



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

Bojan Karlaš (Harvard University), Babak Salimi (UC San Diego), Sebastian Schelter (BIFOLD & TU Berlin)



navigating-data-errors.github.io

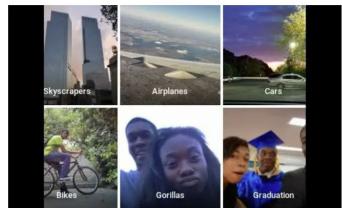


Background: ML apps often behave in unintended ways

Wrong

Google apologises for Photos app's racist blunder

© 1 July 2015



Source: BBC

Biased

MIT
Technology
Review

Featured Topics Newsletters Events Audio

SIGN IN

SUBSCRIBE

SILICON VALLEY

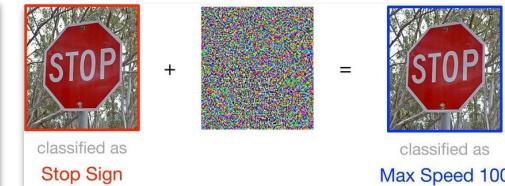
Amazon ditched AI recruitment software because it was biased against women

By Erin Winick

October 10, 2018

Source: MIT Technology Review

Unstable



Source: Xiong et al. ACM Comput. Surv. 2023.



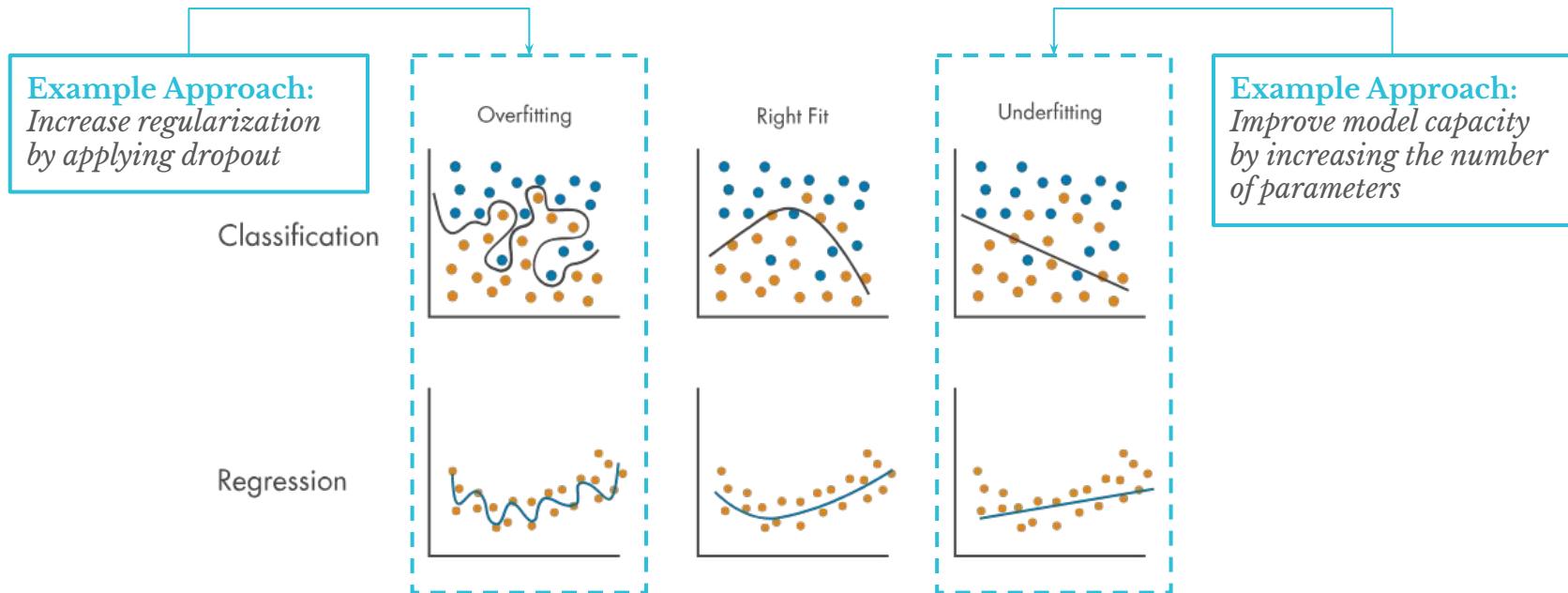
Tesla Autopilot feature was involved in 13 fatal crashes, US regulator says

Federal transportation agency finds Tesla's claims about feature don't match their findings and opens second investigation



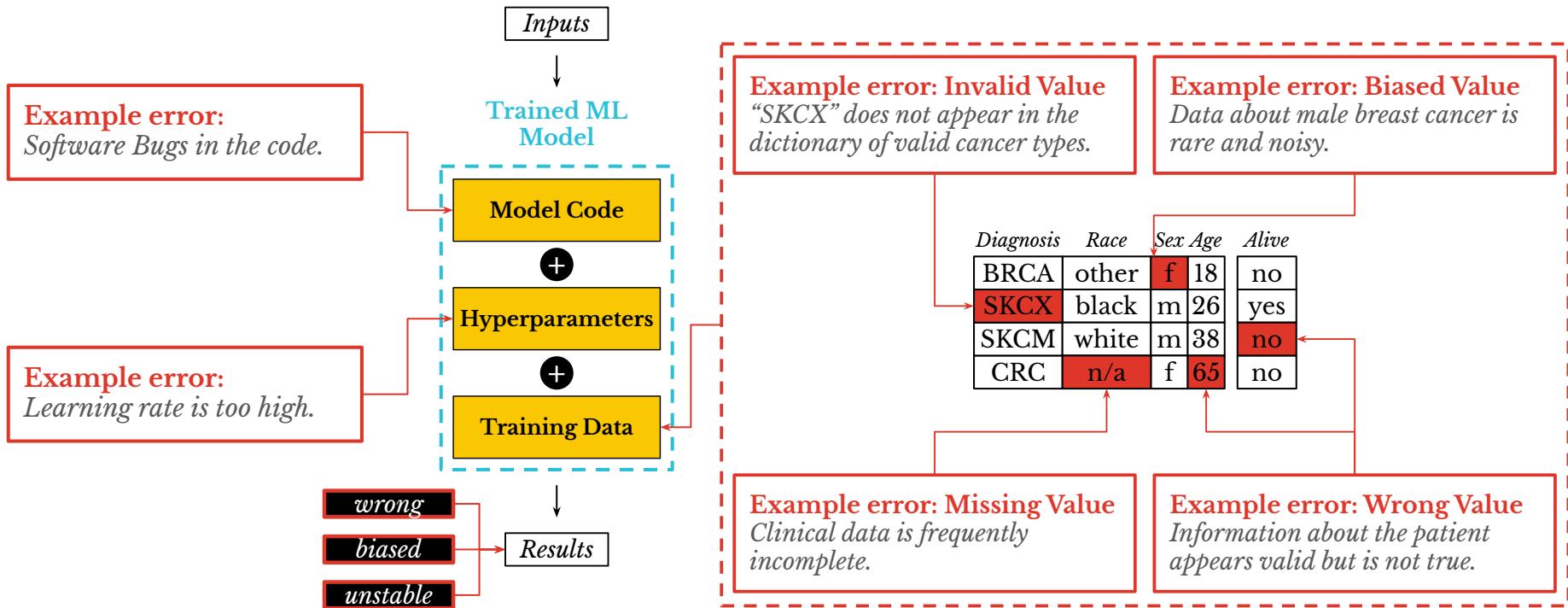
Source: The Guardian

Primary approach: Focus on improving the model



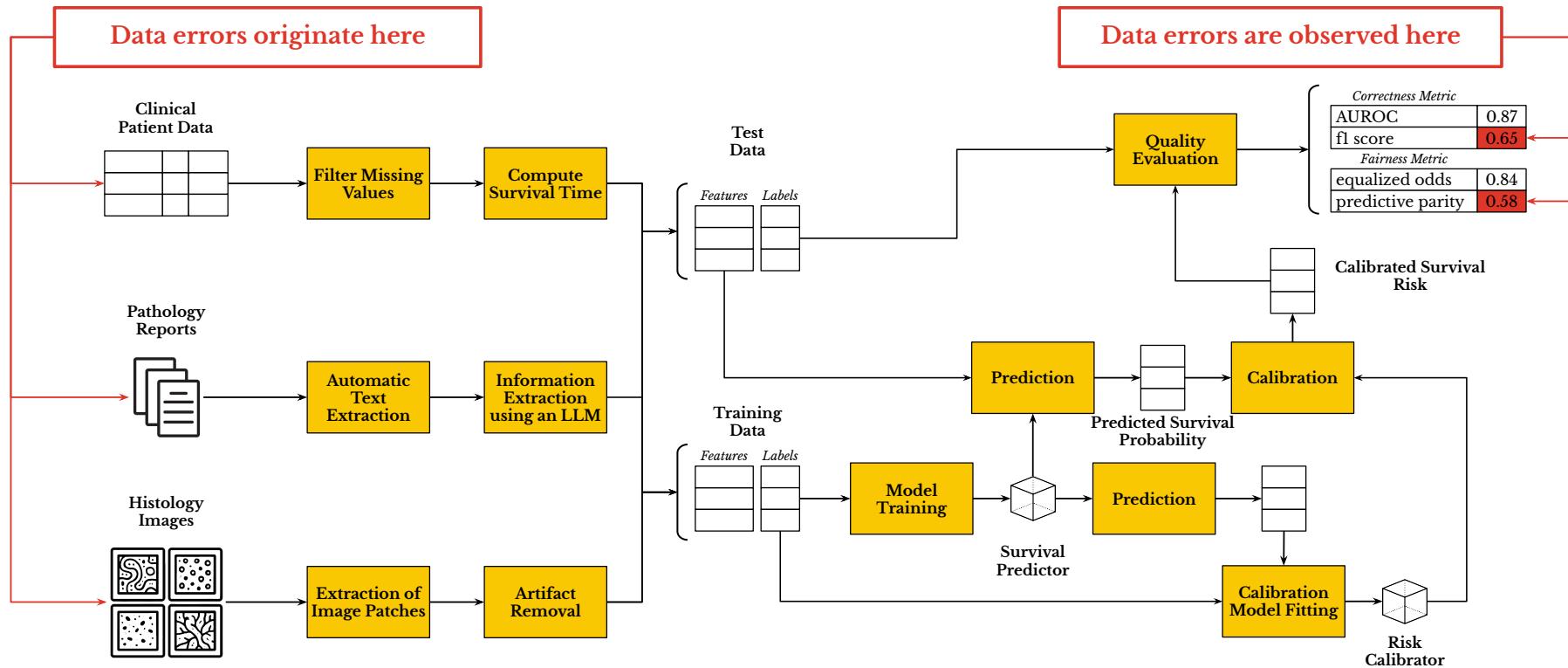
Problem: This is only one piece of the puzzle!

Observation 1: Data is a crucial piece of the puzzle



Challenge 1: Can we identify the most important data errors?

Observation 2: ML apps are built by complex pipelines



Challenge 2: Can we trace data errors as they pass through the pipeline?

Observation 3: Not all data errors are meant to be fixed

For each data error, we can choose to perform one of the following actions:

Discard



Remove the faulty data from the training set.

Repair



Perform manual quality control which might include repeating the data acquisition process.

Ignore



Let the faulty data remain in the training set.

Benefits:

Easy to Perform

Data Quality Improves

No Labor Required

Shortcomings:

Loss of Useful Data

Often Labor-intensive

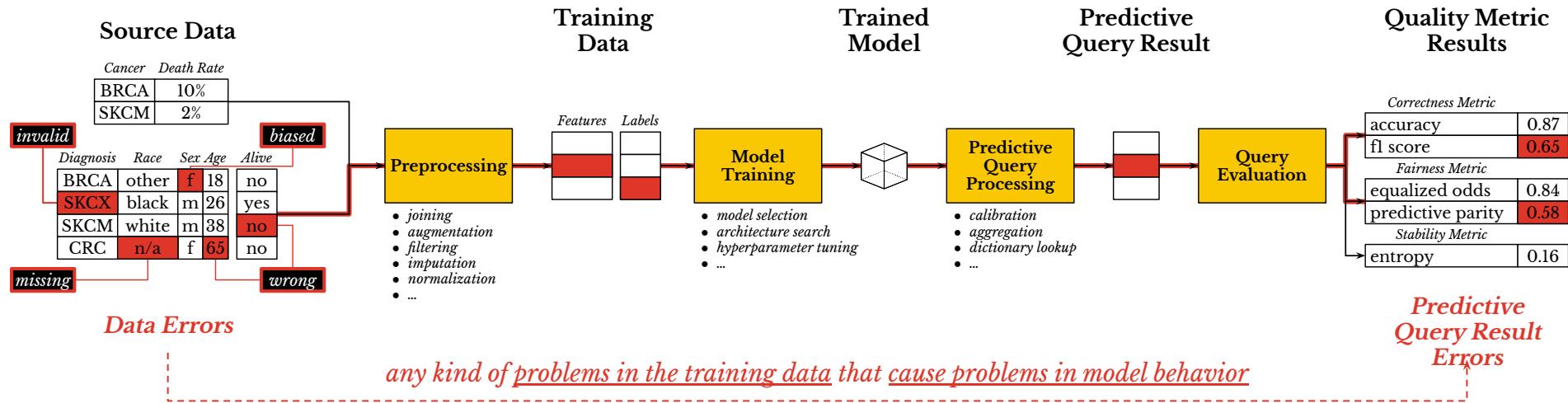
Risk Hurting Model Quality

Optimal trade-off:

Discard or Repair the Portion of Data that will Bring the Highest Model Quality Increase

Challenge 3: Can we ensure reliable model performance after (partial) data repairs?

Tutorial Overview: Data Errors in ML pipelines



Part I: Data Importance for Data Error Detection

What are good approaches for identifying data errors?

Part II: Data Debugging in ML Pipelines

What are practical challenges when debugging complex ML pipelines?

Part III: Learning from Uncertain and Incomplete Data

When we cannot repair all errors, can we still have reliable models?

Opportunities for the Data Management Community

- (1) Data quality is an established discipline in data management, but most practitioners still rely on **manual effort**.
- (2) ML pipelines are data processing pipelines. Models are learned data transformation operators. Many systems have been developed, but most practitioners still rely on **rudimentary scripts for crunching data**.
- (3) Many promising methods for handling data errors suffer from **scalability issues**.

Main Goal: *Present the current state of the art and inspire novel research.*

Part I: Data Importance for Data Error Detection

Bojan Karlaš



- 1) Introducing the Concept of Data Importance**
- 2) Examples of Data Attribution Functions**
- 3) Case Study of Shapley Value as a Measure of Importance**
- 4) Applications of Data Importance**

How can we identify data errors?

Trivial

Solution approach:

Apply a rule-based validation function that performs a dictionary lookup.

invalid

Diagnosis	Race	Sex	Age	Alive
BRCA	other	f	18	no
SKCX	black	m	26	yes
SKCM	white	m	38	no
CRC	n/a	f	65	no

Solution approach:

Check if the value is marked as missing.

missing

Not So Trivial

Solution approach:

Measure the impact of the value on model quality.

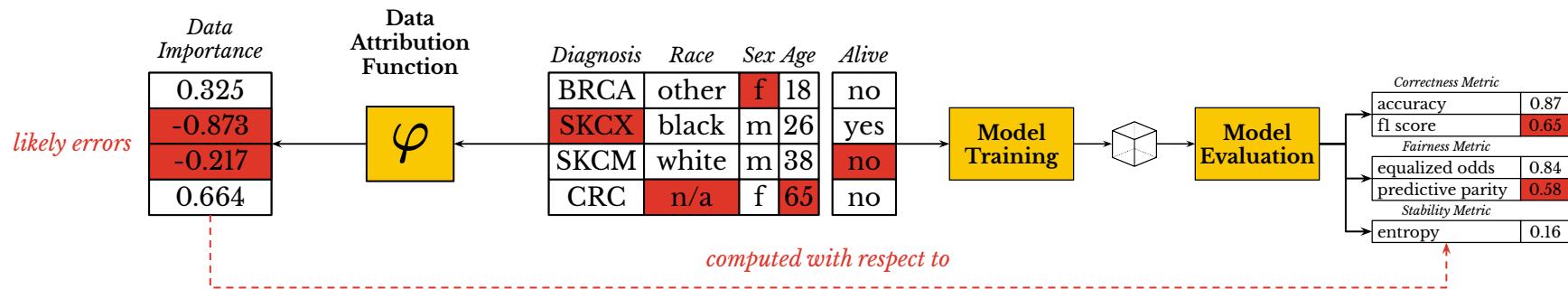
How do we measure this?

That is the main topic of this part of the tutorial.

Recall: Data errors are any kind of problem in the training data that cause problems in model behavior.

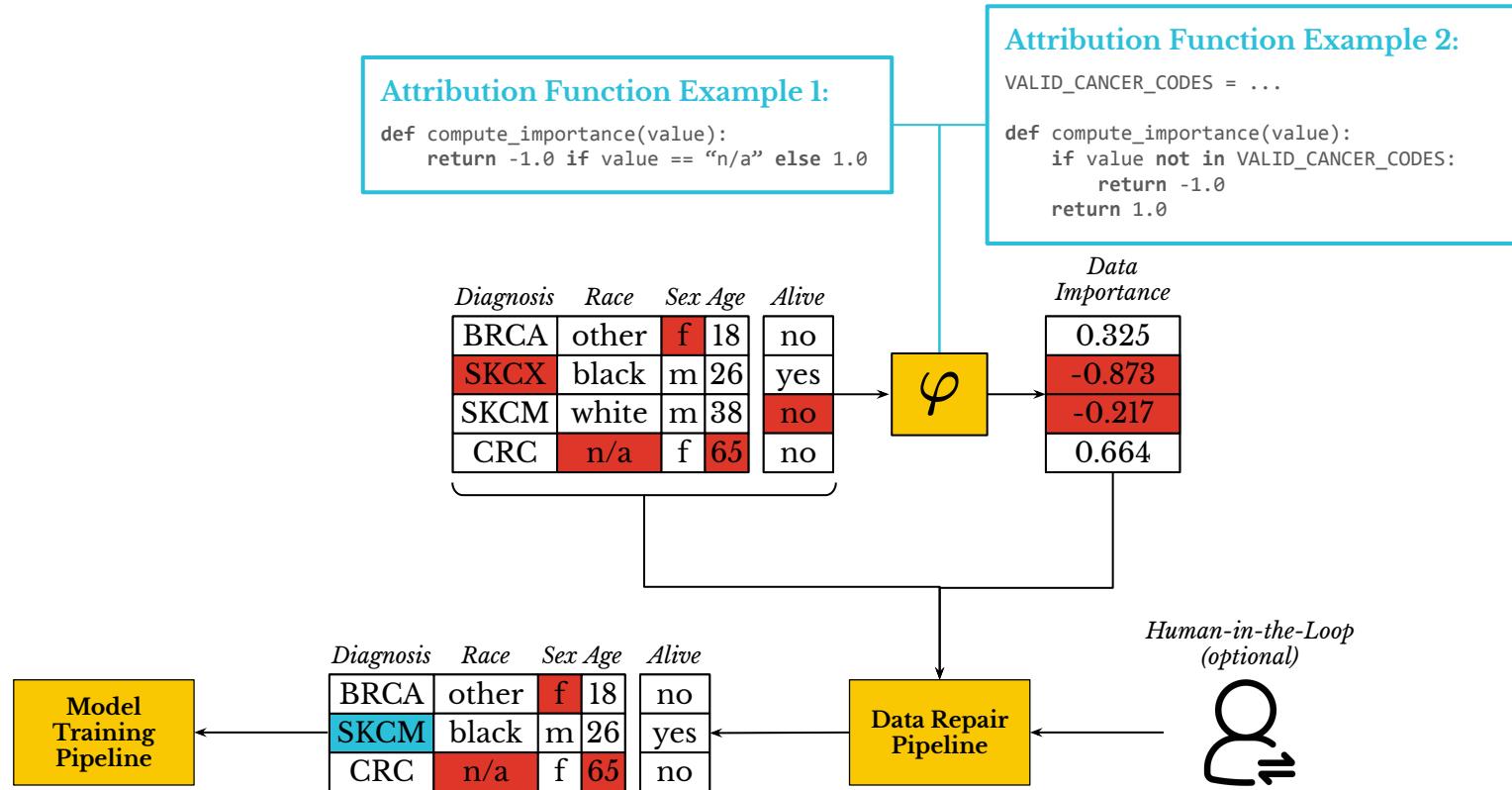
Challenge: Can we define a unified way to think about identifying data errors?

We can define a data attribution function



Recall: Data errors are any kind of problem in the training data that cause problems in model behavior.

How do we use importance to detect data errors?



What makes a good attribution function?

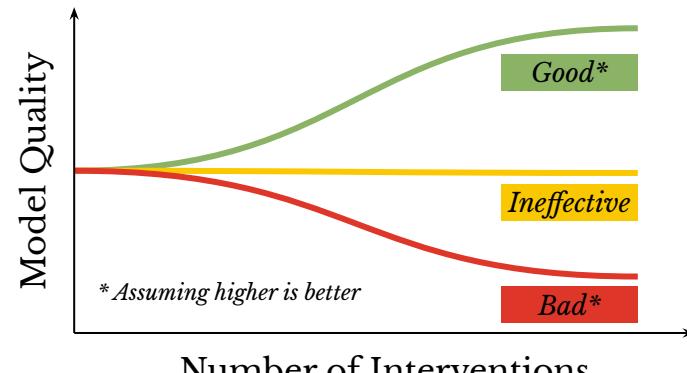
Design Consideration 1

Which model quality metric do we care about improving?

Correctness Metric
accuracy
f1 score
Fairness Metric
equalized odds
predictive parity
Stability Metric
entropy

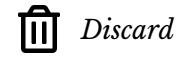
Recall:

Data errors are any kind of problem in the training data that cause problems in model behavior.



Design Consideration 2

What kind of intervention do we intend to apply?



Discard



Repair



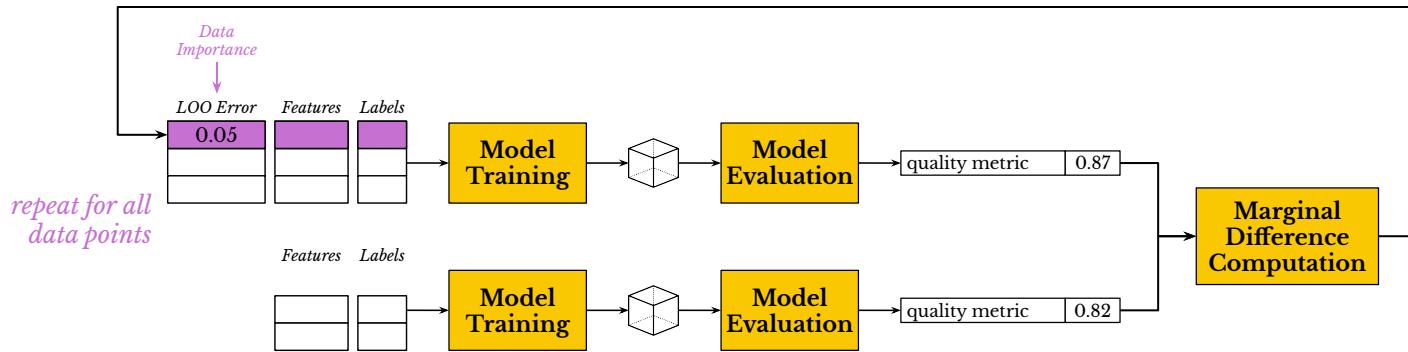
Something Else

Challenge: How do we define an effective attribution function?

- 1) Introducing the Concept of Data Importance
- 2) Examples of Data Attribution Functions**
- 3) Case Study of Shapley Value as a Measure of Importance
- 4) Applications of Data Importance

Leave-one-Out Error

[Approach: Marginal Contribution]



Insights:

- Removing important data points affects model quality.

Approach:

- Remove a data point from the training set, train and evaluate the model again
- Interpret the difference in model quality as data importance.

Benefits:

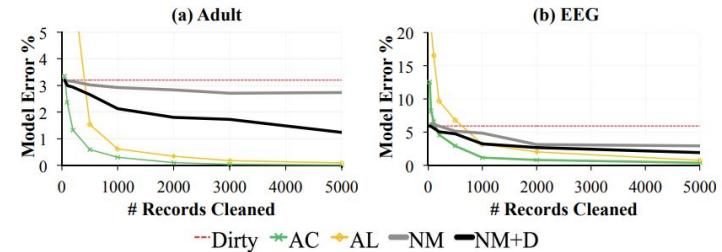
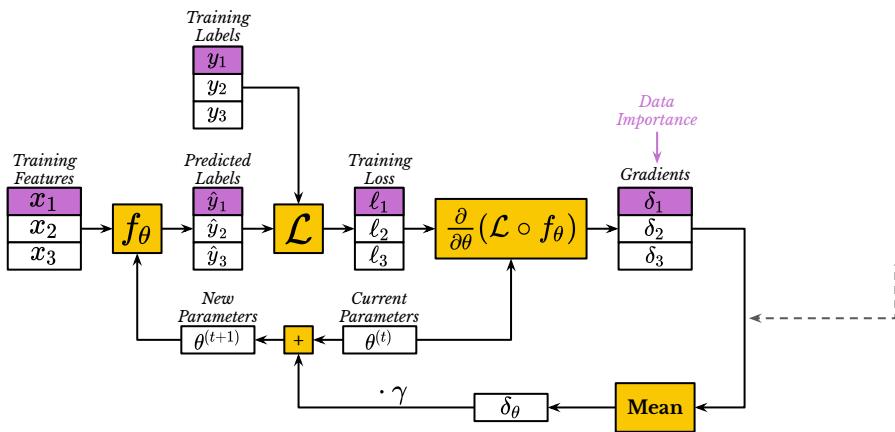
- Very simple to implement.

Shortcomings:

- Requires re-training the model once for each data point.
- Treats data points independently.

Error Gradient

[Approach: Gradient]



Insights:

- Data points vary in their contribution to the gradients that update the model.

Approach:

- Importance is proportional to the magnitude of the gradient.

Benefits:

- Simple to compute.

Shortcomings:

- Treats data points independently.

ActiveClean: Interactive Data Cleaning For Statistical Modeling

Sanjay Krishnan, Jianan Wang*, Eugene Wu*, Michael J. Franklin, Ken Goldberg
UC Berkeley, Simon Fraser University, *Columbia University
sanjaykrishnan, franklin, goldberg@berkeley.edu, jwang@sfu.ca, www.cs.columbia.edu

ABSTRACT

Analyzing often-clean data frequently cleaning some data, executing the analysis, and then cleaning more data based on the results. This is a common pattern in machine learning and statistical modeling, which is an increasingly popular form of data science. However, this pattern can lead to poor performance if one uses and then reuses the same data to both clean it and then train a model. We propose ActiveClean, an iterative framework for cleaning data that can handle mixed dirty and clean data. We evaluate ActiveClean on five real-world datasets and show that it outperforms state-of-the-art data cleaning tools with both real and synthetic errors. The results show that our proposed framework can significantly reduce the number of samples of only clean data needed to correctly train a model for the same amount of data cleaned. Furthermore, for a fixed clean budget, ActiveClean can produce models that are more accurate than either sampling or Active Learning.

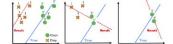


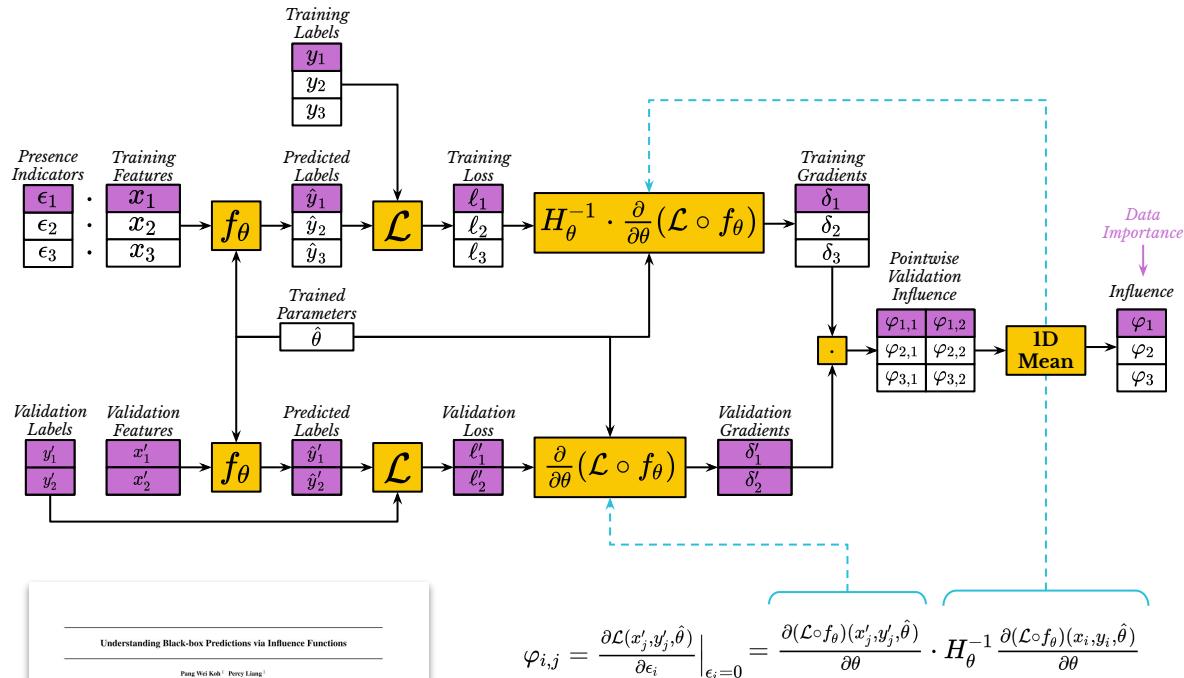
Figure 1(a) Symmetric corruption in one variable can lead to a local minimum. The plot shows a surface with two minima. One path starts at a point labeled '1' and goes down to a local minimum labeled '2'. Another path starts at point '1' and goes up to a local maximum labeled '3'. From there, it goes down to a global minimum labeled '4'.

Krishnan VLDB'16

Krishnan, Sanjay, et al. "Activeclean: Interactive data cleaning for statistical modeling." Proceedings of the VLDB Endowment 9.12 (2016): 948-959. [\[Paper\]](#) [\[Website\]](#)

Influence Function

[Approach: Marginal Contribution, Gradient]



Understanding Black-box Predictions via Influence Functions

Pang Wei Koh · Percy Liang¹

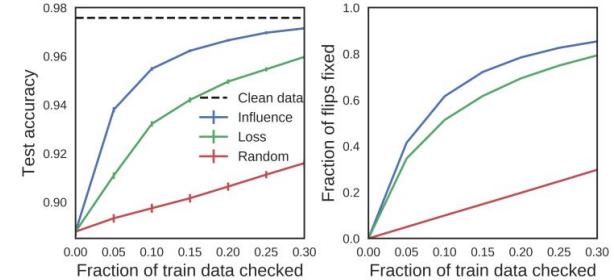
Abstract
 How can we explain the predictions of a black-box model? In this paper, we use influence functions — a classic technique in statistics — to trace a model’s prediction through the training algorithm back to its training data, thereby identifying training data most responsible for a given prediction. To set up influence functions for black-box models, we first introduce presence indicator variables for each data point. We develop a simple, efficient implementation that requires only gradient calculations and Hessian-vector products. We show that even on complex neural networks, influence functions are the theory break down, approximations to influence functions can still provide valuable information. On the other hand, for simple linear models, we demonstrate that influence functions capture the model’s local behavior, debugging models, de-

scribing them, and understanding model behavior.

paper [Koh et al., 2016] by pointing the user to see the theory [Datta et al., 2016; Mitra et al., 2017; Li et al., 2016; Datta et al., 2016; Mitra et al., 2017]. These papers also provide a detailed discussion of how we can explain where the model come from?

In this paper, we take this approach to train the model’s behavior through the lens of influence functions. In particular, we focus on influence functions for a prediction, we ask the counterfactual: what would happen if one of the values of this training point were changed slightly?

Answering this question by perturbing the data and retraining the model is a standard approach to understanding this problem, we are influence functions, a classic technique from robust statistics [Cook & Weisberg, 1994] and machine learning [Koh & Liang, 2017], we can do this without retraining the model. As we highlight a testing point as an influential amount. This allows us to understand the local behavior of the model, the local



Insights:

- The marginal contribution of a single data point can be approximated with gradients.

Approach:

- Introduce presence indicator variables ϵ for each data point and compute the gradient w.r.t. ϵ .

Benefits:

- Easily applicable to arbitrarily complex (twice) differentiable machine learning models.

Shortcomings:

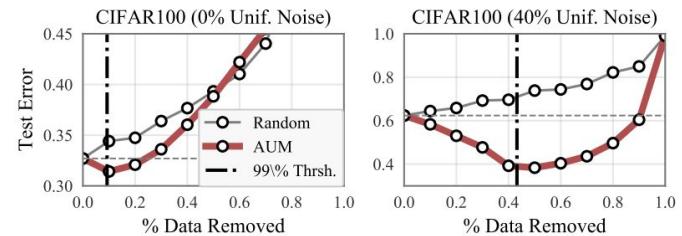
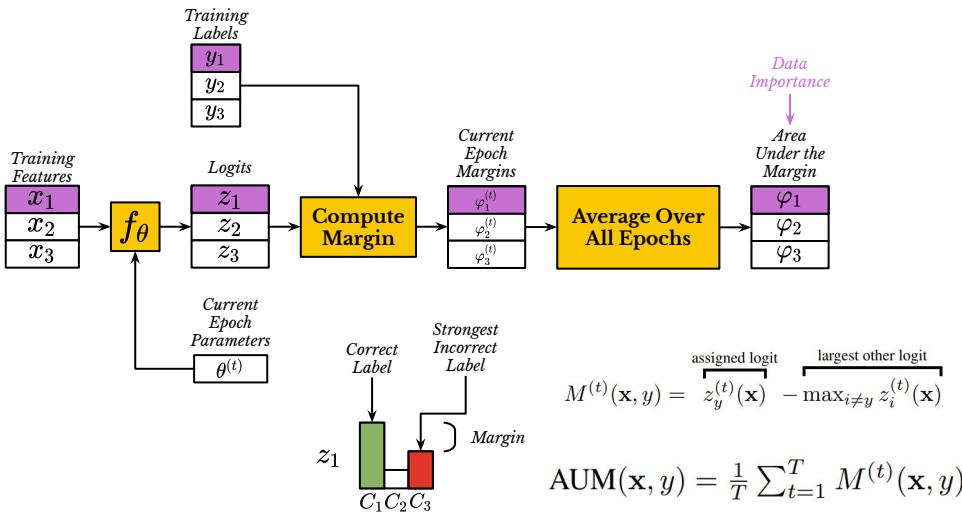
- Treats data points independently.

[Koh ICML ‘17]

Koh, Pang Wei, and Percy Liang. "Understanding black-box predictions via influence functions." International conference on machine learning. PMLR, 2017. [\[Paper\]](#) [\[Code\]](#)

Area Under the Margin

[Approach: Uncertainty Analysis]



Insights:

- If similar samples have the same label, the model will learn to activate only the correct logit.
- In the presence of mislabeled samples, the model will learn to activate alternative logits.

Approach:

- The importance of a data point is proportional to its margin averaged across all training epochs.

Benefits:

- Very simple to implement in a wide array of models.
- Does not rely on a separate clean dataset.

Shortcomings:

- Focuses only on label noise.

Identifying Mislabeled Data using the Area Under the Margin Ranking

Geoff Pleiss¹
Columbia University
gpp12@columbia.edu

Tianqi Zhang²
Stanford University
tianqi.zhang@stanford.edu

Ethan Fetaya³
ASAPP
ethanfetaya@asapp.com

Abstract

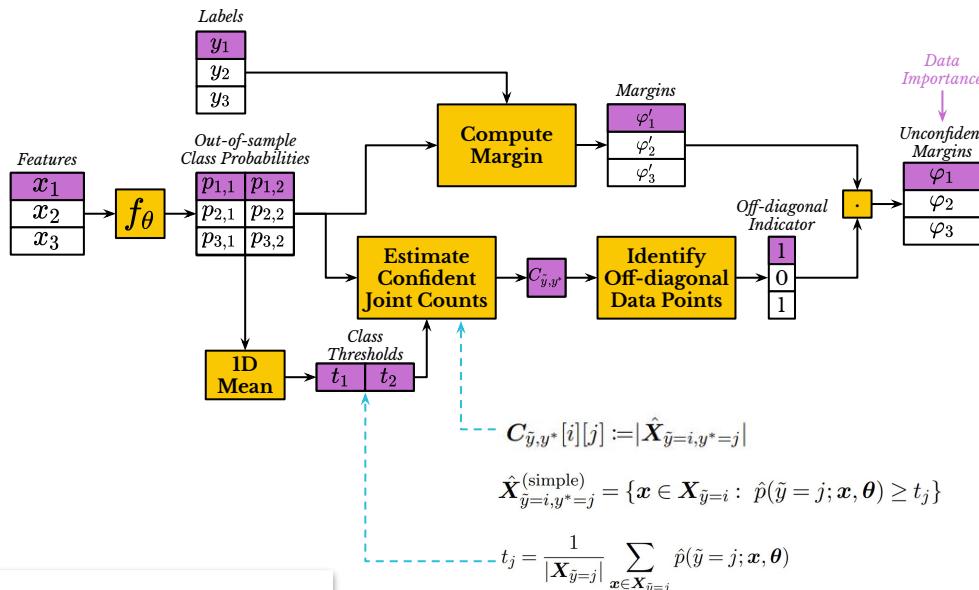
No all data is perfect training or helps generalization; some can be overly ambiguous or completely mislabeled. This paper introduces a new method to identify such data points. Our approach is based on the Area Under the Margin (AUM) statistic, which measures the importance of a data point. We propose a simple procedure—adding a new class populated with personally mislabeled samples—so that the AUM statistic can be computed. This approach consistently improves upon prior work on synthetic and real-world datasets.

[Pleiss NeurIPS '20]

Pleiss, Geoff, et al. "Identifying mislabeled data using the area under the margin ranking." Advances in Neural Information Processing Systems 33 (2020): 17044-17056. [\[Paper\]](#) [\[Blog\]](#) [\[Code\]](#)

Unconfident Margins

[Approach: Uncertainty Analysis]



Journal of Artificial Intelligence Research 70 (2021) 1373-1411
Submitted 06/2020; published 01/2021

Confident Learning: Estimating Uncertainty in Dataset Labels

Curtis G. Northcutt
Massachusetts Institute of Technology
Department of ECECS, Cambridge, MA, USA

Lu Jiang
Google Pixel, Mountain View, CA, USA

Isaac E. Chuang
Massachusetts Institute of Technology
Department of ECECS, Department of Physics, Cambridge, MA, USA

CONF@MIT.EDU

LJUAN@GOOGLE.COM

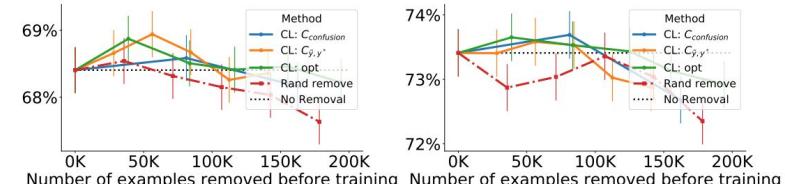
IUCHANG@MIT.EDU

Abstract

Learning rules in the context of datasets of softmaxes typically focus on model predictions, not label quality. Confident learning (CL) is an alternative approach which focuses on the uncertainty of labels. CL is a nonparametric method for datasets, based on the principles of pruning noisy data, using probabilistic thresholds to estimate noise, and ranking examples to train with confidence. Whereas numerous studies have shown that CL can be competitive with state-of-the-art methods, most rely on the assumption of a class-conditional noise process to directly estimate the joint distribution between labels and features. This paper shows that CL can be applied to datasets without this assumption.

[Northcutt JAIR '21]

Northcutt, Curtis, Lu Jiang, and Isaac Chuang. "Confident learning: Estimating uncertainty in dataset labels." Journal of Artificial Intelligence Research 70 (2021): 1373-1411. [\[Paper\]](#) [\[Blog\]](#) [\[Code\]](#)



Insights:

- Given a data point, if a model assigns a higher than average probability to some specific class, it is likely because most similar data points have the same class label. This is likely to be the true label of that data point.

Approach:

- Identify likely mislabeled data points and assign negative importance using the margin. Remaining data points get zero importance.

Benefits:

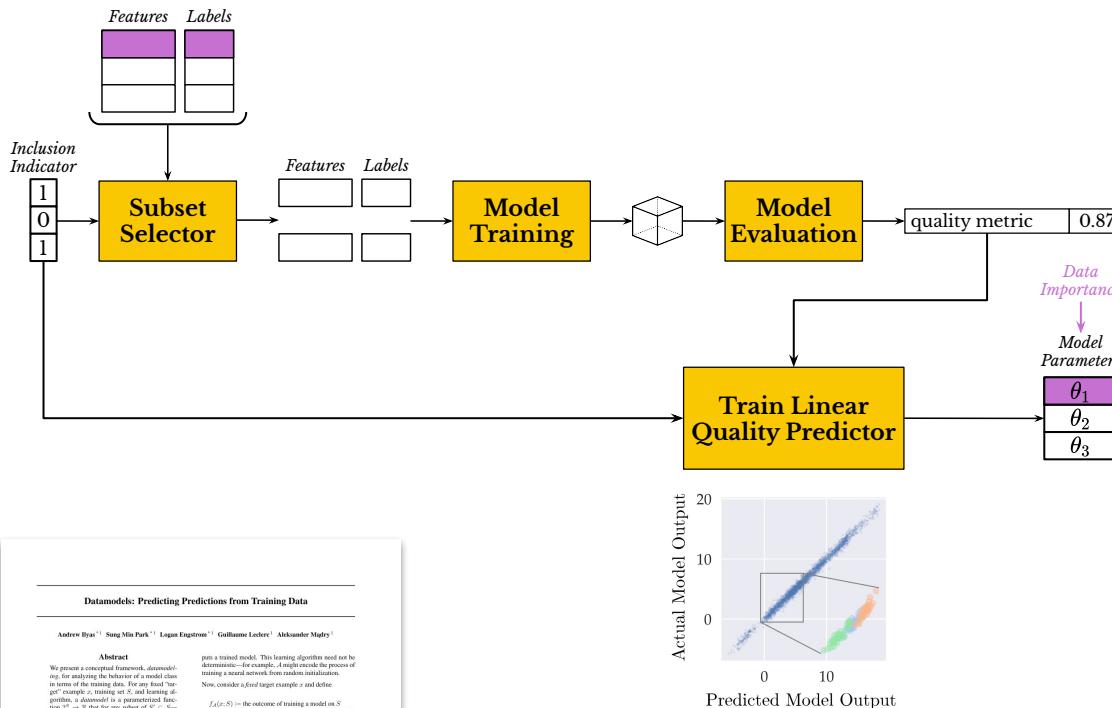
- Very simple to implement in a wide array of models.
- Does not rely on a separate clean dataset.

Shortcomings:

- Focuses only on label noise.
- Relies on having an adequately powerful model.

Model Training Outcome

[Approach: Surrogate Data Model]



[Ilyas ICML '22]

Ilyas, Andrew, et al. "Datamodels: Predicting Predictions from Training Data." Proceedings of the 39th International Conference on Machine Learning. 2022. [\[Paper\]](#) [\[Blog\]](#) [\[Code\]](#)

Insights:

- A linear model can be good at predicting the quality of a model trained on an arbitrary subset of the training data and tested on a single test example.

Approach:

- Train a linear quality predictor and interpret its parameters as data importance.

Benefits:

- Conceptually simple yet powerful framework for analyzing datasets.

Shortcomings:

- The original method requires retraining the model many times.

- 1) Introducing the Concept of Data Importance
- 2) Examples of Data Attribution Functions
- 3) Case Study of Shapley Value as a Measure of Importance**
- 4) Applications of Data Importance

Improving Upon the Marginal Contribution Methods

Recall

Marginal contribution methods treat data points independently, ignoring any interactions that might exist.

Consequence

Let there be a data point that has high importance. If we make two copies of that data point, their individual marginal contribution to the dataset as a whole will be zero.

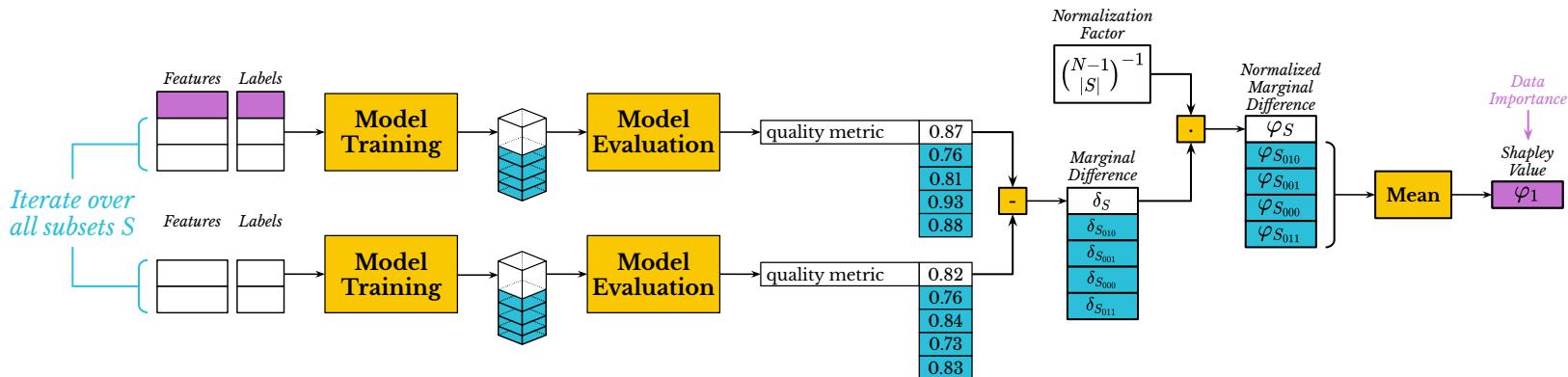
Shapley value

A standard method from game theory for distributing surplus among a coalition of players.

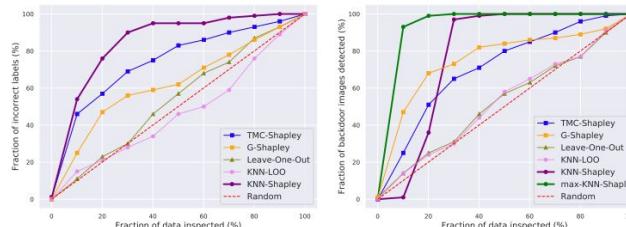
$$\varphi_i = \frac{1}{N} \sum_{S \subseteq X \setminus \{i\}} \binom{N-1}{|S|}^{-1} (u(S \cup \{i\}) - u(S))$$

Approach

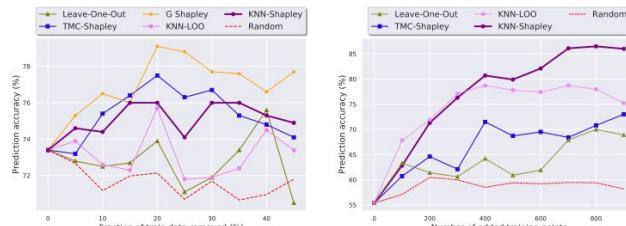
We should measure marginal contribution over all subsets.



Effectiveness at Data Debugging



(a) Noisy labels detection



(c) Data summarization

(d) Data acquisition

Figure 2: The experiment result of (a) noisy label detection on fashion-MNIST dataset; (b) instance-based watermark removal on MNIST dataset; (c) data summarization on UCI Adult Census dataset [15]; (d) data acquisition on MNIST dataset with injected noise. In (a)-(b) the “random” line shows the results of random guess; while in (c)-(d), the “random” line corresponds to the empirical results of the random baseline introduced in Section 4.1.

Table 2: Domain adaptation between MNIST and USPS.

Method	MNIST → USPS	USPS → MNIST
	→	→
KNN-Shapley	31.70% → 47.00%	23.35% → 29.80%
KNN-LOO	31.70% → 37.40%	23.35% → 24.50%
TMC-Shapley	31.70% → 44.90%	23.35% → 29.55%
LOO	31.70% → 29.40%	23.35% → 23.53%

This CVPR 2021 paper is the Open Access version, provided by the Computer Vision Foundation.
Except for the watermark, it is identical to the accepted version.
The final published version of the proceedings is available on IEEE Xplore.

GyF

Scalability vs. Utility: Do We Have to Sacrifice One for the Other in Data Importance Quantification?

Ruoxi Jia¹ Fan Wu^{2*} Xudong Sun³ Jiacen Xu⁴ David Dao⁵
¹Virginia Tech ²UIC ³Shantou Jiaotong University ⁴UC Irvine ⁵ETH Zurich
¹ruoxi.jia.cs@vt.edu ²fanwu@uic.edu ³xudong.su@stu.edu.cn ⁴jxu@math.uci.edu ⁵daod@math.ETHZ.ch

Abstract

I. Introduction

Understanding the importance of each training point is a fundamental problem in machine learning and the related research areas have been progressing to make a range of data workflow tools such as data summarization and data acquisition. In this paper, we propose the leave-one-out error of each training point to indicate its importance quantitatively. This metric is called Shapley value, as it defines a unique value distribution when

[Jia CVPR '21]

Jia, Ruoxi, et al. "Scalability vs. utility: Do we have to sacrifice one for the other in data importance quantification?" Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021. [\[Paper\]](#) [\[Code\]](#)

Benefits and Challenges

Beneficial Properties of the Shapley Value

Symmetry

If two data points have the same contribution to every subset, their value should be the same.

Efficiency

The sum of importances of all data points should equal the marginal contribution of the entire set over an empty set.

Linearity

If the utility function can be expressed as a sum of two other functions, then the importance of a data point using the combined function should equal the sum of importances computed using the individual functions.

Null Player

If a data point has a zero marginal contribution to every single subset, its importance should be zero.

Key Challenge

The number of subsets to enumerate is exponential, making it intractable to compute the exact Shapley value for an arbitrary model.

$$\varphi_i = \frac{1}{N} \sum_{S \subseteq X \setminus \{i\}} \binom{N-1}{|S|}^{-1} (u(S \cup \{i\}) - u(S))$$

Approximation: Monte Carlo Sampling

Challenge

Computing Shapley values is intractable.

Insight

Since Shapley value can be seen as a statistic over exponentially many subsets, we can estimate it using Monte Carlo sampling.

Approach

Use the permutation-based definition of the Shapley value and sample permutations.

$$\varphi_i(v) = \frac{1}{n!} \sum_R [v(P_i^R \cup \{i\}) - v(P_i^R)]$$

$$\phi_i = \mathbb{E}_{\pi \sim \Pi}[V(S_\pi^i \cup \{i\}) - V(S_\pi^i)]$$



Abstract
data valuation is to quantify the contribution of each training datum to the model's performance.

Data Shapley: Equitable Valuation of Data for Machine Learning

Amirata Ghorbani¹ James Zou²

¹note of the market place, similar to labor capital (Ghosh, 2020). It has been suggested that certain data carry individual property, and as such individuals should be compensated for the use of their data. Like labor and capital, a fundamental question is how to equitably value individual's data.

²note of the market place, similar to labor capital (Ghosh, 2020). It has been suggested that certain data carry individual property, and as such individuals should be compensated for the use of their data. Like labor and capital, a fundamental question is how to equitably value individual's data.

[Kwon AISTATS '22]

Kwon, Yongchan, and James Zou. "Beta Shapley: a Unified and Noise-reduced Data Valuation Framework for Machine Learning." International Conference on AI and Statistics. 2022. [\[Paper\]](#) [\[Code\]](#)

[Ghorbani ICML '19]

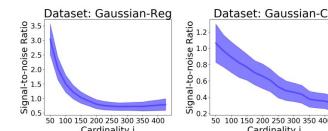
Ghorbani, Amirata, and James Zou. "Data shapley: Equitable valuation of data for machine learning." International conference on machine learning. PMLR, 2019. [\[Paper\]](#) [\[Code\]](#)

Challenge

We need many Monte Carlo samples to produce good estimates.

Insight

When estimating the marginal contribution of a data point to a subset, we empirically observe that larger subsets incur a slower signal-to-noise ratio.

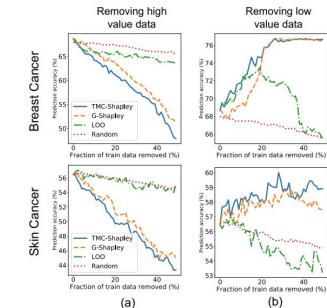


Approach

Leverage the importance sampling strategy and apply a larger weight to smaller subsets, based on the beta distribution.

Benefits

Estimating the Shapley value becomes tractable and is shown to be effective at identifying important data points.



Shortcomings

Each Monte Carlo sample relies on retraining the model from scratch, which is expensive for large models.

Approximation: K-Nearest Neighbor Surrogate Model

Challenge

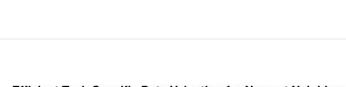
To get good Shapley value estimates, we need to retrain the model many times.

Insight

The simple KNN classifier can make it easy to design efficient and exact algorithms.

Approach

Use the KNN model as a proxy to develop an exact Shapley computation algorithm with polynomial time complexity.



Efficient Task-Specific Data Valuation for Nearest Neighbor Algorithms

Ruoxi Jia¹ David Dao² Boxin Wang³ Frances Ann Hubis⁴ Nezihie Merve Gurel⁵

Bo Li⁶ Ce Zhang⁷ Costas Spanos⁸ Dawn Song⁹

¹UC Berkeley, ²ETH Zurich, ³Zhejiang University, ⁴UIUC

ruoxi@berkeley.edu, danieladao@mit.edu, boxin.wang@zjhu.edu.cn, hubis@illinois.edu, nezihie.gurel@ethz.ch,

bo.li@berkeley.edu, ce.zhang@ethz.ch, costas.spanos@uiuc.edu, song.dawn@berkeley.edu

ABSTRACT

Given a dataset D containing millions of data points and a data point x , it is challenging to query for the k -th nearest neighbor (NN) of x . In this paper, we propose a task-specific valuation framework for data points in D based on their Shapley values, which is a uniquely promising proxy for data quality. We show that our proposed framework is able to handle data points with different dimensions and cardinalities. For general bounded utility functions, the Shapley values of all data points change to zero if a specific data point is removed. Thus, the valuation of a data point x is the sum of its Shapley values for all k data points that are neighbors of x .

In this paper, we focus on the popular family of ML models named KNN classifiers. We propose a fast approximation for KNN classifiers and regressions. The Shapley value of a data point x can be computed in $\mathcal{O}(N \log N)$ time, where N is the number of training samples. Our experiments show that our proposed framework is competitive with state-of-the-art methods.

Figure 1: Shapley Example of Data Valuation.

1. INTRODUCTION

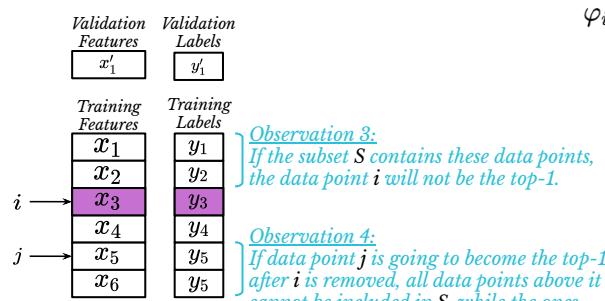
Jia, Ruoxi, et al. "Efficient task-specific data valuation for nearest neighbor algorithms." Proceedings of the VLDB Endowment 12.11 (2019): 1610-1623. [\[Paper\]](#) [\[Code\]](#)

Example Situation

- We are computing the Shapley value of data point i
- Data is sorted by similarity to the validation data point

Observation 1:

Since $K=1$, for any subset S , the top-1 data point will determine the model prediction.



Observation 3:

If the subset S contains these data points, the data point i will not be the top-1.

Observation 4:

If data point j is going to become the top-1 after i is removed, all data points above it cannot be included in S , while the ones below may or may not be included in S .

Starting point: Shapley value definition

$$\varphi_i = \frac{1}{N} \sum_{S \subseteq X \setminus \{i\}} \binom{N-1}{|S|}^{-1} (u(S \cup \{i\}) - u(S))$$

Observation 2:

If data point i is not in the top-1, this term will be zero.

Dynamic Programming

$$\varphi_i(t) = \frac{1}{N} \sum_{j=i+1}^N \sum_{a=1}^{n-j} \binom{N-1}{a}^{-1} (u(\{i\}) - u(\{j\})) \binom{N-j}{a}$$

Final Simplification

$$\varphi_i(t) = \frac{1}{N} \sum_{j=i+1}^N (u(\{i\}) - u(\{j\})) \binom{N-j}{j+1}$$

Result:

After sorting the data, we can compute exact Shapley values in a single pass. Final computational complexity is

$$\mathcal{O}(N \log N)$$

Approximation: Taylor Expansion

Challenge

If we are using a large and complex model, retraining will be extremely slow (preventing Monte Carlo approaches), and the KNN approximation will be biased.

Insight

Models trained with stochastic gradient descent (SGD) compute the loss function many times, over many random subsets of the training dataset. Furthermore, the changes in the model quality metric that are small enough to be effectively approximated with Taylor expansion.

Approach

Redefine the utility function to measure the cumulative impact of a training data point on the validation loss across gradient update steps.

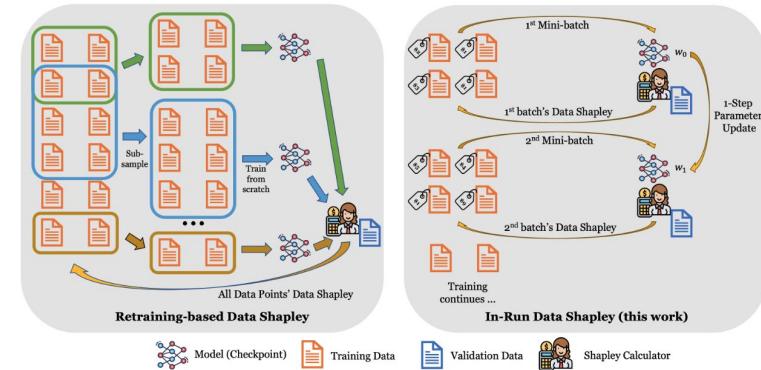
Redefined “local utility function” of subset S of a single SGD minibatch:

$$U^{(t)}(S; z^{(\text{val})}) := \underbrace{\ell(\tilde{w}_{t+1}(S), z^{(\text{val})}) - \ell(w_t, z^{(\text{val})})}_{\text{Model updated only using data from } S} - \underbrace{\ell(w_t, z^{(\text{val})})}_{\text{Model at SGD step } t}$$

$$\tilde{w}_{t+1}(S) := w_t - \eta_t \sum_{z \in S} \nabla \ell(w_t, z)$$

Redefined “global utility function” of subset S over the entire SGD run:

$$U(S) = \sum_{t=0}^{T-1} U^{(t)}(S)$$



Published as a conference paper at ICLR 2025.

DATA SHAPLEY IN ONE TRAINING RUN

Jiachen T. Wang
Princeton University
Prateek Mittal
Princeton University
David Song
UC Berkeley
Rong Jin
Virginia Tech

ABSTRACT

Data Shapley offers a principled framework for attributing the contribution of individual training data points to the final model performance. However, the computation of Shapley requires re-training the entire model from scratch, which becomes computationally infeasible for large-scale models. Additionally, this retraining-based approach is not suitable for large-scale training datasets with many gradient steps, which may often be of interest in practice. This paper introduces a novel approach for calculating Data Shapley in one training run. Our method is specifically designed for assessing data contribution for a particular model during its training process. It does not require re-training the entire model after each update iteration and accumulates these values throughout the training process. We demonstrate that our approach is able to calculate Data Shapley values in linear time to the size of foundation models. In its most optimized implementation, our method adds only a negligible overhead compared to standard training. This allows us to efficiently calculate overall contribution of training data points after the foundation model pre-training stage. We present several case studies that illustrate the effectiveness of our approach and their implications for copyright in generative AI and protecting data privacy.

[Wang ICLR '25]

Wang, Jiachen T., et al. "Data Shapley in One Training Run." The Thirteenth International Conference on Learning Representations. [\[Paper\]](#) [\[Blog\]](#)

- 1) Introducing the Concept of Data Importance
- 2) Examples of Data Attribution Functions
- 3) Case Study of Shapley Value as a Measure of Importance
- 4) Applications of Data Importance**

Influence Function for Explaining Fairness Errors

Challenge

Data attribution gives us an ordered list of data points that impact model quality, but it does not explain what makes these data points impactful.

Insight

If we group important data points based on common predicates, we can derive more powerful conclusions about factors that cause models to underperform.

Approach

First, use influence functions to compute data importance with respect to fairness metrics. Second, use lattice-based search to identify combinations of predicates that define data subsets that are both small and impactful.

Session 4 Responsible Data Management and Fairness

SIGMOD '22, June 12–17, 2022, Philadelphia, PA, USA

Interpretable Data-Based Explanations for Fairness Debugging

Romila Pradhan¹
Purdue University
USA
rpradhan@purdue.edu

Jiongli Zhu²
University of California,
San Diego
USA
jzhu14@ucsd.edu

Boris Glavic³
Illustrative Institute of
Technology
USA
bglavic@iit.edu

Fahad Salimi⁴
University of California,
San Diego
USA
fahadsalimi@cs.sdsu.edu

ABSTRACT

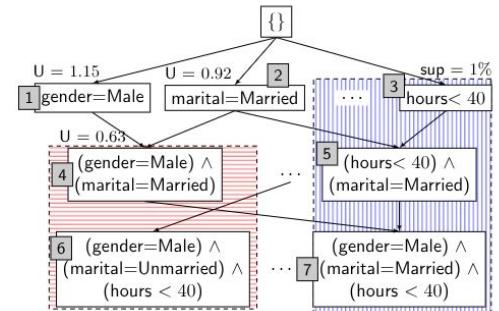
A wide variety of fairness metrics have been proposed to determine if the behavior of machine learning models are used in critical real-world applications in a fair manner. While there has been significant work on generating explanations using existing XAI techniques in mathematically formalized ways, there has been little work on generating compact, interpretable, and causal explanations for fairness metrics. In this paper, we introduce Gourami, a system that produces compact, interpretable, and causal explanations for fairness metrics. Gourami identifies the most impactful subset of the training data that are used for the fairness metric. Specifically, it introduces the concept of causal responsibility which measures the causal effect of a feature on the fairness metric. This concept has been used in causal inference to identify the most impactful features that are responsible for the behavior of ML algorithms – an approach known as causal blame assignment (CBA). We also introduce a lattice-based search for identifying the most impactful predicates that define the most impactful subsets of the training data that are used for the fairness metric. (1) Gourami identifies the most impactful features that are used for the fairness metric. (2) Gourami identifies the most impactful subsets of the training data that are used for the fairness metric. (3) Gourami identifies the most impactful predicates that define the most impactful subsets of the training data that are used for the fairness metric.

Most XAI research to date has focused on generating **feature-based explanations**, which quantify the impact of individual features on the output of a model. In contrast, Gourami generates **causal-based explanations**, which quantify the impact of causal relationships between features on the output of a model. Gourami uses a combination of causal-based methods [26, 35, 36] and logic-based approaches [24, 40, 41] to generate causal-based explanations. Gourami is the first system for generating interpretable explanations for fairness metrics.

Data points ordered by importance

age	education	marital	...	gender	income
39	Bachelors	Never-married	...	Male	$\leq 50K$
53	11th	Never-married	...	Male	$\leq 50K$
28	Bachelors	Married-civ-spouse	...	Female	$\leq 50K$
37	Masters	Married-civ-spouse	...	Female	$\leq 50K$

Lattice-based search identifies predicates that select the most impactful training data subsets



Combinations of predicates that explain model behavior

1	(Gender = Female) \wedge (Relationship = Not married) \wedge (Education = Associate-voc)
2	(Gender = Male) \wedge (Relationship = Spouse) \wedge (Hours < 40) \wedge (Workclass = Federal-gov)
3	(Gender = Male) \wedge (Education = Prof-school)

[Zhu SIGMOD '22]

Pradhan, Romila, et al. "Interpretable data-based explanations for fairness debugging." Proceedings of the 2022 international conference on management of data. 2022. [\[Paper\]](#)

Debugging the LLM Retrieval Corpus

Challenge

Retrieval augmented generation (RAG) is a widely used technique for providing pre-trained large language models (LLMs) with task-specific context. Data errors in the retrieval corpus have a negative impact on model quality.

Insight

The role of a retrieval corpus to an LLM is similar to the role of a training dataset to a classical ML model.

Approach

Define a data attribution function that will compute the importance of data points in the retrieval corpus. Use this to identify and debug data errors.

Improving Retrieval-Augmented Large Language Models via Data Importance Learning

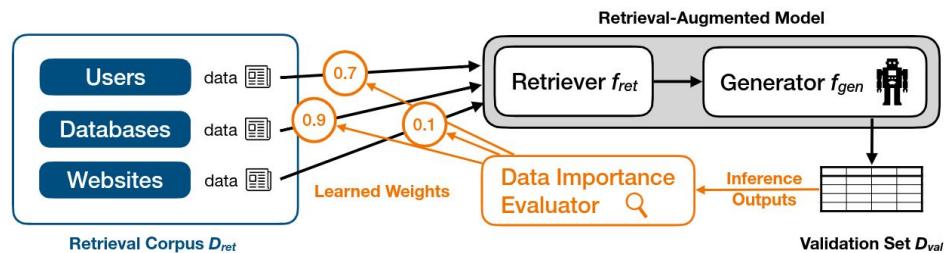
Xiaozhong Lyu¹ Stefan Graepel² Samantha Biagioli³ Shaoqiang Wei⁴
¹Erlangen University, ²ETH Zurich, ³University of Amsterdam, ⁴Apple

Abstract

Retrieval augmentation enables large language models to take advantage of external knowledge bases to improve their performance on downstream tasks. However, the performance of such retrieval-augmented models is limited by the data quality in the retrieval corpus. In this paper, we propose a data importance learning based on multi-linear extension for evaluating the data importance of retrieved documents. We also propose a data pruning algorithm that removes low-importance data from the retrieval corpus. The proposed data importance learning and one key contribution of this paper is a polynomial-time algorithm that computes the data importance of data points in the retrieval corpus using the multi-linear extension. We evaluate our proposed data importance learning on a validation set, the data importance of data points in the retrieval corpus using the proposed multi-linear extension is more accurate than that using the state-of-the-art more efficient (ϵ - δ) approximation algorithm. Our experimental results illustrate that the proposed data importance learning can significantly improve the performance of retrieval-augmented large language models.

[Lyu arXiv '23]

Lyu, Xiaozhong, et al. "Improving retrieval-augmented large language models via data importance learning." arXiv preprint arXiv:2307.03027 (2023). [\[Paper\]](#) [\[Code\]](#)



$$U(f_{gen}, f_{ret}, \mathcal{D}_{val}, \mathcal{D}_{ret}) := \sum_{x_i \subseteq \mathcal{D}_{val}} U(f_{gen}(x_i, f_{ret}(x_i, \mathcal{D}_{ret})))$$

$$\tilde{U}(w_1, \dots, w_M) := \sum_{S \subseteq \mathcal{D}_{ret}} U(S) \underbrace{\prod_{d_i \in S} w_i \prod_{d_i \notin S} (1 - w_i)}_{P[S]}$$

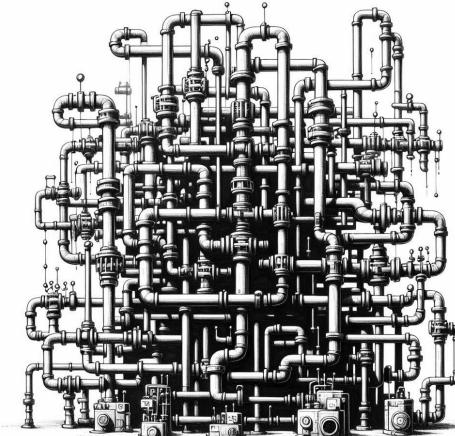
DATASET	GPT-JT (6B)	GPT-JT (6B) W/ RETRIEVAL				GPT-3.5 (175B)
		VANILLA	+LOO	+REWEIGHT	+PRUNE	
BUY	0.102	0.789	0.808	0.815	0.813	0.764
RESTAURANT	0.030	0.746	0.756	<u>0.760</u>	0.761	0.463

Key Takeaways of Part I

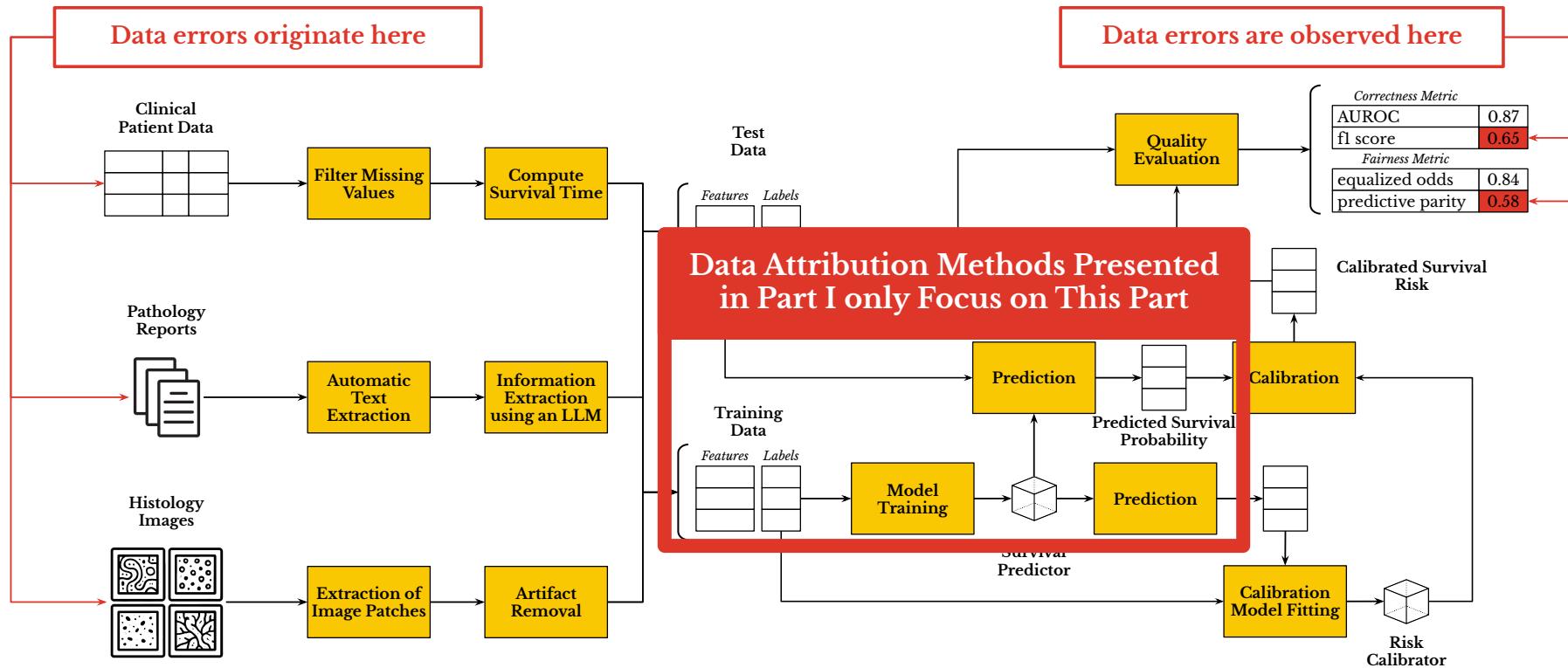
- Data attribution is a useful powerful framework for approaching the problem of data error detection.
- There are many existing data attribution methods with various strengths and shortcomings.
- The most powerful methods face scalability issues that have been tackled by existing research with many opportunities for future improvements.

Part II: Data Debugging in ML Pipelines

Sebastian Schelter



Gap between Attribution Methods and ML Pipelines



1) Gap between Attribution Methods and ML Pipelines

2) Libraries and Systems for ML Pipelines

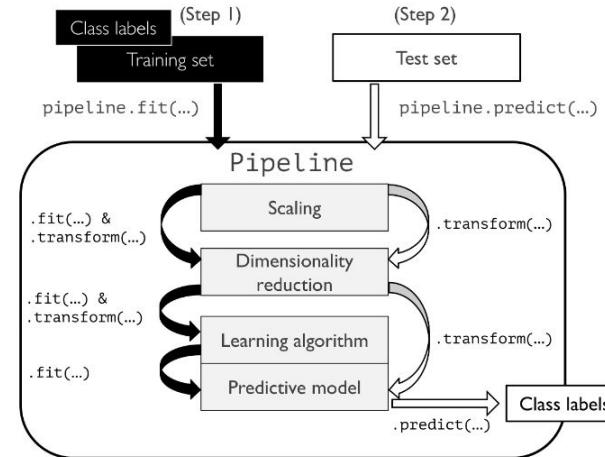
3) Characteristics of Real World ML Pipelines

4) Methods for Debugging ML Pipelines

Scikit-Learn

Highlights

- Among the most popular data science Python libraries
- Has implementations of many machine learning models, as well as feature encoding operators
- Introduced the **estimator/transformer abstraction** for composing complex, nested pipelines
 - **Transformer:** tuple-at-a time transformation
 - **Estimator:** create a data-specific transformer via a global aggregation over the data



Source: <https://vitalflux.com/scikit-machine-learning-pipeline-python-example/>



[Pedregosa JMLR '11]

Pedregosa, Fabian, et al. "Scikit-learn: Machine learning in Python." the Journal of machine Learning research 12 (2011): 2825-2830. [Paper] [Website] [Code]

Tensorflow Extended (TFX)



Highlights

- *End-to-end platform for production ML pipelines*
- *Built on TensorFlow and optimized for scalability, strong emphasis on model validation and monitoring*
- *Includes reusable components for pipelines, inspired by estimator/transformer paradigm*
- *Apache Beam for dataflow operations, Tensorflow for numerical operations*

KDD 2017 Applied Data Science Paper

KDD '17, August 13–17, 2017, Halifax, NS, Canada

TFX: A TensorFlow-Based Production-Scale Machine Learning Platform

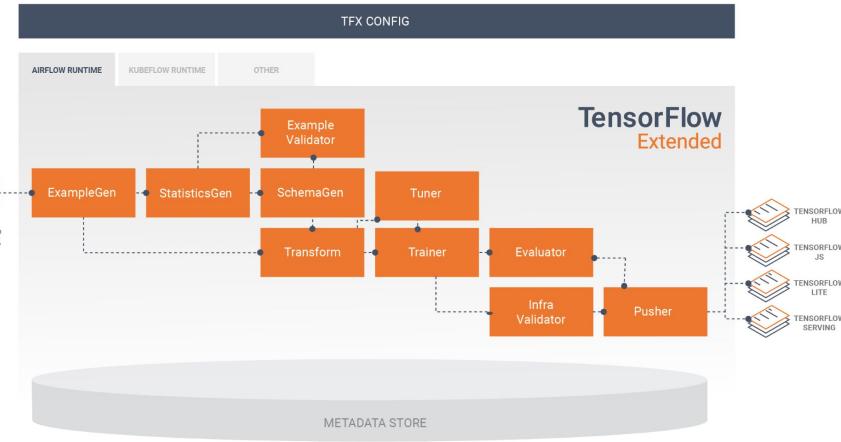
Denis Baylor, Eric Brooks, Hong-Tai Cheng, Noah Frey, Clinton Yu Guo, Zesheng Huang,
 Salma Hritzel, Manisha Iyer, Vipin Jain, Ming Tang, Ravi Kothiyal, Chin-Yan Kuo, Michael Law,
 Chirag Mehta, Akshay Narendran, Nadeem Polycott, Sujith Ramadas, Sudip Roy,
 Steven Enjaoui Whang, Martin Wicke, Jacek Wilkiewicz, Xu Zhang, Martin Zinkevich
 Google Inc.

Abstract machine learning as a tool to gain knowledge from data across a broad spectrum of use cases and products, ranging from recommendation engines to real-time bidding and prediction for advertising [11, 12] and even the protection of endangered species [13].

The challenge of applying machine learning to a specific use case is simple: at the training phase, a single feature vector is used to represent the data. In the inference phase, the model takes features as input and makes predictions. However, the challenge becomes more complex when machine learning needs to be deployed in a production environment. This is because it is often required that, together with the learned model, complete infrastructure is provided to support the system. This can be a challenge when dealing with a diverse range of features that can be represented in different ways, such as categorical variables, sparse vectors, and tensors. In addition, it is often necessary to deal with a wide range of failures that can happen during the execution of the pipeline, such as missing or corrupted data, network issues, and so on. To address these challenges, we have developed a platform called TensorFlow Extended (TFX), a TensorFlow-based general-purpose machine learning platform implemented in Go. In this paper, we show how we used TFX to standardize the components of our machine learning pipeline and how we were able to build a highly reliable and scalable system.

[Baylor SIGKDD '17]

Baylor, Denis, et al. "Tfx: A tensorflow-based production-scale machine learning platform." Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining. 2017. [\[Paper\]](#) [\[Website\]](#) [\[Code\]](#)



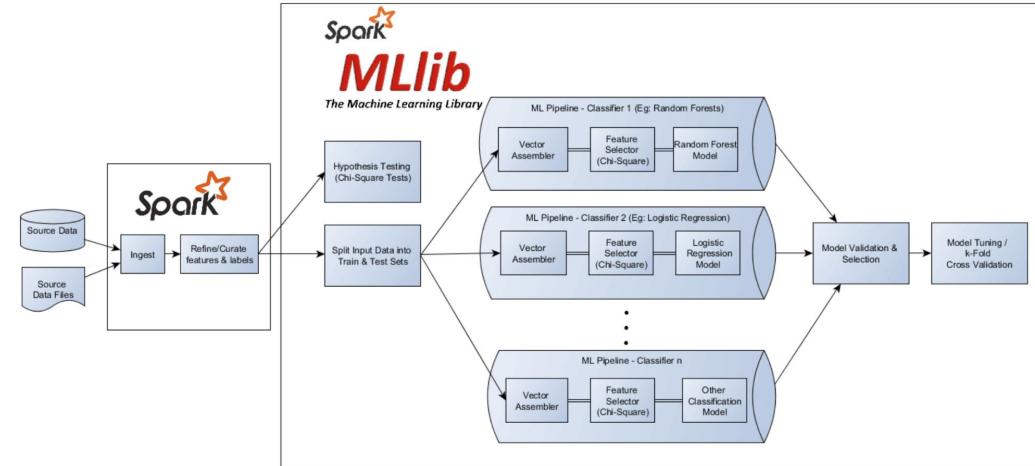
Source: <https://www.tensorflow.org/tfx/guide>

Spark MLLib



Highlights

- Built on top of Apache Spark
- Includes implementations for classification, regression, clustering, collaborative filtering, and dimensionality reduction
- Works natively with Spark DataFrames, SQL, and streaming data
- Adoption of estimator/transformer paradigm from scikit-learn



Source: <https://www.qubole.com/developers/spark-getting-started-guide/workflow>

Journal of Machine Learning Research 17 (2016) 1-7
Submitted 5/15, Published 1/18

MLlib: Machine Learning in Apache Spark

Xiangru Meng¹
Databricks, 160 Spear Street, 10th Floor, San Francisco, CA 94105
meng@dataricks.com

Joseph Bradley²
Databricks, 160 Spear Street, 10th Floor, San Francisco, CA 94105
joseph@dataricks.com

Burak Yavuz³
Databricks, 160 Spear Street, 10th Floor, San Francisco, CA 94105
burak@dataricks.com

Evan Bonney⁴
UC Berkeley, 46 Soda Hall, Berkeley, CA 94720
evan@berkeley.edu

Shivaram Venkataraman⁵
Berkeley, 46 Soda Hall, Berkeley, CA 94720
shiva@cs.berkeley.edu

Darren Liu⁶
Databricks, 160 Spear Street, 10th Floor, San Francisco, CA 94105
dliu@dataricks.com

Jeremy Freudenthal⁷
RISELab Research Center, 18005 Helia Dr, Ashburn, VA 20147
jeremy@rise-lab.org

DB Tsai⁸
Netflix, 999 Market Street, Los Gatos, CA 95032
dbs@netflix.com

Manish Srivastava⁹
OpenSense Labs, 1531 Oracle Street, Menlo Park, CA 94025
manish@opensense.com

[Meng JMLR '16]

Meng, Xiangru, et al. "Mllib: Machine learning in apache spark." Journal of Machine Learning Research 17.34 (2016): 1-7. [Paper] [Website] [Code]

Apache SystemDS



Highlights

- Designed for scalable and efficient execution on both single-node and distributed environments
- Offers a high-level scripting language for expressing ML algorithms and workflows with a declarative R-like language
- Performs cost-based optimization and automatic operator selection for efficient execution across different hardware endpoints
- Optimised feature encoders based on estimator/transformer paradigm



Journal of Machine Learning Research 21 (2020) 1-7

Submitted 5/10; Published 4/10

MLlib: Machine Learning in Apache Spark

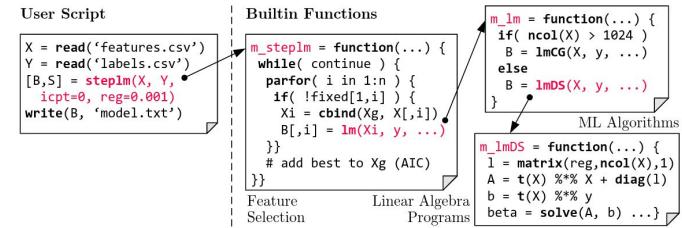
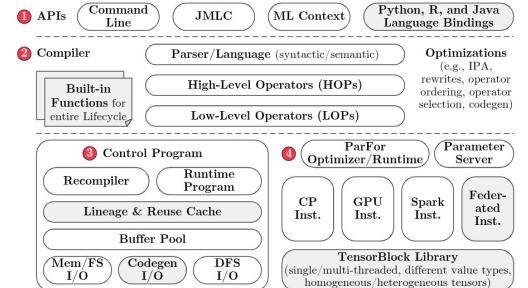
Xiangru Meng¹, JIAO YUAN², HENGJIANG CHEN³¹ 160 Spear Street, 10th Floor, San Francisco, CA 94105² 160 Spear Street, 10th Floor, San Francisco, CA 94105³ 160 Spear Street, 10th Floor, San Francisco, CA 94105Evan Sparks⁴, BRIAN K. HAN⁵, DAVID M. BLAKE⁶⁴ UC Berkeley, 160 Soda Hall, Berkeley, CA 94720⁵ SPARKS@EECS.BERKELEY.EDU⁶ SPARKS@EECS.BERKELEY.EDU

[Boehm CIDR '20]

Boehm, Matthias, et al. "SystemDS: A Declarative Machine Learning System for the End-to-End Data Science Lifecycle." 10th Conference on Innovative Data Systems Research. 2020. [\[Paper\]](#) [\[Website\]](#) [\[Code\]](#)

[Phani VLDB '20]

Phani, Arnab, et al. "UPLIFT: parallelization strategies for feature transformations in machine learning workloads." Proceedings of the VLDB Endowment, Volume 15, Issue 11, 2020. [\[Paper\]](#)



ML Pipelines in the Cloud



Netflix Metaflow

[\[Website\]](#) [\[Documentation\]](#)

Highlights

- Notebook based development environment
- Storing and tracking of code, data and models
- Scaling from local execution to the cloud



Amazon SageMaker

Amazon SageMaker Pipelines

[\[Website\]](#) [\[Documentation\]](#)

Highlights

- Define, automate, and manage end-to-end ML workflows
- Automatically tracks pipeline artifacts
- Leverages AWS Cloud infrastructure



Azure Machine Learning

Azure Machine Learning Pipelines

[\[Website\]](#) [\[Documentation\]](#)

Highlights

- Orchestration of ML workflows with reusable, modular pipeline components
- Versioning, monitoring, and CI/CD integration



Vertex.ai

Vertex AI Pipelines

[\[Website\]](#) [\[Documentation\]](#)

Highlights

- Connects with Vertex AI services
- Tracks pipeline steps, metadata, and artifacts
- Orchestrates ML workflows on Google Cloud

- 1) Gap between Attribution Methods and ML Pipelines
- 2) Libraries and Systems for ML Pipelines
- 3) Characteristics of Real World ML Pipelines**
- 4) Methods for Debugging ML Pipelines

Study of Pipelines at Google

Highlights

- Study of 3000 production pipelines with over 450K models trained over a 4 month period
 - About half the pipelines studied used data- and model-validation operators
 - Input data typically has up to 100 features, but can have over 10K in extreme cases
 - 53% of features were categorical, often with very large domains (averaging over 10M unique values)
 - Training accounts for only 20% of the total runtime cost, over 30% is for model validation and 20% for data ingestion
 - About 1/4 of model training runs results in model deployment
 - Deep learning models account for 60% of pipelines

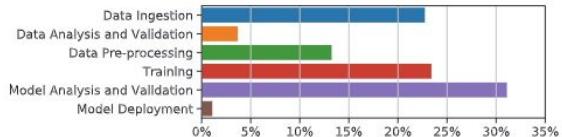
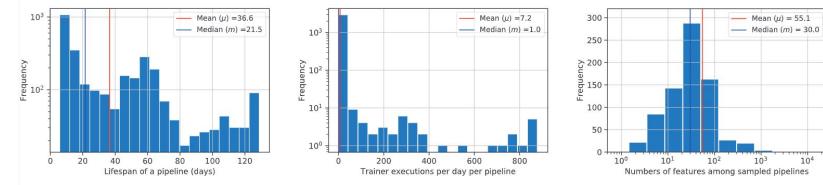


Figure 7: Compute cost of different operators



(a) Distribution of pipeline span.

(b) Distribution of trained models per day

(c) Distribution of the number of features

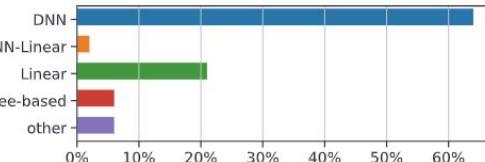
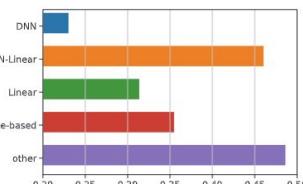


Figure 5: Percentage of Trainer runs with each model type



(f) Model type vs. likelihood of pushes

[Xin SIGMOD '21]

Xin, Doris, et al. "Production machine learning pipelines: Empirical analysis and optimization opportunities." Proceedings of the 2021 international conference on management of data, 2021. [Paper]

Study of Pipelines at Microsoft

Highlights

- Study of over 8M public Jupyter notebooks on GitHub (from 2017, 2019, and 2020), and 2M enterprise pipelines developed with ML.NET*
- Python is emerging as the de-facto standard language for data science (81% of notebooks in 2017 and 91% in 2020)*
- Around 80% cells were linear (no conditional statements) and 76% were completely linear (no conditionals, classes, or functions)*
- Libraries like numpy, matplotlib, pandas, and scikit-learn are used very frequently (e.g., numpy in >60% of notebooks)*
- Few highly used libraries have significant coverage (e.g., top-10 cover ~40% of notebooks, top-100 cover ~75%), but there is a long tail*
- Explicit ML pipelines (defined with sklearn.pipeline) are gaining traction but there are still 5 times more implicit pipelines in GitHub notebooks*
- There is a large number of distinct operators, and a significant portion are user-defined (especially in ML.NET and implicit GitHub pipelines)*

Data Science Through the Looking Glass: Analysis of Millions of GitHub Notebooks and ML.NET Pipelines

Fotis Psallidas, Yizhen Zhu, Bojun Kang*, Jordan Redell, Mattia Ippolito, Suresh Subrahmanyam, Ming Tang, Fengyu Wu, Ce Zhang*, Matias Werner, Aritra Bhattacharya, Kartikasari Kartasasmita, and Michael J. Franklin
Microsoft Research, Seattle, USA

ABSTRACT

The recent success of machine learning (ML) has led to an explosive growth of systems and applications built by researchers, data scientists, and practitioners (data science (DS) practitioners). This quickly shifting paradigm, however, is challenging for system builders and practitioners to keep up. In this paper, we present a novel approach to capturing this paradigm through a wide-angle lens, performing an in-depth analysis of millions of GitHub notebooks and ML.NET pipelines. We analyze the evolution of GitHub notebooks, publicly shared on GitHub—(a) ML notebooks, publicly shared on GitHub—(b) ML.NET pipelines, privately authored in ML.NET within Microsoft, DS encompasses a wide range of activities, including (i) collecting data, (ii) feature engineering, (iii) training models, (iv) model evaluation, (v) ML collaboration, and collecting datasets representative of the field. The results of our analysis show that while the tools we use here are a first step towards the end, and skewed distributions of usage, we can see some interesting findings. Over the past few years, we have used the results

[Psallidas SIGMOD Record '22]

Psallidas, Fotis, et al. "Data science through the looking glass: Analysis of millions of github notebooks and ml. net pipelines." ACM SIGMOD Record 51.2 (2022): 30-37. [Paper]

Dimension	Metric	GH17	GH19	GH20
Notebooks	Total	1.23M	4.6M	8.7M
	Deduped	66.0%	65.5%	65.7%
	Linear	26.4%	29.1%	30.3%
	Completely Linear	21.2%	23.3%	24.6%
Languages	Python	81.7%	91.7%	91.1%
	Other	18.3%	8.3%	8.9%
Cells	Total	34.6M	143.1M	261.2M
Code Cells	Total	64.5%	66.4%	66.9%
	Deduped	41.0%	38.6%	38.5%
	Linear	72.1%	80.2%	79.3%
	Completely Linear	68.3%	76.1%	75.6%
Users	Total	100K	400K	697K

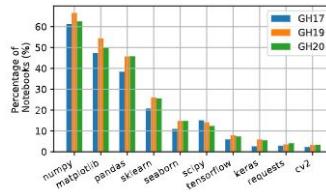


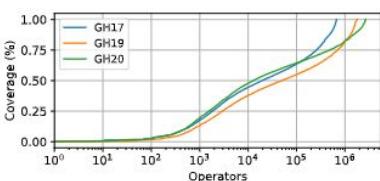
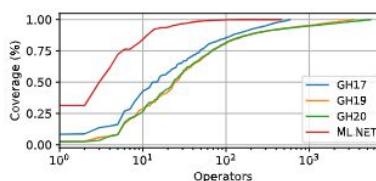
Figure 2: Top-10 used libraries.



Figure 3: DL libraries usage percentages.

	GH17	GH19	GH20	ML.NET
#Pipelines	Implicit	164K	415K	1.4M
	Explicit	10K	129K	252K

	GH17	GH19	GH20	ML.NET
#Distinct Ops	Implicit	668K	1.8M	2.6M
	Explicit	584	3.4K	5.5K

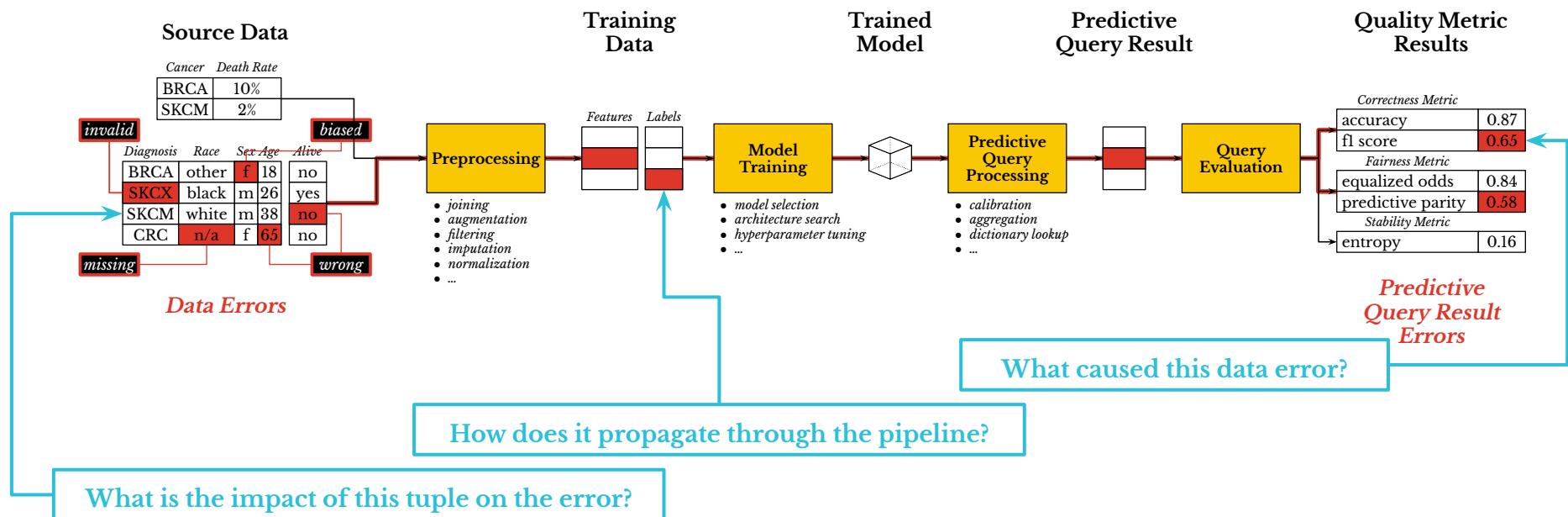


Observations

- No universal way to express ML pipelines, design often prioritises flexibility and ease-of-use
- Many pipelines combine relational / dataflow operators with ML-specific operators based on estimator/transformer abstraction
- Pipelines often executed via multiple runtimes
- Lack of algebraic operator semantics
- Lack of fine-grained data provenance

- 1) Gap between Attribution Methods and ML Pipelines
- 2) Libraries and Systems for ML Pipelines
- 3) Characteristics of Real World ML Pipelines
- 4) Methods for Debugging ML Pipelines**

How should we reason about pipelines?



Modeling ML Pipelines with “Logical Query Plans”

Challenge

Understanding of the semantics of operations and the flow of data required to reason about ML pipelines

Insight

Many common pipeline abstractions offer declarative operations, enables the extraction and definition of “logical query plans” modeling their operations

Approach

Instrument functions of Python data science libraries to extract query plan, enable annotation propagation through operators. Apply rule-based approaches to determine if an error has occurred.

Potential issues in preprocessing pipeline:

- ① Join might change proportions of groups in data
- ② Column ‘age_group’ projected out, but required for fairness
- ③ Selection might change proportions of groups in data
- ④ Imputation might change proportions of groups in data
- ⑤ ‘race’ as a feature might be illegal!
- ⑥ Embedding vectors may not be available for rare names!

Python script for preprocessing, written exclusively with native pandas and sklearn constructs

```
# load input data sources, join to single table
patients = pandas.read_csv(_)
histories = pandas.read_csv(_)
data = pandas.merge([patients, histories], on=['ssn'])

# compute mean complications per age group, append as column
complications = data.groupby('age_group')
    .agg(mean_complications=('complications', 'mean'))
data = data.merge(complications, on='age_group')

# Target variable: people with frequent complications
data['label'] = data['complications'] >
    1.2 * data['mean_complications']

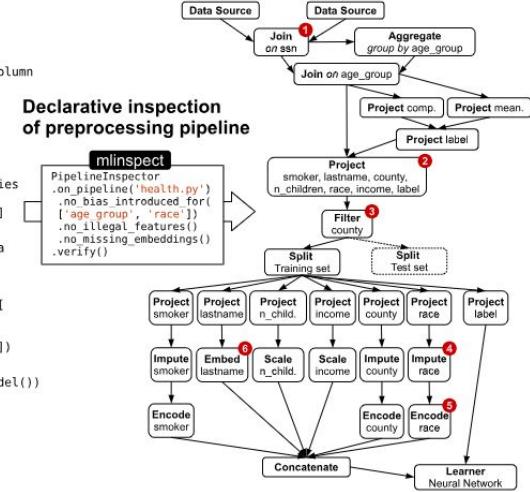
# Project data to subset of attributes, filter by counties
data = data[['smoker', 'last_name', 'county',
    'num_children', 'race', 'income', 'label']]
data = data[data['county'].isin(counties_of_interest)]

# Define a nested feature encoding pipeline for the data
impute_and_encode = sklearn.Pipeline([
    (sklearn.SimpleImputer(strategy='most_frequent')),
    (sklearn.OneHotEncoder())])
featureisation = sklearn.ColumnTransformer(transformers=[
    ('impute_and_encode', ['smoker', 'county', 'race']),
    ('Word2VecTransformer()', 'last_name'),
    (sklearn.StandardScaler(), ['num_children', 'income'])])

# Define the training pipeline for the model
neural_net = sklearn.KerasClassifier(build_fn=create_model())
pipeline = sklearn.Pipeline([
    ('features', featureisation),
    ('learning_algorithm', neural_net)])

# Train-test split, model training and evaluation
train_data, test_data = train_test_split(data)
model_pipeline.fit(train_data, train_data.label)
print(model.score(test_data, test_data.label))
```

Corresponding dataflow DAG for instrumentation, extracted by mlinspect



Data distribution debugging in machine learning pipelines

Stefan Graßberger¹ · Paul Groth¹ · Julia Stoyanovich² · Sebastian Schelter¹

Received: 27 January 2021 / Revised: 9 September 2021 / Accepted: 1 December 2021 / Published online: 31 January 2022
© The Author(s), under exclusive license to Springer Nature Switzerland AG part of Springer Nature 2022

Abstract
Machine learning (ML) is increasingly used to inform important decisions, and the risks arising from this widespread use are prompting attention from policy makers, scientists, and the media. ML applications are often brittle with respect to their input data, which leads to concerns about their correctness, reliability, and fairness. In this paper, we describe mlinspect, a library for inspecting ML pipelines. mlinspect provides a declarative abstraction of ML pipelines and a way to reason about them to address problems collectively as *data distribution bugs*. The key idea is to extract a directed acyclic graph representation of the dataflow from a preprocessing pipeline and to use this representation to automatically instrument the code with predicted annotations. These annotations can then be used to reason about the dataflow and detect errors such as missing or incorrect lineage information from operators to operator. In contrast to existing work, mlinspect operates on declarative abstractions of preprocessing pipelines, which makes it amenable to reasoning and does not require manual code instrumentation. We discuss the design and implementation of the mlinspect library and give a comprehensive end-to-end example that illustrates its functionality.

[Graßberger VLDBJ '22]

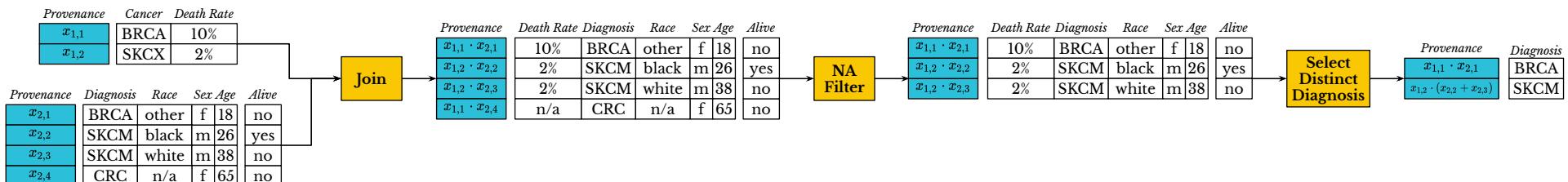
Graßberger, Stefan, et al. "Data distribution debugging in machine learning pipelines." The VLDB Journal 31.5 (2022): 1103-1126. [Paper] [Code]

Leveraging the Provenance Semiring Framework

Highlights

- *Theoretical framework analyzing the relationship between input and output tuples of relational queries*
 - *Allows us to determine the presence of an output tuple as a function of the presence of the input tuples*
 - *Easy to adapt for ML pipelines once logical query plan with “relational-like” operations is known*

Application to an Example Pipeline



Provenance Semantics

Todd J. Green Grigoris Karvounarakis Val Tannen
tjgreen@cis.upenn.edu gkarvoun@cis.upenn.edu val@cis.upenn.edu
Department of Computer and Information Science
University of Pennsylvania
Philadelphia, PA 19104-6389

ABSTRACT
We show that relational algebra calculations for incomplete databases, probabilistic databases, log semantics and why-procedures are particular cases of the same general algorithmic framework. This further suggests a comprehensive probability representation that uses settings of probabilities as parameters. We extend these constructions to datalog and extensions of logic programs. We also show how to implement probability calculations as well as database evaluation for incomplete and probabilistic databases. Finally, we analyze

We observe that in all four cases, the calculations with annotations are strikingly similar. This suggests looking for an algebraic structure on annotations that captures the above particular cases. We propose using commutative semirings for this purpose. In fact, we can show that the laws of commutative semirings are forced by certain expected identities in \mathbb{R}^A . Having identified commutative semirings as the right algebraic structure, we argue that a symbolic representation of executing calculations is just what is needed to record, document, and track \mathcal{R} -queries from input to output.

[Creep SICMOD '07]

Green, Todd J., Grigorios Karvounarakis, and Val Tannen. "Provenance semirings." Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems, 2007 [Paper].

Debugging Preprocessing Pipelines with Datascope

[Attribution Function: Shapley Value]

Challenge

KNN Proxy methods not directly applicable to arbitrary pipelines. Presence of a single source data point does not map directly to a single data point fed to the model.

Insight

Observe three canonical types of pipelines based on shape of produced provenance polynomials. Possible to develop efficient PTIME algorithms for computing the Shapley value for them.

Approach

Compile provenance polynomials to Additive Decision Diagrams and use them to compute Shapley values in PTIME.

Published as a conference paper at ICLR 2024

DATA DEBUGGING WITH SHAPLEY IMPORTANCE OVER MACHINE LEARNING PIPELINES

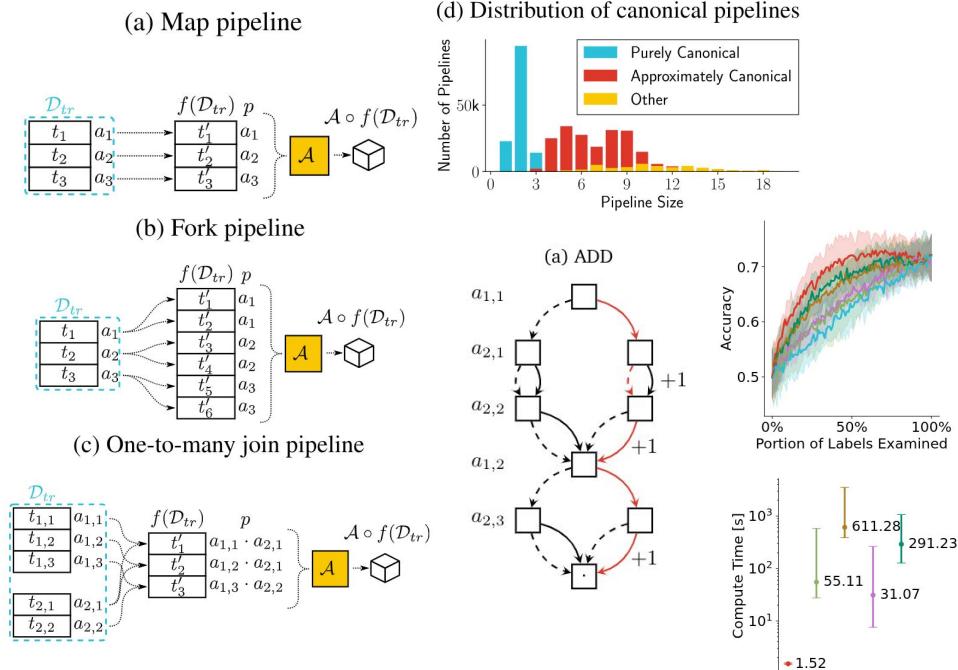
Bojan Karlaš¹, David Basav, Matteo Interlenghi², Sebastian Schelter¹, Wentao Wu¹, Ce Zhang³
¹Harvard University, ²ETH Zurich, ³Microsoft, ⁴University of Amsterdam, ⁵University of Chicago
*Data debugging, harvard.edu

ABSTRACT

A machine learning (ML) model exhibits poor quality (e.g., poor accuracy or fairness), the problem can often be traced back to errors in the training data. Being able to discover the data errors is thus critical for ML system maintenance. Data debugging has received a lot of attention recently. One promising way to measure “data importance” with respect to model quality is to use Shapley values. However, Shapley values are NP-hard to compute, ML model in isolation, without considering the broader ML pipeline for data preparation and feature extraction, which appears to be the majority of real-world ML code. This presents a major challenge for data debugging. In this paper, we propose Datascope, a method for efficiently computing Shapley-based data importance over ML pipelines. Datascope uses a novel approach to decompose the Shapley computation into terms of computational speed. Finally, our experimental evaluation demonstrates that our methods are significantly faster than state-of-the-art methods for data debugging, and in some cases even outperform them. We release our code as an open-source data debugging library available at github.com/czhanglab/Datascope.

[Karlaš ICLR '24]

Karlaš, Bojan, et al. "Data Debugging with Shapley Importance over Machine Learning Pipelines." The Twelfth International Conference on Learning Representations. 2024. [\[Paper\]](#) [\[Website\]](#) [\[Code\]](#)



Debugging Predictive Queries with Rain

[Attribution Function: Influence]

Challenge

Model inference often part of a larger predictive query.

Influence-based attribution methods must account for structure of query.

Insight

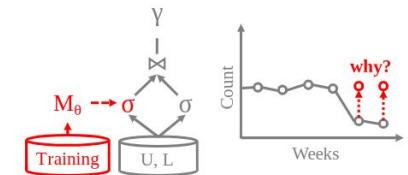
Provenance polynomials for tracking lineage starting from training tuples all the way to predictive query outputs allows us to make the entire expression differentiable.

Approach

User complaints on query outputs (e.g. what-if-queries) used to identify errors. Make the entire query differentiable using provenance polynomials and run the influence framework to identify errors in the training dataset.

Q

```
SELECT COUNT(*)
  FROM Users U JOIN Logins L
    ON U.ID = L.ID
   WHERE L.active_last_month AND
        Mθ.predict(U.*) = "Churn"
```



Complaint-driven Training Data Debugging for Query 2.0

Weiyuan Wu
Simon Fraser University
Burnaby, BC, Canada
wuyw@sfu.ca

Lempros Fidas
Columbia University
New York, NY
lfidas@cs.columbia.edu

Eugene Wu
Columbia University
New York, NY
ewwu@columbia.edu

Jianxin Wang
Simon Fraser University
Burnaby, BC, Canada
jwang@sfu.ca

Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data, June 14–19, 2020, Portland, OR, USA
© 2020, ACM, New York, NY, USA, 18 pages. <https://doi.org/10.1145/3351051.3397000>

ABSTRACT

As the use of machine learning (ML) increases rapidly across all industry sectors, there is a significant interest among practitioners in how to debug ML pipelines. Query 2.0, which integrates model inference into SQL queries. Debugging Query 2.0 is very challenging since an unexpected query result may be caused by many factors (e.g., wrong labels, corrupted features). In response, we propose Rain, a complaint-driven training data debugging framework. Rain allows users to specify complaints on the query's output, and it can quickly identify the training data that are responsible for the complaints.

[Wu SIGMOD '20]

Wu, Weiyuan, et al. "Complaint-driven training data debugging for query 2.0." Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data. 2020. [\[Paper\]](#)

ArgusEyes - Continuous Integration for ML Pipelines

Challenge

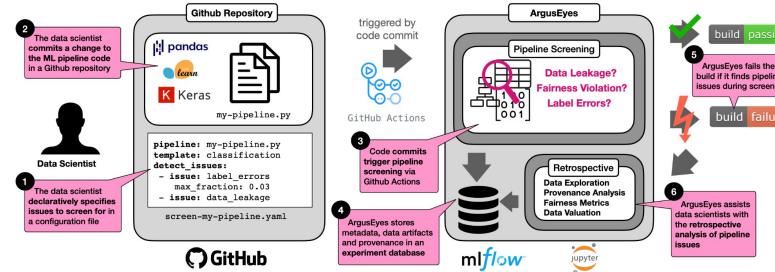
ML systems lack sophisticated testing infrastructure developed for classical software engineering. Many data-related problems only become apparent in production.

Insight

Logical query plans for ML pipelines combined with data debugging techniques enable ML-specific CI infrastructure

Approach

Instrument, execute and screen ML pipelines for declaratively specified pipeline issues, and analyze data artifacts and their provenance to catch potential problems early before deployment to production.



Proactively Screening Machine Learning Pipelines with ArcusEyes

Sebastian Schelter
University of Amsterdam
Amsterdam, The Netherlands
s.schelter@uva.nl

Stefan Graßberger
AIRLab, University of Amsterdam
Amsterdam, The Netherlands
s.grassberger@uva.nl

Shabha Guha
AIRLab, University of Amsterdam
Amsterdam, The Netherlands
s.guha@uva.nl

Bojan Karlaš
Harvard University
Boston, United States

Ce Zhang
ETH Zurich
Zürich, Switzerland

ACT
bejan_karlaaghm.harvard.edu
e_change@mit.edu

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.10.443603>; this version posted May 10, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [aCC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

Data-related issues in ML pipelines. ML pipelines in real-world applications often suffer from a variety of data-related issues. These include missing values, outliers, and noise in training data; missing or noisy test data and label errors [1] in training samples, which negatively impact the accuracy of a pipeline; as well as overfitting to training data [2]. In this paper, we propose a modular approach to mitigate these challenges. Specifically, we propose a modular pipeline, *Anisotropic*, which consists of three main stages: *Feature Extraction*, *Model Selection*, and *Feature Selection*. Anisotropic is designed to be highly modular, allowing users to easily add or remove components, and to be highly extensible, allowing users to easily add new components. We evaluate *Anisotropic* on a variety of datasets, including a variety of continuous integration systems, and show that it outperforms state-of-the-art ML pipelines for dedicatedly specified pipelines. We also show that *Anisotropic* is competitive with state-of-the-art ML pipelines that have been specifically designed for each dataset. We believe that *Anisotropic* will be a valuable tool for data scientists who are already familiar with machine learning. We also believe that *Anisotropic* will be a valuable tool for data scientists who are not yet familiar with machine learning.

[Schelter SIGMOD Demo '23]

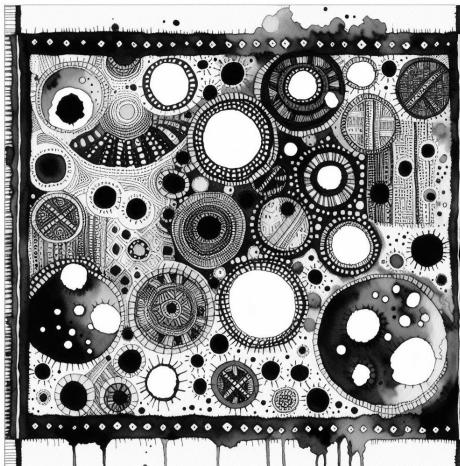
Schelter, et al.: "Proactively Screening Machine Learning Pipelines with ArgusEyes." Proceedings of the 2023 ACM SIGMOD International Conference on Management of Data (demo), 2023. [Paper]

Key Takeaways of Part II

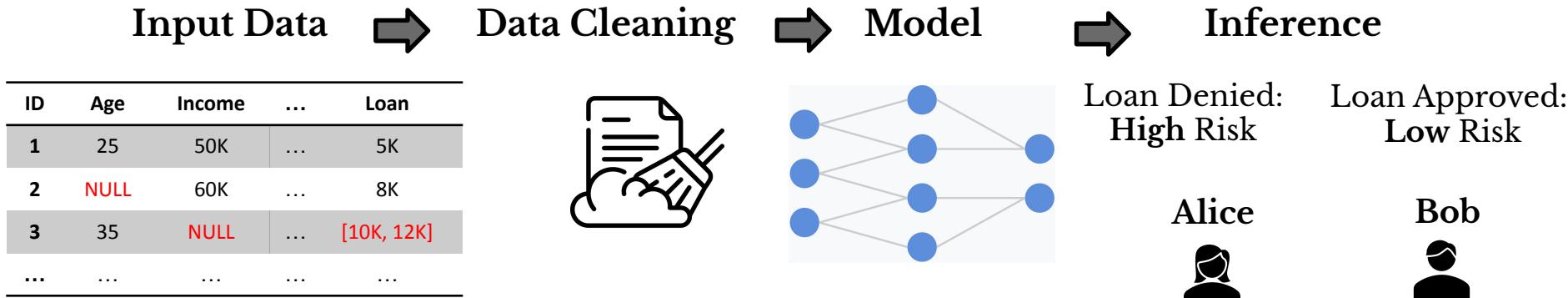
- Attribution methods presented in Part I assume models are trained with source data.
- ML pipelines are complex and present many opportunities for methods development.
- Logical query plans combined with data provenance offer a powerful framework for analyzing ML pipelines.

Part III: Learning from Uncertain and Incomplete Data

Babak Salimi



The Standard ML Pipeline



⚠ Common Assumption: once we “clean” the data, the pipeline consumes accurate and unbiased inputs.

✗ Reality: cleaning/pre-processing yields one reconstruction, driven by heuristic choices & domain assumptions → it can embed hidden bias and hide genuine uncertainty.

→ Key insight for Part III: even after best-effort cleaning, *real-world data remains incomplete and uncertain*. Our models—and the theory behind them—must make that uncertainty explicit rather than ignore it.

Why “Fixing” Data Errors Is Impossible in Principle

Missing values (  / )

Irrecoverable uncertainty: any imputation is just a guess; the true value is unobservable.

Unverifiable assumption: “missing at random,” parametric model of the data, etc.

[Pearl & Mohan, AAAI 2014], [Mohan, Pearl & Tian, NeurIPS 2013]

Measurement / annotation bias ( sentiment,  diagnoses)

Systematic distortion: recorded values can be consistently wrong.

Unverifiable assumption: symmetric, independent label-noise model.

[Pearl, UAI 2010], [Zhang & Yu, IJCAI 2015]

Why “Fixing” Data Errors Is Impossible in Principle

Selection bias & missing counterfactuals (⚠ rejected-loan applicants, excluded patients)

Unknown outcomes: whole sub-populations are never seen.

Finite-sample limits: re-weighting needs the true selection mechanism—which we can’t test.

[Bareinboim, Tian & Pearl, AAAI 2014] [Cortes et al., ALT2008],
[Heckman, Econometrica 1979]

Schema / integration mismatch (⚠ inconsistent units, ✖ fuzzy entity resolution)

Ambiguous merges: no ground-truth correspondences.

Pre-processing bias: heuristics distort original distributions; matching is probabilistic.

[Dong, Halevy & Madhavan, VLDB 2009],
[Getoor & Machanavajjhala, ACM 2012]

Challenges with Traditional Data Pipelines

Input Data



Data Cleaning



Model



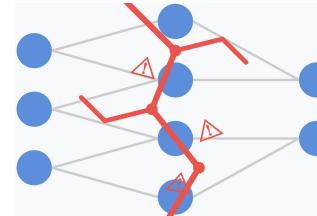
Inference

Loan Denied:
High Risk

Alice

Bob

Loan Approved:
Low Risk



Generalization Failure – Models trained on “repaired” data collapse under real-world shifts.

✗ High-Stakes Mis-decisions – Hidden bias drives flawed credit, medical, and justice outcomes.



⚠ Broken Uncertainty – Bayesian & conformal intervals lose calibration when data are incomplete.

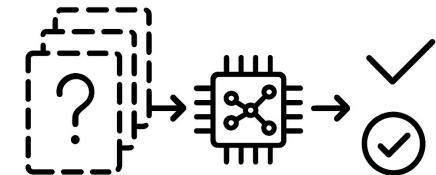
Learning from Incomplete Databases

Perfect cleaning is a myth. Even with best-effort repairs, many plausible datasets remain

Hidden uncertainty \Rightarrow hidden risk. A model trained on one arbitrary repair can look accurate yet flip decisions on another equally valid repair.

Needed: an explicit uncertainty framework.

- capture what is *unknown* in the data,
- propagate that uncertainty through training,
- surface it at inference time.



Practical pay-off.

- Robustness check: see when all admissible models agree (safe to act).
 - Guardrail: abstain or seek more data when predictions diverge.
- Targeted cleaning: focus effort on the cells that actually shrink uncertainty.

Incomplete Databases

Formalism from databases & AI to handle uncertainty by modeling all plausible data interpretations. (*Rooted in modal logic & philosophy*)

Dataset with Quality Issues

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	NULL	60K	...	8K
3	35	NULL	...	[10K, 12K]
...

Q : What is the total income?

Possible Worlds Semantics

Inference:

- All repairs agree \rightarrow Certain answer
 $\text{Range} \leq \tau \rightarrow$ Robust interval (e.g., [5 k – 6 k])
- Range $> \tau \rightarrow$ Uncertain \rightarrow warn / seek more cleaning

Dataset with Quality Issues

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	NULL	60K	...	8K
3	35	NULL	...	[10K, 12K]
...

Q : What is the total income?

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	30	60K	...	8K
3	35	55K	...	7K
...

$$Q(D_1) = 6k$$

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	35	60K	...	8K
3	35	60K	...	8K
...

$$Q(D_2) = 9k$$

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	35	60K	...	8K
3	35	60K	...	8K
...

$$Q(D_3) = 5k$$

Min/Max query result across all possible database repairs.

Range consistent answers:
 $[0.5 - 0.3]$

...

Representing Uncertainty in Databases

C-Tables/M-Tables: Compactly represent multiple possible worlds using variables and conditions.

[Imieliński & Lipski, JACM 1984], [Sundarmurthy et al., ICDT 2017]

Probabilistic Databases: Assign probabilities to possible worlds, quantifying their likelihood.

[Suciu, Olteanu, Ré & Koch, Book 2022]

Answering queries across possible worlds is computationally expensive, often NP-hard or exponential.



ML from Possible Repairs

Inference

- All models ($h_{D_i}^*$) concur \rightarrow **Certain** prediction (e.g., payout = 3 K)
- disagree \rightarrow **Range** prediction (e.g., payout $\in [2 \text{ K}, 4 \text{ K}]$)



Dataset with Quality Issues

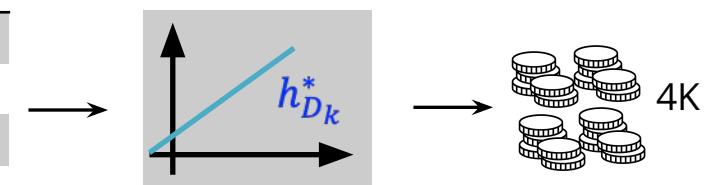
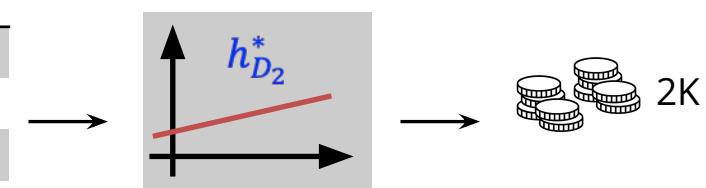
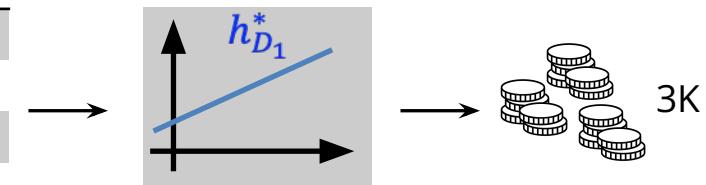
ID	Age	Income	...	Loan
1	25	50K	...	5K
2	NULL	60K	...	8K
3	35	NULL	...	[10K, 12K]
...

Machine-learning analogue of
Consistent Query Answering:
 swap the SQL query Q for a training
 routine T —e.g., gradient descent,
 decision-tree induction, SVM fitting.

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	30	60K	...	8K
3	35	55K	...	7K
...

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	35	60K	...	8K
3	35	60K	...	8K
...

ID	Age	Income	...	Loan
1	25	50K	...	5K
2	35	60K	...	8K
3	35	60K	...	8K
...



KNN Classifiers over Incomplete Information

[Approach: “Certain-kNN” → returns a label only when it is guaranteed across all completions of the missing values]

Insights:

- Missing attributes can flip k-NN labels; intersecting votes across **all** imputations yields a *guaranteed* label.

Approach:

- Model each incomplete record as a value set (hyper-rectangle).
- Two polynomial-time tests (SS, MM) decide if a test point is “certain” without enumerating possible worlds.

Benefits:

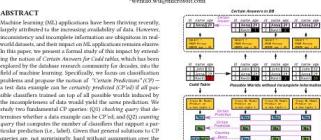
- 100 % precision on “certain” points – i.e., points whose prediction is certain across every imputation.**
- CPClean add-on** ranks the missing cells whose repair would turn “uncertain” points into certain ones, guiding targeted data cleaning.

Shortcomings:

- Guarantees apply only to **numeric-feature k-NN**

Nearest Neighbor Classifiers over Incomplete Information: From Certain Answers to Certain Predictions

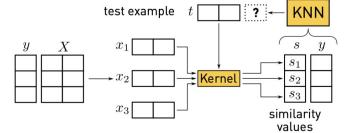
Bojan Karlaš^{1*}, Peng Li², Renchi Wei¹, Nisrine Merve Gürsel¹, Xu Chu¹, Wentao Wu¹, Ce Zhang¹
¹Fudan University, ²ETH Zurich, Georgia Institute of Technology, Microsoft Research
^{*Bojan.karlas, renchi.wei, nisrine.gursel, chuxu, wentao.wu, ce.zhang@fudan.edu.cn, peng.li@inf.ethz.ch, pengli@cs.uga.edu, jwzhang@fudan.edu.cn}



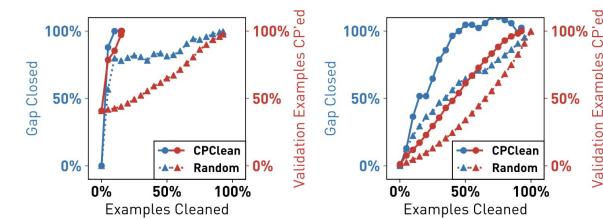
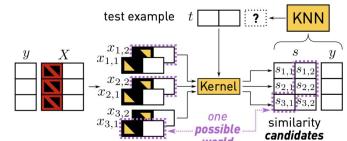
[Karlaš VLDB '20]

Karlaš, Bojan, et al. "Nearest neighbor classifiers over incomplete information: from certain answers to certain predictions." Proceedings of the VLDB Endowment 14.3 (2020): 255-267. [\[Paper\]](#)

a KNN classification over a regular training dataset

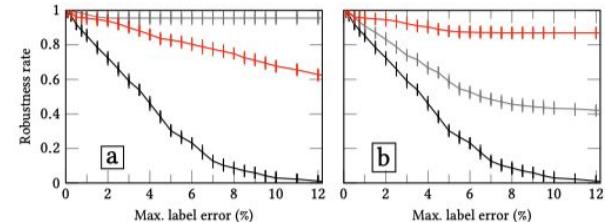
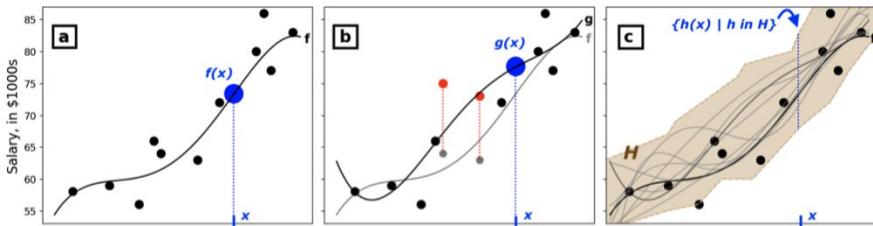


b KNN classification over a training dataset with incomplete information



The Dataset Multiplicity Problem

[Approach: bound model risk across every dataset consistent with the errors]



Insights:

- Introduces a risk interval: the tightest possible lower/upper bound on test error that any admissible dataset can induce for a fixed linear model.

Approach:

- Derive closed-form formulas for the worst- and best-case hinge / logistic loss of any linear classifier under those rules, avoiding enumeration.

Benefits:

- Gives practitioners a numeric certificate of how much reported accuracy can deteriorate.

Shortcomings:

- Theory currently limited to linear models and label-noise rules; deep nets need looser convex relaxations.

The Dataset Multiplicity Problem: How Unreliable Data Impacts Predictions

Anna Meyer
ameyer@wisc.edu
University of Wisconsin - Madison
Madison, USA

Aws Alabargouthi
aws@cs.wisc.edu
University of Wisconsin - Madison
Madison, USA

Loris D'Antoni
loris@cs.wisc.edu
University of Wisconsin - Madison
Madison, USA

ABSTRACT

We introduce dataset multiplicity, a way to study how linear models, under certain assumptions, can fit many datasets with the same test-time predictions. The dataset multiplicity framework asks a counterfactual question of what the set of real-world models (and associated test-time predictions) can be induced by a dataset, given a hypothesis, unlabeled version of the dataset. We discuss how to use this framework to encapsulate various sources of uncertainty in datasets, including missing data, noisy features, missing feature predictors, and noisy labels or features. We show how to exactly analyze the impacts of dataset multiplicity for a specific model architecture, and how to bound the test error of a linear model.

Our empirical analysis shows that real-world datasets, under reasonable assumptions, can have many different samples whose predictions are explained by dataset multiplicity. The degree of domain-specific dataset multiplicity definition determines what samples are considered to be admissible. The number of possible predictions of the models trained on such data and frame this analysis as a ‘multiplicity problem’.

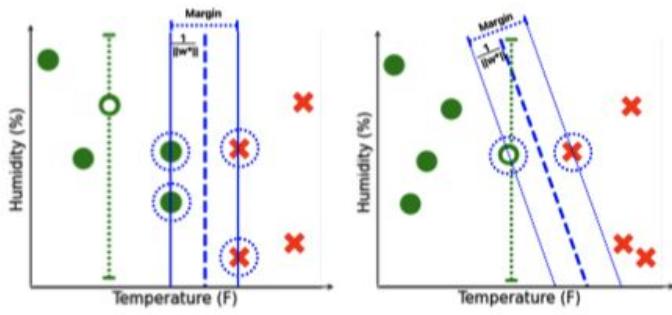
Note that there can be multiple evaluations for

[Meyer FAccT'23]

Meyer, A. P.; Alabargouthi, A.; D'Antoni, L. “The Dataset Multiplicity Problem: How Unreliable Data Impacts Predictions. [Paper]

Certain & Approximately Certain Models for Statistical Learning

[Approach: Fast “certainty test” that lets you skip imputation whenever the missing cells don’t affect the optimum]



Certain and Approximately Certain Models for Statistical Learning

Cheng Zhen
Oregon State University
Corvallis, Oregon
zhen@oregonstate.edu

Nischal Aryal
Oregon State University
Corvallis, Oregon
aryan@oregonstate.edu

Arafa Tariqzadeh
Oregon State University
Corvallis, Oregon
tarizadeh@oregonstate.edu

ABSTRACT

Real-world data is often incomplete and contains missing values. To train accurate models over real-world datasets, users need to spend a substantial amount of time on reasoning, imputing and finding the right model for their data. In this paper, we demonstrate that it is possible to train accurate models directly from incomplete data without any manual intervention. We propose a unified approach for checking the necessity of data imputation to learn accurate models across various widely-used machine learning models. Our approach is based on theoretical guarantees to check this necessity and return accurate models. We empirically show that our proposed algorithm can reduce the amount of time and effort needed for data imputation without loss of accuracy. Our experiments also show that our approach can significantly reduce the number of steps of finding a misspecified mechanism, constructing an accurate

To address the challenge of training accurate models over incomplete data, users usually replace each missing value with a noise, i.e., data imputation, and train their models over the resulting *repaired* data. To repair incomplete data, users must figure out the mechanism of data corruption and choose an appropriate method or base on observed values of some features [38]. Based on this mechanism, researchers have proposed several methods for data repair, such as repair missing values e.g., tree-based or linear regression [21]. Researchers have shown that the desired imputation method may not always be able to find the best solution [21]. This motivates us to challenge to find a model of data missings that results in an accurate model. In this paper, we propose a unified approach for finding a misspecified mechanism, constructing an accurate

[Zhen SIGMOD'24]

Zhen, C. et al. "Certain and Approximately Certain Models for Statistical Learning. [\[Paper\]](#)

Insights:

- Not every example with missing values requires cleaning.
- If the missing cells lie in directions that do not change the model’s optimum, we can train directly on the incomplete data—with full guarantee.

Approach:

- Provide fast algebraic tests (no world enumeration) that decide certainty for linear regression, linear SVM, and two kernel SVMs. When tests pass → output the **certain model** (exactly optimal).
- When tests fail → compute an ϵ -certain model whose loss is within ϵ of the global optimum.

Benefits:

- Skips imputation** for datasets that pass the test, saving cleaning effort and avoiding imputation bias.
- Same code works across several common model families.

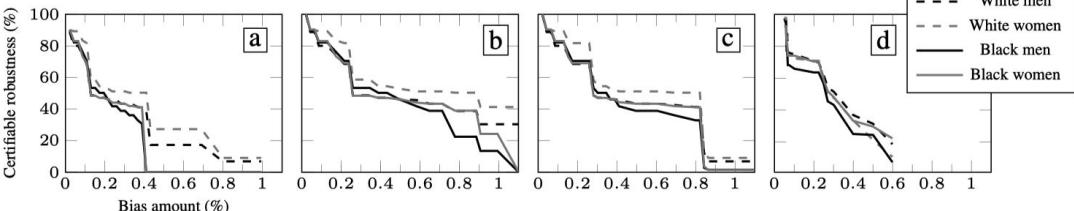
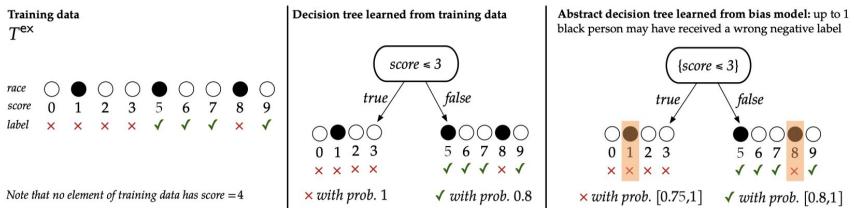
Shortcomings:

- Certainty rarely holds under heavy missingness.
Guarantees limited to the studied linear & kernel models; deep nets need other methods.

Learning from Possible Repairs

Certifying Robustness to Programmable Data Bias in Decision Trees

[Approach — ProgBiasCert: encode “tree + bias program” in SMT to prove the label never flips]



Insights:

- Treat data bias as a **user-written program** (e.g., *age ± 2, race swap, income × 0.9–1.1*).
- A tree is **robust** if its prediction is invariant under **all** transformations allowed by that program.

Approach:

- Translate each path of the decision tree and the bias constraints into a single SMT formula.

Benefits:

- Exact guarantees—no sampling; works with real & categorical features and generates independently checkable proofs

Shortcomings:

- Does not yet handle ensembles or probabilistic bias distributions.

Certifying Robustness to Programmable Data Bias in Decision Trees

Anna P. Meyer, Aws Albarghouthi, and Lorin D’Antoni
Department of Computer Sciences
University of Wisconsin-Madison
Madison, WI 53706
(annameyer, aws, lorin)@cs.wisc.edu

Abstract

Datasets can be biased due to societal inequities, human biases, under-representation of minorities, etc. Our goal is to verify that models produced by a learning algorithm are *pointwise-robust* to potential dataset biases. This is a challenging problem because it entails learning models for a large, or even infinite, number of possible environments, depending on the type of bias. We propose decision-tree learning due to the interpretable nature of the models. Our approach allows programmatically specifying the model’s behavior under different transformations (e.g., a new data feature) and computing types of bias, and targeting bias towards a specific group. To certify robustness, we use a novel symbolic technique to evaluate the model’s behavior under all possible transformations, ensuring that each and every dataset produces the same prediction for a specific test point. We evaluate our approach on datasets that are commonly used in the fairness literature, and demonstrate our approach’s viability on a range of bias models.

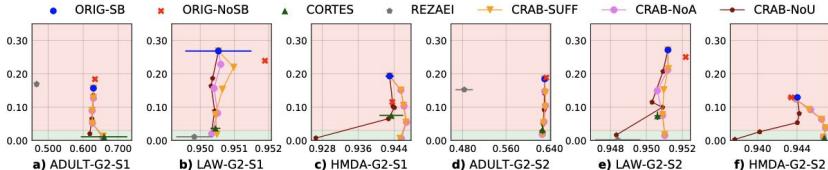
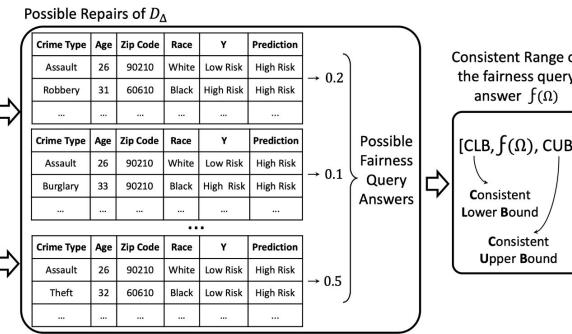
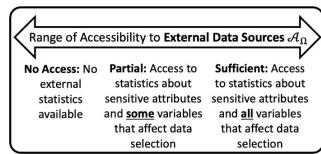
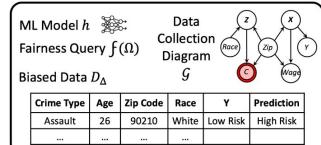
[Meyer NeurIPS'21]

Zhen, C.; Aryal, N.; Termehchy, A.; Chabada, A. S. “Certifying Robustness to Programmable Data Bias in Decision Trees.” [\[Paper\]](#)

Consistent Range Approximation for Fair Predictive Modeling

[Approach: Fair-aware prediction ranges: bound each score so it stays fair under every repair of noisy / missing sensitive attributes]

Input Components



Insights:

- With selection bias we don't know the target-population fairness.
- Treat fairness evaluation as a **query over incomplete data**; answer with a *range* that is guaranteed to contain the truth.

Approach:

- Derive a closed-form range for fairness aggregates.
- Train a classifier that minimises risk while keeping the worst-case value inside the acceptable fairness range.

Benefits:

- Certifies fairness without unbiased samples; needs only the biased data + background knowledge.

Shortcomings:

- Relies on correct causal diagram; ranges may be wide if knowledge is weak.



Consistent Range Approximation for Fair Predictive Modeling

Jiongli Zhu
University of California,
San Diego
jz14@ucsd.edu

Sainyam Galhotra
Cornell University
sgc@cs.cornell.edu

Nazanin Sabri
University of California,
San Diego
nsabri@ucsd.edu

Babak Salimi
University of California,
San Diego
bsalimi@ucsd.edu

ABSTRACT
This paper proposes a novel framework for certifying the fairness of predictive models trained on biased data. It draws from query answering for incomplete data and fairness queries to formulate the problem of consistent range approximation (CRA). The framework employs background knowledge of the data collection process and biased data, working with or without limited statistics about the target population, to derive closed-form ranges of fairness queries. Using CRA, the framework builds predictive models that are certified fair on the target population, regardless of the availability of external data during training. The framework's efficacy is demonstrated through evaluations on real data, showing substantial improvement over existing state-of-the-art methods.

Example 1.1. Consider the dataset in Table 1, which represents

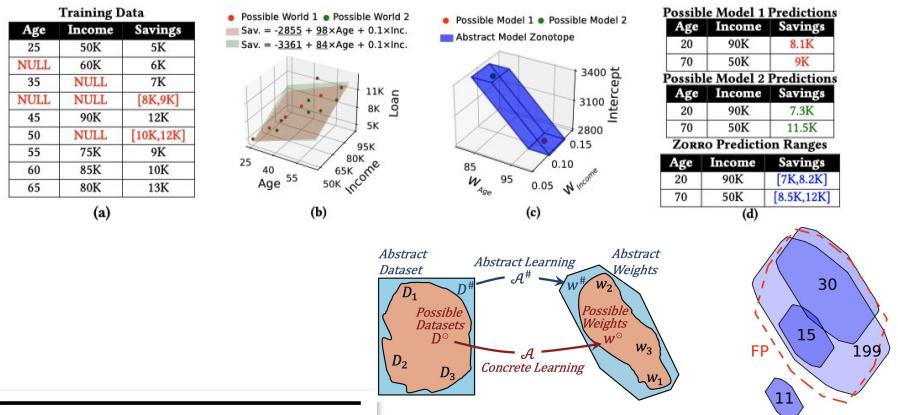
[Zhu VLDB '23]

Consistent Range Approximation for Fair Predictive Modeling. [Paper]

Learning from Possible Repairs

Learning from Uncertain Data: From Possible Worlds to Possible Models

[Approach: Abstract interpretation + zonotopes: train once on a single convex polytope that encodes every possible repair



Learning from Uncertain Data: From Possible Worlds to Possible Models

Jiongli Zhu¹ Su Feng² Boris Glavic³ Babak Salimi¹

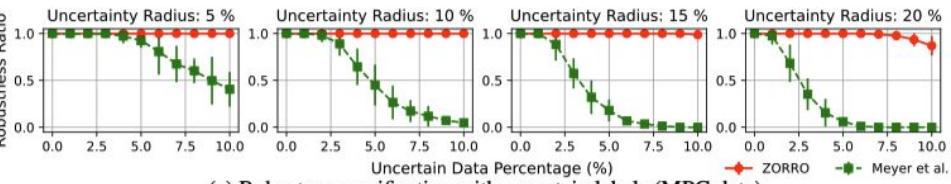
¹University of California, San Diego ²Nanjing Tech University ³University of Illinois, Chicago

Abstract

We introduce an efficient method for learning linear models from uncertain data, where uncertainty is represented as a set of possible variations in the data, leading to predictive multiplicity. Our approach leverages abstract interpretation and zonotopes, a type of convex polytope, to compactly represent these dataset variations. Unlike the symbolic execution of gradient descent on all possible weight vectors simultaneously, we develop techniques to ensure that this process converges to a fixed point and derive closed-form solutions for this fixed point. Our method provides sound over-approximations of all possible optimal models and viable prediction ranges. We validate the theoretical guarantees through theoretical and empirical analysis, highlighting its potential to reason about model and prediction uncertainty due to data quality issues in training data.

[Zhu NeurIPS'24]

Zhu, J.; Feng, S.; Glavic, B.; Salimi, B. "Learning from Uncertain Data: From Possible Worlds to Possible Models. [Paper]



Insights:

- Zonotope = all repairs in a compact affine form.
- Training on the zonotope gives one weight-box that subsumes every per-repair model.

Approach:

- Map each uncertain record to an affine form; the full dataset becomes **one zonotope**. Run gradient descent **symbolically**. Output is a convex box of model weights; any concrete repair yields weights inside this box.

Benefits:

- **Guaranteed intervals for weights & predictions—true model always inside.**

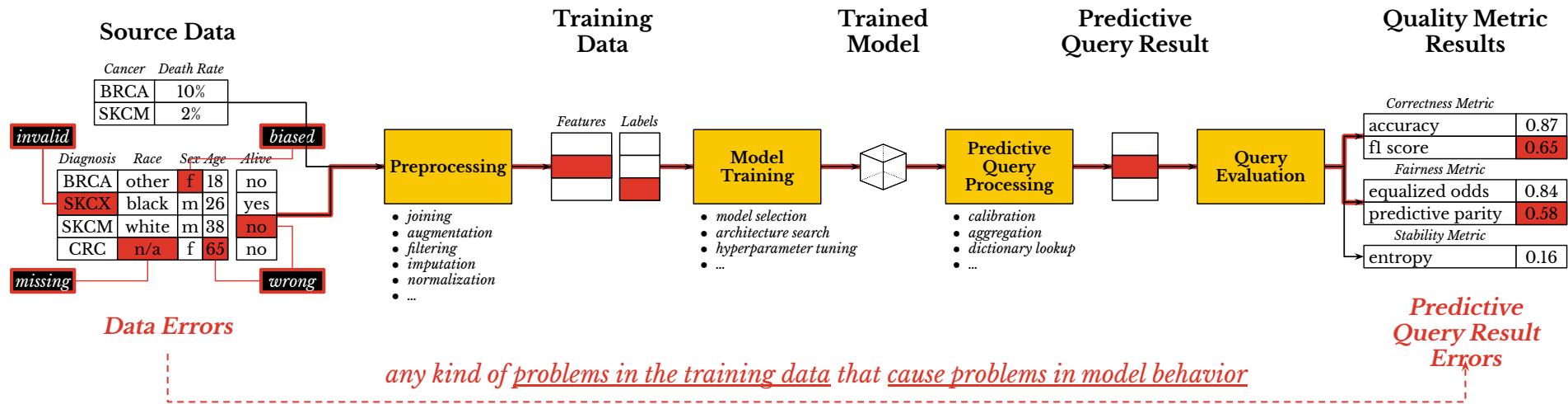
Shortcomings:

- Supports linear models only.

Key Takeaways of Part III

- Residual data uncertainty is inevitable. Cleaning produces at best one plausible version; we must reason over the space of possibilities.
- Guarantee \leftrightarrow coverage trade-off. Certainty methods (Certain-kNN, CRA, ProgBiasCert) give perfect precision or fairness—but may abstain widely.
- Targeted cleaning beats blanket imputation. Algorithms like CPClean and OTClean identify the few cells whose repair actually widens certified coverage.
- Model-side defences matter. Dataset Multiplicity, Certain/Approx-Certain Models, and Zorro show how to train / audit over the whole uncertainty set—returning intervals, ensembles, or risk bounds.
- Certification $>$ best-guess. When stakes are high, prefer guaranteed ranges or proofs of robustness to a single point prediction from a guessed-clean dataset.
- Open frontiers: extend guarantees to deep nets & categorical features, tighten bounds under heavy missingness, and scale zonotope / SMT methods to larger models.

Conclusion: How should we navigate data errors?



Error Detection:
Compute Data Importance

ML Pipeline Debugging:
Leverage Data Provenance

Learning from Uncertain Data:
Apply Possible Worlds Semantics