

# ENHANCING PREDICTION PERFORMANCE THROUGH INFLUENCE MEASURE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In the field of machine learning, the pursuit of accurate models is ongoing. A key aspect of improving prediction performance lies in identifying which data points in the training set should be excluded and which high-quality, potentially unlabeled data points outside the training set should be incorporated to improve the model’s performance on unseen data. To accomplish this, an effective metric is needed to evaluate the contribution of each data point toward enhancing overall model performance. This paper proposes the use of an influence measure as a metric to assess the impact of training data on test set performance. Additionally, we introduce a data selection method to optimize the training set as well as a dynamic active learning algorithm driven by the influence measure. The effectiveness of these methods is demonstrated through extensive simulations and real-world datasets.

## 1 INTRODUCTION

To build effective machine learning models, the significance of individual training data points cannot be overstated. Each data point in the training set contributes uniquely to the model’s learning process, shaping its performance, generalization, and resilience to various challenges (Blum & Langley, 1997). Simply evaluating the model’s performance on the provided data is insufficient; understanding the influence of individual training examples and making informed decisions about their inclusion or exclusion is critical for developing effective and reliable models.

Recent advancements in machine learning have highlighted the importance of strategic data selection and management during training. Techniques such as active learning (Settles, 2009) promote the iterative selection of the most uncertain unlabeled data points for labeling and inclusion in the training set. However, estimating uncertainty in deep neural networks (DNNs) is challenging due to their tendency to exhibit overconfidence (Ren et al., 2021). To address this, methods like deep Bayesian approaches (Gal et al., 2017), query-by-committee (Gorriz et al., 2017), Variational Auto-Encoders (Sinha et al., 2019), adversarial learning (Ducoffe & Precioso, 2018; Mayer & Timofte, 2020), graph convolutional networks (Caramalau et al., 2021), and noise stability (Li et al., 2024) have been proposed to improve uncertainty estimates. However, these approaches assume a well-trained model and fail to consider how model parameters might evolve when the training data is modified.

To bridge this gap, recent studies have focused on quantifying the impact of individual training examples on model behavior, with the main challenge being the identification of an appropriate evaluation metric. The Shapley value has emerged as a promising solution, inspiring a number of Shapley-value-based approaches (Ghorbani & Zou, 2019; Jia et al., 2019b;a; Ghorbani et al., 2020; Kwon & Zou, 2022; Wang & Jia, 2023). However, these methods often require multiple model retrainings and evaluations, making them computationally expensive.

Techniques such as influence functions (Koh & Liang, 2017; Pruthi et al., 2020; Yeh et al., 2018; Chen et al., 2021) offer insights into the effect of individual data points on model predictions, helping to identify and mitigate harmful or overly influential examples. For instance, Chhabra et al. (2024) apply influence functions to measure the impact of training data on a specific model, improving performance by pruning detrimental data points. Their method is highly efficient as it avoids the need for model retraining. However, these algorithms may still fail in certain cases, as demonstrated in the following simple example.

**Binary Classification via Logistic Regression.** We consider a logistic regression model  $p(x) = (1 + e^{-x\beta})^{-1}$ , where  $\beta \in \mathbb{R}$  represents the coefficients. For a training set  $\mathcal{Z} = \{(x_i, y_i)\}_{i=1}^n$  and a validation set  $\mathcal{V}$ , the influence function in Chhabra et al. (2024) for example is defined as

$$\mathcal{I}(-x_i) = \sum_{(x,y) \in \mathcal{V}} \partial_\beta L(y, x; \hat{\beta}) \left[ \sum_{i=1}^n \partial_\beta^2 L(y_i, x_i; \hat{\beta}) \right]^{-1} \partial_\beta L(y_i, x_i; \hat{\beta}),$$

where  $\partial_\beta = \partial/\partial\beta$ ,  $x \in \mathbb{R}$ ,  $y \in \{0, 1\}$  is a binary classification label,  $L$  represents the cross-entropy loss, and  $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}} n^{-1} \sum_{i=1}^n L(y_i, x_i; \beta)$ . According to Chhabra et al. (2024), samples  $x_i$  with negative influence function values negatively impact the model’s performance on the validation set and should be removed. For all training samples  $x_i$ , the term  $\sum_{(x,y) \in \mathcal{V}} \partial_\beta L(y, x; \hat{\beta}) \left[ \sum_{i=1}^n \partial_\beta^2 L(y_i, x_i; \hat{\beta}) \right]^{-1}$  remains constant, with only  $\partial_\beta L(y_i, x_i; \hat{\beta})$  varying.

In the simple logistic regression scenario, we have  $\partial_\beta L(y_i, x_i; \hat{\beta}) = [(1 + e^{-x_i \hat{\beta}})^{-1} - y_i] x_i$ . Assuming  $\sum_{(x,y) \in \mathcal{V}} \partial_\beta L(y, x; \hat{\beta}) \left[ \sum_{i=1}^n \partial_\beta^2 L(y_i, x_i; \hat{\beta}) \right]^{-1} > 0$  and  $x_i > 0$ , the training sample influence is negative if and only if  $y_i = 1$ . However, it is clearly incorrect to solely remove data points from one class, as outliers in the other class may also negatively influence the model’s performance. For higher-dimensional  $x_i$ , as shown in Section 6 using both simulated and real-world data, our method consistently outperforms the approach proposed by Chhabra et al. (2024) and others.

Given the limitations of existing methods, we draw inspiration from local influence measures developed in the statistical community (Zhu et al., 2007; 2011; Shu & Zhu, 2019; Sui et al., 2023) to propose a novel metric in this article. Our approach differs from previous ones that typically evaluate the impact of perturbations on data samples or model parameters with a fixed model assumption. Instead, our measure directly assesses the effect of minor perturbations to training samples on the model’s performance on a validation set, allowing the model to adapt to these changes. To implement this, we have developed a new perturbation manifold and expanded the local influence framework. Using this innovative approach, we introduce two key metrics: one for data trimming, aimed at identifying and removing training set anomalies that compromise model stability on the test set, and another for active learning, which focuses on selecting the most impactful unlabeled data to enhance the prediction performance. Moreover, acknowledging the challenges of slow computation and high memory usage inherent in calculating exact local influence measures, we propose two approximation methods. Our experiments on real-world datasets demonstrate that these approximations achieve performance comparable to exact calculations while significantly reducing computational overhead.

We summarize our contributions as follows,

- (i) Unlike existing local influence measures, we propose a new metric that evaluates the impact of perturbations to training samples on the model’s performance on validation data.
- (ii) The proposed metric is applicable to both data trimming and active learning. For data trimming, it evaluates the effects of minor perturbations to each sample, offering deeper insights into how individual samples impact model performance. In the context of active learning, the method captures the relationship between training samples, unlabeled data, and parameter updates.
- (iii) We propose two approximation methods to alleviate the high computational cost of calculating local influence measures. These algorithms significantly reduce computational overhead while maintaining better performance than other methods, as demonstrated in our experiments.

## 2 A NEW INFLUENCE MEASURE

In this section, we begin by presenting essential background information, including the Perturbation Manifold, before formally defining the proposed metric.

**Perturbation Manifold.** Our definition of the perturbation manifold closely follows that of Shu & Zhu (2019). Given an input sample  $z = (x, y)$  in the training set  $\mathcal{Z} = \{z_i\}_{i=1}^n$  and a machine learning model with an estimated parameter vector  $\hat{\theta}$ , which is trained on  $\mathcal{Z}$ , the prediction probability for class  $c \in \{1, \dots, K\}$  is denoted as  $P(c|x, \hat{\theta})$ . Let  $\omega = (\omega_1, \dots, \omega_p)^\top$  be a perturbation vector that varies within an open subset  $\Omega \subset \mathbb{R}^p$ . The perturbation  $\omega$  is applied to  $x$ , thereby affecting the learning of the parameter vector  $\hat{\theta}$ . We denote the parameter vector obtained by the model after

108 perturbing the training sample  $\mathbf{x}$  with  $\boldsymbol{\omega}$  as  $\hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega})$  with  $\hat{\boldsymbol{\theta}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}$ . We define  $P(c|\mathbf{x} + \boldsymbol{\omega}, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}))$   
 109 as the prediction probability under the perturbation  $\boldsymbol{\omega}$  such that  $\sum_{c=1}^K P(c|\mathbf{x} + \boldsymbol{\omega}, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega})) = 1$ . It  
 110 is assumed that there exists a  $\boldsymbol{\omega}_0 \in \Omega$  such that  $P(c|\mathbf{x} + \boldsymbol{\omega}_0, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)) = P(c|\mathbf{x}, \hat{\boldsymbol{\theta}})$ . Additionally,  
 111 we assume that  $\{P(c|\mathbf{x} + \boldsymbol{\omega}, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}))\}_{c=1}^K$  is positive and sufficiently smooth for all  $\boldsymbol{\omega} \in \Omega$ .  
 112

113 Following the development in (Zhu et al., 2007; 2011), we define  $\mathcal{M} = \{P(c|\mathbf{x} + \boldsymbol{\omega}, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega})) : \boldsymbol{\omega} \in \Omega\}$   
 114 as a perturbation manifold. The tangent space of  $\mathcal{M}$  at  $\boldsymbol{\omega}$  is denoted by  $T_{\boldsymbol{\omega}}$ , which is spanned  
 115 by  $\{\partial l(\boldsymbol{\omega}|c, \mathbf{x}, \hat{\boldsymbol{\theta}}(\mathbf{x})) / \partial \omega_i\}_{i=1}^p$ , where  $l(\boldsymbol{\omega}|c, \mathbf{x}, \hat{\boldsymbol{\theta}}(\mathbf{x})) = \log P(c|\mathbf{x} + \boldsymbol{\omega}, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}))$ . Let  $\mathbf{G}_z(\boldsymbol{\omega}) =$   
 116  $\sum_{c=1}^K \partial_{\boldsymbol{\omega}}^{\top} l(\boldsymbol{\omega}|c, \mathbf{x}, \hat{\boldsymbol{\theta}}(\mathbf{x})) \partial_{\boldsymbol{\omega}} l(\boldsymbol{\omega}|c, \mathbf{x}, \hat{\boldsymbol{\theta}}(\mathbf{x})) P(c|\mathbf{x} + \boldsymbol{\omega}, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}))$  with  $\partial_{\boldsymbol{\omega}} = (\partial / \partial \omega_1, \dots, \partial / \partial \omega_p)$ .  
 117 If  $\mathbf{G}_z(\boldsymbol{\omega})$  is positive definite, then  $\mathcal{M}$  is a Riemannian manifold (Shu & Zhu, 2019) with  $\mathbf{G}_z(\boldsymbol{\omega})$   
 118 serving as the Riemannian metric tensor (Amari, 2012; Amari & Nagaoka, 2000).  
 119

120 Although  $\mathbf{G}_z(\boldsymbol{\omega})$  is often not positive definite in classification problems, we can still reduce the  
 121 dimensionality of the perturbations and reconstruct a Riemannian manifold (Shu & Zhu, 2019).

122 **The Influence Measure.** Let  $L(y', \mathbf{x}'; \boldsymbol{\theta})$  denote the loss function of the model with parameter  $\boldsymbol{\theta}$  on  
 123  $\mathbf{z}' = (\mathbf{x}', y') \notin \mathcal{Z}$ , we can get the expression for the (first-order) influence measure:  
 124

$$125 \text{FI}(\mathbf{z}', \mathbf{z}) = \partial_{\boldsymbol{\omega}} L(y', \mathbf{x}'; \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)) \mathbf{G}_z^{\dagger}(\boldsymbol{\omega}_0) \partial_{\boldsymbol{\omega}}^{\top} L(y', \mathbf{x}'; \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)), \quad (2.1)$$

126 where  $\mathbf{G}_z^{\dagger}(\boldsymbol{\omega}_0)$  is the pseudoinverse of  $\mathbf{G}_z(\boldsymbol{\omega}_0)$ .  
 127

128 We consider a linear perturbation approach where applying perturbation  $\boldsymbol{\omega}$  to the training sample  
 129  $\mathbf{z}$  transforms  $(\mathbf{x}, y)$  into  $(\mathbf{x} + \boldsymbol{\omega}, y)$ . Consequently, the parameter vector updates to  $\hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}) :=$   
 130  $n^{-1} \arg \min_{\boldsymbol{\theta} \in \Theta} \{\sum_{i=1}^n L(y_i, \mathbf{x}_i; \boldsymbol{\theta}) + L(y, \mathbf{x} + \boldsymbol{\omega}; \boldsymbol{\theta}) - L(y, \mathbf{x}; \boldsymbol{\theta})\}$ , with  $\boldsymbol{\omega}_0 = \mathbf{0}$ . Using the chain  
 131 rule, we can derive the expression for  $\partial_{\boldsymbol{\omega}} L(y', \mathbf{x}'; \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0))$ :  
 132

$$133 \partial_{\boldsymbol{\omega}} L(y', \mathbf{x}'; \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)) = \partial_{\boldsymbol{\theta}} L(y', \mathbf{x}'; \hat{\boldsymbol{\theta}}(\mathbf{x})) \partial_{\boldsymbol{\omega}} \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega})|_{\boldsymbol{\omega}=\boldsymbol{\omega}_0}. \quad (2.2)$$

134 Here,  $\partial_{\boldsymbol{\omega}} \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega})|_{\boldsymbol{\omega}=\boldsymbol{\omega}_0} \approx n^{-1} \mathbf{H}_{\hat{\boldsymbol{\theta}}}^{-1} \partial_{\mathbf{x}} \partial_{\boldsymbol{\theta}} L(y, \mathbf{x}; \hat{\boldsymbol{\theta}}(\mathbf{x}))$ , whose derivation is detailed in Appendix  
 135 A), where  $\mathbf{H}_{\hat{\boldsymbol{\theta}}} := n^{-1} \sum_{i=1}^n \partial_{\boldsymbol{\theta}}^2 L(y_i, \mathbf{x}_i; \hat{\boldsymbol{\theta}}(\mathbf{x}))$ . The term  $\partial_{\mathbf{x}} \partial_{\boldsymbol{\theta}} L(y, \mathbf{x}; \hat{\boldsymbol{\theta}}(\mathbf{x}))$  represents the gradient  
 136 of the loss function  $L$  first taken with respect to the model parameters  $\boldsymbol{\theta}$  and then with respect to  $\mathbf{x}$ ,  
 137 evaluated at the perturbed training sample  $(\mathbf{x}, y)$ .  
 138

139 For the computation of  $\mathbf{G}_z(\boldsymbol{\omega}_0)$ ,  $P(c|\mathbf{x} + \boldsymbol{\omega}_0, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0))$  can be directly obtained using the learned  
 140 model parameter vector  $\hat{\boldsymbol{\theta}}$  and the unperturbed sample point  $\mathbf{x}$ . The calculation of  $\partial_{\boldsymbol{\omega}} l(\boldsymbol{\omega}_0|c, \mathbf{x}, \hat{\boldsymbol{\theta}}(\mathbf{x}))$   
 141 requires the application of the chain rule:  
 142

$$143 \begin{aligned} 144 \partial_{\boldsymbol{\omega}} l(\boldsymbol{\omega}_0|c, \mathbf{x}, \hat{\boldsymbol{\theta}}(\mathbf{x})) &= \partial_{\boldsymbol{\omega}} \log P(c|\mathbf{x} + \boldsymbol{\omega}_0, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)) \\ 145 &= \partial_{\boldsymbol{\theta}} \log P(c|\mathbf{x} + \boldsymbol{\omega}_0, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)) \cdot \partial_{\boldsymbol{\omega}} \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega})|_{\boldsymbol{\omega}=\boldsymbol{\omega}_0} \\ 146 &\quad + \partial_{\mathbf{x}} \log P(c|\mathbf{x} + \boldsymbol{\omega}_0, \hat{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\omega}_0)). \end{aligned} \quad (2.3)$$

147 All differentiation operations can be easily computed using backpropagation (Goodfellow et al., 2016)  
 148 in deep learning libraries such as TensorFlow (Abadi et al., 2016) and PyTorch (Paszke et al., 2017).  
 149 This entire process is efficient and does not require retraining the model parameters.  
 150

151 **Theorem 1.** If  $\varphi$  represents a diffeomorphism of  $\boldsymbol{\omega}$ , then  $\text{FI}(\mathbf{z}', \mathbf{z})$  is invariant under any reparamete-  
 152 rization associated with  $\varphi$ .  
 153

154 Compared to widely used measures in Euclidean spaces, such as the Jacobian norm (Novak et al.,  
 155 2018) and Cook’s local influence measure (Cook, 1986), Theorem 1 demonstrates that  $\text{FI}(\mathbf{z}', \mathbf{z})$   
 156 remains invariant under any diffeomorphic transformation (e.g., scaling) of the perturbation vector  $\boldsymbol{\omega}$ .  
 157 The proof of Theorem 1 can be found in Shu & Zhu (2019).  
 158

159 The significance of Theorem 1 is especially pronounced when there are scale differences among the  
 160 dimensions of  $\mathbf{x}$ . For instance, if certain dimensions have significantly larger values than others, the  
 161 contribution of perturbations to those dimensions may appear exaggerated. However, our  $\text{FI}(\mathbf{z}', \mathbf{z})$   
 mitigates this scaling issue by employing the metric tensor of the perturbation manifold instead of  
 that of the standard Euclidean space.

### 3 FI FOR DATA TRIMMING.

The primary goal of data trimming is to eliminate training samples that may compromise the model’s performance on datasets beyond the training set. Since our proposed influence measure (FI) quantifies the impact of each training sample on the model’s performance on validation sets, it serves as a natural tool for data trimming.

From a model robustness perspective, if a small perturbation in a training sample leads to a significant effect on the model’s performance on the validation set, that sample should be excluded, which aligns with the principle of our proposed metric. Given the inherent challenge of ensuring that all samples in the training set are entirely accurate, a sample with excessive influence could severely degrade the model’s overall performance if it contains any contamination. Therefore, to enhance prediction performance on an unseen data set, such samples should be removed from the training set.

Following the setup of Chhabra et al. (2024), we introduce a training set  $\mathcal{Z}$ , a validation set  $\mathcal{V}$ , and a base model  $\mathcal{F}$ . In this context, we assess the impact of each training sample on the model’s performance by computing the FI for each training tuple  $z_i \in \mathcal{Z}$  with respect to  $\mathcal{V}$ . To achieve this, we extend Equation 2.1 to encompass the entire validation set. This involves replacing the loss function for an individual validation sample with the mean loss function across the entire validation set, as follows:

$$\text{FI}^{util}(z) = \partial_{\omega} L(\mathcal{V}; \hat{\theta}(\mathbf{x} + \omega_0)) \mathbf{G}_{z_i}^{\dagger}(\omega_0) \partial_{\omega}^{\top} L(\mathcal{V}; \hat{\theta}(\mathbf{x} + \omega_0)), \quad (3.1)$$

where  $\partial_{\omega} L(\mathcal{V}; \hat{\theta}(\mathbf{x} + \omega_0)) := \frac{1}{|\mathcal{V}|} \sum_{(\mathbf{x}', y') \in \mathcal{V}} \partial_{\omega} L(y', \mathbf{x}'; \hat{\theta}(\mathbf{x} + \omega_0))$ .

The algorithm for computing  $\text{FI}^{util}$  is outlined in Algorithm 1. After computing these values, we can sort all data points in the training set in descending order based on their  $\text{FI}^{util}$  values and remove the top  $b$  points to enhance the model’s performance on the test set. For the detailed data trimming algorithm, please refer to Algorithm 3.

---

#### Algorithm 1 Calculation of $\text{FI}^{util}$

---

**Input:** Training set  $\mathcal{Z}$ , Validation set  $\mathcal{V}$ , Base model  $\mathcal{F}$

**Output:** Influence measure vector  $\text{FI}^{util} \in \mathbb{R}^{|\mathcal{Z}| \times 1}$

```

1: procedure  $\text{FI}^{util}$ -CALCULATION( $\mathcal{Z}, \mathcal{V}, \mathcal{F}$ )
2:   Train  $\mathcal{F}$  with  $\mathcal{Z}$ , and obtain the parameter vector  $\hat{\theta}$ 
3:   Generate an empty vector  $\text{FI}^{util}$  of size  $|\mathcal{Z}| \times 1$ 
4:   Calculate  $\frac{1}{|\mathcal{V}|} \sum_{(\mathbf{x}', y') \in \mathcal{V}} \partial_{\theta} L(y', \mathbf{x}'; \hat{\theta})$  and  $\mathbf{H}_{\hat{\theta}}$ 
5:   for every  $z_i$  in  $\mathcal{Z}$  do
6:     Calculate  $\mathbf{G}_{z_i}(\omega_0)$  and  $\partial_{\mathbf{x}} \partial_{\theta} L(y_i, \mathbf{x}_i; \hat{\theta})$ 
7:      $\partial_{\omega} L(\mathcal{V}; \hat{\theta}) \leftarrow \frac{1}{|\mathcal{V}|} \sum_{(\mathbf{x}', y') \in \mathcal{V}} \partial_{\theta} L(y', \mathbf{x}'; \hat{\theta}) \mathbf{H}_{\hat{\theta}}^{-1} \partial_{\mathbf{x}} \partial_{\theta} L(y_i, \mathbf{x}_i; \hat{\theta})$ 
8:      $\text{FI}^{util}[i] \leftarrow \partial_{\omega} L(\mathcal{V}; \hat{\theta}) \mathbf{G}_{z_i}^{\dagger}(\omega_0) \partial_{\omega}^{\top} L(\mathcal{V}; \hat{\theta})$ 
9:   end for
10:  return  $\text{FI}^{util}$ 
11: end procedure

```

---

### 4 FI FOR ACTIVE LEARNING.

In active learning, the primary objective is to identify the most uncertain or informative samples from an unlabeled pool for annotation, typically in sequential batches. After each round of annotation, the newly labeled data are combined with the existing labeled set to retrain the model and improve its performance. Li et al. (2024) argue that if a small perturbation to the model parameters leads to significant changes in the prediction for a given sample, this indicates high uncertainty for that sample under the current model, suggesting that it should be labeled and added to the training set. Our approach is more fundamental: since the model parameters are derived from the training data, we directly assess how perturbations to the training samples influence the model’s predictions for the unlabeled samples. If slight perturbations to most training samples substantially alter the model’s

prediction for a given sample, it likely contains missing information from the training set and should therefore be included.

We now introduce a method for active learning using the proposed FI. For an unlabeled sample  $\mathbf{x}_{\text{unlabel}}$ , we first assign it a predicted label and treat it as a validation sample. Next, we calculate the influence measure of  $\mathbf{x}_{\text{unlabel}}$  with respect to each point in the training set. The overall influence measure  $\text{FI}^{\text{active}}(\mathbf{x}_{\text{unlabel}})$  is derived by aggregating these individual measures either by averaging or using specific quantiles of these FI values. The  $\text{FI}^{\text{active}}(\mathbf{x}_{\text{unlabel}})$  is defined as follows:

$$\text{FI}^{\text{active}}(\mathbf{x}_{\text{unlabel}}) = g(\{\text{FI}(\mathbf{z}_{\text{unlabel}}, \mathbf{z}_i)\}_{i=1}^n), \quad (4.1)$$

where  $g$  represents the aggregation function,  $\mathbf{z}_{\text{unlabel}} = (\mathbf{x}_{\text{unlabel}}, y_{\text{pred}})$ , and  $y_{\text{pred}}$  is the predicted label assigned by the current model for  $\mathbf{x}_{\text{unlabel}}$ .

During each round of active learning, we begin by applying Algorithm 2 to compute  $\text{FI}^{\text{active}}$ . We then sort all the samples in the unlabeled pool in descending order based on their  $\text{FI}^{\text{active}}$  values. The top-ranked samples are labeled and added to the training set. In the subsequent round, the model is retrained,  $\text{FI}^{\text{active}}$  is recalculated, and the process is repeated. For the detailed active learning algorithm, please refer to Algorithm 4.

---

#### Algorithm 2 Calculation of $\text{FI}^{\text{active}}$

---

**Input:** Labeled pool of training data  $\mathcal{L}$ , Unlabeled pool of training data  $\mathcal{U}$ , Base model  $\mathcal{F}$ , Aggregation function  $g$

**Output:** Influence measure vector  $\text{FI}^{\text{active}} \in \mathbb{R}^{|\mathcal{U}| \times 1}$

```

1: procedure  $\text{FI}^{\text{active}}$ -CALCULATION( $\mathcal{L}, \mathcal{U}, \mathcal{F}, g$ )
2:   Train  $\mathcal{F}$  with  $\mathcal{L}$ , and obtain the parameter vector  $\hat{\theta}$ 
3:   Generate an empty vector  $\text{FI}^{\text{active}}$  of size  $|\mathcal{U}| \times 1$ 
4:   Calculate  $\mathbf{H}_{\hat{\theta}}$ 
5:   for every  $\mathbf{x}_i^{\mathcal{U}}$  in  $\mathcal{U}$  do
6:     Obtain an estimated label  $\hat{y}_i$  with  $\mathcal{F}$ 
7:     Calculate  $\partial_{\theta} L(\hat{y}_i, \mathbf{x}_i^{\mathcal{U}}; \hat{\theta})$ 
8:      $\mathcal{J} \leftarrow \emptyset$ 
9:     for every  $\mathbf{z}_j^{\mathcal{L}}$  in  $\mathcal{L}$  do
10:      Calculate  $\mathbf{G}_{\mathbf{z}_j^{\mathcal{L}}}(\omega_0)$  and  $\partial_{\mathbf{x}} \partial_{\theta} L(y_j^{\mathcal{L}}, \mathbf{x}_j^{\mathcal{L}}; \hat{\theta})$ 
11:       $\partial_{\omega} L(\hat{y}_i, \mathbf{x}_i^{\mathcal{U}}; \hat{\theta}) \leftarrow \partial_{\theta} L(\hat{y}_i, \mathbf{x}_i^{\mathcal{U}}; \hat{\theta}) \mathbf{H}_{\hat{\theta}}^{-1} \partial_{\mathbf{x}} \partial_{\theta} L(y_j^{\mathcal{L}}, \mathbf{x}_j^{\mathcal{L}}; \hat{\theta})$ 
12:       $\mathcal{J} \leftarrow \mathcal{J} \cup \{\partial_{\omega} L(\hat{y}_i, \mathbf{x}_i^{\mathcal{U}}; \hat{\theta}) \mathbf{G}_{\mathbf{z}_j^{\mathcal{L}}}^{\dagger}(\omega_0) \partial_{\omega}^{\top} L(\hat{y}_i, \mathbf{x}_i^{\mathcal{U}}; \hat{\theta})\}$ 
13:     end for
14:      $\text{FI}^{\text{active}}[i] \leftarrow g(\mathcal{J})$ 
15:   end for
16:   return  $\text{FI}^{\text{active}}$ 
17: end procedure

```

---

## 5 APPROXIMATION METHODS

The data selection processes described in the previous two sections have two main drawbacks. First, computing FI requires substantial storage for second-order derivatives, particularly in models with numerous parameters or high-dimensional data. To address this issue, we propose the KFSVD approximation method, which combines K-FAC and Truncated-SVD to effectively reduce storage requirements. Second, the necessity to compute the  $\text{FI}^{\text{util}}$  for all data points in the data trimming algorithm, as well as the need to recalculate  $\text{FI}^{\text{active}}$  for all unlabeled data in each round of the active learning algorithm, significantly decreases computational efficiency. To mitigate this, we implement a subsampling approximation that enhances overall performance. The details of these two approximation methods are as follows.

**KFSVD approximation**, which combines the Kronecker-factored (K-FAC) approximation (Martens & Grosse, 2020; Nickl et al., 2023) with Truncated Singular Value Decomposition (Truncated-SVD)

approximation (Golub & Reinsch, 1971) to mitigate memory consumption associated with the Hessian matrix  $\mathbf{H}_{\hat{\theta}}$  and the second-order partial derivatives  $\partial_{\mathbf{x}}\partial_{\theta}L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))$ . The K-FAC algorithm is a widely recognized technique for approximating the Hessian matrix, which not only accelerates computations but also significantly reduces storage requirements. Meanwhile, the Truncated-SVD approximation employs power iteration and related techniques to compute the top- $k$  eigenvalues and their corresponding eigenvectors of  $\partial_{\mathbf{x}}\partial_{\theta}L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))$ , thereby providing an effective approximation of these second-order partial derivatives. Assuming the dimensionality of the sample covariates is  $d$  and the number of model parameters is  $p$ , the Truncated-SVD approximation enables the decomposition of the second-order partial derivatives as  $\partial_{\mathbf{x}}\partial_{\theta}L(y, \mathbf{x}; \hat{\theta}(\mathbf{x})) = U_{p \times k} \Lambda_{k \times k} V_{d \times k}^{\top}$ , where  $\Lambda_{k \times k}$  is a diagonal matrix containing the top- $k$  eigenvalues of  $\partial_{\mathbf{x}}\partial_{\theta}L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))$  and  $U_{p \times k}$  and  $V_{d \times k}$  are comprised of  $k$  orthogonal vectors. Algorithm 5 and 7 provide the detailed procedures for calculating  $\text{FI}^{\text{util}}$  and  $\text{FI}^{\text{active}}$  using the KFSVD approximation.

**Subsampling approximation**, which utilizes subsampling and random forest techniques to enhance computational efficiency. Specifically, we extract a small subset of samples (e.g., 20%) and compute their FI using Algorithm 5 or 7. These computed values are then used to sort the samples. The features of each sample, along with their corresponding ranks, serve as covariates and target variables to create a new dataset. Subsequently, we train a regression model using random forests on this dataset and leverage the trained model to predict the ranks of other samples. Selections are based on these predicted ranks. Since the computation of FI is the most time-consuming part of the workflow, the overall acceleration is directly correlated with the proportion of samples for which we choose to compute FI accurately. For instance, selecting 20% of the samples can reduce the total processing time to  $\frac{1}{5}$  of the original duration when the dataset is sufficiently large.

In practice, we integrate both approximation methods to develop the comprehensive  $\text{FI}^{\text{util}}$ -based data trimming algorithm and  $\text{FI}^{\text{active}}$ -based active learning algorithm. These algorithms effectively address the storage and computational efficiency challenges inherent in calculating FI. The detailed procedures are outlined in Algorithm 6 and 8.

To better illustrate the computational advantages of the proposed two approximation methods, we analyze their complexity in terms of both memory usage and computational speed. This analysis is also compared with the complexity of *Influence Value (IV)* as introduced by Chhabra et al. (2024).

**Time Complexity.** The following are the time complexities of FI and IV:

- $\text{FI}^{\text{util}}$  without approximation methods:  $\mathcal{O}(p^3 + nd^3 + ndp)$ . Here,  $n$  is the maximum value of the sample sizes of the training set and the validation set,  $d$  is the dimension of the covariates, and  $p$  is the dimension of the model parameters. The term  $\mathcal{O}(p^3)$  comes from computing the inverse of the Hessian matrix. The  $\mathcal{O}(nd^3)$  term is due to inverting  $G$ , repeated  $n$  times. The  $\mathcal{O}(ndp)$  term corresponds to computing  $\partial_{\mathbf{x}}\partial_{\theta}L$  for  $n$  times.
- $\text{FI}^{\text{util}}$  with approximation methods:  $\mathcal{O}(\alpha ndT \log(\alpha n) + \alpha n(p+d)ks + \alpha nd^3 + \sum_{i=1}^L p_i^3)$ . Here,  $\alpha \in (0, 1]$  is the proportion of data for which we compute FI accurately during subsampling.  $T$  is the number of decision trees in the random forest,  $k$  is the Truncated-SVD parameter,  $s$  is the number of iterations for computing each eigenvalue during power iteration, and  $\{p_1, p_2, \dots, p_L\}$  represent the number of parameters in each layer of an  $L$ -layer neural network.
- $\text{FI}^{\text{active}}$  without approximation methods:  $\mathcal{O}(p^3 + n^2d^3 + n^2dp)$ . Here,  $n$  is the maximum value of the sample sizes of the labeled dataset and the unlabeled dataset.
- $\text{FI}^{\text{active}}$  with approximation methods:  $\mathcal{O}(\alpha ndT \log(\alpha n) + \alpha n^2(p+d)ks + \alpha n^2d^3 + \sum_{i=1}^L p_i^3)$ .
- IV:  $\mathcal{O}(p^3 + np)$ .

By employing two approximation techniques, we can significantly reduce the computational gap between our method and the compared methods such as IV. When the dataset is large enough, setting  $\alpha n$  as a constant allows us to achieve better computational complexity than the IV method. For complex models (e.g., neural networks) where  $p$  is sufficiently large, the  $p^3$  term dominates the computational complexity of IV. In such scenarios, the computational costs of our method could be lower than IV.

**Memory Usage.** Below is a discussion of the space complexities associated with FI and IV:

- In the computation of FI, the highest storage requirements are for the Hessian matrix and  $\partial_x \partial_\theta L$ , which together occupy a space of  $\mathcal{O}(p^2 + dp)$ .
- After using the KFSVD approximation method, the space requirement of FI is reduced to  $\mathcal{O}(\sum_{i=1}^L p_i^2 + (p+d)k)$ .
- In the computation of IV, it is necessary to store sample information and the Hessian matrix, resulting in a space complexity of  $\mathcal{O}(p^2 + d)$ .

After applying approximation techniques, our algorithm achieves a space complexity similar to that of the IV. In fact, when  $p$  is large, our method offers an advantage.

## 6 EXPERIMENTAL RESULTS

In this section, we present experimental results that illustrate how our newly proposed metrics,  $FI^{util}$  and  $FI^{active}$ , contribute to enhancing model prediction performance. We compare our algorithms with state-of-the-art strategies in both data trimming and active learning scenarios, thereby validating the effectiveness of our approach on both simulated and real-world datasets.

### 6.1 DATA TRIMMING

In this subsection, we conduct simulations using both linear and nonlinear models to demonstrate how our algorithm enhances data trimming efficiency and evaluate the effectiveness of the proposed FI on real-world datasets. The latest data trimming method, IV, serves as the primary baseline for comparison in these experiments.

#### Validation on 2D Linear Model.

Logistic regression is utilized for this binary classification task. We begin by generating several datasets by sampling from two isotropic 2D Gaussian distributions. Each dataset comprises 150 training samples, 100 validation samples, and 600 test samples. The experimental settings for this

**Table 1: Comparison of two methods on linear model.** Number of cases where FI outperforms IV across 30 random seeds, along with performance improvements. *Acc\_FI*: the mean accuracy by FI, and *Acc\_IV*: by IV.

# of deleted points	# of better case	Acc_FI(%)	Acc_IV(%)
5	23	96.22±0.65	95.77±0.76
10	28	96.20±0.65	94.84±1.07
20	30	96.23±0.64	93.36±2.21

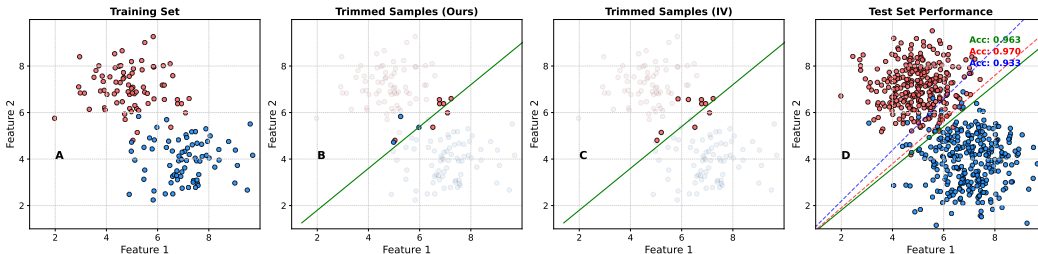
scenario are consistent with those of the study by Chhabra et al. (2024). To account for the randomness inherent in the sampling process, we analyze our method on datasets generated under the same distribution but with different random seeds. As shown in Table 1, our method consistently enhances model performance compared to theirs in most cases, particularly when trimming 5, 10, and 20 samples, with thirty different random seeds employed each time. Moreover, it is evident from Figure 1.C that IV tends to trim samples from a specific class under certain conditions. In this context, Figure 1.D clearly demonstrates that in some scenarios IV fails, while our method continues to perform effectively.

#### Validation on 2D Nonlinear Model.

After demonstrating the effectiveness of our method in linear scenarios, we now extend our examination to nonlinear cases. To achieve this, we construct a binary classification dataset that is non-linearly separable, with each class represented by a crescent-shaped region. A more intuitive understanding can be obtained from Figure 9. We employ a neural network that includes an input layer, two hidden layers with ReLU activation functions, and an output layer with a sigmoid activation function. Similar to the linear case, we conduct repeated experiments in the nonlinear scenario. Each dataset comprises 500 training samples, 250 validation samples, and 250 test samples. As shown in Table 2, our method achieves a higher average accuracy and outperforms the other methods in most cases across the 20 repetitions. Figure 9 illustrates instances where utilizing IV to remove training points can lead to a deterioration in model performance under certain conditions.

#### Validation on the Real-World Datasets.

Here, we evaluate the effectiveness of our proposed methods using four real-world datasets: two tabular datasets, *Adult* (Kohavi, 1996) and *Bank* (Moro



**Figure 1: Performance under Linear Model.** Different colored points represent different classes. **A** shows the training set. **B** and **C** respectively denote the samples to be trimmed by FI and IV. **D** denotes test set. **Green** line: boundary without trimming; **Red** line: boundary after FI trimming; and **Blue** line: boundary after IV trimming.

et al., 2014); a visual dataset, *CelebA* (Liu et al., 2015); and a textual dataset, *Jigsaw Toxicity* (Noever, 2018). Additional details regarding the datasets and experiments can be found in Appendices C & D.1.

Both  $FI^{util}$  and IV are evaluated on the validation set, with Logistic Regression serving as the base model. The results on the test sets of these datasets are presented in Figure 2. Ablation experiments are also conducted in the Appendix E.1 to assess the impact of perturbations on model performance.

**Table 2: Comparison of two methods on nonlinear model.** Number of cases where FI outperform IV across 20 random seeds, along with performance improvements.  $Acc_{FI}$ : the mean accuracy by FI, and  $Acc_{IV}$ : by IV.

# of deleted points	# of better case	Acc <sub>FI</sub> (%)	Acc <sub>IV</sub> (%)
5	17	89.90±2.16	87.50±2.46
10	17	90.24±2.05	88.02±2.62
20	15	90.32±1.67	87.80±2.54

As illustrated in Figure 2, our findings indicate that under a limited budget  $b$ , our  $FI^{util}$  data trimming method consistently outperforms two baseline models, IV and Random Trimming. Among the four datasets examined, Random Trimming demonstrates no improvement in model performance. Although IV exhibits a notable enhancement on *Adult* and *Bank*, it tends to remove important data points on *CelebA*, resulting in decreased performance, and shows no improvement on *Jigsaw Toxicity* compared to Random Trimming. This suggests that IV may fail in certain scenarios. In contrast, our method consistently achieves the maximum improvements across all datasets, particularly on the *Bank* dataset, where accuracy increases by more than 10%. To ensure a more rigorous comparison between FI and IV, we perform K-fold cross-validation on the *Adult* and *Bank* datasets (Appendix E.5). The results demonstrate that our method consistently outperforms IV, highlighting the robustness and stability of our approach.

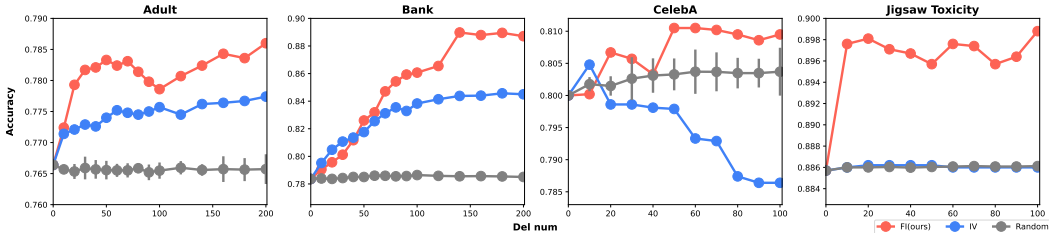
The experimental results above show that our method performs better on real-world datasets than on simulated ones. This is largely due to the simplicity of the simulated datasets, which are two-dimensional with clear boundaries effectively separating the classes and contain no erroneous data points. In such cases, a sufficiently large dataset allows the model to easily find the optimal boundary, minimizing the advantages of data trimming. In contrast, real-world datasets are typically high-dimensional and more complex, often containing errors. Here, the benefits of removing potentially erroneous high-influence points become more evident. To support this, we introduce noise into the real-world datasets and conduct further experiments; details of the noise method are in Appendix D.1, and results are shown in Figure 8. The findings clearly demonstrate that both FI and IV outperform Random Trimming, with our method retaining significant advantages over IV, further confirming its effectiveness.

## 6.2 ACTIVE LEARNING

In this section, we conduct numerical experiments on three simple tabular datasets and three complex image datasets to demonstrate the effectiveness of the proposed FI metric in active learning. We compare our method with two state-of-the-art active learning baselines including IV (Chhabra et al., 2024) and BALD (?Kirsch et al., 2019; 2023). Random Selection is also included as a baseline. All reported results are averaged over three runs to ensure a reliable evaluation. Besides, we present



432  
433  
434  
435  
436  
437  
438  
439



440  
441

**Figure 2:** Accuracy curves of three data trimming methods on test sets of *Adult*, *Bank*, *CelebA* and *Jigsaw Toxicity*.

442  
443

ablation experiments in Appendix E.1, to get more insights into the trade offs between the potential reduction in our method’s performance vs. the quantified computational cost reduction.

444  
445  
446  
447  
448  
449

Note that, IV, as proposed by Chhabra et al. (2024), originally calculates influence values only once during the initial selection, which can hinder overall performance since these values are not updated with subsequent labelings. To ensure fairness in our comparisons, we modified IV to recalculate influence values in each round of selection, and this updated version is used in our experiments.

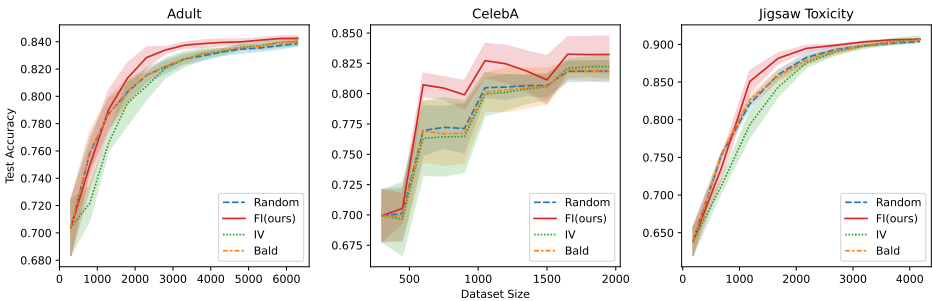
450  
451  
452  
453  
454

For our method, the aggregation function  $g(\{FI(z_{unlabel}, z_i)\}_{i=1}^n)$  in Equation 4.1 denotes the operation of calculating the mean of the data points from the set  $\{FI(z_{unlabel}, z_i)\}_{i=1}^n$  that fall between the 10th and 90th percentiles, sorted in descending order. In other words, we compute the mean of the middle 80% of the data after sorting. This approach effectively excludes extreme values, thereby enhancing the robustness of our final result.

455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465

**Tabular Datasets.** For the tabular datasets, we employed a Logistic Regression model. The total number of annotation rounds for each experiment, along with the settings for the unlabeled pool size and acquisition size, are detailed in Table 4. As shown in Figure 3, our method outperforms the other approaches on these simple tabular datasets. Random Selection performs reasonably well initially, but as data volume increases, improvements in model performance diminish. BALD shows comparable performance to Random Selection on tabular datasets. IV initially underperforms compared to Random Selection, but as more data is added, its performance improves and eventually becomes comparable to BALD. Our method starts similarly to BALD, but its focus on more challenging samples allows it to quickly address missing information once the model achieves a certain accuracy. As a result, in the latter stages of each graph, our method distinctly diverges from the others.

466  
467  
468  
469  
470  
471  
472  
473  
474  
475



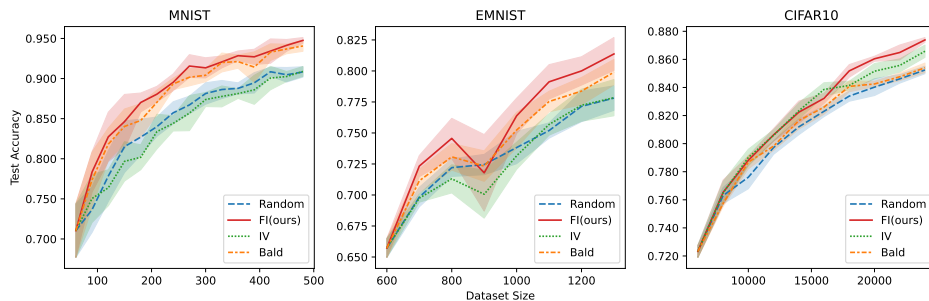
476  
477  
478

**Figure 3:** Classification performance due to the different active learning methods on *Adult*, *CelebA*, and *Jigsaw Toxicity*.

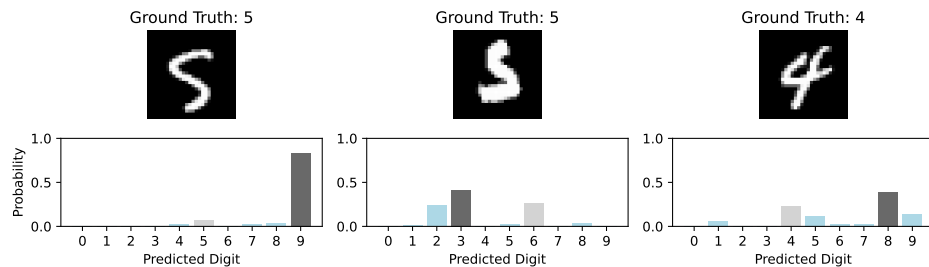
479  
480  
481  
482  
483  
484  
485

**Image Classification.** The data in the tabular datasets have been processed, resulting in high test accuracy with logistic regression. In contrast, for the image datasets—*MNIST* (Lecun et al., 1998), *EMNIST* (Cohen et al., 2017), and *CIFAR-10* (Krizhevsky & Hinton, 2009)—the input consists of raw images, and we employ a Convolutional Neural Network (CNN) as the classifier. To simplify the computation of  $FI^{active}$ , we utilize the outputs from the last layer of the neural network and focus on the parameters of that layer. However, during each training iteration, all parameters of the neural

network are updated, not just those of the final layer. The settings for the unlabeled pool size and acquisition size for each image dataset are provided in Table 5. Experimental results indicate that our method maintains comparable time consumption to other approaches, even with complex image data (see Table 9 for details). Our primary focus in active learning is enhancing model accuracy rather than speed. We developed CNNs tailored for *MNIST* and *EMNIST*, with specifications in Table 6. For *CIFAR-10*, we adopted a model architecture from Trockman & Kolter (2023). As shown in Figure 4, both BALD and IV methods are suited to different scenarios, but our method consistently outperforms others across all datasets, particularly in more complex situations. The relatively straightforward *MNIST* dataset does not fully demonstrate our method’s advantages, so we created more challenging unbalanced and redundant *MNIST* datasets. Comparisons on these datasets reveal our method as the most effective, with advantages even more pronounced than on the original *MNIST*, as illustrated in Figure 10. Additionally, Figure 5 visually represents the selection preferences of  $FI^{active}$ , highlighting its tendency to identify points that are more challenging for the current model to distinguish.



**Figure 4:** Classification performance due to the different active learning methods on *MNIST*, *EMNIST*, and *CIFAR-10*.



**Figure 5:** Images with the highest  $FI^{active}$  selected in the first round of active learning for the *MNIST* problem, along with their corresponding prediction probability distributions.

## 7 CONCLUSION

In this paper, we introduce a novel local influence metric that evaluates the impact of perturbations to training samples on model performance concerning validation samples. This metric is applicable in both data trimming and active learning, offering valuable insights into the contributions of individual samples and their relationships with unlabeled data.

Furthermore, we propose two approximation methods to mitigate the computational costs associated with calculating local influence measures. Our experimental results demonstrate that these algorithms effectively reduce costs while outperforming other state-of-the-art methods.

## REFERENCES

Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. {TensorFlow}: a system for

- 540 {Large-Scale} machine learning. In *12th USENIX symposium on operating systems design and*  
541 *implementation (OSDI 16)*, pp. 265–283, 2016.
- 542
- 543 Shun-ichi Amari. *Differential-geometrical methods in statistics*, volume 28. Springer Science &  
544 Business Media, 2012.
- 545 Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. American  
546 Mathematical Soc., 2000.
- 547 Avrim L Blum and Pat Langley. Selection of relevant features and examples in machine learning.  
548 *Artificial intelligence*, 97(1-2):245–271, 1997.
- 549
- 550 Razvan Caramalau, Binod Bhattarai, and Tae-Kyun Kim. Sequential graph convolutional network  
551 for active learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
552 *recognition*, pp. 9583–9592, 2021.
- 553 Yuanyuan Chen, Boyang Li, Han Yu, Pengcheng Wu, and Chunyan Miao. Hydra: Hypergradient  
554 data relevance analysis for interpreting deep neural networks. *Proceedings of the AAAI Conference*  
555 *on Artificial Intelligence*, 35(8):7081–7089, May 2021. doi: 10.1609/aaai.v35i8.16871. URL  
556 <https://ojs.aaai.org/index.php/AAAI/article/view/16871>.
- 557 Anshuman Chhabra, Peizhao Li, Prasant Mohapatra, and Hongfu Liu. "what data benefits my  
558 classifier?" enhancing model performance and interpretability through influence-based data  
559 selection. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
560 <https://openreview.net/forum?id=HE9eUQ1Avo>.
- 561
- 562 Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist  
563 to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pp.  
564 2921–2926. IEEE, 2017.
- 565 R Dennis Cook. Assessment of local influence. *Journal of the Royal Statistical Society Series B:*  
566 *Statistical Methodology*, 48(2):133–155, 1986.
- 567
- 568 Melanie Ducoffe and Frederic Precioso. Adversarial active learning for deep networks: a margin  
569 based approach. *arXiv preprint arXiv:1802.09841*, 2018.
- 570 Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep Bayesian active learning with image data.  
571 In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference*  
572 *on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1183–  
573 1192. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/gal17a.](https://proceedings.mlr.press/v70/gal17a.html)  
574 [html](https://proceedings.mlr.press/v70/gal17a.html).
- 575 Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning.  
576 In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International*  
577 *Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*,  
578 pp. 2242–2251. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.press/v97/](https://proceedings.mlr.press/v97/ghorbani19c.html)  
579 [ghorbani19c.html](https://proceedings.mlr.press/v97/ghorbani19c.html).
- 580 Amirata Ghorbani, Michael Kim, and James Zou. A distributional framework for data valuation. In  
581 Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine*  
582 *Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 3535–3544. PMLR, 13–  
583 18 Jul 2020. URL <https://proceedings.mlr.press/v119/ghorbani20a.html>.
- 584
- 585 Gene H Golub and Christian Reinsch. Singular value decomposition and least squares solutions. In  
586 *Handbook for Automatic Computation: Volume II: Linear Algebra*, pp. 134–151. Springer, 1971.
- 587 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 588
- 589 Marc Gorriz, Axel Carlier, Emmanuel Faure, and Xavier Giro i Nieto. Cost-effective active learning  
590 for melanoma segmentation, 2017. URL <https://arxiv.org/abs/1711.09168>.
- 591 Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nezihe Merve Gurel, Bo Li, Ce Zhang,  
592 Costas Spanos, and Dawn Song. Efficient task-specific data valuation for nearest neighbor  
593 algorithms. *Proc. VLDB Endow.*, 12(11):1610–1623, July 2019a. ISSN 2150-8097. doi: 10.14778/  
3342263.3342637. URL <https://doi.org/10.14778/3342263.3342637>.

- 594 Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li,  
595 Ce Zhang, Dawn Song, and Costas J Spanos. Towards efficient data valuation based on the  
596 shapley value. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp.  
597 1167–1176. PMLR, 2019b.
- 598  
599 Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch  
600 acquisition for deep bayesian active learning. In Hanna M. Wallach, Hugo Larochelle, Alina  
601 Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in  
602 Neural Information Processing Systems 32: Annual Conference on Neural Information Pro-  
603 cessing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp.  
604 7024–7035, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/  
605 95323660ed2124450caaac2c46b5ed90-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/95323660ed2124450caaac2c46b5ed90-Abstract.html).
- 606 Andreas Kirsch, Sebastian Farquhar, Parmida Atighehchian, Andrew Jesson, Frédéric Branchaud-  
607 Charron, and Yarin Gal. Stochastic batch acquisition: A simple baseline for deep active learning.  
608 *Trans. Mach. Learn. Res.*, 2023, 2023. URL [https://openreview.net/forum?id=  
609 vCwQyNBjW](https://openreview.net/forum?id=vCwQyNBjW).
- 610 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions.  
611 In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference  
612 on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1885–  
613 1894. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/koh17a.  
614 html](https://proceedings.mlr.press/v70/koh17a.html).
- 615  
616 Ron Kohavi. Scaling up the accuracy of naive-bayes classifiers: a decision-tree hybrid. In *Proceedings  
617 of the Second International Conference on Knowledge Discovery and Data Mining, KDD’96*, pp.  
618 202–207. AAAI Press, 1996.
- 619 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.  
620 Technical Report 0, University of Toronto, Toronto, Ontario, 2009. URL [https://www.cs.  
621 toronto.edu/~kriz/learning-features-2009-TR.pdf](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf).
- 622  
623 Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation framework  
624 for machine learning. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.),  
625 *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume  
626 151 of *Proceedings of Machine Learning Research*, pp. 8780–8802. PMLR, 28–30 Mar 2022. URL  
627 <https://proceedings.mlr.press/v151/kwon22a.html>.
- 628 Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document  
629 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.
- 630  
631 Xingjian Li, Pengkun Yang, Yangcheng Gu, Xueying Zhan, Tianyang Wang, Min Xu, and  
632 Chengzhong Xu. Deep active learning with noise stability. In Michael J. Wooldridge, Jennifer G. Dy,  
633 and Sriraam Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024,  
634 Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth  
635 Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024,  
636 Vancouver, Canada*, pp. 13655–13663. AAAI Press, 2024. doi: 10.1609/AAAI.V38I12.29270.  
637 URL <https://doi.org/10.1609/aaai.v38i12.29270>.
- 638 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.  
639 In *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 3730–3738, 2015. doi:  
640 10.1109/ICCV.2015.425.
- 641  
642 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate  
643 curvature, 2020. URL <https://arxiv.org/abs/1503.05671>.
- 644 Christoph Mayer and Radu Timofte. Adversarial sampling for active learning. In *Proceedings of the  
645 IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3071–3079, 2020.
- 646  
647 Sérgio Moro, Paulo Cortez, and Paulo Rita. A data-driven approach to predict the success of bank  
telemarketing. *Decision Support Systems*, 62:22–31, 2014.

- 648 Peter Nickl, Lu Xu, Dharmesh Tailor, Thomas Möllenhoff, and Mohammad Emtiyaz Khan. The  
649 memory-perturbation equation: Understanding model's sensitivity to data. In A. Oh, T. Naumann,  
650 A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information  
651 Processing Systems*, volume 36, pp. 26923–26949. Curran Associates, Inc., 2023.
- 652 David Noever. Machine learning suites for online toxicity detection. *arXiv preprint arXiv:1810.01869*,  
653 2018.
- 654 Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-  
655 Dickstein. Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint  
656 arXiv:1802.08760*, 2018.
- 657 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zach DeVito, Zeming  
658 Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.  
659 URL <https://api.semanticscholar.org/CorpusID:40027675>.
- 660 Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data  
661 influence by tracing gradient descent. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and  
662 H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 19920–19930.  
663 Curran Associates, Inc., 2020.
- 664 Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen,  
665 and Xin Wang. A survey of deep active learning. *ACM computing surveys (CSUR)*, 54(9):1–40,  
666 2021.
- 667 Burr Settles. Active learning literature survey. 2009.
- 668 Hai Shu and Hongtu Zhu. Sensitivity analysis of deep neural networks. In *AAAI*, pp. 4943–4950.  
669 AAAI Press, 2019.
- 670 Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In  
671 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5972–5981, 2019.
- 672 Yang Sui, Yukun Huang, Hongtu Zhu, and Fan Zhou. Adversarial learning of distributional rein-  
673 forcement learning. In *International Conference on Machine Learning*, pp. 32783–32796. PMLR,  
674 2023.
- 675 Asher Trockman and J. Zico Kolter. Patches are all you need? *Trans. Mach. Learn. Res.*, 2023, 2023.  
676 URL <https://openreview.net/forum?id=rAnB7JSMXL>.
- 677 Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine  
678 learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 6388–6421.  
679 PMLR, 2023.
- 680 Jiachen T. Wang, Tianji Yang, James Zou, Yongchan Kwon, and Ruoxi Jia. Rethinking data shapley  
681 for data selection tasks: Misleads and merits. In Ruslan Salakhutdinov, Zico Kolter, Katherine  
682 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings  
683 of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine  
684 Learning Research*, pp. 52033–52063. PMLR, 21–27 Jul 2024.
- 685 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-  
686 attention distillation for task-agnostic compression of pre-trained transformers. In H. Larochelle,  
687 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing  
688 Systems*, volume 33, pp. 5776–5788. Curran Associates, Inc., 2020.
- 689 Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection  
690 for explaining deep neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,  
691 N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*,  
692 volume 31. Curran Associates, Inc., 2018.
- 693 Hongtu Zhu, Joseph G. Ibrahim, Sikyung Lee, and Heping Zhang. Perturbation selection and influence  
694 measures in local influence analysis. *The Annals of Statistics*, 35(6):2565–2588, 2007. ISSN  
695 00905364. URL <http://www.jstor.org/stable/25464601>.
- 696 Hongtu Zhu, Joseph G Ibrahim, and Niansheng Tang. Bayesian influence analysis: a geometric  
697 approach. *Biometrika*, 98(2):307–323, 2011.

## A DERIVATION OF EXPRESSION FOR $\partial_{\omega}\hat{\theta}(\mathbf{x} + \omega)|_{\omega=\omega_0}$

Recall that  $\hat{\theta}(\mathbf{x})$  minimizes the empirical risk:  $R(\theta) := \frac{1}{n} \sum_{i=1}^n L(y_i, \mathbf{x}_i; \theta)$ . We assume that  $R$  is twice-differentiable and strongly convex in  $\theta$ , i.e.,  $\mathbf{H}_{\hat{\theta}} := \partial_{\theta}^2 R(\hat{\theta}) = \frac{1}{n} \sum_{i=1}^n \partial_{\theta}^2 L(y_i, \mathbf{x}_i; \hat{\theta})$  exists and is positive definite. This ensures the availability of  $\mathbf{H}_{\hat{\theta}}^{-1}$ , which will be utilized in the subsequent derivation. The perturbed parameter vector  $\hat{\theta}(\mathbf{x} + \omega)$  can be written as

$$\hat{\theta}(\mathbf{x} + \omega) = \arg \min_{\theta \in \Theta} \left\{ R(\theta) + \frac{1}{n} L(y, \mathbf{x} + \omega; \theta) - \frac{1}{n} L(y, \mathbf{x}; \theta) \right\}. \quad (\text{A.1})$$

Define the parameter change  $\delta(\omega) = \hat{\theta}(\mathbf{x} + \omega) - \hat{\theta}(\mathbf{x})$ , and note that, since  $\hat{\theta}(\mathbf{x})$  does not depend on  $\omega$ , the quantity we aim to calculate can be expressed in terms of it:  $\partial_{\omega}\hat{\theta}(\mathbf{x} + \omega) = \partial_{\omega}\delta(\omega)$ .

According to the definition of  $\hat{\theta}(\mathbf{x} + \omega)$ , we know

$$\mathbf{0} = \partial_{\theta} R(\hat{\theta}(\mathbf{x} + \omega)) + \frac{1}{n} \partial_{\theta} L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x} + \omega)) - \frac{1}{n} \partial_{\theta} L(y, \mathbf{x}; \hat{\theta}(\mathbf{x} + \omega)). \quad (\text{A.2})$$

Next, since  $\hat{\theta}(\mathbf{x} + \omega) \rightarrow \hat{\theta}(\mathbf{x})$  as  $\omega \rightarrow \mathbf{0}$ , we perform a Taylor expansion of the right-hand side:

$$\mathbf{0} = \partial_{\theta} R(\hat{\theta}(\mathbf{x} + \omega)) + \frac{1}{n} \partial_{\theta} L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x} + \omega)) - \frac{1}{n} \partial_{\theta} L(y, \mathbf{x}; \hat{\theta}(\mathbf{x} + \omega)), \quad (\text{A.3})$$

$$\begin{aligned} &\approx [\partial_{\theta} R(\hat{\theta}(\mathbf{x})) + \frac{1}{n} \partial_{\theta} L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x})) - \frac{1}{n} \partial_{\theta} L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))] \\ &\quad + [\partial_{\theta}^2 R(\hat{\theta}(\mathbf{x})) + \frac{1}{n} \partial_{\theta}^2 L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x})) - \frac{1}{n} \partial_{\theta}^2 L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))] \delta(\omega). \end{aligned} \quad (\text{A.4})$$

Solving for  $\delta(\omega)$ , we get

$$\begin{aligned} \delta(\omega) &\approx - [\partial_{\theta}^2 R(\hat{\theta}(\mathbf{x})) + \frac{1}{n} \partial_{\theta}^2 L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x})) - \frac{1}{n} \partial_{\theta}^2 L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))]^{-1} \\ &\quad \cdot [\partial_{\theta} R(\hat{\theta}(\mathbf{x})) + \frac{1}{n} \partial_{\theta} L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x})) - \frac{1}{n} \partial_{\theta} L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))]. \end{aligned} \quad (\text{A.5})$$

Since  $\hat{\theta}(\mathbf{x})$  minimizes  $R$ , we have  $\partial_{\theta} R(\hat{\theta}) = 0$ . We further assume that  $\partial_{\theta}^2 L(y, \mathbf{x}; \hat{\theta}(\mathbf{x}))$  is continuous on  $\mathbf{x}$ , then we have

$$\delta(\omega) \approx - \frac{1}{n} [\partial_{\theta}^2 R(\hat{\theta}(\mathbf{x}))]^{-1} \left[ \frac{1}{n} \partial_{\theta} L(y, \mathbf{x} + \omega; \hat{\theta}(\mathbf{x})) - \frac{1}{n} \partial_{\theta} L(y, \mathbf{x}; \hat{\theta}(\mathbf{x})) \right]. \quad (\text{A.6})$$

After differentiation, the final expression can be obtained,

$$\partial_{\omega}\hat{\theta}(\mathbf{x} + \omega)|_{\omega=\omega_0} = \partial_{\omega}\delta(\omega)|_{\omega=\omega_0} \approx - \frac{1}{n} \mathbf{H}_{\hat{\theta}}^{-1} \partial_{\mathbf{x}} \partial_{\theta} L(y, \mathbf{x}; \hat{\theta}(\mathbf{x})). \quad (\text{A.7})$$

## B ALGORITHMS

---

### Algorithm 3 Data Trimming

---

**Input:** Training set  $\mathcal{Z}$ , Validation set  $\mathcal{V}$ , Base model  $\mathcal{F}$ , Budget  $b$

**Output:** Trimmed Dataset  $\mathcal{Z}'$

- 1: **procedure** DATATRIMMING( $\mathcal{Z}, \mathcal{V}, \mathcal{F}, b$ )
  - 2:     Call the algorithm FI<sup>util</sup>-CALCULATION( $\mathcal{Z}, \mathcal{V}, \mathcal{F}$ ) to obtain FI<sup>util</sup>  $\in \mathbb{R}^{|\mathcal{Z}| \times 1}$
  - 3:     Select  $b$  samples  $\mathbf{z}_i \in \mathcal{Z}$  as  $\{Z_b\}$ , whose FI<sup>util</sup>[ $i$ ] rank in the top  $b$
  - 4:      $\mathcal{Z}' \leftarrow \mathcal{Z} \setminus \{Z_b\}$
  - 5:     **return**  $\mathcal{Z}'$
  - 6: **end procedure**
-

**Algorithm 4** Active Learning

**Input:** Labeled pool of training data  $\mathcal{L}$ , Unlabeled pool of training data  $\mathcal{U}$ , Base model  $\mathcal{F}$ , The number of samples for annotation (per round)  $N$ , Aggregation function  $g$

**Output:** Updated labeled pool  $\mathcal{L}$

```

1: procedure ACTIVELEARNING( $\mathcal{L}, \mathcal{U}, \mathcal{F}, N, g$ )
2:   for  $j \leftarrow 1$  to  $NUM\_rounds$  do
3:     Call the algorithm  $FI^{active}$ -CALCULATION( $\mathcal{L}, \mathcal{U}, \mathcal{F}, g$ ) to obtain  $FI^{active} \in \mathbb{R}^{|\mathcal{U}| \times 1}$ 
4:     Select  $N$  samples  $x_i^{\mathcal{U}}$  as  $\{X_N\}$ , whose  $FI^{active}[i]$  ranks in the top  $N$ 
5:     Take  $\{X_N\}$  out of  $\mathcal{U}$ , and query their labels  $\{Y_N\}$ 
6:     Update  $\mathcal{L} \leftarrow \mathcal{L} \cup \{X_N, Y_N\}$ 
7:   end for
8:   return  $\mathcal{L}$ 
9: end procedure

```

**Algorithm 5** Calculation of  $FI_{approx}^{util}$ 

**Input:** Training set  $\mathcal{Z}$ , Validation set  $\mathcal{V}$ , Base model  $\mathcal{F}$ , Truncated-SVD parameter  $k$

**Output:** Influence measure vector  $FI_{approx}^{util} \in \mathbb{R}^{|\mathcal{Z}| \times 1}$

```

1: procedure  $FI_{approx}^{util}$ -CALCULATION( $\mathcal{Z}, \mathcal{V}, \mathcal{F}, k$ )
2:   Train  $\mathcal{F}$  with  $\mathcal{Z}$ , and obtain the parameter vector  $\hat{\theta}$ 
3:   Generate an empty vector  $FI_{approx}^{util}$  of size  $|\mathcal{Z}| \times 1$ 
4:   Calculate  $H_{\hat{\theta}}$  with K-FAC approximation
5:   Calculate  $\frac{1}{|\mathcal{V}|} \sum_{(x', y') \in \mathcal{V}} \partial_{\theta} L(y', x'; \hat{\theta})$ 
6:   for every  $z_i$  in  $\mathcal{Z}$  do
7:     Calculate  $G_{z_i}(\omega_0)$ 
8:     Calculate  $\Lambda_{k \times k}, U_{p \times k}$  and  $V_{d \times k}$  with power iteration
9:      $\partial_x \partial_{\theta} L(y_i, x_i; \hat{\theta}) \leftarrow U_{p \times k} \Lambda_{k \times k} V_{d \times k}^{\top}$ 
10:     $\partial_{\omega} L(\mathcal{V}; \hat{\theta}) \leftarrow \frac{1}{|\mathcal{V}|} \sum_{(x', y') \in \mathcal{V}} \partial_{\theta} L(y', x'; \hat{\theta}) H_{\hat{\theta}}^{-1} \partial_x \partial_{\theta} L(y_i, x_i; \hat{\theta})$ 
11:     $FI_{approx}^{util}[i] \leftarrow \partial_{\omega} L(\mathcal{V}; \hat{\theta}) G_{z_i}^{\dagger}(\omega_0) \partial_{\omega}^{\top} L(\mathcal{V}; \hat{\theta})$ 
12:   end for
13:   return  $FI_{approx}^{util}$ 
14: end procedure

```

**Algorithm 6** Data Trimming with Approximation

**Input:** Training set  $\mathcal{Z}$ , Validation set  $\mathcal{V}$ , Base model  $\mathcal{F}$ , Budget  $b$ , Truncated-SVD parameter  $k$ , Subsampling rate  $\alpha \in (0, 1]$

**Output:** Trimmed Dataset  $\mathcal{Z}'$

```

1: procedure DATATRIMMINGAPPROX( $\mathcal{Z}, \mathcal{V}, \mathcal{F}, b, k, \alpha$ )
2:   Subsample  $\mathcal{Z}_{\alpha} \subseteq \mathcal{Z}$  s.t.  $z \in \mathcal{Z}_{\alpha}$  w.p.  $\alpha$  for all  $z \in \mathcal{Z}$ 
3:   Call the algorithm  $FI_{approx}^{util}$ -CALCULATION( $\mathcal{Z}_{\alpha}, \mathcal{V}, \mathcal{F}, k$ ) to obtain  $FI_{approx}^{util} \in \mathbb{R}^{|\mathcal{Z}_{\alpha}| \times 1}$ 
4:   for  $z_{i'} = (x_{i'}, y_{i'}) \in \mathcal{Z}_{\alpha}$  do
5:      $r_{i'} \leftarrow$  the rank of  $FI_{approx}^{util}[i']$  sorted in ascending order
6:   end for
7:   Train a random forest  $h$  with  $\{(x_{i'}, r_{i'})\}_{i'=1}^{|\mathcal{Z}_{\alpha}|}$ 
8:   for  $z_i = (x_i, y_i) \in \mathcal{Z}$  do
9:      $\hat{r}_i \leftarrow h(x_i)$ 
10:  end for
11:  Select  $b$  samples  $z_i \in \mathcal{Z}$ , whose  $\hat{r}_i$  rank in the top  $b$ 
12:   $\mathcal{Z}' \leftarrow \mathcal{Z} \setminus \{Z_b\}$ 
13:  return  $\mathcal{Z}'$ 
14: end procedure

```

**Algorithm 7** Calculation of  $\text{FI}_{approx}^{active}$ 

**Input:** Labeled pool of training data  $\mathcal{L}$ , Unlabeled pool of training data  $\mathcal{U}$ , Base model  $\mathcal{F}$ , Aggregation function  $g$ , Truncated-SVD parameter  $k$

**Output:** Influence measure vector  $\text{FI}_{approx}^{active} \in \mathbb{R}^{|\mathcal{U}| \times 1}$

```

1: procedure  $\text{FI}_{approx}^{active}$ -CALCULATION( $\mathcal{L}, \mathcal{U}, \mathcal{F}, g, k$ )
2:   Train  $\mathcal{F}$  with  $\mathcal{L}$ , and obtain the parameter vector  $\hat{\theta}$ 
3:   Generate an empty vector  $\text{FI}_{approx}^{active}$  of size  $|\mathcal{U}| \times 1$ 
4:   Calculate  $\mathbf{H}_{\hat{\theta}}$  with K-FAC approximation
5:   for every  $x_i^{\mathcal{U}}$  in  $\mathcal{U}$  do
6:     Obtain an estimated label  $\hat{y}_i$  with  $\mathcal{F}$ 
7:     Calculate  $\partial_{\theta} L(\hat{y}_i, x_i^{\mathcal{U}}; \hat{\theta})$ 
8:      $\mathcal{J} \leftarrow \emptyset$ 
9:     for every  $z_k^{\mathcal{L}}$  in  $\mathcal{L}$  do
10:      Calculate  $\mathbf{G}_{z_k^{\mathcal{L}}}(\omega_0)$ 
11:      Calculate  $\Lambda_{k \times k}, U_{p \times k}$  and  $V_{d \times k}$  with power iteration
12:       $\partial_x \partial_{\theta} L(y_k^{\mathcal{L}}, x_k^{\mathcal{L}}; \hat{\theta}) \leftarrow U_{p \times k} \Lambda_{k \times k} V_{d \times k}^{\top}$ 
13:       $\partial_{\omega} L(\hat{y}_i, x_i^{\mathcal{U}}; \hat{\theta}) \leftarrow \partial_{\theta} L(\hat{y}_i, x_i^{\mathcal{U}}; \hat{\theta}) \mathbf{H}_{\hat{\theta}}^{-1} \partial_x \partial_{\theta} L(y_k^{\mathcal{L}}, x_k^{\mathcal{L}}; \hat{\theta})$ 
14:       $\mathcal{J} \leftarrow \mathcal{J} \cup \{\partial_{\omega} L(\hat{y}_i, x_i^{\mathcal{U}}; \hat{\theta}) \mathbf{G}_{z_k^{\mathcal{L}}}^{\dagger}(\omega_0) \partial_{\omega}^{\top} L(\hat{y}_i, x_i^{\mathcal{U}}; \hat{\theta})\}$ 
15:     end for
16:      $\text{FI}_{approx}^{active}[i] \leftarrow g(\mathcal{J})$ 
17:   end for
18:   return  $\text{FI}_{approx}^{active}$ 
19: end procedure

```

**Algorithm 8** Active Learning with Approximation

**Input:** Labeled pool of training data  $\mathcal{L}$ , Unlabeled pool of training data  $\mathcal{U}$ , Base model  $\mathcal{F}$ , Number of samples for annotation (per round)  $N$ , Aggregation function  $g$ , Truncated-SVD parameter  $k$ , Subsampling rate  $\alpha \in (0, 1]$

**Output:** Updated labeled pool  $\mathcal{L}$

```

1: procedure ACTIVELEARNINGAPPROX( $\mathcal{L}, \mathcal{U}, \mathcal{F}, N, g, k, \alpha$ )
2:   for  $j \leftarrow 1$  to  $NUM\_rounds$  do
3:     Subsample  $\mathcal{U}_{\alpha} \subseteq \mathcal{U}$  s.t.  $x \in \mathcal{U}_{\alpha}$  w.p.  $\alpha$  for all  $x \in \mathcal{U}$ 
4:     Call the algorithm  $\text{FI}_{approx}^{active}$ -CALCULATION( $\mathcal{L}, \mathcal{U}_{\alpha}, \mathcal{F}, g, k$ ) to obtain  $\text{FI}_{approx}^{active} \in \mathbb{R}^{|\mathcal{U}_{\alpha}| \times 1}$ 
5:     for  $x_{i'} \in \mathcal{U}_{\alpha}$  do
6:        $r_{i'} \leftarrow$  the rank of  $\text{FI}_{approx}^{active}[i']$  sorted in ascending order
7:     end for
8:     Train a random forest  $h$  with  $\{(x_{i'}, r_{i'})\}_{i'=1}^{|\mathcal{U}_{\alpha}|}$ 
9:     for  $x_i \in \mathcal{U}$  do
10:       $\hat{r}_i \leftarrow h(x_i)$ 
11:    end for
12:    Select  $N$  samples  $x_i$  as  $\{X_N\}$ , whose  $\hat{r}_i$  ranks in the top  $N$ 
13:    Take  $\{X_N\}$  out of  $\mathcal{U}$ , and query their labels  $\{Y_N\}$ 
14:    Update  $\mathcal{L} \leftarrow \mathcal{L} \cup \{X_N, Y_N\}$ 
15:  end for
16:  return  $\mathcal{L}$ 
17: end procedure

```

## C DATA SOURCES

**Adult.** The Adult dataset consists of 48,842 instances and 14 features, including categorical and integer types. Extracted from the 1994 Census database by Barry Becker, it focuses on predicting whether an individual’s annual income exceeds \$50,000. Records were filtered based on criteria such



864 as age, gross income, and work hours, making this dataset a valuable resource for classification tasks  
865 in social science and *Access Link: Adult Database*

866 **Bank.** The Bank dataset pertains to direct marketing campaigns conducted by a Portuguese banking  
867 institution, focusing on phone call outreach. The primary objective of this dataset is to classify  
868 whether a client will subscribe to a term deposit, indicated by the binary variable (yes/no). This  
869 dataset is derived from multiple marketing campaigns, which often required several contacts with the  
870 same client to ascertain their interest in the product. *Access Link: Bank Database*

871 **CelebA.** The CelebA dataset is a large-scale facial attribute dataset containing over 200,000 celebrity  
872 images, each annotated with 40 attribute labels. This dataset serves as a valuable resource for tasks  
873 such as facial recognition, attribute prediction, and generative modeling. *Access Link: CelebA*  
874 *Database*

876 **Jigsaw Toxicity.** The Jigsaw dataset of Wikipedia consists of comments from online platforms  
877 that have been labeled for toxicity. It contains a large number of comments, with 28 features of  
878 syntax, sentiment, emotion and outlier word dictionaries. The dataset is commonly used for training  
879 and evaluating machine learning models aimed at detecting harmful or inappropriate content in  
880 user-generated text. *Access Link: Jigsaw Toxicity Database*

881 **MNIST.** The MNIST database is a large collection of handwritten digits that is widely used for training  
882 and testing in the field of machine learning. This dataset contains 70,000 images of handwritten digits  
883 (0-9), each of which is a  $28 \times 28$  pixel grayscale image. *Access Link: MNIST Database*

884 **EMNIST.** The Extended MNIST database consists of handwritten character digits sourced from the  
885 NIST Special Database 19, formatted as  $28 \times 28$  pixel images to align with the structure of the MNIST  
886 dataset. There are six different splits provided in this dataset, and we use the EMNIST Letters with  
887 145,600 characters. *Access Link: EMNIST Database*

888 **CIFAR10.** The CIFAR10 database, developed by the Canadian Institute for Advanced Research, is a  
889 widely utilized collection of images for training machine learning and computer vision algorithms. It  
890 consists of 60,000 color images, each measuring  $32 \times 32$  pixels, categorized into 10 classes: airplanes,  
891 cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks. Each class contains 6,000 images. *Access*  
892 *Link: CIFAR10 Database*

## 894 D EXPERIMENT DETAILS

### 896 D.1 EXPERIMENT DETAILS IN 6.1

898 **Data Construction for Real-World Datasets.** In this experiment, the preprocessing methods for  
899 all real-world datasets are consistent with those used in Chhabra et al. (2024). For *CelebA*, we use the  
900 extracted features provided by the authors Liu et al. (2015), and for *Jigsaw Toxicity*, we obtain text  
901 embeddings using the MiniLM transformer model (Wang et al., 2020). Further details are provided  
902 below.

- 903 • **Adult.** This dataset contains 37,692 samples, with 30,162 for training and 7,530 for testing.  
904 There are 102 features, and the target is to predict if income exceeds \$50k (yes) or not (no).
- 905 • **Bank.** This dataset consists of 30,490 samples, divided into 18,292 training samples and  
906 12,198 test samples. There are 50 features, and the target is to predict if the client will  
907 subscribe a term deposit (yes/no).
- 908 • **CelebA.** This dataset includes 104,163 samples, with 62,497 for training and 41,666 for  
909 testing. There are 39 features, and the aim is to predict whether a person is smiling (yes) or  
910 not (no).
- 911 • **Jigsaw Toxicity.** This dataset consists of 30,000 samples, split into 18,000 training samples  
912 and 12,000 test samples. There are 385 features, and the target is to determine if a tweet is  
913 toxic (yes) or not (no).

915 Given that the initial test accuracy of the model on the original datasets generally exceeds 90%, the  
916 impact of data trimming is minimal. Accordingly, we randomly sample 5,000 instances from the  
917 original training set to serve as the new training set, and 4,200 instances from the test set to serve as  
the new test set.

**Data Construction for Noisy Real-World Datasets.** Consider the additional experiments on data trimming presented in Figure 8. For *Adult*, *Bank*, and *Jigsaw Toxicity*, our construction method follows the same approach as in the Real-World Datasets experiment, with the additional step of randomly sampling 500 instances from the training sets and adding white noise with a variance of 0.1. For *CelebA*, since the variables are binary (-1, 1), we introduced noise by randomly selecting six feature columns and flipping their values (transforming -1 to 1 and 1 to -1). This approach generates a noisy dataset for supplementary experiments.

**Random Trimming.** We randomly remove  $b$  (budget) data points from the training sets under five random seeds, and evaluate the average performance on the test sets in each iteration.

**Experimental Procedure and Parameter Settings.** First, a model is trained on the initial dataset. Based on this trained model, we then implement three data trimming strategies, removing  $b$  data points. Finally, we retrain the model to evaluate the effectiveness of the different trimming strategies. The experimental procedure is illustrated in Figure 6, highlighting the key steps involved in our study. Additionally, the parameter settings are detailed in Table 3.

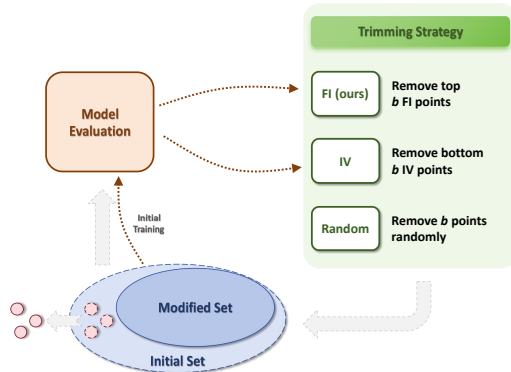


Figure 6: Flowchart of data trimming.

Table 3: Parameter Settings of Data Trimming.

	Adult	Bank	CelebA	Jigsaw Toxicity
optimizer	SGD	SGD	Adam	Adam
learning_rate	1e-2	1e-2	1e-4	1e-2
weight_decay	1e-2	1e-2	1e-6	1e-2
	Adult +noise	Bank +noise	CelebA +noise	Jigsaw Toxicity +noise
optimizer	SGD	SGD	SGD	Adam
learning_rate	1e-2	1e-2	1e-2	1e-1
weight_decay	1e-2	1e-4	1e-6	1e-4

## D.2 EXPERIMENT DETAILS IN 6.2

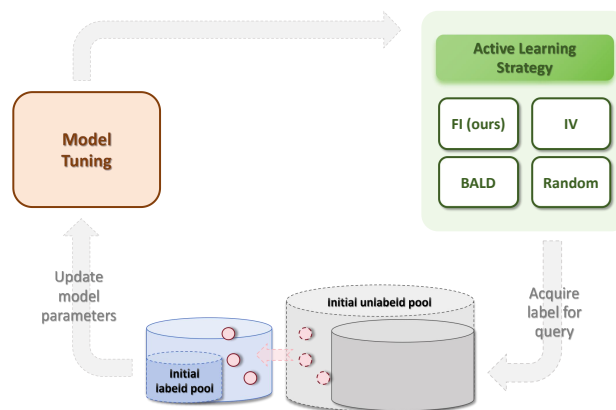


Figure 7: Flowchart of active learning.

**Data Construction for Unbalanced and Redundant MNIST.** For *Unbalanced MNIST*, we set the total sample size to 27,500. Categories 0 to 4 are each allocated an equal sample size, representing 1/55 of the total sample. Similarly, categories 5 to 9 are assigned equal sample sizes, with each constituting 10/55 of the total sample. This allocation strategy ensures a deliberate imbalance among

the classes. Samples are systematically drawn based on these ratios from the original dataset to create this new unbalanced dataset. For **Redundant MNIST**, the task is delineated to classify solely the digits 1 and 7, presented in equal proportions. If the acquisition function selects an input from any class other than 1 or 7, the labeling function designates a “neither” category. This setup leads to a three-way classification scheme during training, categorized as 1 vs. 7 vs. neither. This design allows us to explore the effectiveness of the learning model in dealing with class imbalance and partial class information, critical aspects in real-world applications where similar conditions are often encountered.

**Active Learning Experiment.** Tables 4 and 5 detail the parameter settings for active learning with tabular and image data, respectively. Table 6 describes the neural network architectures employed for *MNIST* and *EMNIST*.

**Table 4:** Active Learning Experiment Configuration for Tabular Datasets

Attribute	Adult	CelebA	Jigsaw Toxicity
Number of Classes	2	2	2
Rounds	12	11	8
Initial Pool	300	300	180
Unlabeled Pool Size	5000	3000	5000
Acquisition Size	500	150	500

**Table 5:** Active Learning Experiment Configuration for Image Datasets

Attribute	MNIST	EMNIST	CIFAR10
Number of Classes	10	37	10
Rounds	14	10	9
Initial Pool	60	600	6000
Unlabeled Pool Size	420	2000	10000
Acquisition Size	30	100	2000

**Table 6:** Architecture of MNIST CNN

Layer Type	Activation	Output Dimensions (incl. Padding)
Conv2d	ReLU	(32, 14, 14), P=1
Conv2d	ReLU	(64, 7, 7), P=1
Dropout	-	-
Linear	ReLU	256
Linear	-	num_classes

## E ADDITIONAL EXPERIMENTS

### E.1 ABLATION EXPERIMENTS

Here, we supplement ablation experiments to evaluate the impact of perturbations on model performance, taking the *Bank* dataset as a case study. Specifically, we introduce Gaussian noise to  $b$  samples in the training set with the highest and the lowest FI values, respectively, where  $b$  takes values from [50, 100, 150, 200, 250, 300]. After introducing the perturbations, we evaluate the model’s performance on the test set. The results are presented in Table 7 below. It can be observed that samples with higher FI values are more sensitive to perturbations, as their accuracy decreases more significantly compared to samples with lower FI values.

To address the trade-off between computational cost and performance, we add ablation experiments on two datasets. As shown in the table, the approximate method is over three times faster than

**Table 7: Impact of Perturbations on Model Accuracy for Data Trimming.**

$b$	Acc. after Perturbations (Top $b$ -FI Samples)	Change in Acc.	Acc. after Perturbations (Bottom $b$ -FI Samples)	Change in Acc.
0	78.36%	/	78.36%	/
50	77.59%	-0.77%	78.49%	0.14%
100	77.19%	-1.17%	78.33%	-0.03%
150	77.07%	-1.29%	77.45%	-0.91%
200	75.47%	-2.89%	77.42%	-0.93%
250	75.42%	-2.93%	77.19%	-1.17%
300	75.45%	-2.91%	77.04%	-1.31%

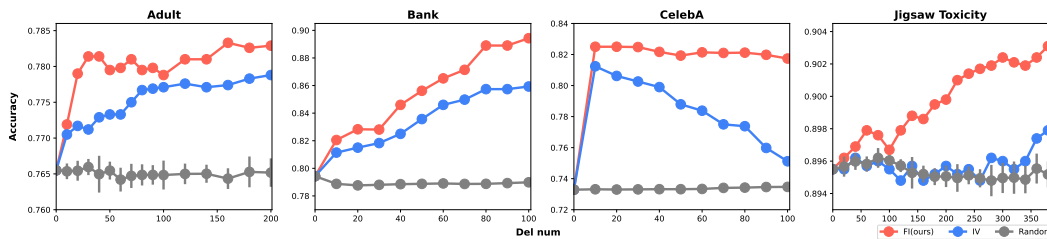
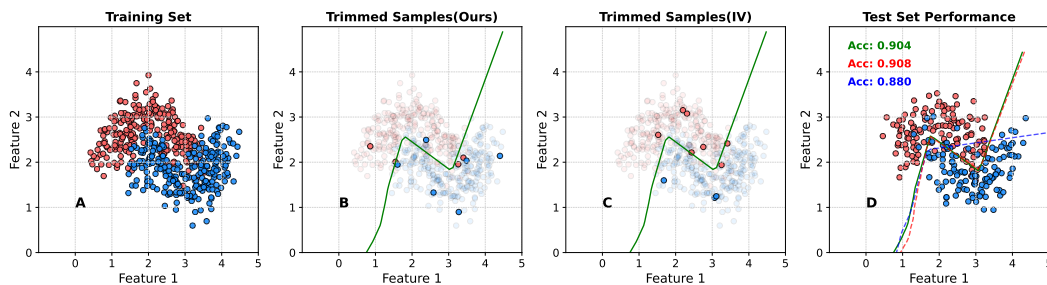
the exact computation, with minimal impact on performance. It even achieves better results on the CelebA dataset. Additionally, as the dataset size increases, the proportion of time spent on training and prediction with the random forest decreases, making the speedup even more significant.

**Table 8: Comparison of Approximate and Exact Methods on Different Datasets**

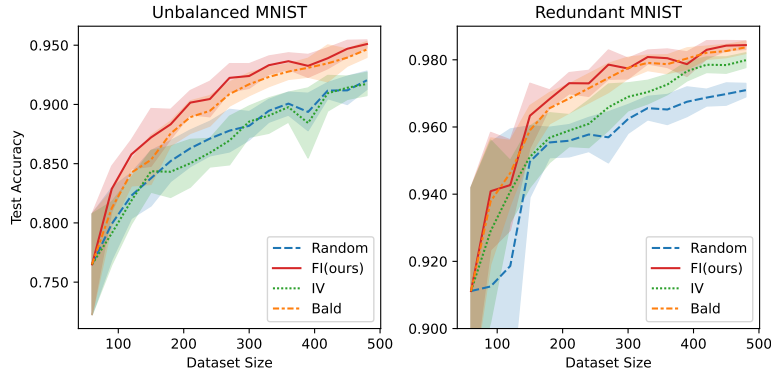
Dataset	Approx. Accuracy (%)	Exact Accuracy (%)	Approx. Time (s)	Exact Time (s)
MNIST	93.8	94.0	393	1228
CelebA	80.9	78.1	269	1058

## E.2 ADDITIONAL EXPERIMENTS FOR DATA TRIMMING

Here, we present experimental results on four real-world datasets with noise in Figure 8, and nonlinear experimental results in Figure 9.

**Figure 8: Accuracy curves of three data trimming methods on test sets of *Adult*, *Bank*, *CelebA* and *Jigsaw Toxicity*.****Figure 9: Performance under Nonlinear Model.** Different colored points represent different classes. **A** shows the training set. **B** and **C** respectively denote the samples to be trimmed by FI method and IV method. **D** denotes test set. **Green** line: boundary without trimming. **Red** line: boundary after FI trimming. **Blue** line: boundary after IV trimming.

## E.3 ADDITIONAL EXPERIMENTS FOR ACTIVE LEARNING



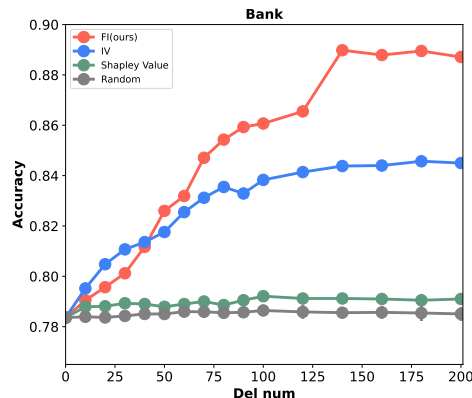
**Figure 10:** Classification performance due to the different active learning methods on *Unbalanced MNIST*, and *Redundant MNIST*.

**Table 9:** Time (seconds) for active learning algorithms over toy case on *MNIST*, *EMNIST*, and *CIFAR-10*.

Methods	MNIST	EMNIST	CIFAR-10
FI(ours)	19	45	372
IV	13	15	204
BALD	9	15	214

## E.4 ADDITIONAL EXPERIMENTS FOR COMPARISON WITH DATA SHAPLEY

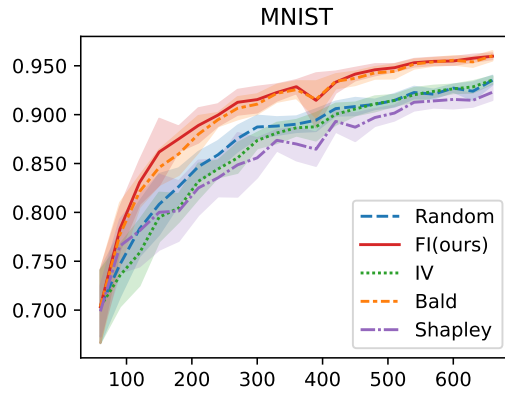
**Data Shapley** is indeed a well-known metric for data valuation (Ghorbani & Zou, 2019) and is applicable in this context. However, based on our experimental results, Shapley value-based methods do not demonstrate any advantage over our approach or IV. Wang et al. (2024) share a similar perspective on applying Data Shapley to data selection tasks, noting that its performance can be comparable to random selection. Additionally, their high computational cost presents a significant drawback. Specifically:



**Figure 11:** Performance comparison of different methods on *Bank*, including the newly added Shapley method.

- In our experiments on data trimming with *Bank*, we observed that the Data Shapley method led to only a slight improvement in model performance, with results showing marginal differences compared to Random selection (see Figure 11).

- In our experiments on active learning with *MNIST*, we found that Data Shapley performed the worst, even slightly underperforming compared to Random selection (see Figure 12).
- In terms of time consumption, the Data Shapley method takes approximately 860 times longer than FI.

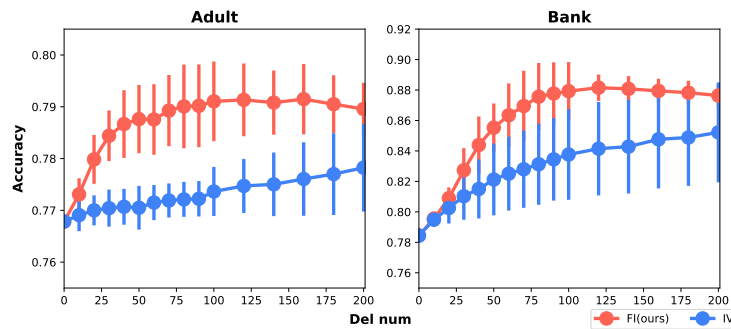


**Figure 12:** Performance comparison of different methods on *MNIST*, including the newly added Shapley method.

### E.5 K-FOLD CROSS-VALIDATION FOR DATA TRIMMING

Although all the datasets used in our experiments have clearly defined training and test sets, we further conduct 10-fold cross-validation on Data Trimming to verify the robustness of our method across different splits of the data.

In this experimental setup, the initial training set was randomly divided into 10 equally-sized subsets. In each iteration, one subset was held out as the validation set, while the remaining 9 subsets were used as the training set. The test set was fixed and remained the same as in the previous experiments, ensuring consistency across all evaluations. The results of 10-fold experiments are detailed in Figure 13, where we compare the overall performance and stability of our approach FI against the IV method.



**Figure 13:** K-Fold Cross-Validation for Data Trimming: Performance Comparison of FI and IV on *Adult* and *Bank*.