# 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)
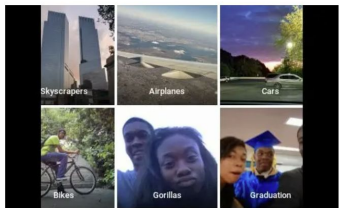
navigating-data-errors.github.io

# Background: ML apps often behave in unintended ways

## Wrong

Google apologises for Photos app's racist blunder
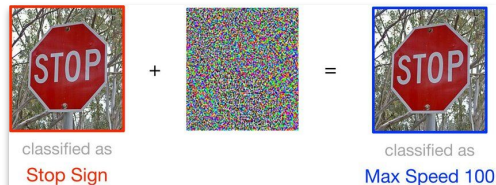
*Source: BBC*

## Biased

Amazon ditched AI recruitment software because it was biased against women

By Erin Winick — October 10, 2018

*Source: MIT Technology Review*

## Unstable

classified as
Stop Sign

classified as
Max Speed 100

*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

**Example Approach:** *Increase regularization by applying dropout*

**Example Approach:** *Improve model capacity by increasing the number of parameters*

Overfitting

Right Fit

Underfitting

Classification

Regression

*Source: MathWorks*

**Problem:** *This is only one piece of the puzzle!*

3

# Observation 1: Data is a crucial piece of the puzzle

Inputs

Trained ML Model

Model Code

Hyperparameters

Training Data

Results

wrong

biased

unstable

**Example error:** *Software Bugs in the code.*

**Example error:** *Learning rate is too high.*

**Example error: Invalid Value** *"SKCX" does not appear in the dictionary of valid cancer types.*

**Example error: Biased Value** *Data about male breast cancer is rare and noisy.*

**Example error: Missing Value** *Clinical data is frequently incomplete.*

**Example error: Wrong Value** *Information about the patient appears valid but is not true.*

| 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 |

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

4

# Observation 2: ML apps are built by complex pipelines



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

5

# 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 🗑 | Repair 🔨 | Ignore 🥂 |
|---|---|---|
| *Remove the faulty data from the training set.* | *Perform manual quality control which might include repeating the data acquisition process.* | *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

**Source Data**

| Cancer | Death Rate |
|--------|-----------|
| BRCA | 10% |
| SKCM | 2% |

*invalid*

*biased*

| 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 |

*missing*

*wrong*

*Data Errors*

**Training Data**

**Trained Model**

**Predictive Query Result**

**Quality Metric Results**

**Preprocessing**
- *joining*
- *augmentation*
- *filtering*
- *imputation*
- *normalization*
- *...*

*Features*  *Labels*

**Model Training**
- *model selection*
- *architecture search*
- *hyperparameter tuning*
- *...*

**Predictive Query Processing**
- *calibration*
- *aggregation*
- *dictionary lookup*
- *...*

**Query Evaluation**

*Correctness Metric*

| accuracy | 0.87 |
|----------|------|
| f1 score | 0.65 |

*Fairness Metric*

| equalized odds | 0.84 |
|----------------|------|
| predictive parity | 0.58 |

*Stability Metric*

| entropy | 0.16 |
|---------|------|

*Predictive Query Result Errors*

*any kind of <u>problems in the training data</u> that <u>cause problems in model behavior</u>*

**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*                                                                                                    *Not So 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 |

**biased**

**Solution approach:**
*Measure the impact of the value on model quality.*

**Solution approach:**
*Check if the value is marked as missing.*

**missing**

**wrong**

**How do we measure this?**
*That is the main topic of this part of the tutorial.*

**Recall:** *Data errors are any kind of* <u>*problem in the training data*</u> *that cause* <u>*problems in model behavior*</u>.

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

11

# We can define a data attribution function



**Recall:** *Data errors are any kind of <u>problem in the training data</u> that cause <u>problems in model behavior</u>.*

# How do we use importance to detect data errors?

**Attribution Function Example 1:**

```python
def compute_importance(value):
    return -1.0 if value == "n/a" else 1.0
```

**Attribution Function Example 2:**

```python
VALID_CANCER_CODES = ...

def compute_importance(value):
    if value not in VALID_CANCER_CODES:
        return -1.0
    return 1.0
```

| 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 |

$\varphi$

*Data Importance*

| |
|---|
| 0.325 |
| -0.873 |
| -0.217 |
| 0.664 |

*Human-in-the-Loop (optional)*

**Model Training Pipeline**

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

**Data Repair Pipeline**

# 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 <u>problem in the training data</u> that cause <u>problems in model behavior</u>.*

## Design Consideration 2

*What kind of intervention do we intend to apply?*

*Discard*

*Repair*

*Something Else*

Model Quality

*Good\**

*Ineffective*

*Bad\**

*\* Assuming higher is better*

Number of Interventions

**Challenge:** *How do we define an effective attribution function?*

# 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.

16

# 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.

**[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]

17

# Influence Function

[Approach: **Marginal Contribution, Gradient**]

*Training Labels*
$y_1$
$y_2$
$y_3$

*Presence Indicators*
$\epsilon_1$
$\epsilon_2$
$\epsilon_3$

*Training Features*
$x_1$
$x_2$
$x_3$

$f_\theta$

*Predicted Labels*
$\hat{y}_1$
$\hat{y}_2$
$\hat{y}_3$

$\mathcal{L}$

*Training Loss*
$\ell_1$
$\ell_2$
$\ell_3$

$H_\theta^{-1} \cdot \frac{\partial}{\partial \theta}(\mathcal{L} \circ f_\theta)$

*Training Gradients*
$\delta_1$
$\delta_2$
$\delta_3$

*Trained Parameters*
$\hat{\theta}$

*Data Importance*

*Pointwise Validation Influence*
$\varphi_{1,1}$ $\varphi_{1,2}$
$\varphi_{2,1}$ $\varphi_{2,2}$
$\varphi_{3,1}$ $\varphi_{3,2}$

**1D Mean**

*Influence*
$\varphi_1$
$\varphi_2$
$\varphi_3$

*Validation Labels*
$y'_1$
$y'_2$

*Validation Features*
$x'_1$
$x'_2$

$f_\theta$

*Predicted Labels*
$\hat{y}'_1$
$\hat{y}'_2$

$\mathcal{L}$

*Validation Loss*
$\ell'_1$
$\ell'_2$

$\frac{\partial}{\partial \theta}(\mathcal{L} \circ f_\theta)$

*Validation Gradients*
$\delta'_1$
$\delta'_2$

$$\varphi_{i,j} = \frac{\partial \mathcal{L}(x'_j, y'_j, \hat{\theta})}{\partial \epsilon_i}\bigg|_{\epsilon_i = 0} = \frac{\partial(\mathcal{L} \circ f_\theta)(x'_j, y'_j, \hat{\theta})}{\partial \theta} \cdot H_\theta^{-1} \frac{\partial(\mathcal{L} \circ f_\theta)(x_i, y_i, \hat{\theta})}{\partial \theta}$$



**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]

18

# Area Under the Margin

[Approach: **Uncertainty Analysis**]



Training Labels: $y_1$, $y_2$, $y_3$

Training Features: $x_1$, $x_2$, $x_3$

$f_\theta$

Logits: $z_1$, $z_2$, $z_3$

**Compute Margin**

Current Epoch Margins: $\varphi_1^{(t)}$, $\varphi_2^{(t)}$, $\varphi_3^{(t)}$

**Average Over All Epochs**

Area Under the Margin: $\varphi_1$, $\varphi_2$, $\varphi_3$ — *Data Importance*

Current Epoch Parameters: $\theta^{(t)}$

Correct Label / Strongest Incorrect Label — Margin

$$M^{(t)}(\mathbf{x}, y) = \overbrace{z_y^{(t)}(\mathbf{x})}^{\text{assigned logit}} - \overbrace{\max_{i \neq y} z_i^{(t)}(\mathbf{x})}^{\text{largest other logit}}$$

$$\text{AUM}(\mathbf{x}, y) = \frac{1}{T} \sum_{t=1}^{T} M^{(t)}(\mathbf{x}, y)$$



CIFAR100 (0% Unif. Noise) / CIFAR100 (40% Unif. Noise)
- Random
- AUM
- 99\% Thrsh.

**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**

19

# Unconfident Margins

[Approach: **Uncertainty Analysis**]



Labels
$y_1$
$y_2$
$y_3$

Features
$x_1$
$x_2$
$x_3$

$f_\theta$

Out-of-sample Class Probabilities
$p_{1,1}$ $p_{1,2}$
$p_{2,1}$ $p_{2,2}$
$p_{3,1}$ $p_{3,2}$

**Compute Margin**

Margins
$\varphi'_1$
$\varphi'_2$
$\varphi'_3$

*Data Importance*

Unconfident Margins
$\varphi_1$
$\varphi_2$
$\varphi_3$

**Estimate Confident Joint Counts**

$C_{\tilde{y},y^*}$

**Identify Off-diagonal Data Points**

Off-diagonal Indicator
1
0
1

1D Mean

Class Thresholds
$t_1$ $t_2$

$$C_{\tilde{y},y^*}[i][j] := |\hat{X}_{\tilde{y}=i,y^*=j}|$$

$$\hat{X}^{(\text{simple})}_{\tilde{y}=i,y^*=j} = \{x \in X_{\tilde{y}=i}: \ \hat{p}(\tilde{y}=j; x, \theta) \geq t_j\}$$

$$t_j = \frac{1}{|X_{\tilde{y}=j}|} \sum_{x \in X_{\tilde{y}=j}} \hat{p}(\tilde{y}=j; x, \theta)$$

(a) ResNet18 Validation Accuracy    (b) ResNet50 Validation Accuracy

## 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.

[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]

# Model Training Outcome

[Approach: **Surrogate Data Model**]



**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.

[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]

# 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.*

**Approach**

*We should measure marginal contribution over all subsets.*

**Shapley value**

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

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

# Effectiveness at Data Debugging



(a) Noisy labels detection

(b) Watermark removal

(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 | |
|---|---|---|---|---|
| *K*NN-Shapley | **31.70%** → **47.00%** | | **23.35%** → **29.80%** | |
| *K*NN-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% | |

[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 <u>exponential</u>, making it intractable to compute the exact Shapley value for an arbitrary model.*

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

# 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.*

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

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

**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.*

**[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]

26

# 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.*

*Example Situation*

- *We are computing the Shapley value of data point i*
- *Data is <u>sorted by similarity</u> to the validation data point*

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

| Validation Features | Validation Labels |
|---|---|
| $x'_1$ | $y'_1$ |

| Training Features | Training Labels |
|---|---|
| $x_1$ | $y_1$ |
| $x_2$ | $y_2$ |
| $x_3$ | $y_3$ |
| $x_4$ | $y_4$ |
| $x_5$ | $y_5$ |
| $x_6$ | $y_5$ |

$i \rightarrow x_3$

$j \rightarrow x_5$

*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*

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

*Observation 2:*
*If data point i is not in the top-1, this term will be zero.*

*Dynamic Programming*

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

*Final Simplification*

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

**Result:**
*After sorting the data, we can compute exact Shapley values in a single pass. Final computational complexity is*
$$\mathcal{O}(N \log N)$$

[ Jia VLDB '19]
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]

# 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:*

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

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

*Redefined "global utility function" of subset S over the entire SGD run:*

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



**[Wang ICLR '25]**
Wang, Jiachen T., et al. "Data Shapley in One Training Run." The Thirteenth International Conference on Learning Representations. [Paper] [Blog]

# 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.*

*Data points ordered by importance*



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



*Combinations of predicates that explain model behavior*





**[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]

30

# 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.*



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

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

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

[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]

# 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



**Challenge:** *How should we debug 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 data processing operators*
- *Characterized by the fit/transform and estimator/transformer abstractions for building pipelines*



Source: https://vitalflux.com/sklearn-machine-learning-pipeline-python-example/

# Tensorflow Extended (TFX)

## Highlights

- *End-to-end platform for production ML pipelines*
- *Built on TensorFlow and optimized for scalability*
- *Includes reusable components such as ExampleGen, Transform, Trainer, Evaluator, and Pusher for building robust ML pipelines*
- *Supports orchestration with Airflow, Kubeflow, and Vertex AI*
- *Strong emphasis on model validation and monitoring*



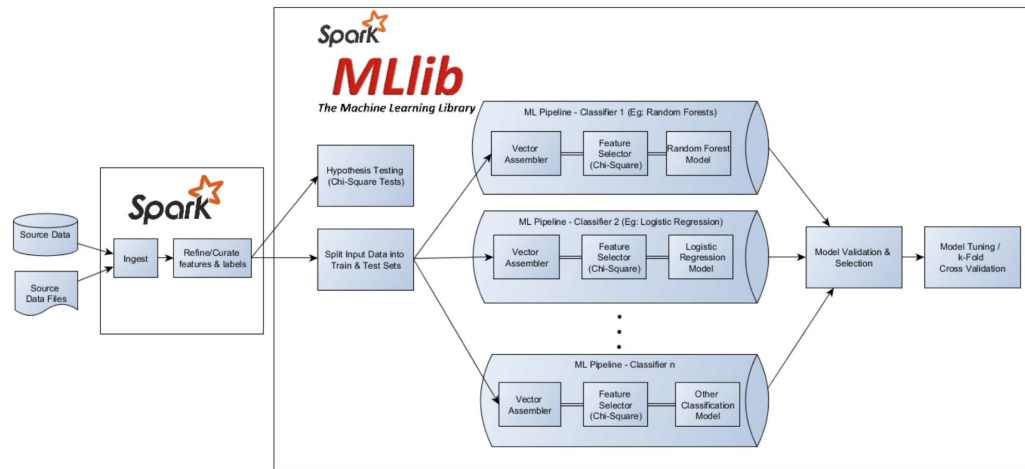Source: https://www.tensorflow.org/tfx/guide

**[TFX]**
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]

# 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*
- *Provides a high-level API for constructing, tuning, and evaluating machine learning pipelines using transformers and estimators*



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

Journal of Machine Learning Research 17 (2016) 1-7     Submitted 5/15; Published 4/16

MLlib: Machine Learning in Apache Spark

Xiangrui Meng†     MENG@DATABRICKS.COM
*Databricks, 160 Spear Street, 13th Floor, San Francisco, CA 94105*
Joseph Bradley     JOSEPH@DATABRICKS.COM
*Databricks, 160 Spear Street, 13th Floor, San Francisco, CA 94105*
Burak Yavuz     BURAK@DATABRICKS.COM
*Databricks, 160 Spear Street, 13th Floor, San Francisco, CA 94105*
Evan Sparks     SPARKS@CS.BERKELEY.EDU
*UC Berkeley, 465 Soda Hall, Berkeley, CA 94720*
Shivaram Venkataraman     SHIVARAM@EECS.BERKELEY.EDU
*UC Berkeley, 465 Soda Hall, Berkeley, CA 94720*
Davies Liu     DAVIES@DATABRICKS.COM
*Databricks, 160 Spear Street, 13th Floor, San Francisco, CA 94105*
Jeremy Freeman     FREEMANJ11@JANELIA.HHMI.ORG
*HHMI Janelia Research Campus, 19005 Helix Dr, Ashburn, VA 20147*
DB Tsai     DBT@NETFLIX.COM
*Netflix, 970 University Ave, Los Gatos, CA 95032*
Manish Amde     MANISH@ORIGAMILOGIC.COM
*Origami Logic, 1134 Crane Street, Menlo Park, CA 94025*

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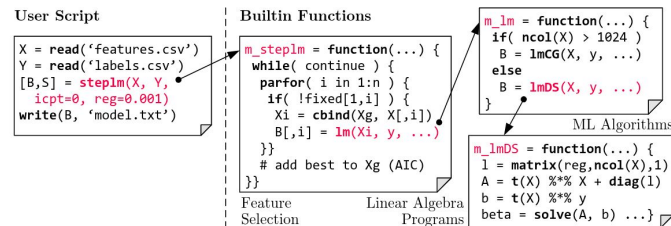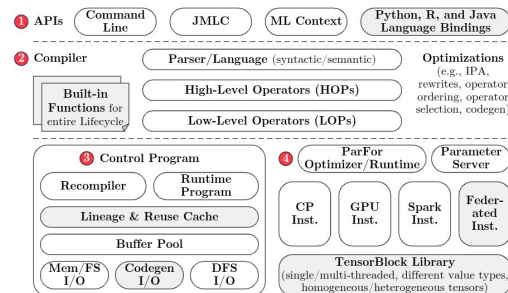Meng, Xiangrui, et al. "Mllib: Machine learning in apache spark." Journal of Machine Learning Research 17.34 (2016): 1-7. [Paper] [Website] [Code]

# Apache SystemDS



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*

- *Provides tools for lineage tracing, intermediate result inspection, and performance analysis to aid in model development and debugging*

**[SystemDS]**
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]

# 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 Pipelines**

[Website] [Documentation]

Highlights

- *Define, automate, and manage end-to-end ML workflows*
- *Automatically tracks data, code, parameters, and model artifacts*
- *Leverages AWS Cloud infrastructure*

**Azure Machine Learning Pipelines**

[Website] [Documentation]

Highlights

- *Orchestration of ML workflows with reusable, modular pipeline components*
- *Versioning, monitoring, and CI/CD integration*
- *Runs pipelines on scalable Azure compute targets*

**Vertex AI Pipelines**

[Website] [Documentation]

Highlights

- *Connects with Vertex AI services like training, hyperparameter tuning, and model deployment*
- *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*

- *Deep learning models account for 60% of pipelines*

- *Pipelines often have a large lifespan, averaging 36 days*

- *About 1/4 model training runs results in model deployment*



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.
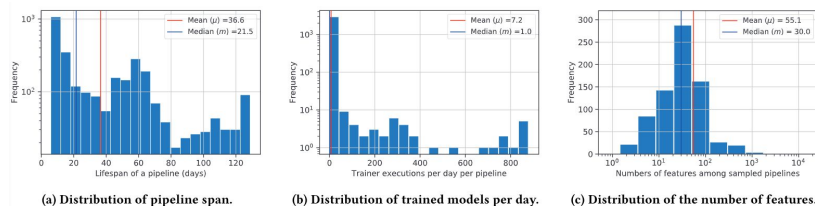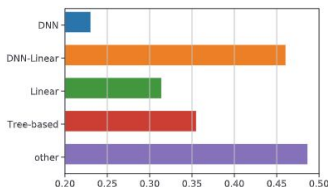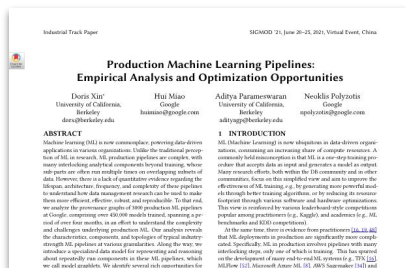


Figure 5: Percentage of Trainer runs with each model type



(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]

42

# 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)*

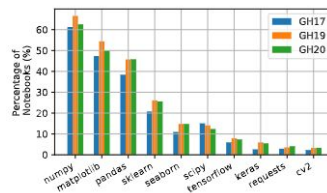| Dimension | Metric | GH17 | GH19 | GH20 |
|---|---|---|---|---|
| Notebooks | Total | 1.23M | 4.6M | 8.7M |
| | Deduped | 66.0% | 65.5% | 65.7% |
| | Linear | 26.4% | 29.1% | 30.3% |
| | Completely Linear | 21.2% | 23.3% | 24.6% |
| Languages | Python | 81.7% | 91.7% | 91.1% |
| | Other | 18.3% | 8.3% | 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 | N/A |
| | Explicit | 10K | 129K | 252K | 29.7M |
| #Distinct Ops | Implicit | 668K | 1.8M | 2.6M | N/A |
| | Explicit | 584 | 3.4K | 5.5K | 23.5K |

[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]

43

# How should we reason about pipelines?

### Source Data

| Cancer | Death Rate |
|--------|-----------|
| BRCA | 10% |
| SKCM | 2% |

*invalid*

*biased*

| 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 |

*missing*

*wrong*

***Data Errors***

### Training Data

*Features*   *Labels*

**Preprocessing**

- *joining*
- *augmentation*
- *filtering*
- *imputation*
- *normalization*
- *...*

### Trained Model

**Model Training**

- *model selection*
- *architecture search*
- *hyperparameter tuning*
- *...*

**Predictive Query Processing**

- *calibration*
- *aggregation*
- *dictionary lookup*
- *...*

### Predictive Query Result

**Query Evaluation**

### Quality Metric Results

*Correctness Metric*

| accuracy | 0.87 |
|----------|------|
| f1 score | 0.65 |

*Fairness Metric*

| equalized odds | 0.84 |
|----------------|------|
| predictive parity | 0.58 |

*Stability Metric*

| entropy | 0.16 |
|---------|------|

***Predictive Query Result Errors***

**What caused this data error?**

**How does it propagate through the pipeline?**

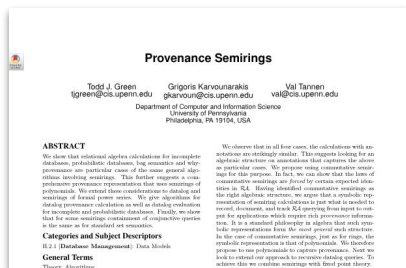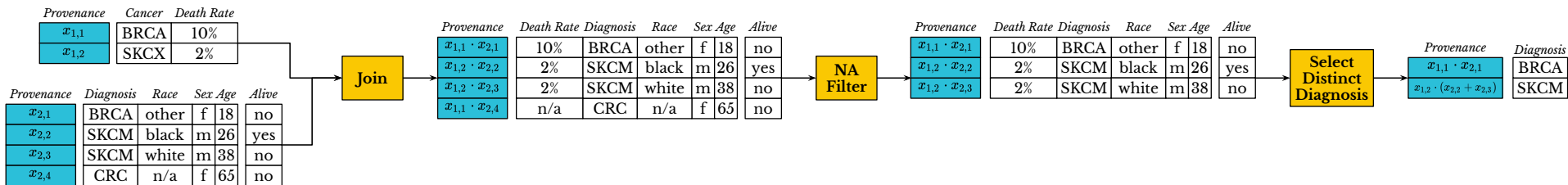**What is the impact of this tuple on the error?**

45

# Leveraging the Provenance Semiring Framework

## Highlights

- *Theoretical framework analyzing the relationship between input and output tuples of relational queries*
- *It allows us to determine the presence of an output tuple as a function of the presence of an input tuples*

## Application to an Example Pipeline



[Green SIGMOD '07]
Green, Todd J., Grigoris 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

*Computing the Shapley value using the KNN proxy method assumes that the presence of a single source data point maps directly to a single data point fed to the model. Hence, the results are not directly applicable to arbitrary pipelines.*

## Insight

*We can use the provenance framework to analyze pipelines and develop PTIME algorithms for computing the Shapley value. We notice that there are three canonical types of pipelines that are both representative of real-world pipelines, and lend themselves to efficient Shapley value computation.*

## Approach

*Compile provenance polynomials to Additive Decision Diagrams and use them to compute Shapley values in PTIME.*
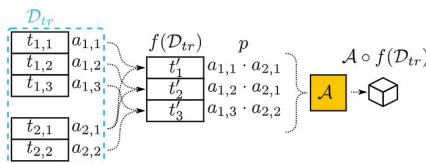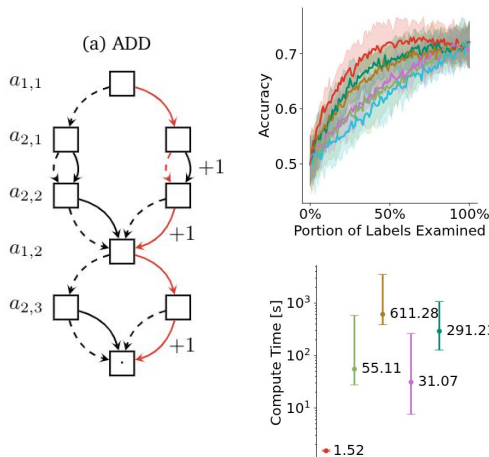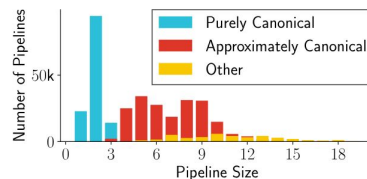
(a) Map pipeline

(b) Fork pipeline

(c) One-to-many join pipeline

(d) Distribution of canonical pipelines

(a) ADD

47

# Debugging Predictive Queries with Rain

[Attribution Function: Influence]

## Challenge

*The existing influence-based attribution methods assume that the model predictions are directly used for computing model quality. However, model inference is often part of a larger predictive query.*

## Insight

*Using provenance polynomials to track 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) are 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"
```



Research 15: Machine Learning for Cleaning, Integration, and Search    SIGMOD '20, June 14–19, 2020, Portland, OR, USA

**Complaint-driven Training Data Debugging for Query 2.0**

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

Lampros Flokas
Columbia University
New York, NY
lamflokas@cs.columbia.edu

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

Jiannan Wang
Simon Fraser University
Burnaby, BC, Canada
jnwang@sfu.ca

**ABSTRACT**

As the need for machine learning (ML) increases rapidly across all industry sectors, there is a significant interest among commercial database providers to support "Query 2.0", which integrate model inference into SQL queries. Debugging Query 2.0 is very challenging since an unexpected query result may be caused by the bugs in training data (e.g., wrong labels, corrupted features). In response, we propose Rain, a complaint-driven training data debugging system. Rain allows users to specify complaints over the query's

**1 INTRODUCTION**

Database researchers have long advocated the value of integrating model inference within the DBMS: data used for model inference is already in the DBMS, it brings the code (model) to the data, and it provides a familiar relational user

# Debugging Data Distributions with MLinspect

**Challenge**

*Some data errors are not necessarily caused by values in source data, but rather by the pipeline itself.*

**Insight**

*Detecting such errors requires on-the-fly analysis of the distribution of data as it passes through the pipeline.*

**Approach**

*Instrument functions of Python data science libraries, track lineage of operators and measure changes in data distribution. Apply rule-based approaches to determine if an error has occurred (e.g. if a bias against a sensitive group has been introduced).*



**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**

```python
# 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())])
featurisation = 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', featurisation),
    ('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))
```
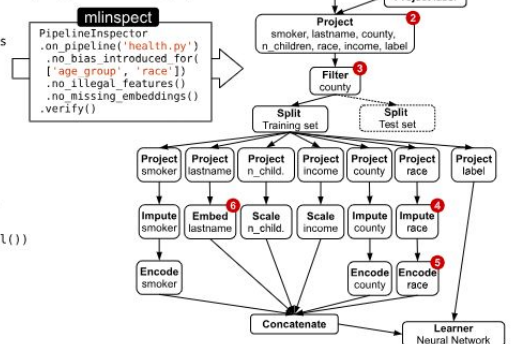
**Corresponding dataflow DAG for instrumentation, extracted by *mlinspect***

**Declarative inspection of preprocessing pipeline**

```
mlinspect
PipelineInspector
.on_pipeline('health.py')
.no_bias_introduced_for(
    ['age_group', 'race'])
.no_illegal_features()
.no_missing_embeddings()
.verify()
```

**[Grafberger VLDBJ '22]**
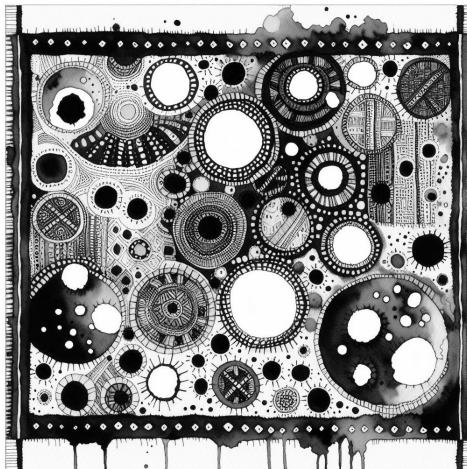Grafberger, Stefan, et al. "Data distribution debugging in machine learning pipelines." The VLDB Journal 31.5 (2022): 1103-1126. [Paper] [Code]

49

# 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**
- **Data provenance is a powerful framework for analyzing ML pipelines**

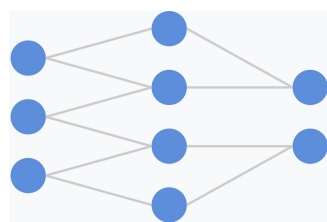# Part III:
# Learning from Uncertain and Incomplete Data

Babak Salimi

# The Standard ML Pipeline

**Input Data** ➡️ **Data Cleaning** ➡️ **Model** ➡️ **Inference**

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

Loan Denied: **High** Risk

Loan Approved: **Low** Risk

**Alice**

**Bob**

⚠️ **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**

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

Loan Denied:
**High** Risk

Loan Approved:
**Low** Risk

**Alice**

**Bob**

📉 **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 ⇒ 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 → Certain answer**
  **Range ≤ τ → Robust interval** (e.g., [5 k – 6 k])
- **Range > τ → Uncertain → warn / seek more cleaning**

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

$$Q(D_1) = 6k$$

### 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 | 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 |
| ... | ... | ... | ... | ... |

$$Q(D_3) = 5k$$

Range consistent answers:
**[0.5 - 0.3]**

Min/Max query result across all possible database repairs.

# 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 → *Certain* prediction**   (e.g., payout = 3 K)
- **disagree → *Range* prediction**   (e.g., payout ∈ [2 K , 4 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 |
| ... | ... | ... | ... | ... |

$h^*_{D_1}$ → 3K

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

$h^*_{D_2}$ → 2K

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

$h^*_{D_k}$ → 4K

# 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**

# 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.

62

# 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]



(a) Data cleaning is not needed   (b) Data cleaning is needed

**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 and resources imputing and finding proper values for missing data items. In this paper, we demonstrate that it is possible to learn accurate models directly from data with missing values for certain training data and target models. We propose a unified approach for checking the necessity of data imputation to learn accurate models across various widely-used machine learning paradigms. We build efficient algorithms with theoretical guarantees to check this necessity and return accurate models in cases where imputation is unnecessary. Our extensive experiments indicate that our proposed algorithms significantly reduce the amount of time and effort needed for data imputation without imposing considerable computational overhead.

To address the problem of training over incomplete data, users usually replace each missing data item with a value, i.e., data imputation, and train their models over the resulting *repaired* data. To repair incomplete data, users must figure out the mechanisms and causes of data missingness, e.g., completely at random or based on observed values of some features [28]. Based on this mechanism, they build a (statistical) model for missing data and replace the missing values with some measurements defined over this model, e.g., mean. Users may also leverage a variety of ML models to repair missing data, e.g., tree-based or linear regression [21]. Researchers have shown that the desired imputation method may vary depending on the downstream ML task [20]. Hence, it is often challenging to find a model of data missingness that results in an accurate ML model for a downstream task [20]. The aforementioned steps of finding a missingness mechanism, constructing an accurate

**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.
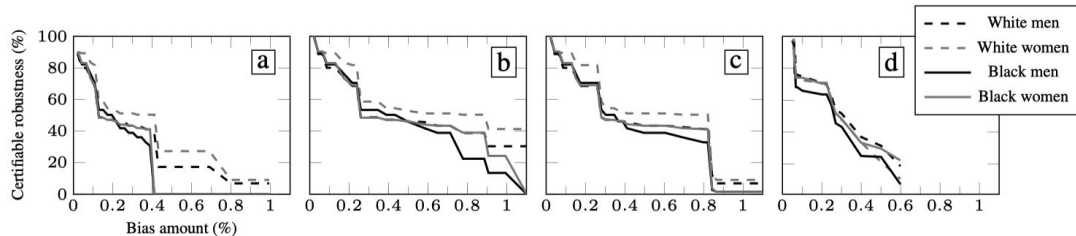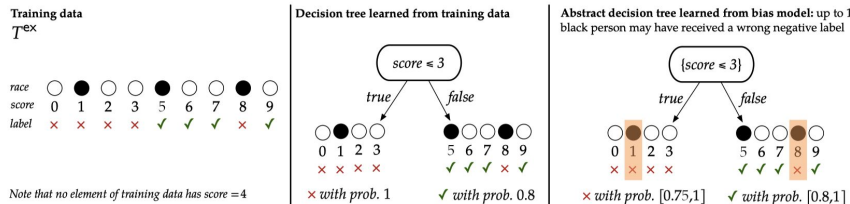- Same code works across several common model families.

**Shortcomings:**

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

**[Zhen SIGMOD'24]**
Zhen, C. et al. "Certain and Approximately Certain Models for Statistical Learning. [Paper]

# Certifying Robustness to Programmable Data Bias in Decision Trees

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





Note that no element of training data has score = 4

## 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.
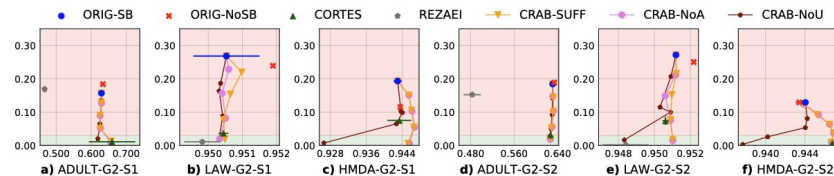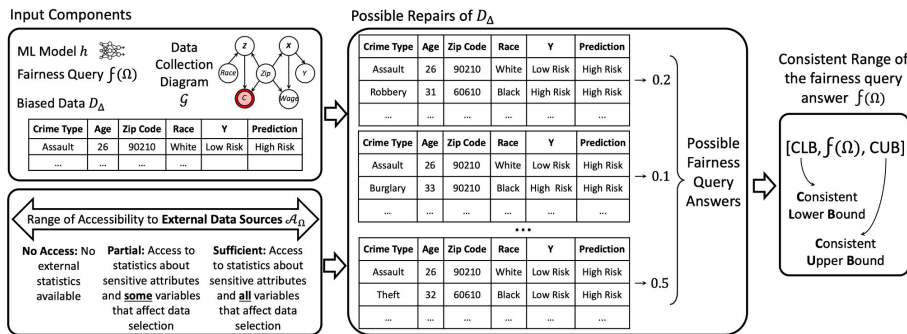
## Benefits:

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

## Shortcomings:

- Does not yet handle ensembles or probabilistic bias distributions.

[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**

ML Model $h$ · Data Collection Diagram $\mathcal{G}$

Fairness Query $f(\Omega)$

Biased Data $D_\Delta$

| Crime Type | Age | Zip Code | Race | Y | Prediction |
|---|---|---|---|---|---|
| Assault | 26 | 90210 | White | Low Risk | High Risk |
| ... | ... | ... | ... | ... | ... |

Range of Accessibility to **External Data Sources** $\mathcal{A}_\Omega$

**No Access:** No external statistics available

**Partial:** Access to statistics about sensitive attributes and **some** variables that affect data selection

**Sufficient:** Access to statistics about sensitive attributes and **all** variables that affect data selection

**Possible Repairs of $D_\Delta$**

| Crime Type | Age | Zip Code | Race | Y | Prediction |
|---|---|---|---|---|---|
| Assault | 26 | 90210 | White | Low Risk | High Risk |
| Robbery | 31 | 60610 | Black | High Risk | High Risk |
| ... | ... | ... | ... | ... | ... |

$\rightarrow 0.2$

| Crime Type | Age | Zip Code | Race | Y | Prediction |
|---|---|---|---|---|---|
| Assault | 26 | 90210 | White | Low Risk | High Risk |
| Burglary | 33 | 90210 | Black | High Risk | High Risk |
| ... | ... | ... | ... | ... | ... |

$\rightarrow 0.1$

| Crime Type | Age | Zip Code | Race | Y | Prediction |
|---|---|---|---|---|---|
| Assault | 26 | 90210 | White | Low Risk | High Risk |
| Theft | 32 | 60610 | Black | Low Risk | High Risk |
| ... | ... | ... | ... | ... | ... |

$\rightarrow 0.5$

Possible Fairness Query Answers

Consistent Range of the fairness query answer $f(\Omega)$

$[CLB, f(\Omega), CUB]$

**C**onsistent **L**ower **B**ound

**C**onsistent **U**pper **B**ound

a) ADULT-G2-S1    b) LAW-G2-S1    c) HMDA-G2-S1    d) ADULT-G2-S2    e) LAW-G2-S2    f) HMDA-G2-S2

ORIG-SB    ORIG-NoSB    CORTES    REZAEI    CRAB-SUFF    CRAB-NoA    CRAB-NoU

## 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.

[Zhu VLDB '23]
Consistent Range Approximation for Fair Predictive Modeling. [Paper]

# Learning from Uncertain Data: From Possible Worlds to Possible Models

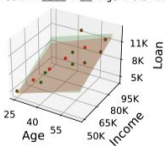[Approach: Abstract interpretation + zonotopes: train once on a single convex polytope that encodes every possible repair



(a) Robustness verification with uncertain labels (MPG data).



**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.

66

# 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 ↔ 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.**