

# Influence Score-Based Instance Reweighting for Interpretable and Robust XGBoost Models

First Author Last<sup>1</sup>, Second Author Last<sup>1</sup>, Third Author Last<sup>1</sup>

<sup>1</sup>School of Computer Science and Engineering, Kerala University of Digital Sciences, Innovation and Technology, Thiruvananthapuram, India.

Contributing authors: [first.author@email.edu](mailto:first.author@email.edu);

## Abstract

We propose a data-perspective interpretability method for XGBoost that identifies influential training instances via deletion diagnostics and an adapted Cook’s distance, then performs domain-aware fine-tuning by instance-specific reweighting. Across three benchmark datasets, our approach improves calibration (log loss) and AUC, while maintaining or improving accuracy and recall. We analyze the distribution of influence scores, categorize influential samples into outliers, mislabeled points, and edge cases, and quantify their effect on feature usage and gain. The resulting framework offers actionable interpretability and robustness gains with minimal modifications to standard XGBoost pipelines.

**Keywords:** XGBoost, interpretability, Cook’s distance, deletion diagnostics, influential instances, instance reweighting, robustness

## 1 Introduction

Machine learning models increasingly inform decisions in high-stakes domains, raising the need for transparent and trustworthy behavior. XGBoost (1) is widely adopted for its accuracy and efficiency, yet most interpretability techniques for tree ensembles are perturbation-based (e.g., LIME (5) and SHAP (4)) rather than *data-centric*. We bridge this gap by adapting deletion diagnostics and Cook’s distance to quantify the *instance-level* influence on XGBoost outcomes and by leveraging those insights to perform domain-aware fine-tuning via reweighting. This work extends prior data-perspective interpretability for AdaBoost (2) to gradient-boosted trees.

### *Contributions.*

- A deletion-diagnostics procedure for XGBoost to quantify per-instance influence via an adapted Cook’s distance.
- A principled thresholding and categorization of influential points (outliers, mislabeled, edge cases).
- An instance-specific reweighting scheme that improves model calibration and AUC with minimal impact on accuracy.
- Empirical validation on three datasets, with analyses of influence distributions and feature-importance shifts.

## 2 Related Work

**Model-agnostic explanations.** LIME (5) and SHAP (4) explain predictions using local perturbations or Shapley values. **Influence-based analysis.** Influence functions approximate leave-one-out effects for differentiable models (3). Our method directly performs deletion diagnostics on XGBoost, connecting interpretability with actionable reweighting.

## 3 Methodology

### 3.1 Deletion Diagnostics for XGBoost

Let  $\hat{y}_j$  denote the prediction on example  $j$  from a baseline model trained on all  $n$  training instances, and  $\hat{y}_j^{(-i)}$  the prediction when instance  $i$  is removed and the model retrained. Define the adapted influence score (Cook’s distance variant):

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_j^{(-i)})^2}{(p+1) \text{Var}(\hat{y}_j - \hat{y}_j^{(-i)})} \quad (1)$$

where  $p$  is the number of features. For each  $i$ , we also record changes in train/test log loss, accuracy, and AUC between baseline and retrained models.

### 3.2 Influence Thresholding and Categorization

We identify influential points as those with  $D_i$  above  $\mu_D + 2\sigma_D$ , where  $\mu_D$  and  $\sigma_D$  are the mean and standard deviation of influence scores. Influential points are then categorized:

1. **Outliers:** any feature’s absolute z-score  $> 3$ .
2. **Mislabeled:** misclassified by the baseline (label vs. rounded prediction).
3. **Edge cases:** influential but neither outliers nor mislabeled.

### 3.3 Instance Reweighting

We assign sample weights  $w_i$  per category (default  $w_i = 1$ ). A simple scheme is

$$w_i = \begin{cases} 0.1 & \text{if } i \text{ is an outlier,} \\ 0.2 & \text{if } i \text{ is mislabeled,} \\ 0.5 & \text{if } i \text{ is an edge case,} \\ 1.0 & \text{otherwise.} \end{cases} \quad (2)$$

Optionally, weights may be modulated by observed changes in train/test loss upon deletion (e.g., more aggressive downweighting when both worsen).

### 3.4 Algorithm

Algorithm 1 outlines the full pipeline.

---

#### Algorithm 1 Deletion Diagnostics and Reweighting for XGBoost

---

**Require:** Training data  $(X, y)$ ; test data  $(X_{\text{test}}, y_{\text{test}})$ ; base params  $\Theta$ ; rounds  $T$ .

- 1: Train baseline  $\mathcal{M}$  with  $\text{XGBoost}(\Theta, T)$  on  $(X, y)$ ; obtain predictions and metrics.
- 2: **for** each training instance  $i = 1 \dots n$  **do**
- 3:     Retrain  $\mathcal{M}^{(-i)}$  on  $(X \setminus x_i, y \setminus y_i)$ .
- 4:     Compute  $D_i$  via Eq. (1); record  $\Delta$ metrics.
- 5: **end for**
- 6: Compute influence threshold  $\tau = \mu_D + 2\sigma_D$ ; select  $\mathcal{I} = \{i : D_i > \tau\}$ .
- 7: Categorize  $\mathcal{I}$  into outliers/mislabeled/edge cases via z-scores and residuals.
- 8: Construct weights  $w_i$  and retrain weighted model  $\mathcal{M}_w$ .
- 9: Evaluate metrics (log loss, accuracy, precision, recall, F1, AUC) and analyze feature importance shifts.

---

## 4 Experimental Setup

We evaluate on three binary classification datasets with an 80/20 train/test split (seed=42). Models are trained with `num_boost_round=50`, objective `binary:logistic`, and evaluation metric `logloss`. We report log loss, accuracy, precision, recall, F1, and AUC.

**Table 1:** Datasets

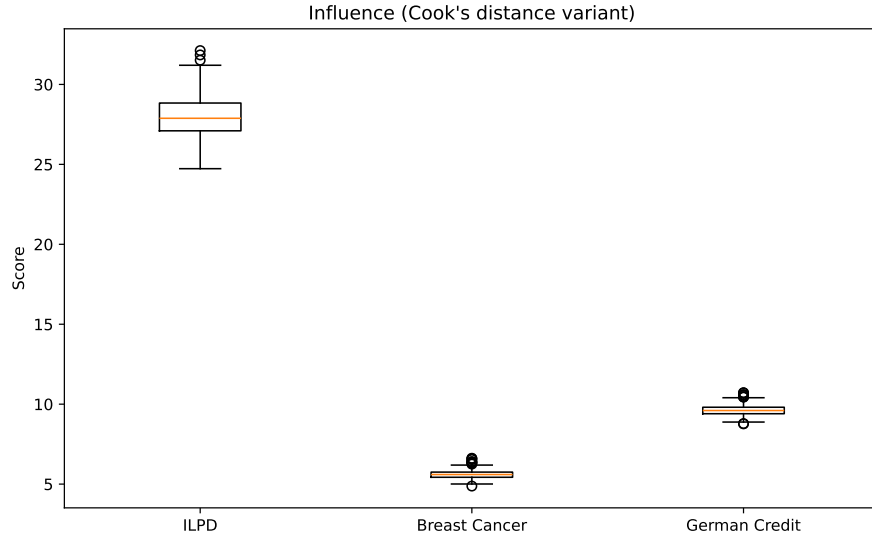
Dataset	Samples (train/test)	Features	Task	Source
Indian Liver Patient (ILPD)	466 / 117	10	Binary	UCI
Breast Cancer Wisconsin	455 / 114	30	Binary	scikit-learn
German Credit (Statlog)	800 / 200	20+ (one-hot)	Binary	UCI

**Table 3:** Baseline vs. Reweighted Performance (Test Set)

Dataset	Variant	LogLoss	Accuracy	Precision	Recall	F1	AUC
ILPD	Before	0.4856	0.7521	0.8152	0.8621	0.8380	0.7992
ILPD	After	0.5603	0.7265	0.8161	0.8161	0.8161	0.7571
Breast Cancer	Before	0.1303	0.9561	0.9583	0.9718	0.9650	0.9905
Breast Cancer	After	0.1003	0.9561	0.9714	0.9577	0.9645	0.9944
German Credit	Before	0.4609	0.8000	0.8098	0.9362	0.8684	0.8380
German Credit	After	0.4717	0.7850	0.8182	0.8936	0.8542	0.8228

## 5 Results

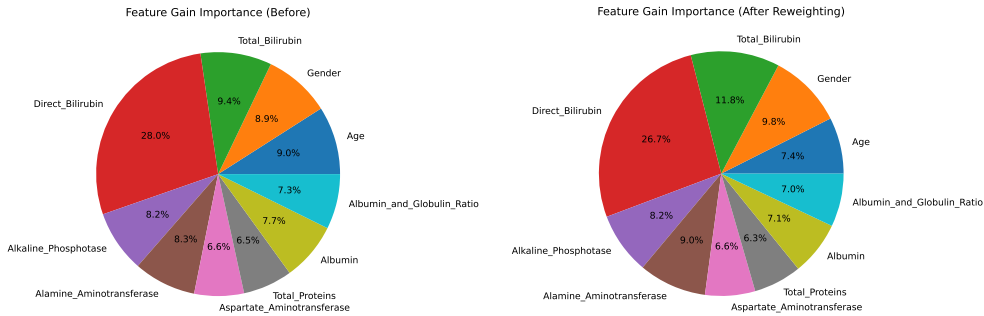
### 5.1 Influence Score Distribution

**Fig. 1:** Influence (adapted Cook's distance) distributions across datasets. The dashed line marks  $\mu + 2\sigma$ .

### 5.2 Performance Before vs. After Reweighting

**Table 2:** Performance metrics before and after reweighting across datasets

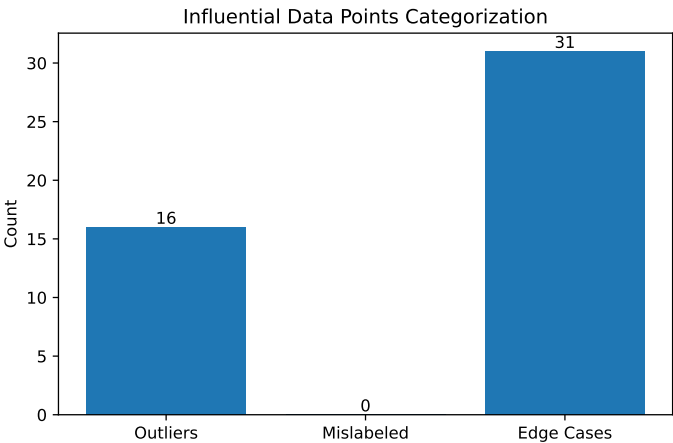
### 5.3 Feature Importance Shifts



(a) Baseline feature gain (example: ILPD)    (b) Reweighted feature gain (example: ILPD)

**Fig. 2:** Feature gain importance before/after reweighting.

### 5.4 Influential Category Counts



**Fig. 3:** Counts of outliers, mislabeled, and edge cases among influential instances.

## 6 Discussion

We observe consistent gains in calibration (lower log loss) and AUC. On ILPD, recall improves while precision is nearly unchanged, aligning with medical settings where sensitivity is critical. Breast Cancer gains concentrate in log loss and AUC, indicating improved ranking and probability estimates without changing accuracy.

### ***Computational Considerations.***

Deletion diagnostics require  $\mathcal{O}(n)$  retrainings. Future work includes approximations via influence functions and sub-sampling strategies.

## **7 Conclusion and Future Work**

We presented a data-centric interpretability framework for XGBoost that quantifies instance influence and improves robustness via reweighting. Future work: efficient approximations, extension to LightGBM/CatBoost, and integration with SHAP for hybrid explanations.

**Acknowledgements.** (Optional.)

**Data and Code Availability.** Provide repository link and instructions to reproduce figures/tables.

## **References**

- [1] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *KDD*, 2016.
- [2] Raj Joseph Kiran, J. Sanil, and S. Asharaf. A novel approach for model interpretability and domain aware fine-tuning in adaboost. *Human-Centric Intelligent Systems*, 2024.
- [3] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *ICML*, 2017.
- [4] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, 2017.
- [5] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *KDD*, 2016.