

# Navigating Data Errors in Machine Learning Pipelines: Identify, Debug, and Learn

Bojan Karlaš

Harvard University  
bkarlas@mgh.harvard.edu

Babak Salimi

University of California, San Diego  
bsalimi@ucsd.edu

Sebastian Schelter

BIFOLD & TU Berlin  
schelter@tu-berlin.de

## Abstract

Addressing data errors—such as missing, incorrect, noisy, biased, or out-of-distribution values—is essential to building reliable machine learning (ML) systems. Traditional methods often focus on refining the training process to minimize error symptoms or repairing data errors indiscriminately, without addressing their root causes. These isolated approaches ignore how errors originate and propagate through the interconnected stages of ML pipelines—data preprocessing, model training, and prediction—resulting in superficial fixes and suboptimal solutions. Consequently, they miss the opportunity to understand how data errors impact downstream tasks and to implement targeted, effective interventions.

In recent years, the research community has made significant progress in developing holistic approaches to identify the most harmful data errors, prioritize impactful repairs, and reason about their effects when errors cannot be fully resolved. This tutorial surveys prominent work in this area and introduces practical tools designed to address data quality issues across the ML lifecycle. By combining theoretical insights with hands-on demonstrations, attendees will gain actionable strategies to diagnose, repair, and manage data errors, enhancing the reliability, fairness, and transparency of ML systems in real-world applications.

## ACM Reference Format:

Bojan Karlaš, Babak Salimi, and Sebastian Schelter. 2024. Navigating Data Errors in Machine Learning Pipelines: Identify, Debug, and Learn. In *Proceedings of SIGMOD (SIGMOD’25)*. ACM, New York, NY, USA, 8 pages.

## 1 Introduction

ML systems are increasingly deployed in high-stakes domains such as healthcare, finance, and law enforcement, where their decisions profoundly impact individuals and communities. To ensure trust in these systems, it is essential to focus on their accuracy, fairness, robustness, and reliability [30, 50, 78]. Developing trustworthy ML systems involves navigating a complex multi-stage pipeline—spanning data preparation, model training, predictive queries, and evaluation—where failures at any stage can lead to significant performance degradation. A key observation is that many such failures are directly caused by data errors, including missing, incorrect, invalid, biased, and out-of-distribution values. These errors propagate through ML pipelines, compromising outcomes even when models and algorithms are otherwise well-refined. Debugging and

mitigating these errors consumes considerable time and effort for developers, highlighting the need to systematically address data quality issues. Existing methods for addressing bias in ML models primarily focus on algorithm-specific solutions, which often treat the symptoms of poor data quality rather than tackling its root causes [5, 51]. Traditional data cleaning techniques [40, 42, 48, 68] produce a single “best-guess” version of cleaned data but offer no guarantees of unbiasedness or representativeness. Similarly, explainability methods, while valuable for interpreting model predictions, often analyze models in isolation and overlook the broader ML pipeline, where errors from upstream stages can propagate and amplify (Figure 1). In real-world ML applications, training data is typically derived from multiple heterogeneous source datasets through complex ML pipelines [4, 6, 67, 84]. These pipelines integrate, transform, and encode raw data into features, exacerbating data quality challenges such as inconsistencies, noise, and bias [28]. To overcome these limitations, a holistic, data-centric approach that scrutinizes and refines the entire ML pipeline is necessary to ensure robust, fair, and reliable systems in deployment [26].

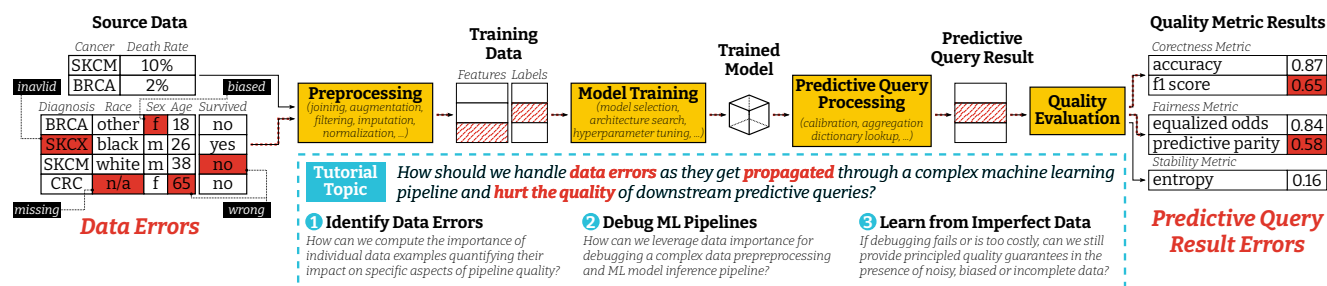
**Outline of the Tutorial.** Addressing these challenges requires a shift in perspective: understanding ML pipelines as interconnected workflows where data quality issues must be tackled holistically across all stages. This tutorial focuses on recent work that highlights the importance of quantifying task-specific data contributions, such as Shapley values for identifying problematic data points [21, 34], reasoning about end-to-end pipelines to trace error propagation [23, 24, 72], and providing quality guarantees to enable robust learning under uncertainty, incomplete data, and distributional inconsistencies [55, 62, 93]. These techniques collectively provide the foundation for understanding how data errors propagate through ML pipelines and how targeted interventions can mitigate their downstream impact.

**Structure and Main Takeaways.** This 3-hour tutorial equips participants with actionable tools to identify, debug, and reason about data quality challenges while enhancing accuracy, fairness, reliability, and robustness in ML systems. Structured into two 90-minute parts, the initial survey session will introduce methods for identifying data errors, debugging ML pipelines, and learning from imperfect data. The second part of the tutorial consists of a hands-on session, where attendees will learn to apply practical tools to real-world tasks, such as prioritizing problematic data points, tracing error propagation, and implementing robust learning strategies. By combining these perspectives, participants will gain a holistic understanding of how to address data quality issues, enabling them to build robust, transparent, and trustworthy ML systems for real-world challenges.

---

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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGMOD’25, June 22–27, 2025, Berlin, Germany  
© 2024 Copyright held by the owner/author(s).



**Figure 1: Data errors are often the root cause of many quality issues in modern ML pipelines. Handling such errors as they traverse complex pipeline steps is a key challenge for practitioners. This tutorial covers some significant recent developments, presents novel tools, and explores opportunities for future work in this space.**

**Target Audience.** This tutorial will cover currently available methods and tools for navigating data errors from both a theoretical and a practical perspective. This will provide participants with a deeper understanding of the errors’ impact on data science workflows, as well as actionable steps for solving many real-world issues. We envision that the tutorial will be valuable for several groups of participants: (1) *practitioners* – data scientists and ML engineers who would learn about some new tools that could be added to their arsenal, (2) *researchers* – members of the data management and ML systems communities who are interested in this space would get acquainted with the current progress of the field as well as some open challenges, and (3) *system builders* – engineers and system architects who are working on tools for ML development would hear about methods that could potentially result in new features of their systems.

**Materials.** We make the slides and code for this tutorial available under an open license at: <https://navigating-data-errors.github.io>

## 2 Outline of the Survey

In the first part of the tutorial, we will present a survey of relevant works covering the notion of data importance as a framework for identifying data errors, the application of these methods for debugging end-to-end ML pipelines, and the approaches for providing quality guarantees despite the presence of data errors (Figure 1).

### 2.1 Data Importance for Data Error Detection

Repairing data is often a very costly process requiring a lot of human effort. Therefore, identifying data with the most significant negative impact on the downstream ML model could allow practitioners to prioritize their efforts. The general strategy pursued by recent approaches is to define some method of measuring the importance of individual units of data with respect to their impact on the downstream ML model. In this part of the survey, we will explore various methods for quantifying data importance, as well as some practical approaches for overcoming computational challenges that arise when applying these methods to real datasets.

**Quantifying Data Importance.** We will start by covering a simple way to measure data importance using the *leave-one-out* (LOO) score. We will discuss several generalizations of this approach including the Shapley value [21, 34], Banzhaf value [80], Beta

Shapley [43], and others. We will also cover gradient-based methods [41, 42], as well as some uncertainty-based methods [59, 63]. Finally, we will cover methods geared toward specific aspects of model quality such as fairness [66], as well as methods specialized for retrieval augmented generation used in inference based on large language models [47].

*Take-away:* Attendees will get acquainted with the notion of data importance in the context of data debugging, as well as various methods for quantifying it. They will also develop a sense of the strengths and weaknesses of various methods, allowing them to reason about selecting the best method for their own data debugging scenario.

**Overcoming Computational Challenges.** Despite their effectiveness, many methods for quantifying data importance come with enormous computational costs, limiting their applicability in real-world data debugging tasks [35]. For example, the Shapley value involves a sum over exponentially many subsets, making it intractable in practical settings. In this part, we will explore approaches for speeding up the computation, including Monte Carlo methods [21], using the K-nearest neighbors as a proxy model [33], and model-based estimation [14].

*Take-away:* This part will introduce the attendees to certain algorithmic approaches for making data importance computationally efficient by leveraging various tools well-known to the data management community. We hope that our overview will inspire some of the attendees to make future contributions.

### 2.2 Data Debugging in ML Pipelines

The discussed techniques are designed for a static, preprocessed training dataset of a given model. However, this assumption ignores the circumstances in real-world ML applications, where models are trained on data that is preprocessed as part of an *ML pipeline* [6, 39, 44, 60, 67, 73, 76, 83–85]. Such a pipeline typically accesses several heterogeneous input datasets, integrates them and encodes them into features to produce the actual training data for the model (Figure 1). This raises challenges for debugging – data errors should be identified in the *source data* of a pipeline, while existing debugging methods are designed for already preprocessed *training data* (the output of the preprocessing step). The second part of our survey will bring these ML pipelines into play.

**Libraries and Systems for ML pipelines.** We will start by giving an overview of the implementation of ML pipelines, which typically combine several systems and libraries. Examples from the open source space include combinations of Pandas [81] with scikit-learn [61, 74], Spark [86] with SparkML [52], Tensorflow Transform [6] with Apache Beam [2], and systems such as Apache SystemDS [9], MLflow Recipes [15] from Databricks or Metaflow [53] from Netflix. Furthermore, we will cover proprietary pipeline abstractions in commercial cloud services such as Amazon SageMaker [3], Microsoft Azure ML [56], or Google’s Vertex AI [22].

*Take-away: Attendees will learn about shared design patterns and abstractions across various libraries, as well as their shortcomings which make debugging more difficult.*

**Characteristics of Real-World ML Pipelines.** Next, we will summarize two large-scale empirical studies on the characteristics of real-world ML pipelines encountered in large companies, code repositories and cloud platforms. In particular, we will focus on a study of thousands of production ML pipelines at Google [84], and on insights from the analysis of millions of GitHub notebooks and ML.Net pipelines from Microsoft [67].

*Take-away: Attendees will learn about detailed usage statistics for common libraries and operators in these pipelines, e.g., that a small number of highly popular libraries and operations dominates the pipelines. At the same time, they will be made aware that there exists a long tail of niche operators, e.g., hundreds of system-provided operators accompanied by thousands of user-defined operators. Moreover, they will be confronted with findings which contradict conventional wisdom, e.g., that data ingestion, data preprocessing and model analysis account for higher compute costs than model training, or that a large proportion of pipelines train traditional non-neural ML models.*

**Inspecting and Debugging Data in ML pipelines.** Finally, we will discuss techniques to inspect pipelines and debug their input and output data. We will start by reviewing techniques for the basic analysis of ML pipelines [24, 57, 64]. Subsequently, we will dive into work on provenance-based data debugging of ML pipelines and discuss approaches which leverage fine-grained provenance information [27] to reason about the input and output data of a pipeline. Examples include a continuous integration system to screen pipelines for issues like data leakage and label errors [72], techniques to efficiently compute data importance over pipelines of various shapes [39], or the identification of training data points whose removal fixes user complaints in prediction queries [20, 83].

*Take-away: Attendees will be given a detailed overview of different pipeline representations, the efficient computation of fine-grained data provenance for a pipeline, and how this enables the adaptation of existing data debugging techniques to the pipeline’s source data. We will furthermore highlight the connection to related areas such as incremental view maintenance of the pipeline outputs based on changes in their inputs.*

### 2.3 Learning from Uncertain and Incomplete Data

While debugging pipelines by identifying and repairing data errors is an important aspect of the ML development lifecycle, it can also quickly become too costly or, in certain scenarios, even impossible.

For example, the information needed to repair missing values could require an expensive data acquisition process or could simply be inaccessible. Before time and resources are spent on data debugging, an equally important question that arises is – *do we even need to debug?* Answering this question depends on the ability to establish guarantees for the quality of predictive queries in spite of the presence of data errors, which is the topic that we will cover in this part of the survey.

Uncertain and incomplete data, arising from data errors, missing values, and biases, are pervasive in real-world ML applications. These imperfections distort the underlying data distribution, compromising the fairness, accuracy, and generalizability of ML models. Traditional approaches, such as fairness-aware learning, data cleaning techniques [10, 38, 48, 68, 87, 88], and methods addressing selection bias [13, 16, 31, 45, 69] or labeling errors [36, 89, 91], often fail to recover a representative dataset, limiting their effectiveness in practice. Furthermore, robust model learning methods aim to ensure resilience against adversarial perturbations [32, 70, 77, 90] and distributional shifts [7, 58], but rely on restrictive assumptions that rarely hold. These challenges are exacerbated in complex ML pipelines, where data imperfections propagate and interact, amplifying performance degradation.

**Quantifying and Handling Incomplete Data.** Recent advancements have shifted focus from perfecting data to reasoning under its inherent uncertainty and incompleteness. Early contributions extended nearest neighbor classifiers to handle incomplete information, ensuring predictions align with the most reliable available data [40]. The dataset multiplicity problem formalized the challenges posed by unreliable or conflicting datasets, emphasizing the need to quantify uncertainty and its impact on predictions [55]. To address noisy and incomplete inputs, frameworks have emerged that provide statistical guarantees by constructing certain and approximately certain models, ensuring robust performance in tasks like linear regression and support vector machines [92]. Another critical development is the possible worlds framework, which trains models across multiple plausible interpretations of uncertain data, enhancing resilience in ambiguous settings [93]. Further, methods addressing fairness concerns, such as consistent range approximation, certify that predictive models remain unbiased despite biases in training data [94]. Robustness to programmable data biases has also been explored, where decision trees are evaluated over biased datasets to ensure consistent predictions and fairness [54].

*Take-away: As part of this overview, attendees will learn about the limitations of traditional approaches for cleaning, debiasing, and learning robust models. The tutorial will highlight recent progress in learning from incomplete, uncertain, and inconsistent data, as well as methods for propagating the impact of these imperfections on ML models. We will also cover applications for quantifying robustness and improving data cleaning. As a result, participants will gain practical insights into building resilient models that can effectively learn and perform under real-world data imperfections.*

### 2.4 Open Challenges & Conclusion

To conclude the survey part, we will highlight some problems that we believe to be relevant to this space but have received little attention from the research community. A major concern is the

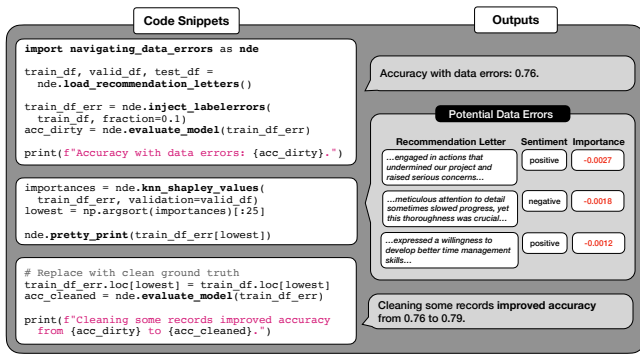


Figure 2: Data importance for data error detection in the hands-on session – Attendees run code snippets to inject synthetic label errors into the data, identify the most strongly affected tuples via data importance, and observe how prioritized cleaning helps recover model performance.

scalability of several of the presented methods for computing data importance and learning from imperfect data [29, 40, 55], and the research community is already actively working on ways to nevertheless apply such techniques at the scale of real-world data [35]. Some methods that we will cover leverage proxy models as a strategy for improving computational efficiency. For example, using the K-nearest neighbor model can provide very efficient solutions for computing Data Shapley values [33, 39], but it may not always give the best results in situations where the inductive bias of the proxy model is incompatible with the actual model being used [37].

Another underexplored direction is the connection between data debugging and low-latency machine unlearning [17] – since several debugging techniques at their core assess the impact of repeatedly removing data points (or groups of data points) from a model [29], the insights from recent research on data debugging could benefit ongoing efforts to design data-driven applications that forget critical data fast [75].

From a long-term perspective, there are questions of how the discipline of data debugging will be impacted by the advent of AI-assisted programming [8, 19, 71], and how the broader AI development lifecycle will be impacted by the recently introduced AI-regulation such as the EU AI Act, GDPR, SCPA, etc. [1, 11, 18, 79].

*Take-away: Attendees will get an overview of opportunities for future research directions in the space of data debugging.*

### 3 Outline of the Hands-On Session

The hands-on session is divided into two parts, each of which will be implemented in a separate Google Colab notebook. We will provide an easy-to-install Python package called `navigating_data_errors` for the code and data used throughout the hands-on-session. The first part will be a fully guided walkthrough of realistic data science scenarios where participants will get acquainted with tools for data error handling. The second part will give willing attendees the opportunity to apply these tools to a real-world situation where they are given a poor-quality dataset and are tasked with performing interventions in order to improve the downstream ML model.

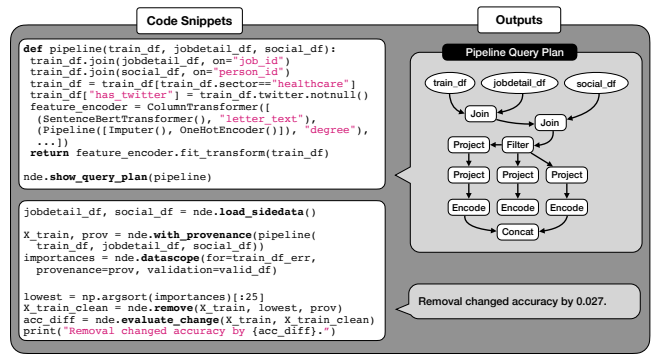


Figure 3: Incorporating preprocessing pipelines into data debugging during the hands-on session – Attendees will define and visualise pipelines for integrating, filtering and encoding data, and learn how to debug and modify a pipeline’s source data via fine-grained provenance information.

### 3.1 Tool Overview

**Structure.** The hands-on session will start with a one-hour introduction to various tools for identifying data errors, computing pipeline provenance, and quantifying uncertainty in model training and predictions. This part will leverage synthetically generated data from a hiring scenario, in particular a set of recommendation letters together with multiple tables of side data such as demographic information and social media details of the applicants. The corresponding ML use case will be to train a classifier to predict the sentiment of a recommendation letter.

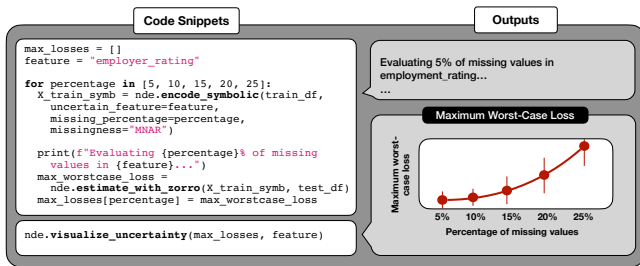
We will walk attendees through various examples of software tools for methods discussed in the survey such as kNN-Shapley [33], Gopher [66], `mlinspect` [25, 72], `Datascopes` [39], and `Zorro` [93].

**Content.** We will introduce the APIs of these tools, provide code snippets with usage examples, as well as a set of short 5-minute programming tasks for attendees, to encourage them to explore the tools themselves.

**Identifying Data Errors.** As sketched in Figure 2, we will start by showcasing how to identify and “recover” from data errors via data debugging. The data for this part consists of a single preprocessed table as training data without any complex features. We will inject synthetic noise such as label errors into the data and show how this negatively impacts the downstream quality metrics of the classifier. We will apply tools from Section 2.1 to identify impactful tuples with data errors, provide them to an “oracle” cleaning function and show how such prioritized cleaning improves quality metrics.

*Task for attendees: Given a “cleaning oracle” function, attendees will be asked to implement an iterative cleaning solution for data with label errors, which will apply data repairs and result in measurable improvements of model quality.*

**Incorporating Pipelines.** Next, we will introduce ML pipelines for data preprocessing into the scenario (as discussed in Section 2.2), which include additional side tables, and preprocess the data with complex operations such as (fuzzy) joins, filters, projections with user-defined functions, as well as costly feature encoders. As sketched



**Figure 4: Learning from imperfect data during the hands-on session – Attendees will inject synthetic missing values into the data to simulate real-world imperfections, and will see how this incomplete data impacts prediction reliability.**

in Figure 3, we will visualise the pipeline and show how to compute fine-grained data provenance for its outputs. Next, attendees will learn to identify the previously injected data errors in the source data of the pipeline based on the provenance information and the introduced tools.

*Task for attendees: The attendees should now extend the code of their iterative cleaning solution from the previous task to make it work on the ML pipeline.*

Reasoning about Uncertainty in the Predictions. We conclude our tool introduction with scenarios discussed in Section 2.3, where data quality issues cannot be fully resolved through cleaning. Here, we demonstrate how to reason about and quantify uncertainty in model training and predictions. As illustrated in Figure 4, using a subset of the data, we inject synthetic missing attributes and uncertain labels to simulate real-world imperfections. We will focus on Zorro [93], a framework that symbolically propagates uncertainty due to missing values through the training process, allowing us to compute prediction ranges for model outputs.

Attendees will observe how incomplete and uncertain data impact prediction reliability and robustness and will visualize the resulting uncertainty ranges for specific test points. We will also demonstrate the application of symbolic processing for influence analysis and data cleaning, showcasing how these techniques help identify and address problematic data points. By comparing a baseline model trained on imputed data to the uncertainty-aware model trained with Zorro, we will highlight how reasoning under uncertainty improves model robustness and enables reliable decision-making in imperfect data environments.

*Task for attendees: Attendees will compute prediction ranges for the data using Zorro and compare these ranges to the predictions of a baseline model trained with simple imputation. They will summarize their observations, focusing on the differences in prediction variability and the reliability of the two models under imperfect data conditions.*

### 3.2 Data Debugging Challenge

In the final half hour of the hands-on session, we will present attendees with a challenging data cleaning task, inspired by recent benchmarks for data-centric AI development [49]. The attendees will be given access to a prepared training dataset with data errors unknown to them, and access to a classifier with a validation set.

Moreover, they will be given an “oracle” function, which allows them to specify a limited set of training tuples to clean (by supplying their identifiers). This oracle function will then evaluate the classifier (retrained on the partially cleaned data) on a hidden test set, and report the metric on this test set to the attendee. This will allow attendees to test their previously acquired knowledge about the data debugging tools in a challenging example scenario. We additionally plan to implement a live “leaderboard” to motivate the more competitive attendees, showcase the submissions that introduced the highest improvements, and foster subsequent discussions among presenters and attendees.

## 4 Prerequisites & Context

**Prerequisites for the survey part.** The target audience for our survey are researchers and practitioners with an interest in the intersection of data management, machine learning, and data quality. We intend to cover both theoretical aspects as well as practical aspects related to the design and deployment of real-world ML applications to appeal to a large audience. The survey assumes a very basic understanding of machine learning.

**Prerequisites for the hands-on session.** For the hands-on session, attendees will need a laptop with internet access, and basic Python and data wrangling skills. We plan to implement the tutorial and tasks in Google Colab notebooks to avoid any local software installation or data downloads.

**Ethics.** Our hands-on tutorial will only use artificial, synthetically generated data and software publicly available under open-source licenses. Several problems and methods discussed in this tutorial are crucial for upcoming regulations such as the EU AI Act, which introduces comprehensive data governance and data compliance requirements for ML applications [18].

**Relationship to Previous Tutorials.** We would like to note that our tutorial shares some overlap with several tutorials presented in recent years at data management venues, with key differences that we highlight here. The tutorials on “Data Cleaning: Overview and Emerging Challenges” at SIGMOD’16 [12] and “Data Collection and Quality Challenges for Deep Learning” at VLDB’20 [82] both provided extensive coverage of data-quality related topics. However, the field has produced some substantial developments over the past years in terms of methods and tools which we will present in our tutorial. The tutorial “Explainable AI: Foundations, Applications, Opportunities for Data Management Research” [65] at SIGMOD’22 also covered data importance methods, among other topics, but they focused on the context of model interpretability, as opposed to data debugging. Similarly, the tutorial “Applications and Computation of the Shapley Value in Databases and Machine Learning” at SIGMOD’24 [46] focuses entirely on the Shapley value and its various applications, while in our case this is presented as one of the various tools available for identifying data errors in ML pipelines.

## Presenters

**Bojan Karlaš** is a postdoctoral research fellow at Harvard University. He works with Eugene Semenov at the Cutaneous Biology Research Center of Massachusetts General Hospital, David Liu at the Department of Medical Oncology of Dana-Farber Cancer Institute, and Kun-Hsing Yu at the Department of Biomedical Informatics of Harvard Medical School. His current research focus is on developing machine learning pipelines for processing biomedical data and extracting clinically meaningful insights. Previously, he did his Ph.D. at the Systems Group of ETH Zürich working with Ce Zhang where he was developing systems for managing the ML development lifecycle with a specific focus on data debugging.

**Babak Salimi** is an Assistant Professor at the Halicioğlu Data Science Institute and an affiliate of the Department of Computer Science and Engineering at the University of California, San Diego. His research focuses on the intersection of data management, machine learning, and responsible data science, with an emphasis on developing tools and methods that promote transparency, fairness, reliability, and robustness in algorithmic decision-making. Specifically, his work addresses key challenges in data management for machine learning, such as data cleaning, integration, and debiasing, while also creating frameworks that ensure interpretability, trust, and ethical use of data in real-world applications. His contributions have been recognized with several prestigious awards, including the NSF CAREER Award (2024), SIGMOD Research Highlight Award (2020), SIGMOD Best Paper Award (2019), and VLDB Best Demonstration Paper Award (2018).

**Sebastian Schelter** is a professor at the Berlin Institute for the Foundations of Learning and Data (BIFOLD) and Technische Universität Berlin. He conducts research at the intersection of data management and machine learning, which addresses data-related problems in ML applications that cause negative economic, societal or scientific impact. This research is accompanied by efficient and scalable open source implementations, many of which are applied in real-world use cases, for example in the Amazon Web Services cloud and in large European e-Commerce platforms. He was previously at the University of Amsterdam, New York University, and Amazon Research. His contributions have been recognized with the SIGMOD Systems Award (2023), SIGMOD Best Demo Runner-Up Award (2023), and Best Paper Runner-Up Award from the Table Representation Learning workshop at NeurIPS (2022).

## References

- [1] 2016. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation). <https://gdpr-info.eu/>. (2016). [Online; accessed 17-Feb-2019].
- [2] Tyler Akidau, Robert Bradshaw, Craig Chambers, Slava Chernyak, Rafael J Fernández-Moctezuma, Reuven Lax, Sam McVeety, Daniel Mills, Frances Perry, Eric Schmidt, et al. 2015. The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing. *Proceedings of the VLDB Endowment* 8, 12 (2015), 1792–1803.
- [3] Amazon Web Services. 2020. SageMaker Pipelines. <https://aws.amazon.com/sagemaker/pipelines/>.
- [4] Rajveer Bachkaniwala, Harshith Lanka, Kexin Rong, and Ada Gavrilovska. 2024. Lotus: Characterization of Machine Learning Preprocessing Pipelines via Framework and Hardware Profiling. In *2024 IEEE International Symposium on Workload Characterization (IISWC)*. IEEE, 30–43.
- [5] Agathe Balayn, Christoph Loh, and Geert-Jan Houben. 2021. Managing bias and unfairness in data for decision support: a survey of machine learning and data engineering approaches to identify and mitigate bias and unfairness within data management and analytics systems. *The VLDB Journal* 30, 5 (2021), 739–768.
- [6] Denis Baylor, Eric Breck, Heng-Tze Cheng, Noah Fiedel, Chuan Yu Foo, Zakaria Haque, Salem Haykal, Mustafa Ispir, Vihan Jain, Levent Koc, et al. 2017. TFx: A tensorflow-based production-scale machine learning platform. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 1387–1395.
- [7] Aharon Ben-Tal, Dick den Hertog, Anja De Waegenaere, Bertrand Melenberg, and Gijs Rennen. 2011. Robust Solutions of Optimization Problems Affected by Uncertain Probabilities. *Advanced Risk & Portfolio Management® Research Paper Series* (2011). <https://api.semanticscholar.org/CorpusID:761793>
- [8] Sahil Bhatia, Jie Qiu, Niranjan Hasabnis, Sanjit A Seshia, and Alvin Cheung. 2024. Verified Code Transpilation with LLMs. *NeurIPS* (2024).
- [9] Matthias Boehm et al. 2020. SystemDS: A Declarative Machine Learning System for the End-to-End Data Science Lifecycle. *CIDR* (2020).
- [10] Toon Calders and Sicco Verwer. 2010. Three naive Bayes approaches for discrimination-free classification. *Data Mining and Knowledge Discovery* 21, 2 (2010), 277–292.
- [11] California Consumer Privacy Act (CCPA). 2024. Requests to Delete. <https://oag.ca.gov/privacy/ccpa#sectiond>.
- [12] Xu Chu, Ihab F. Ilyas, Sanjay Krishnan, and Jiannan Wang. 2016. Data Cleaning: Overview and Emerging Challenges. In *Proceedings of the 2016 International Conference on Management of Data (San Francisco, California, USA) (SIGMOD '16)*. Association for Computing Machinery, New York, NY, USA, 2201–2206. <https://doi.org/10.1145/2882903.2912574>
- [13] Corinna Cortes, Mehryar Mohri, Michael Riley, and Afshin Rostamizadeh. 2008. Sample selection bias correction theory. In *International conference on algorithmic learning theory*. Springer, 38–53.
- [14] Ian Covert, Chanwoo Kim, Su-In Lee, James Zou, and Tatsunori Hashimoto. 2024. Stochastic Amortization: A Unified Approach to Accelerate Feature and Data Attribution. *NeurIPS* (2024).
- [15] Databricks. 2022. Mlflow Recipes. <https://mlflow.org/docs/latest/recipes.html>.
- [16] Wei Du and Xintao Wu. 2021. Robust fairness-aware learning under sample selection bias. *arXiv preprint arXiv:2105.11570* (2021).
- [17] Sebastian Schelter et al. 2021. HedgeCut: Maintaining Randomised Trees for Low-Latency Machine Unlearning. In *SIGMOD*.
- [18] EU AI Act. 2024. Article 10: Data and Data Governance. <https://artificialintelligenceact.eu/article/10/>.
- [19] Raul Castro Fernandez et al. 2023. How large language models will disrupt data management. *VLDB* (2023).
- [20] Lampros Flokas, Weiyuan Wu, Yejia Liu, Jiannan Wang, Nakul Verma, and Eugene Wu. 2022. Complaint-driven training data debugging at interactive speeds. In *Proceedings of the 2022 International Conference on Management of Data*. 369–383.
- [21] Amirata Ghorbani and James Zou. 2019. Data shapley: Equitable valuation of data for machine learning. In *International conference on machine learning*. PMLR, 2242–2251.
- [22] Google. 2021. Vertex AI Pipelines. <https://cloud.google.com/vertex-ai/docs/pipelines>.
- [23] Stefan Grafberger, Paul Groth, and Sebastian Schelter. 2023. Automating and Optimizing Data-Centric What-If Analyses on Native Machine Learning Pipelines. *Proceedings of the ACM on Management of Data* 1, 2 (2023), 1–26.
- [24] Stefan Grafberger, Paul Groth, Julia Stoyanovich, and Sebastian Schelter. 2022. Data distribution debugging in machine learning pipelines. *The VLDB Journal* 31, 5 (2022), 1103–1126.
- [25] Stefan Grafberger, Shubha Guha, Julia Stoyanovich, and Sebastian Schelter. 2021. Mlinspect: A data distribution debugger for machine learning pipelines. In *Proceedings of the 2021 International Conference on Management of Data*. 2736–2739.
- [26] Stefan Grafberger, Zeyu Zhang, Sebastian Schelter, and Ce Zhang. 2024. Red Onions, Soft Cheese and Data: From Food Safety to Data Traceability for Responsible AI. *IEEE Data Eng. Bull.* 47, 1 (2024), 63–81.
- [27] Todd J Green, Grigoris Karvounarakis, and Val Tannen. 2007. Provenance semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. 31–40.
- [28] Shubha Guha, Falaah Arif Khan, Julia Stoyanovich, and Sebastian Schelter. 2022. Automated Data Cleaning Can Hurt Fairness in Machine Learning-based Decision Making. *ICDE* (2022).
- [29] Zayd Hammoudeh and Daniel Lowd. 2024. Training data influence analysis and estimation: A survey. *Machine Learning* 113, 5 (2024), 2351–2403.
- [30] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudik, and Hanna Wallach. 2019. Improving fairness in machine learning systems: What do industry practitioners need?. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–16.
- [31] Jiayuan Huang, Arthur Gretton, Karsten Borgwardt, Bernhard Schölkopf, and Alex Smola. 2006. Correcting sample selection bias by unlabeled data. *Advances in neural information processing systems* 19 (2006).
- [32] Jinyuan Jia, Xiaoyu Cao, and Neil Zhenqiang Gong. 2021. Intrinsic certified robustness of bagging against data poisoning attacks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 7961–7969.

- [33] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nezihe Merve Gürel, Bo Li, Ce Zhang, Costas J Spanos, and Dawn Song. 2019. Efficient task-specific data valuation for nearest neighbor algorithms. *arXiv preprint arXiv:1908.08619* (2019).
- [34] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li, Ce Zhang, Dawn Song, and Costas J Spanos. 2019. Towards efficient data valuation based on the shapley value. In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 1167–1176.
- [35] Ruoxi Jia, Fan Wu, Xuehui Sun, Jiachen Xu, David Dao, Bhavya Kaikhura, Ce Zhang, Bo Li, and Dawn Song. 2021. Scalability vs. utility: Do we have to sacrifice one for the other in data importance quantification?. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8239–8247.
- [36] Heinrich Jiang and Ofir Nachum. 2020. Identifying and correcting label bias in machine learning. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 702–712.
- [37] Kevin Jiang, Weixin Liang, James Y Zou, and Yongchan Kwon. 2023. Opendataval: a unified benchmark for data valuation. *Advances in Neural Information Processing Systems* 36 (2023).
- [38] Faisal Kamiran and Toon Calders. 2012. Data preprocessing techniques for classification without discrimination. *Knowledge and information systems* 33, 1 (2012), 1–33.
- [39] Bojan Karlaš, David Dao, Matteo Interlandi, Sebastian Schelter, Wentao Wu, and Ce Zhang. 2023. Data Debugging with Shapley Importance over Machine Learning Pipelines. In *The Twelfth International Conference on Learning Representations*.
- [40] Bojan Karlaš, Peng Li, Renzhi Wu, Nezihe Merve Gürel, Xu Chu, Wentao Wu, and Ce Zhang. 2020. Nearest Neighbor Classifiers over Incomplete Information: From Certain Answers to Certain Predictions. *Proc. VLDB Endow.* 14, 3 (2020), 255–267. <https://doi.org/10.5555/3430915.3442426>
- [41] Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*. PMLR, 1885–1894.
- [42] Sanjay Krishnan, Jiannan Wang, Eugene Wu, Michael J Franklin, and Ken Goldberg. 2016. Activeclean: Interactive data cleaning for statistical modeling. *Proceedings of the VLDB Endowment* 9, 12 (2016), 948–959.
- [43] Yongchan Kwon and James Zou. 2021. Beta shapley: a unified and noise-reduced data valuation framework for machine learning. *arXiv preprint arXiv:2110.14049* (2021).
- [44] Peng Li, Zhiyi Chen, Xu Chu, and Kexin Rong. 2023. DiffPrep: Differentiable Data Preprocessing Pipeline Search for Learning over Tabular Data. *Proceedings of the ACM on Management of Data* 1, 2 (2023), 1–26.
- [45] Anqi Liu and Brian Ziebart. 2014. Robust classification under sample selection bias. *Advances in neural information processing systems* 27 (2014).
- [46] Xuan Luo and Jian Pei. 2024. Applications and Computation of the Shapley Value in Databases and Machine Learning. In *Companion of the 2024 International Conference on Management of Data (SIGMOD/PODS '24)*. Association for Computing Machinery, New York, NY, USA, 630–635. <https://doi.org/10.1145/3626246.3654680>
- [47] Xiaozhong Lyu, Stefan Grafberger, Samantha Biegel, Shaopeng Wei, Meng Cao, Sebastian Schelter, and Ce Zhang. 2023. Improving retrieval-augmented large language models via data importance learning. *arXiv preprint arXiv:2307.03027* (2023).
- [48] Mohammad Mahdavi and Ziawasch Abedjan. 2020. Baran: Effective error correction via a unified context representation and transfer learning. *Proceedings of the VLDB Endowment* 13, 12 (2020), 1948–1961.
- [49] Mark Mazumder, Colby Banbury, Xiaozhe Yao, Bojan Karlaš, William Gaviria Rojas, Sudnya Damos, Greg Damos, Lynn He, Alicia Parrish, Hannah Rose Kirk, et al. 2024. Dataperf: Benchmarks for data-centric ai development. *Advances in Neural Information Processing Systems* 36 (2024).
- [50] Sean McGregor. 2021. Preventing repeated real world AI failures by cataloging incidents: The AI incident database. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 15458–15463.
- [51] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* 54, 6 (2021), 1–35.
- [52] Xiangrui Meng, Joseph Bradley, Burak Yavuz, Evan Sparks, Shivaram Venkataraman, Davies Liu, Jeremy Freeman, DB Tsai, Manish Amde, Sean Owen, et al. 2016. Mllib: Machine learning in apache spark. *Journal of Machine Learning Research* 17, 34 (2016), 1–7.
- [53] Metaflow.org. 2024. A framework for real-life ML, AI, and data science. <https://metaflow.org>.
- [54] Anna Meyer, Aws Albarghouthi, and Loris D’Antoni. 2021. Certifying Robustness to Programmable Data Bias in Decision Trees. *Advances in Neural Information Processing Systems* 34 (2021).
- [55] Anna P Meyer, Aws Albarghouthi, and Loris D’Antoni. 2023. The dataset multiplicity problem: How unreliable data impacts predictions. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 193–204.
- [56] Microsoft. 2018. Azure Machine Learning Pipelines. <https://learn.microsoft.com/en-us/azure/machine-learning/concept-ml-pipelines>.
- [57] Mohammad Hossein Namaki, Avriela Floratou, Fotis Psallidas, Subru Krishnan, Ashvin Agrawal, Yinghui Wu, Yiwen Zhu, and Markus Weimer. 2020. Vamsa: Automated provenance tracking in data science scripts. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 1542–1551.
- [58] Hongseok Namkoong and John C Duchi. 2016. Stochastic Gradient Methods for Distributionally Robust Optimization with f-divergences. In *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.), Vol. 29. Curran Associates, Inc. [https://proceedings.neurips.cc/paper\\_files/paper/2016/file/4588e674d3f0fa985047d4c3f13ed0d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/4588e674d3f0fa985047d4c3f13ed0d-Paper.pdf)
- [59] Curtis Northcutt, Lu Jiang, and Isaac Chuang. 2021. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research* 70 (2021), 1373–1411.
- [60] Kwanghyun Park, Karla Saur, Dalitso Banda, Rathijit Sen, Matteo Interlandi, and Konstantinos Karanasos. 2022. End-to-end optimization of machine learning prediction queries. In *Proceedings of the 2022 International Conference on Management of Data*. 587–601.
- [61] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* 12 (2011), 2825–2830.
- [62] Alireza Pirhadi, Mohammad Hossein Moslemi, Alexander Cloninger, Mostafa Milani, and Babak Salimi. 2024. Otlean: Data cleaning for conditional independence violations using optimal transport. *Proceedings of the ACM on Management of Data* 2, 3 (2024), 1–26.
- [63] Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. 2020. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems* 33 (2020), 17044–17056.
- [64] Neoklis Polyzotis, Martin Zinkevich, Sudip Roy, Eric Breck, and Steven Whang. 2019. Data validation for machine learning. *Proceedings of machine learning and systems* 1 (2019), 334–347.
- [65] Romila Pradhan, Aditya Lahiri, Sainyam Galhotra, and Babak Salimi. 2022. Explainable ai: Foundations, applications, opportunities for data management research. In *Proceedings of the 2022 International Conference on Management of Data (SIGMOD/PODS '22)*. 2452–2457.
- [66] Romila Pradhan, Jiongli Zhu, Boris Glavic, and Babak Salimi. 2022. Interpretable data-based explanations for fairness debugging. In *Proceedings of the 2022 International Conference on Management of Data*. 247–261.
- [67] Fotis Psallidas, Yiwen Zhu, Bojan Karlaš, Jordan Henkel, Matteo Interlandi, Subru Krishnan, Brian Kroth, Venkatesh Emani, Wentao Wu, Ce Zhang, et al. 2022. Data science through the looking glass: Analysis of millions of github notebooks and ml. net pipelines. *ACM SIGMOD Record* 51, 2 (2022), 30–37.
- [68] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *Proc. VLDB Endow.* 10, 11 (2017), 1190–1201. <https://doi.org/10.14778/3137628.3137631>
- [69] Ashkan Rezaei, Anqi Liu, Omid Memarrast, and Brian Ziebart. 2020. Robust Fairness under Covariate Shift. *arXiv preprint arXiv:2010.05166* (2020).
- [70] Elan Rosenfeld, Ezra Winston, Pradeep Ravikumar, and Zico Kolter. 2020. Certified robustness to label-flipping attacks via randomized smoothing. In *International Conference on Machine Learning*. PMLR, 8230–8241.
- [71] Sebastian Schelter and Stefan Grafberger. 2024. Messy Code Makes Managing ML Pipelines Difficult? Just Let LLMs Rewrite the Code! *arXiv:2409.10081* [cs.DB] <https://arxiv.org/abs/2409.10081>
- [72] Sebastian Schelter, Stefan Grafberger, Shubha Guha, Bojan Karlaš, and Ce Zhang. 2023. Proactively screening machine learning pipelines with arguseyes. In *Companion of the 2023 International Conference on Management of Data*. 91–94.
- [73] Maximilian E Schüle, Luca Scalerandi, Alfons Kemper, and Thomas Neumann. 2023. Blue Elephants Inspecting Pandas: Inspection and Execution of Machine Learning Pipelines in SQL. In *EDBT*. 40–52.
- [74] Scikit-learn. 2024. Pipelines and composite estimators. <https://scikit-learn.org/stable/modules/compose.html>.
- [75] Supreeth Shastri, Vinay Banakar, Melissa Wasserman, Arun Kumar, and Vijay Chidambaram. [n. d.]. Understanding and Benchmarking the Impact of GDPR on Database Systems. *Proceedings of the VLDB Endowment* 13, 7 ([n. d.]).
- [76] Shafaq Siddiqi, Roman Kern, and Matthias Boehm. 2023. SAGA: A Scalable Framework for Optimizing Data Cleaning Pipelines for Machine Learning Applications. *Proceedings of the ACM on Management of Data* 1, 3 (2023), 1–26.
- [77] Jacob Steinhardt, Pang Wei Koh, and Percy Liang. 2017. Certified defenses for data poisoning attacks. In *NeurIPS*.
- [78] Julia Stoyanovich, Serge Abiteboul, Bill Howe, HV Jagadish, and Sebastian Schelter. 2022. Responsible data management. *Commun. ACM* 65, 6 (2022), 64–74.
- [79] DigiChina Stanford University. [n. d.]. Internet Information Service Algorithmic Recommendation Management Provisions. <https://digichina.stanford.edu/work/translation-internet-information-service-algorithmic-recommendation-management-provisions-opinion-seeking-draft/>
- [80] Jiachen T Wang and Ruoxi Jia. 2023. Data banzhaf: A robust data valuation framework for machine learning. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 6388–6421.

- [81] Wes McKinney. 2010. Data Structures for Statistical Computing in Python. In *Proceedings of the 9th Python in Science Conference*, Stéfan van der Walt and Jarrod Millman (Eds.), 56 – 61. <https://doi.org/10.25080/Majora-92bf1922-00a>
- [82] Steven Euijong Whang and Jae-Gil Lee. 2020. Data collection and quality challenges for deep learning. *Proceedings of the VLDB Endowment* 13, 12 (2020), 3429–3432.
- [83] Weiyuan Wu, Lampros Flokas, Eugene Wu, and Jiannan Wang. 2020. Complaint-driven training data debugging for query 2.0. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 1317–1334.
- [84] Doris Xin, Hui Miao, Aditya Parameswaran, and Neoklis Polyzotis. 2021. Production machine learning pipelines: Empirical analysis and optimization opportunities. In *Proceedings of the 2021 international conference on management of data*. 2639–2652.
- [85] Gyeong-In Yu, Saeed Amizadeh, Sehoon Kim, Artidoro Pagnoni, Ce Zhang, Byung-Gon Chun, Markus Weimer, and Matteo Interlandi. 2021. WindTunnel: towards differentiable ML pipelines beyond a single model. *Proceedings of the VLDB Endowment* 15, 1 (2021), 11–20.
- [86] Matei Zaharia, Mosharaf Chowdhury, Michael J Franklin, Scott Shenker, and Ion Stoica. 2010. Spark: Cluster computing with working sets. In *2nd USENIX workshop on hot topics in cloud computing (HotCloud 10)*.
- [87] Richard S. Zemel, Yu Wu, Kevin Swersky, Toniann Pitassi, and Cynthia Dwork. 2013. Learning Fair Representations. In *ICML (3) (JMLR Workshop and Conference Proceedings, Vol. 28)*. JMLR.org, 325–333.
- [88] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. 2018. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. 335–340.
- [89] Hantian Zhang, Ki Hyun Tae, Jaeyoung Park, Xu Chu, and Steven Euijong Whang. 2023. iFlipper: Label Flipping for Individual Fairness. *Proceedings of the ACM on Management of Data* 1, 1 (2023), 1–26.
- [90] Xuezhou Zhang, Xiaojin Zhu, and Stephen Wright. 2018. Training set debugging using trusted items. In *AAAI*.
- [91] Yixuan Zhang, Boyu Li, Zenan Ling, and Feng Zhou. 2023. Mitigating Label Bias in Machine Learning: Fairness through Confident Learning. *arXiv preprint arXiv:2312.08749* (2023).
- [92] Cheng Zhen, Nischal Aryal, Arash Termehchy, and Amandeep Singh Chabada. 2024. Certain and Approximately Certain Models for Statistical Learning. *Proceedings of the ACM on Management of Data* 2, 3 (2024), 1–25.
- [93] Jiongli Zhu, Su Feng, Boris Glavic, and Babak Salimi. 2024. Learning from Uncertain Data: From Possible Worlds to Possible Models. *NeurIPS* (2024).
- [94] Jiongli Zhu, Sainyam Galhotra, Nazanin Sabri, and Babak Salimi. 2023. Consistent range approximation for fair predictive modeling. *VLDB* (2023).