

# DistilKaggle: A Distilled Dataset of Kaggle Jupyter Notebooks

Mojtaba Mostafavi Ghahfarokhi  
Department of Computer Engineering  
Sharif University of Technology  
m.mostafavi@sharif.edu

Mohammad Abolnejadian  
Department of Computer Engineering  
Sharif University of Technology  
mohammad.abolnejadian@sharif.edu

Arash Asgari  
Department of Computer Engineering  
Sharif University of Technology  
arash.asgari@sharif.edu

Abbas Heydarnoori\*  
Department of Computer Engineering  
Sharif University of Technology  
heydarnoori@sharif.edu

## ABSTRACT

Jupyter notebooks have become indispensable tools for data analysis and processing in various domains. However, despite their widespread use, there is a notable research gap in understanding and analyzing the contents and code metrics of these notebooks. This gap is primarily attributed to the absence of datasets that encompass both Jupyter notebooks and extracted their code metrics. To address this limitation, we introduce DistilKaggle, a unique dataset specifically curated to facilitate research on code metrics in Jupyter notebooks, utilizing the Kaggle repository as a prime source. Through an extensive study, we identify thirty-four code metrics that significantly impact Jupyter notebook code quality. These features such as *lines of code cell*, *mean number of words in markdown cells*, *performance tier of developer*, etc., are crucial for understanding and improving the overall effectiveness of computational notebooks. The DistilKaggle dataset which is derived from a vast collection of notebooks constitutes two distinct datasets: (i) *Code Cells and Markdown Cells Dataset* which is presented in two CSV files, allowing for easy integration into researchers' workflows as dataframes. It provides a granular view of the content structure within 542,051 Jupyter notebooks, enabling detailed analysis of code and markdown cells; and (ii) The *Notebook Code Metrics Dataset* focused on the identified code metrics of notebooks. Researchers can leverage this dataset to access Jupyter notebooks with specific code quality characteristics, surpassing the limitations of filters available on the Kaggle website. Furthermore, the reproducibility of the notebooks in our dataset is ensured through the code cells and markdown cells datasets, offering a reliable foundation for researchers to build upon. Given the substantial size of our datasets, it becomes an invaluable resource for the research community, surpassing the capabilities of individual Kaggle users to collect such extensive data. For accessibility and transparency, both the dataset

and the code utilized in crafting this dataset are publicly available at <https://github.com/ISE-Research/DistilKaggle>.

## CCS CONCEPTS

• **Software and its engineering** → **Software libraries and repositories**.

## KEYWORDS

Open dataset, Kaggle, Jupyter notebooks, Code metrics, Code quality

### ACM Reference Format:

Mojtaba Mostafavi Ghahfarokhi, Arash Asgari, Mohammad Abolnejadian, and Abbas Heydarnoori. 2024. DistilKaggle: A Distilled Dataset of Kaggle Jupyter Notebooks. In *21st International Conference on Mining Software Repositories (MSR '24)*, April 15–16, 2024, Lisbon, Portugal. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3643991.3644882>

## 1 INTRODUCTION

Jupyter notebooks have emerged as the predominant coding environment for data scientists [20]. They offer numerous advantages over traditional software development environments, including seamless documentation, code sharing, result analysis, visual and intuitive code development, and the ability to compile, execute, and modify code cells without re-executing the entire notebook. The quality and comprehensibility of implementations in computational notebooks are pivotal for various purposes, such as education for imparting coding best practices to future data scientists [29], notebook reusability [8], and maintainability [6, 31, 33]. Recent studies have concentrated on enhancing the quality of Jupyter notebooks through improved documentation [13, 22, 27, 28, 33], notebook structure [10, 25, 32], and managing notebook variants and revisions [10, 24]. Many of these applications and researches use various datasets of computational notebooks.

Jupyter notebooks possess several features distinguishing them from other programming environments, such as interactive programming, code cells, markdown cells, the possibility of arranging cell executions in different orders within a notebook, and the availability of each cell's output right after the code, and so on. In many cases, these features can overshadow script-based coding styles, such as independent Python scripts. For instance, in Jupyter notebooks, the result of each executed cell is attached to it, motivating developers to use commands to produce outputs. In addition to code cells, Jupyter notebooks include another type of cell called

\*Abbas Heydarnoori is currently affiliated with the Department of Computer Science at Bowling Green State University, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MSR '24, April 15–16, 2024, Lisbon, Portugal

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-0587-8/24/04...\$15.00  
<https://doi.org/10.1145/3643991.3644882>

a markdown cell, where developers can explain each part of their code. This feature incentivizes developers to enhance the quality of their code documentation by utilizing HTML tags in markdown cells.

The widespread adoption of notebooks in recent years, along with their distinctive features, has motivated our commitment to delivering a comprehensive and current collection for researchers to utilize with enhanced efficiency and ease. This dataset comprises notebooks paired with features designed to assess code quality, as elaborated upon in this paper.

Quaranta et al. [21] introduced KGTorrent, a dataset comprising 248,761 Jupyter notebooks, totaling over 180GB in size. Despite its utility for the community, we identified several shortcomings in this dataset, leading us to present DistilKaggle. The motivations behind introducing DistilKaggle are outlined below:

- (1) The latest Jupyter notebook in the KGTorrent dataset dates back to October 2020, and the code provided by the authors to refresh the dataset no longer functions due to changes in Kaggle's notebook downloading policies. In contrast, the notebooks in this paper were published on the Kaggle platform between October 2020 and October 2023.
- (2) Downloading notebooks in a quantity comparable to the dataset presented here would take months due to Kaggle API's request rate limit. Our dataset offers a convenient access to a substantial volume of valuable data—Jupyter notebooks.
- (3) In contrast to the substantial 180GB size of the KGTorrent dataset, our dataset has been streamlined to a manageable size of 3GB by organizing notebooks into two dataframes for code cells and markdown cells, enhancing download efficiency and facilitating a swift utilization. Despite this reduction, full reproducibility is maintained through the provided dataframes. Still, the complete notebooks dataset is provided upon request.
- (4) A common step in works on code quality involves extracting code quality metrics. After an extensive review of previous works, we identified 34 static code quality features, like lines of code cell, mean number of words in markdown cells, performance tier of developer. These metrics are extracted from all the notebooks in our dataset and presented in the CSV format. We anticipate that these dataset metrics will significantly save researchers' time and effort. Additionally, for works requiring notebooks with specific characteristics, the features dataframe empowers researchers to create custom subset data by filtering based on their study's requirements.

In this paper, we collected a set of notebooks' features that either were proposed by previous studies for static code analysis or were presented as being effective for code quality by the studies focused on Jupyter notebooks. In the next step, we extracted these features from our notebooks' dataset. In summary, we made the following contributions:

- A dataset of Jupyter notebooks crawled from the Kaggle platform.
- Code for crawling and scrapping the notebooks.
- A dataset of features extracted from each notebook.
- Code for feature extraction.

## 2 DATASET

There are different platforms with large datasets of notebooks that provide a large amount of information about the notebooks and their creators. Each of these platforms also has a different level of popularity among scholars and data scientists, which may increase the validity of the studies that are conducted based on their data. Also, based on the aim of the studies, the researchers need a set of information that is only available on a small number of these platforms.

Kaggle, a subsidiary of Google, is an online community of data scientists and machine learning practitioners. Kaggle allows users to find and publish datasets, explore and build models in a web-based data-science environment, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges. Thus, we chose the Kaggle platform as the source of the notebooks for this study. To the best of our knowledge, DistilKaggle is one of the first Jupyter notebook metrics dataset on this scale and can be used for different data mining and code analysis purposes.

In the following sections, we will explain the metrics, whose names may not inherently convey their nature, in detail and the methodology of obtaining them. All of the metrics, their abbreviations in DistilKaggle, and the papers introduced them are presented in Table 1.

### 2.1 Metrics of Notebooks

In the following, we discuss notebook metrics in two groups: *notebook-based metrics* and *script-based metrics*.

**2.1.1 Notebook-Based Metrics.** The first set of metrics are those that have been specifically introduced for Jupyter notebooks. In recent years, several articles have worked on the code quality of Jupyter notebooks, and various metrics have been introduced in them [5, 9, 14, 16, 25, 25, 27–31, 33, 34]. After extracting all these metrics, we discussed with five notebook developers to identify the ones that are effective in code comprehension. Finally, out of 23 available metrics, 17 of them were selected, and we implemented appropriate algorithms to obtain all of them. These criteria are introduced in Table 1 and some of them are described here:

- **Number of Visualization Data Type:** Considering that it is possible to display the output of each cell in notebooks by executing it, this metric shows the number of visual outputs, including all kinds of images and graphs.
- **Number of Executed Cells:** This attribute specifies the number of cells that have been executed in a notebook and have an execution order number.
- **H1, H2 and H3 headlines:** These metrics show the number of headlines that are written inside the markdowns and are usually used to express titles and subtitles.
- **Distribution of Markdowns:** For calculating the distribution of markdowns throughout a notebook, we suggest the inverse coefficient of variation. The coefficient of variation (CV), also known as relative standard deviation (RSD), is a statistical measure commonly used for comparing diversity in workgroups and is defined as the ratio of the standard deviation to the mean [3].

**Table 1: Features Extracted in this Study**

Feature	Abbreviation	Reference
<b>Script-Based Metrics</b>		
<i>Basic Metrics</i>		
Lines of Code	LOC	[4, 12, 23, 28]
Number of Blank Lines of Code	BLC	[4, 12, 23]
Lines of Comments	LOCom	[4, 12, 23]
Number of Comment Words	CW	[12]
Number of Statements	S	[12, 23]
Number of Parameters	P	[12, 23]
Number of User-Defined Functions	UDF	[9]
<i>Complexity Metrics</i>		
Cyclomatic Complexity	CyC	[4, 9, 12, 23]
Nested Block Depth	NBD	[12, 23]
Kind of Line of Code Identifier Density	KLCID	[12]
<i>Halstead Metrics</i>		
Number of Operands	OPRND	[12, 19, 23]
Number of Operators	OPRAT	[4, 12, 19]
Number of Unique Operands	UOPRND	[12, 19]
Number of Unique Operators	UOPRAT	[12, 19]
<i>Readability Metrics</i>		
Average Line Length of Code	ALLC	[4]
Number of Identifiers	ID	[4]
Average Number of Identifiers	AID	[4]
<b>Notebook-Based Metrics</b>		
Number of Code Cells	CC	[18, 28]
Mean Number of Lines in Code Cells	MeanLCC	[12]
Number of Imports	I	[12, 18]
Mean Number of Words in Markdown Cells	MeanWMC	[28]
Number of H1 tags in the Markdown	H1	[18, 26, 27]
Number of H2 tags in the Markdown	H2	[18, 27, 31]
Number of H3 tags in the Markdown	H3	[18, 27, 31]
Number of Markdown Cells	MC	[14, 18, 28]
Mean Number of Lines in Markdown Cells	MeanLMC	[9]
Number of Markdown Words	MW	[12, 33]
Number of Lines in Markdown Cells	LMC	[23, 28]
Distribution of Markdown Cells	DMC	[18, 28]
Distribution of Imports	DI	[18, 28]
Performance Tier	PT	[2, 23]
External API Popularity	EAP	[5, 23, 34]
Number of Visualization Data Type	NVD	[1, 16]
Number of Executed Cells	EC	[18]

**Table 2: Columns of Dataframes**

Column	Description
Kernel ID	This key represents the unique identifier assigned to the Kaggle project associated with the Jupyter notebook.
Cell Index	The Cell Index indicates the sequential order of the cell within the Jupyter notebook.
Source	The Source key contains the actual text written in the cell.
Output Type	The Output Type key signifies the type of output generated by the execution of a code cell in the Jupyter notebook. It can take on various values, including: stream, display_data, error, execute_result. This column only exists in code cells Dataframe.
Execution Count	The Execution Count helps in tracking the order of code execution and understanding the flow of computations within the notebook. This column only exists in code cells Dataframe.

where:

- $n$ : the size of population
- $x_i$ : each value from the population

To calculate the distribution of markdown, we considered the number of words of each markdown as the size of it. After applying Equation 1 to the notebooks of the final dataset, we tested the results and evaluated the confidence of our outputs. For instance, the highest CV (the most unbalanced distribution) belongs to notebook P<sup>1</sup> and the lowest CV (the most balanced distribution) belongs to notebook Q<sup>2</sup>.

- **Distribution of Imports:** We adopted a similar approach as the distribution of markdowns for assessing the distribution of the imports feature within Jupyter notebooks, leveraging the inverse coefficient of variation as a statistical measure.
- **Performance Tier:** This feature is present in the metadata of the Kaggle repository, which, based on specific rules, determines the level of expertise of notebook developers with one of the levels 0 to 4.
- **External API Popularity:** This metric assigns a number to each notebook that represents to what extent the APIs and libraries used in that notebook are popular [34]. The more its value, the more the APIs used in the notebook are frequently used by other developers which results in better CU based on previous studies [23]. In order to measure this criterion in notebooks: First, we counted the number of times each API was imported in our dataset and assigned a popularity score to each API based on the frequency at which they were used in the whole dataset. Then, for each notebook, we summed up the popularity score of APIs used in that notebook to obtain the External API Popularity[23] score for that notebook.

$$CV = \sigma / \mu \quad (1)$$

where:

$$\sigma = \sqrt{\sum (x_i - \mu)^2 / n} \quad (2)$$

<sup>1</sup>P: <https://www.kaggle.com/code/xiwuhan/jmtc-2nd-place-solution>

<sup>2</sup>Q: <https://www.kaggle.com/code/anokas/two-sigma-time-travel-eda>

**2.1.2 Script-Based Metrics.** Given that each code cell in Jupyter notebooks is a regular Python script, many of the metrics for code scripts proposed by prior studies [4, 7, 12, 15, 17, 23] are also applicable to code cells in Jupyter notebooks. Considering that a small percentage of notebooks use the concept of class and object orientation<sup>3</sup> and have a weaker structure, some metrics were ignored.

Script-based metrics including seven basic metrics, two complexity metrics, four Halstead metrics [19] and three readability metrics as introduced in Table 1 and some of them are explained below:

- **Cyclomatic Complexity:** Cyclomatic complexity assigns a number for code complexity based on code graph. It was first used by Mathias [15] in order to examine the underlying nature of code designed to study the process of program comprehension.
- **KLCID:** Kind of Line of Code Identifier Density (KLCID) is One of the complexity metrics which analyzes the cognitive complexity of program comprehension tasks. The KLCID is calculated by counting the lines of conceptually unique code and calculating the density of identifiers.

## 2.2 Dataset Construction Methodology

Our journey in generating the dataset began by obtaining the paper IDs shared on Meta Kaggle API between October 2020 and October 2023. Due to Kaggle’s platform policy, each machine can download limited number of notebooks per day. To overcome this limitation, we utilized multiple servers with various IP addresses, enabling the download of 293,290 notebooks from the Kaggle platform. The response payload containing the notebooks is stored in JSON format, with a total size exceeding 300GB.

In the subsequent step, we organized the code cells and markdown cells into two distinct DataFrames (CSV files). This not only reduced the overall size of the dataset to 3GB but also enhanced its usability for future data analysis by fellow researchers. These refined DataFrames form an integral part of our final dataset, DistilKaggle. The detailed description of the columns of the Dataframes is presented in Table 2.

Following a comprehensive investigation into code quality metrics utilized in prior studies to evaluate code quality and comprehensibility, we employed the markdowns and code cells dataframes to extract these metrics for each Jupyter notebook. In many instances, utilizing either the code cells or markdown cells proved sufficient for capturing these features. Certain features, such as the performance tier indicating the developer’s expertise level, were extracted using the Meta Kaggle dataset [11].

Finally, as Figure 1 shows, the resulting features dataset is compact at 121MB, yet it offers invaluable insights into the code quality of the notebooks.

## 3 SAMPLE APPLICATION

The sheer magnitude of this dataset renders it versatile for diverse code analysis tasks. Many prior studies have concentrated on utilizing Jupyter notebooks for statistical analysis[5, 18, 29–31] and implementing deep learning models[14, 27, 28, 33], underscoring our dataset’s potential for various research applications.

<sup>3</sup>Based on the statistics of our notebooks dataset, only 3.5% of notebooks defined classes to implement objects.

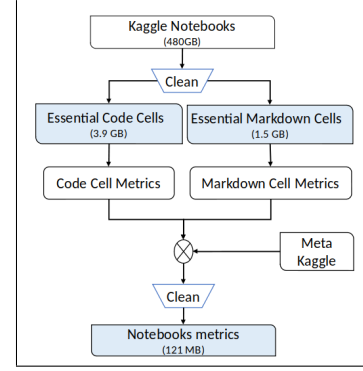


Figure 1: The approach for creating DistilKaggle

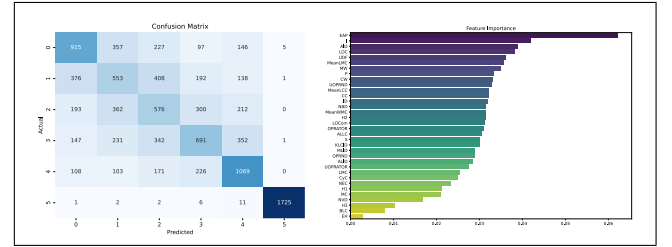


Figure 2: (Left) Confusion matrix results for the Random Forest Classifier trained with 100 trees, utilizing 80% of the dataset for training. (Right) The influence of individual features on the classifier’s prediction output (performance tier). Notably, API popularity emerged as the most decisive factor.

To demonstrate the dataset’s utility, we used a Random Forest classifier to predict developers’ performance tier. We trained the classifier on 80% of the DistilKaggle dataset and tested it on the remaining 20%. The achieved F1 score of 57% significantly outperformed a baseline Random predictor, which had an accuracy of 20%. This stark contrast highlights the distinct coding styles between experienced and novice programmers, suggesting a need for further exploration. Figure 2 showcases the Random Forest classifier’s performance. For example, an interesting result is “the more specialized the developers, the more famous API they use”.

## 4 CONCLUSIONS

In this paper, we introduced DistilKaggle, an extensive dataset comprising 293,290 Jupyter notebooks, each accompanied by its corresponding static code quality metrics. Our vision for DistilKaggle is to serve as a catalyst for future research endeavors in the realm of code analysis. The substantial volume of notebooks positions it as an ideal resource for the development and evaluation of large-scale deep learning models.

Notably, we enriched DistilKaggle by incorporating code quality metrics from the KGTorrent dataset, creating a combined dataset totaling 542,051 (293,290 notebooks from DistilKaggle and 248,761 notebooks from KGTorrent) notebooks. In this process, we extracted code quality metrics of these notebooks, amplifying the depth and richness of our dataset.

## REFERENCES

- [1] Vishakha Agrawal, Yong-Han Lin, and Jinghui Cheng. 2022. Understanding the characteristics of visual contents in open source issue discussions: a case study of jupyter notebook. In *International Conference on Evaluation and Assessment in Software Engineering*. 249–254.
- [2] Eman Abdullah AlOmar, Anthony Peruma, Mohamed Wiem Mkaouer, Christian D Newman, and Ali Ouni. 2021. Behind the scenes: On the relationship between developer experience and refactoring. *Journal of Software: Evolution and Process* (2021).
- [3] Arthur G Bedeian and Kevin W Mossholder. 2000. On the use of the coefficient of variation as a measure of diversity. *Organizational Research Methods* 3, 3 (2000), 285–297.
- [4] Raymond PL Buse and Westley R Weimer. 2009. Learning a metric for code readability. *IEEE Transactions on software engineering* 36, 4 (2009), 546–558.
- [5] Malinda Dilhara, Ameya Ketkar, and Danny Dig. 2021. Understanding Software-2.0: a study of machine learning library usage and evolution. *ACM Transactions on Software Engineering and Methodology* 30, 4 (2021), 1–42.
- [6] Helen Dong, Shurui Zhou, Jin LC Guo, and Christian Kästner. 2021. Splitting, renaming, removing: A study of common cleaning activities in Jupyter notebooks. In *36th IEEE/ACM International Conference on Automated Software Engineering Workshops*. IEEE, 114–119.
- [7] Jonathan Dorn. 2012. A general software readability model. *MCS Thesis available from (<https://web.eecs.umich.edu/~weimerw/students/dorn-mcs-pres.pdf>)* 5 (2012), 11–14.
- [8] Will Epperson, April Yi Wang, Robert DeLine, and Steven M Drucker. 2022. Strategies for reuse and sharing among data scientists in software teams. In *44th International Conference on Software Engineering: Software Engineering in Practice*. 243–252.
- [9] Konstantin Grotov, Sergey Titov, Vladimir Sotnikov, Yaroslav Golubev, and Timofey Bryksin. 2022. A large-scale comparison of Python code in Jupyter notebooks and scripts. In *19th International Conference on Mining Software Repositories*. 353–364.
- [10] Yuan Jiang, Christian Kästner, and Shurui Zhou. 2022. Elevating Jupyter Notebook Maintenance Tooling by Identifying and Extracting Notebook Structures. In *2022 IEEE International Conference on Software Maintenance and Evolution*. IEEE, 399–403.
- [11] Kaggle. 2020. *Kaggle Meta Data*, <https://www.kaggle.com/datasets/kaggle/meta-kaggle>. Retrieved April 5, 2022 from <https://www.kaggle.com/datasets/kaggle/meta-kaggle>
- [12] Nadia Kasto and Jacqueline Whalley. 2013. Measuring the difficulty of code comprehension tasks using software metrics. In *15th Australasian Computing Education Conference-Volume 136*. 59–65.
- [13] Jiali Liu, Nadia Boukhefifa, and James R Eagan. 2019. Understanding the role of alternatives in data analysis practices. *IEEE Transactions on Visualization and Computer Graphics* 26, 1 (2019), 66–76.
- [14] Xuye Liu, Dakuo Wang, April Wang, Yufang Hou, and Lingfei Wu. 2021. HA-ConvGNN: Hierarchical attention based convolutional graph neural network for code documentation generation in Jupyter Notebooks. *arXiv:2104.01002* (2021).
- [15] Karl S Mathias, James H Cross, T Dean Hendrix, and Larry A Barowski. 1999. The role of software measures and metrics in studies of program comprehension. In *37th Annual Southeast Regional Conference*.
- [16] Jorge Piazentin Ono, Juliana Freire, and Claudio T Silva. 2021. Interactive data visualization in jupyter notebooks. *Computing in Science & Engineering* 23, 2 (2021), 99–106.
- [17] JR Parker and Katrin Becker. 2003. Measuring effectiveness of constructivist and behaviourist assignments in CS102. *ACM SIGCSE Bulletin* 35, 3 (2003), 40–44.
- [18] João Felipe Pimentel, Leonardo Murta, Vanessa Braganholo, and Juliana Freire. 2019. A large-scale study about quality and reproducibility of jupyter notebooks. In *IEEE/ACM 16th International Conference on Mining Software Repositories*. 507–517.
- [19] Daryl Posnett, Abram Hindle, and Premkumar Devanbu. 2011. A simpler model of software readability. In *8th International Conference on Mining Software Repositories*. 73–82.
- [20] Fotis Psallidas, Yiwen Zhu, Bojan Karlas, Jordan Henkel, Matteo Interlandi, Subru Krishnan, Brian Kroth, Venkatesh Emani, Wentao Wu, Ce Zhang, and et al. 2022. Data Science through the looking Glass. *ACM SIGMOD Record* 51, 2 (2022), 30–37.
- [21] Luigi Quaranta, Fabio Calefato, and Filippo Lanubile. 2021. KGTorrent: A dataset of python Jupyter notebooks from kaggle. In *18th IEEE/ACM International Conference on Mining Software Repositories*. 550–554.
- [22] Adam Rule, Aurélien Tabard, and James D Hollan. 2018. Exploration and explanation in computational notebooks. In *CHI Conference on Human Factors in Computing Systems*. 1–12.
- [23] Simone Scalabrino, Gabriele Bavota, Christopher Vendome, Mario Linares-Vasquez, Denys Poshyvanyk, and Rocco Oliveto. 2019. Automatically assessing code understandability. *IEEE Transactions on Software Engineering* 47, 3 (2019), 595–613.
- [24] Krishna Subramanian, Ilya Zubarev, Simon Völker, and Jan Borchers. 2019. Supporting data workers to perform exploratory programming. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–6.
- [25] Sergey Titov, Yaroslav Golubev, and Timofey Bryksin. 2022. Resplit: improving the structure of jupyter notebooks by re-splitting their cells. In *29th IEEE International Conference on Software Analysis, Evolution and Reengineering*. 492–496.
- [26] Ashwin Prasad Shivarpatna Venkatesh and Eric Bodden. 2021. Automated cell header generator for Jupyter notebooks. In *1st ACM International Workshop on AI and Software Testing/Analysis*. 17–20.
- [27] Ashwin Prasad Shivarpatna Venkatesh, Jiawei Wang, Li Li, and Eric Bodden. 2023. Enhancing Comprehension and Navigation in Jupyter Notebooks with Static Analysis. In *30th IEEE International Conference on Software Analysis, Evolution and Reengineering*.
- [28] April Yi Wang, Dakuo Wang, Jaimie Drozdal, Michael Muller, Soya Park, Justin D Weisz, Xuye Liu, Lingfei Wu, and Casey Dugan. 2022. Documentation matters: Human-centered AI system to assist data science code documentation in computational notebooks. *ACM Transactions on Computer-Human Interaction* 29, 2 (2022), 1–33.
- [29] Jiawei Wang, Li Li, and Andreas Zeller. 2020. Better code, better sharing: on the need of analyzing jupyter notebooks. In *42nd ACM/IEEE International Conference on Software Engineering: New Ideas and Emerging Results*. 53–56.
- [30] Jiawei Wang, Li Li, and Andreas Zeller. 2021. Restoring execution environments of Jupyter notebooks. In *43rd IEEE/ACM International Conference on Software Engineering*. IEEE, 1622–1633.
- [31] Jiawei Wang, KUO Tzu-Yang, Li Li, and Andreas Zeller. 2020. Assessing and restoring reproducibility of Jupyter notebooks. In *35th IEEE/ACM International Conference on Automated Software Engineering*. IEEE, 138–149.
- [32] John Wenskovitch, Jian Zhao, Scott Carter, Matthew Cooper, and Chris North. 2019. Albireo: An interactive tool for visually summarizing computational notebook structure. In *IEEE Visualization in Data Science*. IEEE, 1–10.
- [33] Chenyang Yang, Shurui Zhou, Jin LC Guo, and Christian Kästner. 2021. Subtle bugs everywhere: Generating documentation for data wrangling code. In *36th IEEE/ACM International Conference on Automated Software Engineering*. IEEE, 304–316.
- [34] Chenguang Zhu, Ripon K Saha, Mukul R Prasad, and Sarfraz Khurshid. 2021. Restoring the Executability of Jupyter Notebooks by Automatic Upgrade of Deprecated APIs. In *36th IEEE/ACM International Conference on Automated Software Engineering*. 240–252.