of 20,000 instances in this dataset
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).*
  - This dataset is a cleaned, processed, and sampled dataset from the original 2019 StackOverflow Developer’s Survey, which contained nearly 90,000 instances
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.*

---

\*Code and data are available at: [LINK](#).

- Each instance contains the responder’s responses to the survey, including demographic, subjective, and factual information regarding their developer’s career, involvement in the StackOverflow online forum, and their knowledge of programming
5. *Is there a label or target associated with each instance? If so, please provide a description.*
    - There is no label or target associated with each response. Each response was only given an index number
  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.*
    - The survey from StackOverflow allowed some questions to be skipped, so there are a few missing entries that were eventually discarded.
  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.*
    - There is no explicit information in the dataset that is in relation to each individual instances
  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.*
    - None
  9. *Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.*
    - Unknown
  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.*
    - This dataset linked to the dataset StackOverflow has published in 2019 on <https://insights.stackoverflow.com/survey/2019#overview>
  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.*
    - No
  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.*
    - No
  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.*
    - Yes, this dataset does identify each instances by age group, gender, and ethnicity.
  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.*
    - No
  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.*
    - No
  16. *Any other comments?*
    - None

## 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.

- The survey was distributed by StackOverflow on their website. They have made the survey accessible through a variety of media such as ads and direct messages.
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?*
    - The survey was an online website hosted by StackOverflow.
  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 dataset was sampled uniformly without replacement from the large dataset.
  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)?*
    - Unknown.
  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 survey was open from January 23rd, 2019 to February 14th, 2019.
  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.*
    - Unknown
  7. *Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (for example, websites)?*
    - The data was from a survey on a website.
  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.*
    - Unknown
  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.*
    - Unknown
  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).*
    - Unknown
  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.*
    - Unknown
  12. *Any other comments?*
    - None

### 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.*
  - Yes, missing values were removed. Neutral responses were also treated as missing values and removed.
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.*
  - Yes, the raw data is accessible on [https://stackoverflow.com/users/prediction-data?\\_ga=2.188917099.1567384284.1650887000-373250572.1643712221](https://stackoverflow.com/users/prediction-data?_ga=2.188917099.1567384284.1650887000-373250572.1643712221)

3. *Is the software that was used to preprocess/clean/label the data available? If so, please provide a link or other access point.*
  - Yes, the software used to preprocess/clean the data is R. It is accessible on <https://www.r-project.org/>
4. *Any other comments?*
  - None

## Uses

1. *Has the dataset been used for any tasks already? If so, please provide a description.*
  - Unknown.
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.*
  - Yes, the repository is [https://github.com/ray0130/job\\_satisfaction](https://github.com/ray0130/job_satisfaction)
3. *What (other) tasks could the dataset be used for?*
  - The dataset can be used to gain more information on the developer's behavior, technological knowledge, and activity on StackOverflow.
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?*
  - Unknown
5. *Are there tasks for which the dataset should not be used? If so, please provide a description.*
  - No
6. *Any other comments?*
  - No

## 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.*
  - No, but the raw data is distributed on StackOverflow's Website
2. *How will the dataset be distributed (for example, tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?*
  - No
3. *When will the dataset be distributed?*
  - Unapplicable
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.*
  - No, but the raw data is under the Open Database License (ODbL)
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.*
  - Unknown
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.*
  - Unknown
7. *Any other comments?*
  - None

## Maintenance

1. *Who will be supporting/hosting/maintaining the dataset?*
  - Ray Wen
2. *How can the owner/curator/manager of the dataset be contacted (for example, email address)?*
  - Ray can be contacted via email on ray.wen@mail.utoronto.ca
3. *Is there an erratum? If so, please provide a link or other access point.*
  - None
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)?*
  - No
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.*
  - Unknown
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.*
  - Unknown
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.*
  - Others who wish to extend/augment/build on/contribute to the dataset can access the data through the repository or download the raw data on StackOverflow
8. *Any other comments?*
  - None

## References

Gebru, Timnit, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. “Datasheets for Datasets.” *Communications of the ACM* 64 (12): 86–92.