# JUST SELECT TWICE: LEVERAGING LOW QUALITY DATA TO IMPROVE DATA SELECTION

Anonymous authors

Paper under double-blind review

## Abstract

Data valuation is crucial for assessing the impact and quality of individual data points, enabling the ranking of data by importance for efficient data collection, storage, and training. Many data valuation methods are sensitive to outliers and require a large level of noise to effectively distinguish low-quality data from highquality data, making them particularly useful for data removal tasks. In particular, optimal transport-based methods exhibit notable performance in outlier detection but show only moderate effectiveness in high-quality data selection, due to their sensitivity to outliers and insensitivity to small variations. To mitigate the issue of insensitivity to high-quality data and facilitate effective data selection, in this paper, we propose a straightforward two-stage approach, JST, that initially does data valuation as usual, but then performs a second-round data selection where the identified low-quality data points are designated as the validation set to perform data valuation again. In this way, high-quality data become outliers with respect to the new validation set and can be naturally identified<sup>1</sup>. We empirically evaluate an instantiation of our framework based on optimal transport method for data selection and data pruning on several standard datasets and our framework demonstrates superior performance compared to pure data valuation, especially under small noise conditions. Additionally, we show the general applicability of our framework to influence function based and reinforcement learning based data valuation methods.

028 029

031

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

022

024

025

026

027

#### 1 INTRODUCTION

Access to large, high-quality datasets is essential in machine learning. However, in real-world data collection and curation pipelines, individual data points often inherit different levels of quality and vary in their importance on the impact of training (Liang & Zou, 2022; Sorscher et al., 2022; Liu et al., 2021). Therefore, it is critical to understand and assess such properties of data and to effectively prioritize highly valuable data sources for subset selection. This can assist practitioners in improving model performance efficiently (Karr et al., 2006; Jiang et al., 2023) and make strategic, cost-effective decisions in data marketplaces and exchanges (Olaleye & Adusei, 2024).

There have been considerable recent efforts to develop various data valuation methods aimed at evaluating individual data points and assigning a value to each one (Jiang et al., 2023; Sim et al., 2022). This can help to quantify differences between points and rank them based on their assigned value, establishing a certain order of quality or importance in the training process. For instance, methods based on optimal transport (Just et al., 2023) and influence functions (Koh & Liang, 2017) employ sensitivity analysis to quantify the impact of individual data points on dataset distance and training outcome, respectively. Additionally, importance weights can be obtained through reinforcement learning (Yoon et al., 2020) to evaluate individual data points.

In essence, these data valuation methods should naturally provide a metric for high-quality data subset selection. However, unfortunately, most data valuation methods have predominantly excelled in scenarios with large noise, including intensely corrupted samples (Just et al., 2023), randomly flipping labels (Just et al., 2023; Koh & Liang, 2017; Yoon et al., 2020), or domain adaption due to

 <sup>&</sup>lt;sup>1</sup>To be clear, in typical data valuation methods, outliers are considered with respect to a clean validation set, with noisy data treated as outliers. In our framework, the validation set is noisy, so high-quality clean data are considered outliers.



Figure 1: Sensitivity comparison between outliers and in-distribution data. We inject white noise into 50% of the CIFAR-10 training set with different noise levels and inspect top k data points for each data valuation method from lowest values for outliers and highest values for in-distribution data, respectively. We evaluate the difference between the normalized outlier detection rate (NDR<sub>outlier</sub>) and the normalized in-distribution data detection rate (NDR<sub>in-distribution</sub>). The outlier data detection rate is much larger than in-distribution data detection rate, reflecting higher sensitivity of these data valuation methods to outliers.

substantial mismatch between training and target domains (Koh & Liang, 2017; Yoon et al., 2020). In
particular, the optimal transport-based method in (Just et al., 2023) demonstrates strong performance
in outlier detection but exhibits only moderate effectiveness in selecting high-quality data due to its
well-known property of sensitivity to outliers and insensitivity to small variations (Villani, 2021).

As demonstrated in Figure 1, the sensitivity of these data valuation methods to outliers is much larger
than the high-quality data. Therefore, in this paper, we address the following question: *How can we improve the sensitivity of data valuation methods to high-quality data and make them better when applied to data selection?*

080 To this end, we propose JST (Just Select Twice), a data subset selection framework that augments 081 existing data valuation methods by leveraging outliers detected by these methods, as illustrated in 082 Figure 2. Our key insight is that incorporating a second-round subset selection, using detected outliers 083 as the new validation set, allows high-quality data to be identified as outliers with respect to the new 084 validation set in this new context. This approach enhances the sensitivity of data valuation methods to high-quality data, establishing a more meaningful order of "importance" or "quality" for subset 085 selection. Even with the inclusion of numerous high-quality data points within the new validation 086 set alongside low-quality data in the second-round subset selection due to the non state-of-the-art 087 performance of data valuation methods, we observe that the signals from such a validation set continue 088 to suffice in achieving superior performance in subset selection compared to pure data valuation 089 methods. 090

Before diving into the details, we summarize our contributions as follows: (1) JST, a novel and straightforward two-stage framework augmenting existing data valuation methods, rendering them more applicable for data subset selection tasks. (2) We empirically evaluate an instantiation of our framework based on optimal transport method for data selection and data pruning on six standard datasets, showing its superior outperformance to pure data valuation, particularly in small noise settings. We also demonstrate its general applicability to influence function and reinforcement learning based data valuation methods.

098 099

100

064

065

066

067

068

069

070 071

#### 2 BACKGROUND AND RELATED WORK

101 We provide a small amount of background and describe related work.

Data Selection. Active learning and coreset construction are two widely used methods for data subset
 selection. These methods aim to identify the most representative training data points. Active learning
 involves iteratively choosing points to label from a large unlabeled dataset based on the model's
 uncertainty or other heuristics, such as the entropy of predicted class probabilities (Sener & Savarese,
 2018; Coleman et al., 2020). Recently proposed data selection methods in continual learning are
 related to these approaches, determining examples to be stored or labeled for ongoing model training
 (Castro et al., 2018; Aljundi et al., 2019).



Figure 2: Illustration of JST framework. In stage 1, we perform data valuation as usual and select relatively low-quality data as the validation set. In stage 2, we perform data valuation on the remaining training samples and the new validation set. In this way, high-quality data stands out as outliers with low value scores.

124 In contrast, coreset construction begins with the entire dataset and tries to select a small subset that 125 encapsulates the essence of the full dataset. The goal is to have the model trained on the subset 126 perform approximately as well as the one trained on the entire dataset (Feldman, 2020; Huang 127 et al., 2021). Many works have tackled this setting to accelerate clustering and learning, proposing 128 coresets in k-means and k-medians clustering (Har-Peled & Kushal, 2005), SVMs (Tsang et al., 129 2005), Bayesian logistic regression (Huggins et al., 2016), and Bayesian inference (Campbell & Broderick, 2019). Recently, newly proposed coreset selection methods have demonstrated efficiency 130 in neural network training (Mirzasoleiman et al., 2020; Killamsetty et al., 2021b;a; Tukan et al., 131 2023). In addition, related coreset techniques have been applied to active learning (Sener & Savarese, 132 2018) and continual learning (Tiwari et al., 2022). 133

134 Data Valuation. Data valuation assesses the contribution of each individual data point to the overall 135 performance of a model, aiming to distribute the validation performance across the training data 136 points (Jiang et al., 2023). Formally, given a training dataset  $D_{tr} = \{z_i\}_{i=1}^N$ , a validation dataset 137  $D_v$ , and a model performance metric *PERF* evaluated on the validation set  $D_v$ , data valuation methods assign a scalar score to each training sample  $z_i$  in  $D_{tr}$  to split *PERF* evaluated on  $D_v$ . For 138 instance, if we define a utility function U over all subsets  $S \subseteq D_{tr}$  of the training data as U(S) :=139  $PERF(\mathcal{A}(S))$  evaluated on the validation set  $D_v$  with a learning algorithm  $\mathcal{A}$ , a straightforward 140 method to evaluate the contribution of a training sample  $z_i$  is to calculate the leave-one-out (LOO) 141 value,  $U(D_{tr}) - U(D_{tr} \setminus \{z_i\})$ . This represents the change in model performance when the point is 142 excluded from the training set (Just et al., 2023). 143

In practice, feasible alternatives to LOO can be used to estimate the impact of the weight change of a 144 data point on the model performance. Two such approaches are optimal transport-based methods 145 (LAVA) (Just et al., 2023) and influence function-based methods (Koh & Liang, 2017). The latter 146 technique employs the Wasserstein distance between training and validation sets as the performance 147 metric, while the influence function-based approach uses the validation loss. Both methods measure 148 the gradient as the scalar value score with respect to the weight change of a data point. In addition, 149 reinforcement learning based-methods (DVRL) (Yoon et al., 2020) learn a weight function to minimize 150 the weighted empirical risk using policy gradients, and so obtain optimal importance weights as the 151 scalar value scores.

152 Indeed, data valuation is intimately related to data selection methods, but it formalizes the data 153 selection problem by aiming to select a subset of the training set that matches a desired target 154 distribution, as represented by the validation set (Just et al., 2023; Yoon et al., 2020), because data 155 valuation methods incorporate validation performance into individual training samples. This approach 156 is essentially different from active learning and coreset construction, which instead seek to compress 157 the full training dataset. Therefore, data valuation can be applied to data selection in the scenario of 158 robust learning with a mismatch between training set and target set (Yoon et al., 2020). However, as demonstrated in Figure 1, data valuation methods are typically more sensitive to outliers rather 159 than high-quality data. This property makes data valuation an ineffective approach when seeking 160 to distinguish high-quality data at a finer-grained scale, especially when there is a small mismatch 161 between the training and target sets. To address this issue, motivated by the property, we propose

a straightforward two-stage framework, called JST, to make data valuation more suitable for data selection. To the best of our knowledge, our work is the first to leverage this property to augment existing data valuation methods.

166 167 3 PRELIMINARIES

168

169

187

191 192 193

194

We briefly detail the components, setup and, notation used to describe our method.

**Data Selection:** We are given a training set  $D_{tr} = \{z_1, \ldots, z_n\}$  containing n data points, where each  $z_i = (x_i, y_i)$  is drawn from a source distribution  $p_{src}(z)$ , as well as a validation set  $D_v = \{z'_1, \ldots, z'_m\}$  with m data points, where each  $z'_i = (x'_i, y'_i)$  is drawn from a target distribution  $p_{trg}(z')$ (typically n > m). Both sets share the same input-output space  $\mathcal{X} \times \mathcal{Y}$  but differ in their underlying distributions, i.e.,  $p_{src}(z) \neq p_{trg}(z')$ .

Given a selection budget k ( $k \le n$ ), the goal of a data selection procedure is to identify a subset  $\hat{D} = {\hat{z}_1, \ldots, \hat{z}_k}$  where  $\hat{D} \subseteq D_{tr}$ , such that the distribution of  $\hat{D}$  closely matches the distribution of the validation set  $D_v$ , to minimize the impact on the learned model. Therefore, the selected subset should approximate the distribution of the validation set, i.e.,  $P(\hat{D}) \approx p_{trg}(z')$ , where  $P(\hat{D})$  is the distribution constructed from  $\hat{D}$ .

**Data Valuation:** Similarly, we are given a training set  $D_{tr} = \{z_1, \ldots, z_n\}$  containing *n* data points  $z_i = (x_i, y_i)$  and a validation set  $D_v = \{z'_1, \ldots, z'_m\}$  with *m* data points  $z'_i = (x'_i, y'_i)$ . As before both sets share the same input-output space  $\mathcal{X} \times \mathcal{Y}$ .

The goal of data valuation is to understand and distribute the validation performance across training data points. To achieve this, we use a function  $\mathcal{V}$  computed over the training set  $D_{tr}$  to find a score vector  $\overline{s} \in \mathbb{R}^n$  that represents the allocation to each data point, described as follows:

$$\overline{s} := \mathcal{V}(D_{tr}, D_v), \text{ where } \overline{s} \in \mathbb{R}^n.$$
 (1)

Given a score vector  $\overline{s} = [s_1, \dots, s_n]$  representing scores  $s_i$  for each data point  $z_i = (x_i, y_i)$  in training set  $D_{tr}$ , we can express the process of ranking (in descending order) and indexing the k points starting from index r as follows:

$$D_{sel} := \mathcal{R}(\overline{s})[r: r+k], \text{ where } |D_{sel}| = k \text{ and } D_{sel} \subseteq D_{tr}.$$
(2)

#### 4 JST: JUST SELECT TWICE

We now present JST, a straightforward two-stage approach to augment existing data valuation methods.
In the first stage of the process, we perform data valuation as usual and select data points with lowest value scores as the validation set. Then, in the second stage of the process, we perform data valuation on remaining training data points and select the data points of low value scores as high-quality data.

# 200 Stage 1 (Data Valuation).

Civen a training set  $D_{tr}$  and a validation set  $D_v$ , let  $D_{sel}$  of size  $|D_v|$  represent the selected data points of lowest value scores obtained by the ranking function  $\mathcal{R}(\cdot)[|D_{tr}| - |D_v| : |D_{tr}|]$  of data valuation  $\mathcal{V}(D_{tr}, D_v)$ .

#### 206 Stage 2 (Data Selection).

Substitute the original validation set  $D_v$  with 207  $D_{sel}$  obtained in the first step and remove 208 data points in  $D_{sel}$  from the training set 209  $D_{tr}$ , denoted as  $D'_{tr} = D_{tr} \setminus D_{sel}$ . Then 210 perform data valuation  $\mathcal{V}(\cdot, \cdot)$  again but in-211 stead on the changed training set  $D'_{tr}$  and 212 validation set  $D_{sel}$  and minus the score vec-213 tor, i.e.,  $\overline{s'} = -\mathcal{V}(D'_{tr}, D_{sel})$ . Finally, the 214 top k high-quality data can be selected by 215  $\mathcal{R}(s')[1:1+k]$ . We summarize our framework in Algorithm 1.

#### Algorithm 1 JST Selection

**Input:** training set  $D_{tr} = \{(x_i, y_i)\}_{i=1}^n$ , validation set  $D_v = \{(x'_i, y'_i)\}_{i=1}^m$ , data valuation  $\mathcal{V}(\cdot, \cdot)$ **Output:** data value score vector  $\overline{s'}$ 

Stage one:

Perform data valuation:

 $\overline{s} \leftarrow \mathcal{V}(D_{tr}, D_v)$ 

Select backward data points:

$$D_{sel} \leftarrow \mathcal{R}(\overline{s})[|D_{tr}| - |D_v| : |D_{tr}|]$$

#### Stage two:

Remove training data points:

 $D'_{tr} \leftarrow D_{tr} \setminus D_{sel}$ 

Perform data valuation again and minus the score vector:

 $\overline{s'} \leftarrow -\mathcal{V}(D'_{tr}, D_{sel})$ 

216 Practical Implementation. In practice, similar to other data valuation and selection methods 217 (Sorscher et al., 2022; Jiang et al., 2023; Just et al., 2023; Xia et al., 2023), we begin by using a 218 neural network trained on the validation set  $D_v$  to extract features into a low-dimensional space. 219 This allows us to leverage feature relationships effectively and compute data value scores efficiently. 220 Given that the manually cleaned validation set  $D_v$  is often small, a ResNet-18 model (He et al., 2016) is an appropriate choice for feature extraction. However, in the second round, prior knowledge from 221 the clean validation set  $D_v$  can erroneously align the new validation set  $D_{sel}$  containing low-quality 222 data with the remaining training samples  $D'_{tr}$ , negatively impacting performance. Given our limited access to noisy data, we find that leveraging a pretrained ResNet-18 model on ImageNet1K (Deng 224 et al., 2009) suffices. 225

To maintain a perfect match between our algorithm and pure data valuation methods, we use the same number  $|D_v|$  of backward training data points  $D_{sel}$  in the second round as the validation set, which is always mixed with a non-obvious proportion of high-quality data. However, we find that behavior does not affect the performance of our algorithm. This is because even a small portion of signals from low-quality data is sufficient to pop out high-quality data with low value scores. Therefore, in practice, we can reduce the size of the validation set  $|D_{sel}|$  to involve more training samples for selection in the second round without degrading performance.



Ablation experiments aimed at these two notions are deferred to Appendix B.

233

255

256

261 262 263

264 265

266 267

268

269

Figure 3: JST framework performance comparison on high-quality data selection, instantiated with optimal transport based method. JST framework shows a notable improvement in selection precision compared with pure data valuation method.



Figure 4: Illustration of the principle of JST framework. In the second-round data selection, highquality data are popped out as outliers with low value scores, demonstrating near-perfect separation.

# 5 EXPERIMENTAL EVALUATION

271 272

In this section, we conduct experiments to verify the effectiveness of our proposed JST framework in 273 augmenting data valuation methods for selecting data that align with the distribution of validation set 274 from training set. We focus on two main tasks: high-quality data selection (Section 5.1) and raw 275 dataset pruning (Section 5.2), using the JST framework instantiated with the optimal transport based 276 data valuation method. Additionally, we demonstrate the general applicability of our JST framework to other data valuation methods, including influence function based and reinforcement learning based approaches (Section 5.4). To cover the complexity of tasks across different scales, we utilize several 278 datasets, including MNIST, Fashion-MNIST (Xiao et al., 2017), CIFAR-10, CIFAR-100 (Krizhevsky 279 et al., 2009), SVHN (Netzer et al., 2011), and Food-101 (Bossard et al., 2014). 280

281 282

283

284

285

286

5.1 HIGH-QUALITY DATA SELECTION

In real-world scenarios, training data often contain corrupted images (Hendrycks & Dietterich, 2019; Li et al., 2020). We evaluate the effectiveness of JST framework in selecting high-quality data in the presence of corrupted images.

Setup. We employ four standard datasets: MNIST, Fashion-MNIST, CIFAR-10, and CIFAR-100. 287 For all four datasets, we inject white noise into 25% of the training set and utilize the clean test set 288 as the validation set. To demonstrate the superior performance of JST framework compared with 289 the pure data valuation method, we add two levels of small noise to the training set as depicted 290 in Figure 3 for noise-level comparative analysis. To achieve the goal of identifying "high-quality" 291 training samples, we rank data points based on their value scores in descending order, selecting a 292 subset with the highest values. For every selection budget, we compute the selection precision, i.e, the 293 percentage of the points that are uncorrupted within the selected points. We compare our framework 294 with pure data valuation method and random selection baseline across these four datasets. 295

**Results.** As depicted in Figure 3, our experimental results show a notable enhancement in selection 296 precision across all four datasets, where our JST framework outperforms the pure data valuation 297 method on nearly all selection budgets. Globally, our framework significantly improves the mean 298 rank of uncorrupted data compared to pure data valuation methods. With a very small noise level of 299 standard deviation 0.6 and 3.0 (Figure 3 A), our framework push forward the mean rank by 3,487 for 300 MNIST, 1,798 for Fashion-MNIST, 700 for CIFAR-10, and 861 for CIFAR-100, from initial values 301 of 24,814, 24,269, 22,212, and 22,106, respectively. This underscores the feasibility of employing 302 data selection using our framework under challenging circumstances with small noise. When facing 303 relatively larger noise with a standard deviation of 1.2 and 6.0 (Figure 3 B), our framework still 304 achieves improvements of 785 for MNIST, 607 for Fashion-MNIST, 843 for CIFAR-10, and 1,167 305 for CIFAR-100, from initial values of 23,064, 23,072, 21,869, and 21,937, respectively. Remarkably, as illustrated in Figure 4, with a noise level of standard deviation 0.6, the data valuation method 306 after being augmented by our framework achieves near-perfect separation between uncorrupted and 307 corrupted data in the simple datasets MNIST and Fashion-MNIST. Furthermore, Figure 4 clearly 308 confirms the feasibility of the principle of our framework, wherein high-quality data are assigned 309 low value scores and popped out as outliers in the second-round data selection, vice versa to the 310 first-round data valuation.

311312313

5.2 RAW DATASET PRUNING

With the help of commercial search engines, several web-crawled large datasets have been curated (Li et al., 2017; Sun et al., 2017). These datasets typically contain heterogeneous noise, such as label noise and out-of-distribution samples, and individual training data points usually vary in quality. We further demonstrate the effectiveness of our JST framework in real-world data collection scenarios.

Setup. We conduct experiments on two web-crawled datasets, SVHN and Food-101, focusing on the task of raw dataset pruning with being provided a small size of manually cleaned validation set.

More specifically, SVHN is a digit classification dataset with 10 categories (from 0 to 9), cropped from pictures of real-world house number plates obtained from Google Street View images. It contains 73,257 digits for training and 26,032 digits for testing. Similarly, but with greater complexity, Food-101 dataset, which was crawled online, consists of 101 food categories with 750 training images and



Figure 5: JST framework performance comparison on raw dataset pruning with different pruning ratio, instantiated with optimal transport based method.

250 test images per category, totaling 101,000 images. The labels for test images have been manually cleaned, while the training set contains some noise, primarily in the form of intense colors, incorrect labels and out-of-distribution samples. We randomly sampled 5,000 digits from the test set of SVHN and 7,500 images from the test set of Food-101 to create their respective validation sets, while leaving the remaining samples to be used as the test sets.

342 Likewise to Section 5.1, we rank training data 343 points based on their value scores in descending 344 order to identify "high-quality" samples. How-345 ever, to mitigate the issue of class imbalance, 346 which negatively impacts the trained model per-347 formance (Johnson & Khoshgoftaar, 2019), we simply rearrange the data points sequentially 348 according to their class one by one after each 349 round of ranking. Then, to evaluate the efficacy 350 of our JST framework, we train a ResNet-18 351 model from scratch on the training set of SVHN 352 and Food-101 with varying pruning ratios, re-353 spectively. Additionally, to explore the effective-354 ness of our framework on pretrained models, we 355 finetune a ResNet-50 model (He et al., 2016), 356 initially pretrained on ImageNet1K, with differ-357 ent pruning ratios. We test all trained neural net-358 works on the reconstructed testing set, excluding samples from the validation set to prevent prior 359 knowledge and maintain a fair accuracy com-360 parison. Still, we compare our JST framework 361 with the pure data valuation method and random 362 selection baseline. 363

For the implementation of neural network training, all experiments are conducted on NVIDIA
H100 GPUs with PyTorch (Paszke et al., 2019).
In training from scratch experiments, we utilize an SGD optimizer with a momentum of 0.9,
weight decay of 5e-4, and an initial learning rate of 0.1 with cosine annealing learning rate decay



Figure 6: Visualization of extremal images for class apple pie in the Food-101 dataset. The top five (most valuable) and bottom five (least valuable) images are selected. For random selection, after shuffling, first five images are shown. Our framework, compared to the pure data valuation and random baseline, selects diverse, relevant images and excludes low-quality, out-of-distribution samples.

strategy over 200 epochs. For experiments on SVHN, we set a batch size of 256; for experiments on Food-101, we set a batch size of 128 and employ data augmentations of random crop and random horizontal flip. In finetuning experiments on Food-101, the experiment setting is the same as training from scratch experiments except the number of epochs, which is set to 40. Across each experiment, we perform three individual runs with different random seeds and report the mean performance and the standard deviation as error bars.

376

334

335 336 337

338

339

340

341

**Results.** As illustrated in Figure 5, in both SVHN and Food-101 experiments, our framework achieves much higher accuracy compared to the pure data valuation method. It also outperforms the

random selection baseline across most pruning ratios especially in cases with large pruning ratios.
Further, in Figure 6, we visualize extremal images of the Food-101 dataset for one class (apple pie)
in the ranking. It demonstrates the efficacy of JST, as our framework selects the most relevant images
to the class, while random selection and the pure data valuation method include monotone image
patterns, low-quality data, or out-of-distribution samples, thereby hurting model performance. More
extremal images are deferred at Appendix F.

384 385

386

5.3 CROSS-DOMAIN RETRIEVAL

In real-world scenarios, datasets often consist of samples from multiple domains with imbalanced
 proportions. Cross-domain retrieval aims to identify and select relevant samples from such a mixed domain dataset to align with the distribution of a given target domain.

Setup. We conduct experiments on one time-series tabular dataset: HHAR (Stisen et al., 2015).
 HHAR dataset is designed to predict human activities within specific time segments, classifying them into six categories: biking, standing, walking, walking upstairs, and walking downstairs.

393 Following the data preprocessing method outlined in previous studies (Ragab et al., 2023; He et al., 394 2023), HHAR dataset is segmented into non-overlapping segments of 128 time steps for classification 395 and then split into training and testing sets as applied in prior research. We further divide HHAR 396 sequential data segments into four domains based on the devices and motion sensors used: phone-397 accelerometer, watch-accelerometer, phone-gyroscope, and watch-gyroscope. To comprehensively 398 evaluate the effectiveness of our JST framework for cross-domain retrieval, we conduct experiments 399 using each of these domains as the validation set, extracted from the testing set. The training set is treated as a mixed-domain dataset. 400

401 Based on previous studies (Kwon & Zou, 2021; Wang & Jia, 2023; Kwon & Zou, 2023), we apply the 402 k-means clustering algorithm to the data values. This approach allows us to partition the training set 403 into two groups: a potential target domain data group and an other data group. Given our expectation 404 that potential target domain data are likely to have high values, we identify the cluster with the higher 405 average data value as the potential target domain data group for retrieval. Following common practice in the literature, we compare the clustering results against the annotated ground truth and calculate 406 the F1-score to evaluate our JST framework. We compare our JST framework with the pure data 407 valuation method. Note that for a fair comparison, we calculate the F1-score using only the data 408 propagated into the second round of our JST framework for the pure data valuation method. 409

410 411

F1-score	Round 1	Round 2
Phone - Accelerometer	$0.458 \pm 0.001$	$0.652 \pm 0.001$
Watch - Accelerometer	$0.212 \pm 0.001$	$0.305 \pm 0.001$
Phone - Gyroscope	$0.347 \pm 0.001$	$0.582 \pm 0.002$
Watch - Gyroscop	$0.511 \pm 0.003$	$\textbf{0.688} \pm \textbf{0.001}$

Table 1: JST framework performance comparison on cross-domain retrieval, instantiated with optimal transport based method. The average and standard error of the F1-score are denoted by 'average±standard error'. All the results are based on 10 repetitions. Boldface numbers denote the best method. JST framework shows a notable improvement in F1-score compared with pure data valuation method.

421

Results. As shown in Table 1, our JST framework consistently outperforms the pure data valuation method, even under conditions of severe domain imbalance. Specifically, in the HHAR dataset, the balance score — calculated as the average proportion of the minority domain to the majority domain when selecting two random domains — is as low as 47%. On average, our JST framework improves the F1-score by 0.175 across four domains in the HHAR dataset.

427 428

429

#### 5.4 APPLYING JST FRAMEWORK TO DIFFERENT DATA VALUATION METHODS

To demonstrate the general applicability of our JST framework, we apply it to other two data valuation
 methods: influence function based and reinforcement learning based approaches. We use dataset
 CIFAR-10 and CIFAR-100 to evaluate the effectiveness of JST. The experimental setup is the same as

460

461

in Section 5.1 with white noise standard deviation of 6.0 and 9.0, except the size of training set and validation set. As these two methods suffer from running time for large dataset, to save computational resources, we randomly sample 10,000 data points from the training set as the training set and 1,000 data points from the test data as the validation set. We use the last layer of ResNet-18 model as their dependent training model. As illustrated in Figure 7, these two methods augmented by our framework, JST, improve selection precision to varying extents, corresponding to higher sensitivity of outliers compared with in-distribution data as shown in Figure 1.



Figure 7: Applying JST framework to different data valuation methods on CIFAR-10 and CIFAR-100. Both methods (A: influence function based, B: reinforcement learning based) augmented by JST framework improve selection precision to varying extents.



Figure 8: A: Evaluating JST framework on marginal contribution based data valuation methods on CIFAR-10. These methods do not improve selection precision compared with the pure data valuation method. B: Sensitivity to outliers compared with in-distribution data for marginal contribution based data valuation methods. Following the experimental setup from Figure 1, we evaluate the difference between NDR<sub>outlier</sub> and NDR<sub>in-distribution</sub> on CIFAR-10. The results show that these methods exhibit similar levels of sensitivity to outliers and in-distribution data, which explains the underperformance of our approach when applied to marginal contribution based methods.

# 486 6 WHERE DOES THE JST FRAMEWORK UNDERPERFORM?

488 Finally, we discuss scenarios where our JST framework underperforms, showing a similar selection 489 precision to the pure data valuation method. Specifically, we observe that marginal contribution based 490 data valuation methods yield unsatisfactory results, as these methods fail to demonstrate a higher 491 detection rate for outliers compared to in-distribution data. We test our framework on three marginal 492 contribution based methods: LOO (Jiang et al., 2023), AME (Lin et al., 2022), and Data Banzhaf (Wang & Jia, 2023). The experimental setup is the same as in Section 5.4 on CIFAR-10 with white 493 noise standard deviation of 10.0, except the reduced size of training set (5,000) and validation set 494 (500) for LOO to save runtime. As shown in Figure 8 A, these three methods augmented by our 495 framework do not improve selection precision. This is because these methods lack the property of 496 higher sensitivity to outliers (Figure 8 B). Based on this underperformance analysis, we propose 497 a simple test framework to validate the expected behavior of our approach: before applying it to 498 real-world applications, we can assess whether the data valuation method exhibits higher sensitivity 499 to outliers on a specific type of data. 500

501 502

# 7 CONCLUSION

503 In this paper, we introduce JST, a straightforward two-stage framework designed to enhance the 504 sensitivity of existing data valuation methods to high-quality data for subset selection. The JST 505 framework identifies low-quality data points as a validation set, thereby allowing high-quality 506 data to stand out as outliers during the data valuation process. Experiments on multiple datasets 507 demonstrate the effectiveness of JST framework in selecting high-quality data and pruning raw 508 datasets, particularly in scenarios with small noise. We successfully apply JST to different data 509 valuation methods, highlighting the general applicability of our framework, thus making it a valuable 510 tool for enhancing data selection in machine learning.

However, our work focuses solely on data valuation and selection based on loss and accuracy, without considering other important aspects such as fairness across subpopulations, which is crucial in data-centric methods. Additionally, adapting our framework to a wider range of data valuation methods and the currently heated LLM model settings remains a practical interest. We leave these two aspects for future work.

516 517 518

527

528

529

530

# References

- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for
   online continual learning. *Advances in neural information processing systems*, 32, 2019.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101–mining discriminative components with random forests. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VI 13*, pp. 446–461. Springer, 2014.
- Trevor Campbell and Tamara Broderick. Automated scalable bayesian inference via hilbert coresets.
   *Journal of Machine Learning Research*, 20(15):1–38, 2019.
  - Francisco M. Castro, Manuel J. Marin-Jimenez, Nicolas Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *Proceedings of the European Conference on Computer Vision* (ECCV), September 2018.
- Cody Coleman, Christopher Yeh, Stephen Mussmann, Baharan Mirzasoleiman, Peter Bailis, Percy
   Liang, Jure Leskovec, and Matei Zaharia. Selection via proxy: Efficient data selection for
   deep learning. In *International Conference on Learning Representations*, 2020. URL https:
   //openreview.net/forum?id=HJg2b0VYDr.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- 539 Dan Feldman. Core-sets: Updated survey. *Sampling techniques for supervised or unsupervised tasks*, pp. 23–44, 2020.

550

551

555

587

- Sariel Har-Peled and Akash Kushal. Smaller coresets for k-median and k-means clustering. In
   *Proceedings of the twenty-first annual symposium on Computational geometry*, pp. 126–134, 2005.
- Huan He, Owen Queen, Teddy Koker, Consuelo Cuevas, Theodoros Tsiligkaridis, and Marinka Zitnik.
   Domain adaptation for time series under feature and label shifts. In *International Conference on Machine Learning*, pp. 12746–12774. PMLR, 2023.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
  - Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- Jiawei Huang, Ruomin Huang, Wenjie Liu, Nikolaos Freris, and Hu Ding. A novel sequential coreset
   method for gradient descent algorithms. In *International Conference on Machine Learning*, pp. 4412–4422. PMLR, 2021.
- Jonathan Huggins, Trevor Campbell, and Tamara Broderick. Coresets for scalable bayesian logistic regression. *Advances in neural information processing systems*, 29, 2016.
- Kevin Fu Jiang, Weixin Liang, James Zou, and Yongchan Kwon. Opendataval: a unified benchmark for data valuation. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL https://openreview.net/forum?id=
   eEK99egXeB.
- Justin M Johnson and Taghi M Khoshgoftaar. Survey on deep learning with class imbalance. *Journal* of Big Data, 6(1):1–54, 2019.
- Hoang Anh Just, Feiyang Kang, Jiachen T Wang, Yi Zeng, Myeongseob Ko, Ming Jin, and Ruoxi Jia.
  Lava: Data valuation without pre-specified learning algorithms. *arXiv preprint arXiv:2305.00054*, 2023.
- Alan F Karr, Ashish P Sanil, and David L Banks. Data quality: A statistical perspective. *Statistical Methodology*, 3(2):137–173, 2006.
- Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer.
   Grad-match: Gradient matching based data subset selection for efficient deep model training. In
   *International Conference on Machine Learning*, pp. 5464–5474. PMLR, 2021a.
- Krishnateja Killamsetty, Xujiang Zhao, Feng Chen, and Rishabh Iyer. Retrieve: Coreset selection for efficient and robust semi-supervised learning. *Advances in neural information processing systems*, 34:14488–14501, 2021b.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In International conference on machine learning, pp. 1885–1894. PMLR, 2017.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation framework
   for machine learning. *arXiv preprint arXiv:2110.14049*, 2021.
- Yongchan Kwon and James Zou. Data-oob: Out-of-bag estimate as a simple and efficient data value. In *International Conference on Machine Learning*, pp. 18135–18152. PMLR, 2023.
  - Junnan Li, Caiming Xiong, and Steven CH Hoi. Mopro: Webly supervised learning with momentum prototypes. *arXiv preprint arXiv:2009.07995*, 2020.
- Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool. Webvision database: Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862*, 2017.
- Weixin Liang and James Zou. Metashift: A dataset of datasets for evaluating contextual distribution shifts and training conflicts. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=MTex8qKavoS.

603

604

605

611

619

634

- Jinkun Lin, Anqi Zhang, Mathias Lécuyer, Jinyang Li, Aurojit Panda, and Siddhartha Sen. Measuring the effect of training data on deep learning predictions via randomized experiments. In *International Conference on Machine Learning*, pp. 13468–13504. PMLR, 2022.
- Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,
  Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training
  group information. In *International Conference on Machine Learning*, pp. 6781–6792. PMLR, 2021.
  - Baharan Mirzasoleiman, Kaidi Cao, and Jure Leskovec. Coresets for robust training of deep neural networks against noisy labels. *Advances in Neural Information Processing Systems*, 33:11465–11477, 2020.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.
   Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, pp. 7. Granada, Spain, 2011.
- Sunday Adewale Olaleye and Akwasi Gyamerah Adusei. The new reality of data economy and
   productization: A conceptual paper. 2024.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
  Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,
  high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Mohamed Ragab, Emadeldeen Eldele, Wee Ling Tan, Chuan-Sheng Foo, Zhenghua Chen, Min Wu,
  Chee-Keong Kwoh, and Xiaoli Li. Adatime: A benchmarking suite for domain adaptation on time
  series data. ACM Transactions on Knowledge Discovery from Data, 17(8):1–18, 2023.
- Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In International Conference on Learning Representations, 2018. URL https://openreview.net/forum?id=HlaIuk-RW.
- Rachael Hwee Ling Sim, Xinyi Xu, and Bryan Kian Hsiang Low. Data valuation in machine
   learning:" ingredients", strategies, and open challenges. In *IJCAI*, pp. 5607–5614, 2022.
- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=UmvSlP-PyV.
- Allan Stisen, Henrik Blunck, Sourav Bhattacharya, Thor Siiger Prentow, Mikkel Baun Kjærgaard,
   Anind Dey, Tobias Sonne, and Mads Møller Jensen. Smart devices are different: Assessing and
   mitigatingmobile sensing heterogeneities for activity recognition. In *Proceedings of the 13th ACM conference on embedded networked sensor systems*, pp. 127–140, 2015.
- Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable
   effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pp. 843–852, 2017.
- Rishabh Tiwari, Krishnateja Killamsetty, Rishabh Iyer, and Pradeep Shenoy. Gcr: Gradient coreset
   based replay buffer selection for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 99–108, 2022.
- Ivor W Tsang, James T Kwok, Pak-Ming Cheung, and Nello Cristianini. Core vector machines: Fast svm training on very large data sets. *Journal of Machine Learning Research*, 6(4), 2005.
- Murad Tukan, Samson Zhou, Alaa Maalouf, Daniela Rus, Vladimir Braverman, and Dan Feldman.
   Provable data subset selection for efficient neural networks training. In *International Conference* on Machine Learning, pp. 34533–34555. PMLR, 2023.
  - Cédric Villani. Topics in optimal transportation, volume 58. American Mathematical Soc., 2021.

- Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 6388–6421.
   PMLR, 2023.
- Kiaobo Xia, Jiale Liu, Jun Yu, Xu Shen, Bo Han, and Tongliang Liu. Moderate coreset: A universal method of data selection for real-world data-efficient deep learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum? id=7D5EECbOaf9.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
   machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Jinsung Yoon, Sercan Arik, and Tomas Pfister. Data valuation using reinforcement learning. In International Conference on Machine Learning, pp. 10842–10851. PMLR, 2020.

APPENDIX

# 704 705

# A DATA VALUATION ALGORITHMS

706 This section provides a detailed explanation of three data valuation algorithms applied in our study: optimal transport based method, influence function based method and reinforcement learning based 708 method. Although we have defined notations in the preliminaries section, a thorough set of notations 709 is presented here for clarity. The input space is denoted by  $\mathcal{X}$  and the output space by  $\mathcal{Y}$ . We denote 710 the training set by  $D_{tr} = \{z_i\}_{i=1}^n$ , where each  $z_i = (x_i, y_i)$  is drawn from a source distribution 711  $p_{src}(z)$ , the validation dataset by  $D_v = \{z'_i\}_{i=1}^m$ , where each  $z'_i = (x'_i, y'_i)$  is drawn from a target 712 distribution  $p_{trg}(z')$ , and a model performance metric *PERF* evaluated on the validation set  $D_v$ . Typically, n > m and  $p_{src}(z)$  is not required to be the same as  $p_{trg}(z')$ , i.e.,  $p_{src}(z) \neq p_{trg}(z')$ . 713 714 The goal of data valuation is to distribute the validation performance across training data points and compute a data score value  $s(z_i)$  for each training data point  $z_i$ . 715

716

718

#### 717 LAVA (OPTIMAL TRANSPORT BASED)

719LAVA (Just et al., 2023) measures the sensitivity of validation performance to changes in the training<br/>data. It examines how the optimal transport cost between the training set  $D_{tr}$  and the validation<br/>set  $D_v$  changes when a particular data point in  $D_{tr}$  is assigned increased weight. The sensitivity is<br/>determined by calculating the gradient of the optimal transport cost with respect to the probability<br/>mass.720721

The optimal transport cost gradient can be calculated as the data value score for the training data point  $z_i$  as follows:

where  $\{t_i^*\}_{i=1}^n$  is the optimal solution of the dual problem, which is expressed as:

 $s(z_i) := t^*[i] - \frac{1}{n-1} \sum_{j \neq i} (t^*[j])$ 

 $t^*, u^* := \arg \max_{t, u \in \mathcal{C}^0(\mathcal{X} \times \mathcal{Y})^2} \left\langle t, \frac{1}{n} \delta_{(x_i, y_i)} \right\rangle + \left\langle u, \frac{1}{m} \delta_{(x'_i, y'_i)} \right\rangle$ 

726

727

728 729

730

731 732

735

736

737

738

745 746 where  $C^0(\mathcal{X} \times \mathcal{Y})$  denotes the set of all continuous functions defined on  $\mathcal{X} \times \mathcal{Y}$ , and  $\delta_{(x,y)}$  represents the Delta measure at  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ . The negative gradient indicates that the loss will decrease when the data point is given more weight, signifying a higher value for the data point, whereas a positive gradient indicates the opposite.

#### 739 740 INFLUENCE FUNCTION

In machine learning, influence function (Koh & Liang, 2017) is used to evaluate the impact of a data point on a model's performance by upweighting a specific training point. We denote the influence of a training data point  $z_i = (x_i, y_i)$  on the loss  $L(z, \theta)$  with respect to the parameters  $\theta \in \Theta$  of a validation data point  $z'_j = (x'_j, y'_j)$  as  $I(z_i, z'_j)$ . This can be obtained by:

$$I(z_i, z'_j) := -\nabla_{\theta} L(z'_j, \hat{\theta})^{\top} \mathcal{H}_{\hat{\theta}}^{-1} \nabla_{\theta} L(z_i, \hat{\theta})$$

where  $\nabla_{\theta} L(z'_j, \hat{\theta})$  represents the gradient of the loss  $L(z'_j, \hat{\theta})$  with respect to optimized parameters  $\hat{\theta}$  evaluated at the validation data point  $z'_j$ . Similarly,  $\nabla_{\theta} L(z_i, \hat{\theta})$  represents the gradient of the loss  $L(z_i, \hat{\theta})$  with respect to optimized parameters  $\hat{\theta}$  evaluated at the training data point  $z_i$ . The term  $\mathcal{H}_{\hat{\theta}}$  denotes the Hessian matrix of empirical risk with respect to the optimized model parameters  $\hat{\theta}$ , defined as  $\frac{1}{n} \sum_{i=1}^{n} \nabla^2_{\theta} L(z_i, \hat{\theta})$ .

754 The negative influence predicts a decrease in loss, while the positive influence predicts an increase. 755 Therefore, higher negative influences correspond to more valuable data points, whereas larger positive influences correspond to less valuable points. To evaluate the impact of each training data point  $z_i$  on the whole valuation set  $D_v$ , we simply sum the influence of the training data point  $z_i$  on each validation data point  $z'_i$  as the data value score  $s(z_i)$ , denoted as:

$$s(z_i) := \sum_{j=1}^m I(z_i, z'_j)$$

DVRL (REINFORCEMENT LEARNING BASED)

DVRL (Yoon et al., 2020) involves using reinforcement learning algorithm to compute the importance weight for each training data point as the data value score. The objective function that DVRL solves in training a model  $g : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ , which maps data points to their importance weight, is expressed as:

$$\min_{g \in G} \mathbb{E}_{(x',y') \sim p_{trg}(z')} [\mathbb{L}(f_g(x'),y')]$$

770 771 772

773 774

775

776

777

778

779

781

782 783

784

785 786

787

788

789

790

791

793

794

796

797

798

799

800

801

s.t.  $f_g = \arg\min_{f \in F} \mathbb{E}_{(x,y) \sim p_{src}(z)}[g(x,y) \cdot \mathbb{L}(f(x),y)]$ 

where  $G := \{g : \mathcal{X} \times \mathcal{Y} \to [0, 1]\}, F := \{f : \mathcal{X} \to \mathcal{Y}\}$  and the loss  $\mathbb{L}$  can be MSE or cross entropy. The data value score of a training data point  $s(z_i)$  is computed as the importance weight  $g(x_i, y_i)$ . The objective can be optimized using policy gradient methods. A large importance weight indicates a more crucial data point for the training process, signifying its high value.

#### **B** PRACTICAL IMPLEMENTATION OF JST

This section we provide the ablation studies of JST framework instantiated by optimal transport based data valuation method on CIFAR-10 dataset with a noise level of standard deviation 6.0.

Performance Comparison on Using Different Feature Extractors 1.00 Round1: trained 0.90 Round2: trained Round2: pretrained 0.85 0.95 0.80 Selection Precision 0.80 0.80 Selection Precision 0.75 0.70 0.65 0.60 Round1: pretrained 0.75 0.55 Round2: pretrained (Round1: pretrained) Round1: trained Round2: pretrained (Round1: trained) 0 70 0 50 50000 20000 30000 ò 10000 20000 30000 40000 ò 10000 40000 50000 Number of Data Selected Number of Data Selected

802 Figure 9: Performance comparison on using different feature extractors in each round. A: In the 803 second-round data selection, we have ResNet-18 model trained on the validation set  $D_{v}$  to extract 804 features compared with a pretrained ResNet-18 model on ImageNet1K. The feature extractor of ResNet-18 model trained on the validation set  $D_v$  hurts the performance in the second-round data 805 selection. B: In both rounds, we have pretrained ResNet-18 model on ImageNet1K to extract features. 806 Our JST framework improves selection precision, though the improvement is less pronounced in the 807 second-round, compared with employing ResNet-18 model trained on the validation set  $D_v$  in the 808 first-round data valuation. This confirms that the improvement in our JST framework is not due to 809 switching to a pretrained ResNet-18 model on ImageNet1K for the second-round data selection.

758 759

760 761 762

763 764

765

766

767

768



Figure 10: Performance comparison on using different sizes of validation set in the second-round data selection. In the second-round, we ablate different sizes of validation set  $|D_{sel}|$  to use, including 5,000, 4,000, 3,000, 2,000, 1,000. The selection precision lines closely align together, indicating the robustness of our JST framework, not requiring many signals from low-quality data to pop out high-quality data.

#### C HOW TO COMPUTE SENSITIVITY COMPARISON?

For in-distribution data detection rate, we rank data points based on their value scores in descending order, selecting a subset with the highest values. For every selection budget, we compute the selection precision, i.e, the percentage of the points that are uncorrupted within the selected points, as indistribution data detection rate. For outlier detection rate, we rank data points based on their value scores in ascending order, selecting a subset with the lowest values. For every selection budget, we compute the selection precision, i.e, the percentage of the points that are corrupted within the selected points, as outlier detection rate. Finally, to account for potential imbalance between corrupted and uncorrupted data, we normalize both detection rates.



Figure 11: Further visualization of JST framework performance comparison on raw dataset pruning. We plot the differences on accuracy between our framework and both the pure data valuation method and the random selection baseline, highlighting the improvements achieved by JST framework.

# E RUNTIME COMPARISON



Figure 12: Runtime comparison on CIFAR-10 with varying sizes of training set. We use the same experimental setup as described in Section 5.1 but vary the dataset size. We observe that our JST framework approximately doubles the runtime across all methods, as it does not require additional overhead. The optimal transport based method stands out for its exceptional efficiency, completing in just a few seconds while exhibiting favorable near-linear time complexity (Just et al., 2023).

JST on Datasets	Time
SVHN - Round 1	9.9 sec.
SVHN - Total	21.6 sec.
Food101 - Round 1	44.4 sec.
Food101 - Total	90.7 sec.

Table 2: Runtime comparison for the raw dataset pruning task. Leveraging the exceptional efficiency of the optimal transport based method, our JST framework achieves data valuation in remarkably short time.

# F EXTREMAL IMAGES





