

# GoodCore: Data-effective and Data-efficient Machine Learning through Coreset Selection over Incomplete Data

CHENGLIANG CHAI, Beijing Institute of Technology, China

Given a dataset with incomplete data (e.g., missing values), training a machine learning model over the incomplete data requires two steps. First, it requires a data-effective step that cleans the data in order to improve the data quality (and the model quality on the cleaned data). Second, it requires a data-efficient step that selects a core subset of the data (called coresets) such that the trained models on the entire data and the coresets have similar model quality, in order to improve the training efficiency. The first-data-effective-then-data-efficient methods are too costly, because they are expensive to clean the whole data; while the first-data-efficient-then-data-effective methods have low model quality, because they cannot select high-quality coresets for incomplete data.

In this paper, we investigate the problem of coresets selection over incomplete data for data-effective and data-efficient machine learning. The essential challenge is how to model the incomplete data for selecting high-quality coresets. To this end, we propose the GoodCore framework towards selecting a good coresset over incomplete data with low cost. To model the unknown complete data, we utilize the combinations of possible repairs as possible worlds of the incomplete data. Based on possible worlds, GoodCore selects an expected optimal coresset through gradient approximation without training ML models. We formally define the expected optimal coreset selection problem, prove its NP-hardness, and propose a greedy algorithm with an approximation ratio. To make GoodCore more efficient, we further propose optimization methods that incorporate human-in-the-loop imputation or automatic imputation method into our framework. Experimental results show the effectiveness and efficiency of our framework with low cost.

CCS Concepts: • Computer systems organization → Embedded systems; Redundancy; Robotics; • Networks → Network reliability.

Additional Key Words and Phrases: datasets, neural networks, gaze detection, text tagging

## ACM Reference Format:

Chengliang Chai. 2018. GoodCore: Data-effective and Data-efficient Machine Learning through Coreset Selection over Incomplete Data. *J. ACM* 37, 4, Article 111 (August 2018), 28 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Data-effective machine learning (ML) (a.k.a. data-centric AI [61]) aims at obtaining high-quality training data to release the value of AI, because it is well-known that dirty data may severely degrade the performance of ML models [20, 60].

Data-efficient ML focuses on making the training process more efficient. A commonly used strategy is to select a core subset of training data (or coresets) [25, 57] to represent the entire dataset such that ML models trained on the coresets can achieve similar performance to the ML models trained on the entire dataset.

---

Author's address: Chengliang Chai, ccl@bit.edu.cn, Beijing Institute of Technology, China.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 0004-5411/2018/8-ART111

<https://doi.org/XXXXXXX.XXXXXXX>

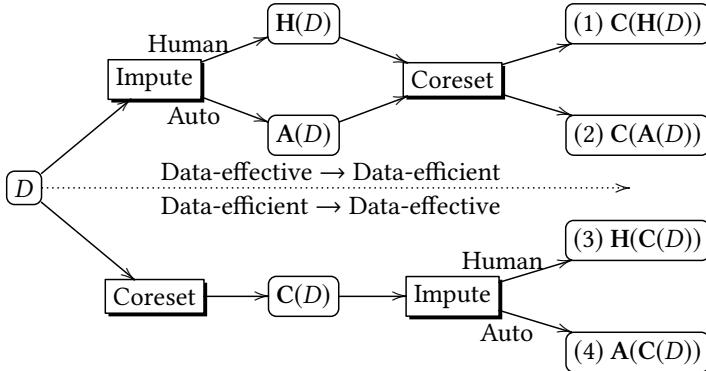


Fig. 1. Sequential methods.

Solution	Accuracy	Human Cost	Machine Cost
(1) C(H(D))	High	High	Low
(2) C(A(D))	Low	None	Low
(3) H(C(D))	Low	Low	Low
(4) A(C(D))	Low	None	Low
Our goal	High	None or Low	Low
(5) H(G(D))	High	Low	High
(6) A(G(D))	Medium	None	High
(7) G(D, ∪ <sup>H</sup> )	High	Low	Low
(8) G(D, ∪ <sup>A</sup> )	Medium	None	Low

Fig. 2. A comparison of different approaches (1–4: sequential methods; 5–8: our solutions).

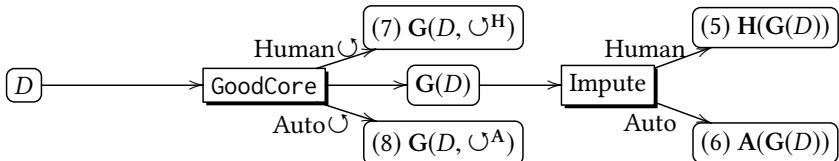


Fig. 3. Our proposal and its variants.

Apparently, users desire both data-effective ML (for training better ML models) and data-efficient ML (for saving training cost). In this work, our main goal is to support both data-effective and data-efficient ML over *incomplete data* where there are many missing values, which is very common in real-world scenarios [20, 52, 72].

**Running data-effective and data-efficient tools sequentially.** Intuitively, we can either run data imputation methods first for data-effective and then run coresnet selection algorithms denoted by  $C(\cdot)$  for data-efficient, or vice versa. Moreover, for data-effective solutions through data cleaning, we generally consider two cases, either human-based solutions denoted by  $H(\cdot)$  or automatic solutions denoted by  $A(\cdot)$ . In summary, we have the following four cases, as shown in Figure 1:

- First data-effective (impute) then data-efficient (coreset):

- (1) Impute-Human:  $H(D) \rightarrow \text{Coreset}: C(H(D))$
- (2) Impute-Auto:  $A(D) \rightarrow \text{Coreset}: C(A(D))$

- First data-efficient (coreset) then data-effective (impute):

- (3) Coreset:  $C(D) \rightarrow$  Impute-Human:  $H(C(D))$
- (4) Coreset:  $C(D) \rightarrow$  Auto-Human:  $A(C(D))$

Next let's discuss the pros and cons of the above approaches.

Case (1) has high human cost, low machine cost, and high accuracy in terms of the trained ML models. Case (2) has zero human cost, low machine cost, but with low accuracy because automatic imputation may not be good enough. Case (3) has low human cost, low machine cost, but with low accuracy because coreset selection over a dirty dataset may not ensure to compute a “good” coreset. Case (4) has no human cost, low machine cost, but with low accuracy with the similar reason as (3). The comparison of the above four methods can be found in Figure 2.

**Our goal.** Clearly, a primary goal is to achieve high accuracy for ML models, where only case (1) can achieve. Case (2) achieves low accuracy because automatic imputation is hard to be accurate. The main obstacle for making (1) practical is its high human cost. Hence, our *main goal* is to achieve high accuracy with no or low human cost, and with low machine cost.

Consider cases (3) and (4), the main reason for them to achieve low accuracy is because they cannot compute a good coreset directly from the dirty data. Intuitively, if we can compute a good coreset directly from the dirty data, we can cheaply clean the coreset to achieve high accuracy, where the “goodness” means that the subset of tuples selected from the dirty data is similar to the subset of tuples selected from the clean data.

**Challenge.** The main challenge of computing a good coreset from dirty data is to accurately estimate the ground truth of each missing value; otherwise, we cannot select a coreset to well represent the clean data. This is a known hard problem because each missing value may have multiple possible repairs. Also, because a coreset selection algorithm is typically iterative that each tuple is selected per iteration [54], selecting a bad tuple may cause cascade amplification to the following iterations, resulting in a bad coreset.

**Our proposal.** To tackle the above challenge, we model the combinations of possible repairs as possible worlds of the original dirty data  $D$ . We then formulate it as an optimization problem for selecting an expected optimal coreset that can represent the possible worlds of  $D$  via gradient approximation without training in advance. We prove this problem to be NP-hard. We propose an approximate algorithm, called GoodCore, denoted by  $G(\cdot)$ , with the main idea to iteratively add a tuple with the highest utility into the coreset. After a good coreset is computed, we can either use human imputation or automatic imputation to impute the data, as shown in Figure 3. We further elaborate these two methods below:

- (5) GoodCore:  $G(D) \rightarrow$  Impute-Human:  $H(G(D))$
- (6) GoodCore:  $G(D) \rightarrow$  Impute-Auto:  $A(G(D))$

However, one main drawback is that modeling possible worlds of  $D$  is computationally expensive, which hinders the practicability of the GoodCore algorithm. To address this high computational cost problem, we further propose **optimization** methods to integrate imputation-in-the-loop (with either humans or automatic methods) into the GoodCore algorithm (see methods 7 and 8 in Figure 3). To this end, the optimized algorithms can significantly reduce the number of possible worlds, thus achieving low computational cost.

- (7) GoodCore with human-in-the-loop imputation:  $G(D, \cup^H)$
- (8) GoodCore with machine-in-the-loop imputation:  $G(D, \cup^A)$

A comparison of methods (5)–(8) is given in Figure 2. Note that method (7) is the best solution because it can achieve a high ML accuracy with low human cost and low machine cost.

**Contributions** We make the following contributions.

- (i) *Two birds with one stone.* We study the problem of solving both data-effective and data-efficient ML in one framework, which is an important but not addressed problem. (Section 3)
- (ii) *NP-hardness and approximate solutions.* We prove the NP-hardness of the problem. We propose a greedy algorithm with an approximate ratio. (Section 4)
- (iii) *Optimizations.* We develop optimization techniques that integrate imputation-in-the-loop into the coresset selection process, to improve the efficiency while achieving high accuracy. We also analyze the convergence rate of our method and theoretically prove that it can converge fast. (Section 5)
- (iv) *Experiments.* We conduct extensive experiment on 6 real-world datasets and compare with 10 baselines to show that GoodCore can select a well-performed coresset to achieve both data-effective and data-efficient ML while consuming a low human cost. (Section 7)

## 2 BACKGROUND OF CORESET SELECTION

In this section, we introduce the background of coresset selection on complete data, denoted by  $D_c$ .

### 2.1 Gradient Descent for Machine Learning

**Gradient descent** [41] is the most typical optimization algorithm to train ML models. At a high level, it tweaks the parameters iteratively to minimize a given convex and differentiable function to its local minimum.

Let  $D_c = \{t_1, t_2, \dots, t_n\}$  be a set of train tuples (without missing values), where  $t_i = (\mathbf{x}_i, \mathbf{y}_i)$ ,  $\mathbf{x}_i \in \mathbb{R}^d$  denotes the vector of features and  $\mathbf{y}_i$  denotes the corresponding label. The goal of training on  $D_c$  is to find the best parameter  $\theta^*$  of a model by minimizing the loss:

$$\theta^* = \arg \min_{\theta \in \vartheta} f(\theta), f(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta, t_i) \quad (1)$$

where  $\vartheta$  is the parameter space. For ease of representation, we abbreviate  $f_i(\theta, t_i)$  as  $f_i(\theta)$  to represent the loss of the  $i$ -th train example. Generally speaking, the gradient descent approach is always applied to find the minimizer of Eq. 1, where the **full gradient** (sum of the gradients over all training tuples), denoted by  $\nabla f(\theta) = \sum_{i=1}^n \nabla f_i(\theta)$ , has to be computed iteratively.

Besides incremental gradient methods like stochastic gradient descent (SGD) that can be leveraged to accelerate the iterative gradient computation, there are other popular and orthogonal methods, such as coresets, which will be discussed next.

### 2.2 Coreset over Complete Data

**Coreset.** To make training more efficient, instead of learning from entire  $D_c$ , one research question is that whether we can compute a small subset  $C(D_c)$  of  $D_c$  such that learning with  $C(D_c)$  can hopefully achieve the same performance as learning with  $D_c$ . This small selected subset is called **coreset** [25, 57]. In the following, we simply write  $C(D_c)$  as  $C$  when it is clear from the context.

The state-of-the-art coresset selection solutions are mostly based on gradient approximation [36, 54]. Suppose that  $\theta$  denotes the parameter of an ML model trained over the full dataset, and  $\theta'$  denotes the parameter of the same model trained over the coresset. Intuitively, the objective of gradient approximation for coresset selection is to make  $\nabla f(\theta')$  as close as possible to  $\nabla f(\theta)$ . To this end, existing solutions focus on *selecting the coresset that minimizes the upper bound of gradient approximation error* ( $\|\nabla f(\theta) - \nabla f(\theta')\|$ ). Next, let's formally define it from scratch.

**Gradient-based coresset selection** is to minimize the **gradient approximation error (GA error)** between the full gradient w.r.t.  $D_c$  and the weighted sum of gradients w.r.t. the coresset  $C$

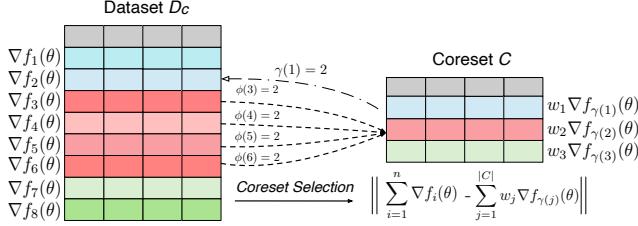


Fig. 4. Example of coresset selection.

(or coresset gradient). Formally, Eq. 2 tries to minimize the GA error by considering all possible parameters  $\theta \in \vartheta$  (*i.e.*,  $\max_{\theta \in \vartheta}$ ), where “ $\|\cdot\|$ ” denotes the normed difference. Next, we introduce the coresset gradient.

$$\begin{aligned}
 C^* = \arg \min_{C \subseteq D_c, w_j \geq 0} & \max_{\theta \in \vartheta} \left\| \underbrace{\sum_{i=1}^n \nabla f_i(\theta)}_{\text{full gradient}} - \underbrace{\sum_{j=1}^{|C|} w_j \nabla f_{\gamma(j)}(\theta)}_{\text{coreset gradient}} \right\|, \\
 \text{s.t. } & |C| \leq K
 \end{aligned} \tag{2}$$

**gradient approximation error**

Because the coresset is a subset of the complete dataset (*i.e.*,  $C \subseteq D_c$ ), we use  $\gamma(j) = i$  (where  $j \in [1, |C|]$ ,  $i \in [1, n]$ ) to denote that the  $j$ -th tuple in  $C$  (denoted by  $c_j$ ) is the  $i$ -th tuple in  $D_c$ , *i.e.*,  $t_i$ . In other words,  $\gamma$  is an index mapping from  $C$  to  $D_c$ .

Recall that the key idea of the coresset is to *use a subset of tuples to represent the entire set*. Eq. 2 potentially contains another important mapping  $\phi$  from  $D_c$  to  $C$  to indicate this, *i.e.*,  $\phi(i) = j$ ,  $i \in [1, n]$ ,  $j \in [1, |C|]$ , which is highly related to the weight. Specifically, let  $\phi(i) = j$  denote that we will assign  $t_i$  to  $c_j$  (use  $c_j$  to represent  $t_i$ ) and use  $\nabla f_{\gamma(j)}$  to represent  $\nabla f_i$ . Each  $t_i$  will be assigned to one and only one  $c_j$ , but each  $c_j$  might be assigned with multiple tuples in  $D_c$ . Based on  $\phi$ ,  $w_j$  is defined as the weight of the  $c_j$ , which is the number of tuples in  $D_c$  assigned to the  $c_j$ , *i.e.*,  $w_j = |\{t_i | \phi(i) = j, i \in [1, n]\}|$  ( $c_j$  is utilized to represent  $w_j$  tuples in  $D_c$ ).

Next let's use an example to better illustrate Eq. 2.

**EXAMPLE 1.** Let's consider a case of the gradients of each tuple, as shown in Figure 4. Suppose that for any  $\theta$ ,  $\nabla f_1(\theta) \approx \nabla f_2(\theta)$ ,  $\nabla f_3(\theta) \approx \nabla f_4(\theta) \approx \nabla f_5(\theta) \approx \nabla f_6(\theta)$  and  $\nabla f_7(\theta) \approx \nabla f_8(\theta)$ . In this case, based on Eq. 2, if one wants to find an optimal coresset with a size of 3, *i.e.*,  $K = 3$ , the solution can be  $C^* = \{t_2, t_5, t_7\}$  ( $\gamma(1) = 2, \gamma(2) = 5$  and  $\gamma(3) = 7$ ), associated with  $w_1 = 2, w_2 = 4, w_3 = 2$  because  $\phi(1) = \phi(2) = 1, \phi(3) = \phi(4) = \phi(5) = \phi(6) = 2$  and  $\phi(7) = \phi(8) = 3$ . In this way,  $C^*$  can be one of the optimal coressets that can well approximate the full gradient because  $\left\| \sum_{i=1}^8 \nabla f_i(\theta) - \sum_{j=1}^3 w_j \nabla f_{\gamma(j)}(\theta) \right\|$  is minimized, which is close to 0.

**Key observation.** We can observe from Example 1 that in order to minimize the GA error, we should set  $\phi(i) = j$ , where  $\nabla f_i$  and  $\nabla f_{\gamma(j)}$  are likely to be close. Therefore, computing the coresset is similar to computing the  $K$  exemplars [64] of the gradients, if all the gradients of tuples can be computed.

**Upper bound minimization of GA error.** We can see from Eq. 2 that to solve the equation, the gradients have to be computed, which have a close relationship with the parameter  $\theta$ . However, the

main bottleneck is that the entire parameter space  $\vartheta$  is too expensive to explore. Hence, a typical solution is to first compute the upper bound of GA error (Eq. 3), then generalize [7, 30, 54] the upper bound computation to the entire parameter space (Eq. 4), and finally select the coresets to minimize the bound. To be specific, using the triangle equation, for any particular  $\theta$ , we have:

$$\left\| \sum_{i=1}^n \nabla f_i(\theta) - \sum_{j=1}^{|C|} w_j \nabla f_{Y(j)}(\theta) \right\| \leq \sum_{i=1}^n \|\nabla f_i(\theta) - \nabla f_{Y(\phi_\theta(i))}(\theta)\| \quad (3)$$

Together with the aforementioned observation, given a coresets  $C$ , the upper bound is minimized when  $\phi$  assigns every tuple  $t_i$  to the tuple in  $C$  with most gradient similarity, i.e.,  $\left\| \sum_{i=1}^n \nabla f_i(\theta) - \sum_{j=1}^{|C|} w_j \nabla f_{Y(j)}(\theta) \right\| \leq \sum_{i=1}^n \min_{c_j \in C} \|\nabla f_i(\theta) - \nabla f_{Y(j)}(\theta)\|$ .

**For the entire space**  $\vartheta$ , it has been proved in recent works [7, 30, 54] that for convex ML problems (corresponding to an optimization problem in which the objective function is a convex function), the normed gradient difference between tuples can be efficiently bounded by:

$$\forall i, j, \max_{\theta \in \vartheta} \|\nabla f_i(\theta) - \nabla f_j(\theta)\| \leq \max_{\theta \in \vartheta} O(\|\theta\|) \cdot \|\mathbf{x}_i - \mathbf{x}_j\| \quad (4)$$

where  $\|\mathbf{x}_i - \mathbf{x}_j\|$  is the Euclidean distance between feature vectors of two tuples, namely *feature distance*, and  $O(\|\theta\|)$  is a constant. Hence, we can conclude that **GA error can be bounded independent of the optimization problem in practice**, i.e., any particular  $\theta$ . Finally, considering Eq. 3 and Eq. 4 together, *the coreset selection problem can be converted to*:

$$C^* = \arg \min_{C \subseteq D_c} \sum_{i=1}^n \min_{c_j \in C} s_{ij}, \text{ s.t. } |C| \leq K \quad (5)$$

where  $s_{ij} = \|\mathbf{x}_i - \mathbf{x}_{Y(j)}\|$  for ease of representation. The above equation indicates that given a train data  $D_c$  and a coresets  $C$ , we use  $S = \sum_{i=1}^n \min_{c_j \in C} s_{ij}$  to score the coresets. The lower the score, the smaller upper bound of the GA error we can get, which indicates a better coresets. To summarize, solving Eq. 5 is to minimize the upper bound of the GA error (i.e., select the coresets with the lowest score) by just considering the feature vectors of the training tuples without training in advance.

Note that Eq. 4 holds for tuples associated with the same label [7, 30]. Therefore, in practice, we respectively select coresets for tuples with different labels and combine them. Suppose that we aim to select a coresets with size  $K$  for a binary classification task (label 1: 60%, label 0: 40%), so we select a coresets with size  $60\%K$  for tuples with label 1 and another one with  $40\%K$  for tuples with label 0.

**Our scope.** In this paper, we focus on the convex problems (*e.g.*, logistic regression, support vector machine, etc.) because for such problems the gradient difference can be well bounded by the difference between feature vectors. Note that, for other ML algorithms such as deep neural networks, they can also be trained using selected coresets to achieve good training accuracy (see Section 7 for our experimental findings).

### 3 CORESET OVER INCOMPLETE DATA

In this section, we will formally define the problem of coresets selection over incomplete data (Section 3.1) and then describe our proposed framework to solve the problem (Section 3.2).

#### 3.1 Problem Definition

As discussed above, we have to compute the coresets score  $S$ , so as to produce a good coresets. To this end, the feature distances can be computed as a pre-processing step, based on which the coresets

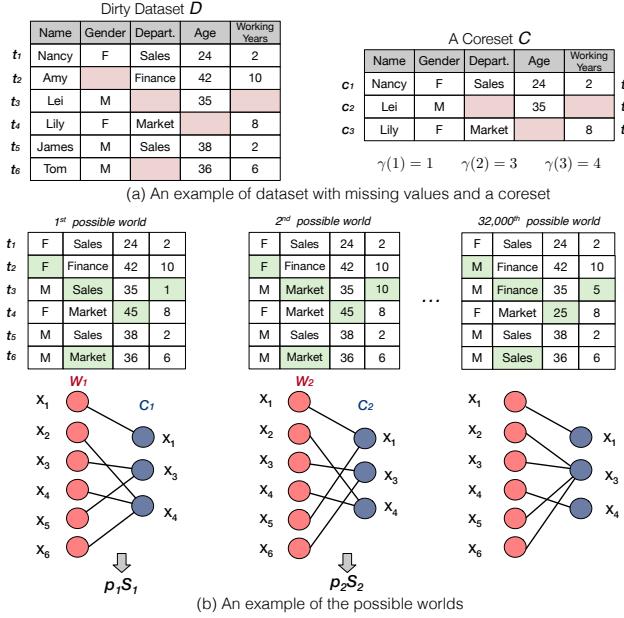


Fig. 5. Example of coresset selection with missing values

score can be computed. However, when there exists incomplete data with missing values, even the feature distances are hard to compute accurately, let alone selecting a proper coresset.

**Incomplete data.** Formally, suppose that  $D$  has  $M$  attributes, denoted by  $\{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_M\}$ . Each attribute  $\mathcal{A}_m, m \in [1, M]$  represents a domain set including the Null, (i.e.,  $\text{Null} \in \mathcal{A}_m$ ), in which each tuple in  $D$  can take value on this attribute.  $|\mathcal{A}_m|$  denotes the domain size. Then, each tuple  $t_i \in \mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_m$ . Let  $t_i[m]$  denote the value of the  $m$ -th attribute of  $t_i$ , i.e.,  $t_i[m] \in \mathcal{A}_m$ .

For a tuple  $t_i \in D$ , if  $\exists t_i[m] = \text{Null}, m \in [1, M]$ ,  $t_i$  is an incomplete tuple, denoted by  $\mathbb{I}[t_i] = 1$ , otherwise  $\mathbb{I}[t_i] = 0$ . Let us better illustrate this using an example.

**EXAMPLE 2.** As shown in Figure 5(a), there are 6 tuples in the table  $D$  with five attributes (an excerpt from a large table). For example,  $\mathcal{A}_2$  is the Gender attribute, i.e.,  $\mathcal{A}_2 = \{M, F, \text{Null}\}$ . Among these tuples,  $t_2, t_3, t_4, t_6$  have missing values, e.g.,  $\mathbb{I}[t_2] = 1, \mathbb{I}[t_1] = 0$ . Given a coresset as shown on the right side, if there are no missing values, we can assign each tuple  $t_i \in D$  to its most similar tuple in  $C$  (compute  $\min_{c_j \in C} s_{ij}$ ), and then sum these feature distances up to compute the coresset score  $S$ . However, given these missing values, the feature distances cannot be computed accurately (e.g.,  $s_{12}, s_{13}, s_{22}$ , etc.), and thus the assignment of tuples in  $D$  cannot be determined precisely. Hence, the coresset score is not precise, and thereby leads to a coresset that cannot well represent the full complete (clean) data.

As discussed above, *imputation before coresset selection* suffers from either large cost (human imputation) or large number of possible repairs (automatic imputation), while *imputation after coresset selection* cannot obtain a good coresset because of the inaccurate feature distance computation (see Example 2).

Therefore, an essential problem is to select a good coresset that can represent the complete dataset  $D_c$ , which relies on accurate coresset score computation given  $D_c$  that is the unknown ground truth.

Fortunately, the possible repairs of  $D$  can be modeled by possible worlds [9–11, 22], based on which we can effectively select the coresets over incomplete data.

**Possible worlds.** Given the incomplete dataset  $D$ ,  $\forall t \in D$  and  $\mathbb{I}[t] = 1$ ,  $\forall t[m] = \text{Null}$ ,  $m \in [1, M]$ , we assign a value in  $\mathcal{A}_m \setminus \{\text{Null}\}$  to  $t[m]$  as an imputation (a.k.a. a possible repair). Thus, we have an assignment for all the missing values in  $D$ , which corresponds to a possible world  $W$ . Since there exist a large number of possible assignments, we define the set of possible worlds as  $\mathcal{I}_W = \{W_k | k \in [1, |\mathcal{I}_W|]\}$ .

Let us better illustrate this using an example.

**EXAMPLE 3.** Given  $D$ , for tuples  $t_2, t_3, t_4, t_6$  with missing values, we have a large number of possible assignment as shown in Figure 5(b), each of which corresponds to a possible world (we omit the Name attribute because there is no missing value on this attribute). Suppose that there are 2 (4/100/10) types of values of the attribute Gender (Department/Age/Working years), there exist 32,000 possible worlds in total.

Note that for numerical attributes, we will bin them into different buckets, such that we can treat them as categorical values and avoid the unlimited number of possible worlds.

Even with possible worlds, the score computation of coresets remains challenging. Each possible world of  $D$  is a complete dataset, and thus given a coreset, the score can be directly computed considering the feature distances, as discussed in Section 2.2. However, the crucial issue is that each possible world could be the ground truth, i.e.,  $D_c$ , but each one leads to a different score.

**EXAMPLE 4.** As shown in Figure 5(b), the two possible worlds  $W_1$  and  $W_2$  are only different in  $t_3$ , leading to a different feature vector  $\mathbf{x}_3$ , which makes the score computation a difference. To be specific, given the same coreset  $C$  with tuples  $t_1(c_1), t_3(c_2)$  and  $t_4(c_3)$ , because of a different  $\mathbf{x}_3$ , the closest feature distance of  $\mathbf{x}_5$  in  $W_2$  becomes  $\mathbf{x}_1$ , rather than  $\mathbf{x}_3$  in  $W_1$ . And the closest feature distance of  $\mathbf{x}_6$  in  $W_2$  becomes  $\mathbf{x}_3$ , rather than  $\mathbf{x}_4$  in  $W_1$ . Therefore, the coreset scores, i.e., the sum of these closest feature distances of tuples are different among possible worlds.

Example 4 shows that different possible worlds make the mapping  $\phi$  different, which leads to different scores. Hence, to get a good coreset without the ground truth, an intuitive solution is to compute the expected coreset score considering all possible worlds. By doing so, although we cannot get the complete data ( $D_c$ ) in advance, we can focus on how to select an informative coreset that can represent the possible worlds of  $D$ .

Next, we formally define the studied problem.

**Expected optimal coreset selection over incomplete data.** Given  $D$ , we have a number of possible worlds  $\mathcal{I}_W = \{W_k\}$ . Then given a subset (coreset)  $C \subset D$ , for different  $W_k$ , we have the corresponding  $C_k$  with the same tuples as  $C$  but probably different imputations. For  $C_k$ , we can compute a score  $S_k = \sum_{i=1}^n \min_{c_j \in C_k} s_{ij}$ , where  $s_{ij} = \|\mathbf{x}_i - \mathbf{x}_{Y(j)}\|$  and both feature vectors are from  $\{W_k\}$ . Then, we have the expectation  $E[C] = \sum_{k=1}^{|\mathcal{I}_W|} p_k S_k$ , where  $p_k$  denotes the probability of the appearance of  $\{W_k\}$ . Finally, our problem becomes how to compute the coreset  $C$  with the lowest expectation of GA error upper-bound. Formally, we have

$$C^* = \arg \min_{C \subset D} E[C], \text{ s.t. } |C| \leq K \quad (6)$$

For example, given  $D$ , the corresponding possible worlds and a coreset  $C$  in Figure 5, we have different  $C_k$  with the same tuples (containing  $t_1, t_3, t_4$ ) but probably different imputations. For each  $C_k$ , we will compute  $S_k$ , and finally compute  $E[C]$ . Solving Eq. 6 can result in an informative coreset with incomplete tuples being selected. After these tuples imputed by a human, i.e., Case (5), or state-of-the-art automatic method, i.e., Case (6), we can derive a good coreset.

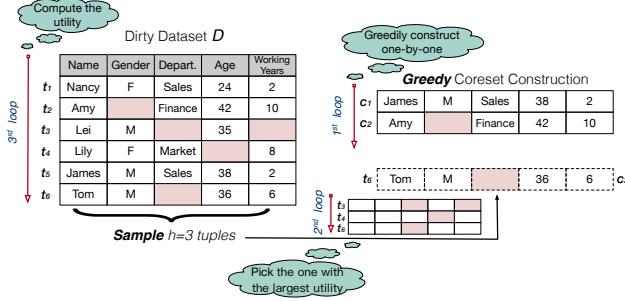


Fig. 6. The GoodCore framework.

**Algorithm 1:** GoodCore Framework

---

**Input:** Incomplete train data  $D$ , coreset size  $K$ , sample size  $h$ .  
**Output:** A coreset  $C \subseteq D$ , weight  $\mathbb{W} = \{w_j\}, |C| = |\mathbb{W}| = K$ .

```

1  $C = \emptyset;$ 
2 while  $|C| < K$  do
3   /*1st loop*/
4   Sample  $h$  tuples as  $T_{sample} \subseteq D \setminus C$ 
5   for each tuple  $t \in T_{sample}$  do
6     /*2nd loop*/
7      $E[t|C] = \text{ComputeUtility}(t, C, D);$  /*3rd loop*/
8      $t^* = \arg \max_{t \in T_{sample}} E[t|C];$ 
9      $C = C \cup \{t^*\};$ 
10  for  $t \in C$  do
11    if  $\mathbb{I}[t] = 1$  then
12      Impute  $t$  by a human or automatic method.
13  for  $j = 1$  to  $|C|$  do
14    for  $i = 1$  to  $n$  do
15      if  $c_j = \arg \min_{c_{j'} \in C} \max_{\theta \in \mathcal{S}} \|\nabla f_i(\theta) - \nabla f_{Y(j')}(\theta)\|$  then
16         $w_j += 1;$ 
17 return  $C, \mathbb{W};$ 

```

---

### 3.2 Goodcore Framework

Next, we will introduce our proposed GoodCore framework to solve Eq. 6, which is non-trivial because it is NP-hard. But fortunately, we prove that it has the sub-modular property (see Section 4). Hence, GoodCore uses a greedy framework with three loops to solve the problem with an approximate ratio.

At a high level, the greedy strategy adds one tuple with the largest “utility” to the coreset iteratively, which can be considered as the first loop. In each iteration, we have to iterate tuples in  $D$  to select the one with the largest utility, which is the second loop. Naturally, we have to compute the utility of each tuple, where all tuples in  $D$  have to be considered, leading to the third loop.

Next, we will further illustrate the framework using Figure 6 and Algorithm 1.

**The first loop (lines 3-17)** of the greedy algorithm is to add the tuple  $t^*$  with the maximum *utility* (*i.e.*,  $E[t|C] = E[C] - E[C \cup \{t\}]$ ) into the coresset iteratively for  $K$  times. To be specific, the “utility” of a tuple  $t$  denotes the reduction of expectation of GA error after adding  $t$  into the coresset  $C$ .

Suppose that  $K = 3$ . Figure 6 (the 1<sup>st</sup> loop part) shows the situation that there already have been 2 tuples in  $C$ , and we are going to add the third tuple into the coresset.

**The second loop (lines 6-7)** computes the utilities of tuples that are not in coresset  $C$ , based on which the best one is picked for the first loop. An ideal solution is to consider all tuples in  $D \setminus C$ , which is prohibitively expensive, so in practice we use an efficient method to accelerate this loop by uniformly sampling  $h$  tuples as  $T_{sample}$  (line 5) and then selecting the best one from  $T_{sample}$  (line 14). The difference is that theoretically, considering all tuples has an approximate ratio  $1 - \frac{1}{e}$  (because of the sub-modular property), while the sampling method holds a  $(1 - \frac{1}{e} - \epsilon)$  ratio [53], where  $\epsilon$  is related to the sampling ratio.

As shown in Figure 6, suppose that  $h = 3$ , and we sample  $\{t_3, t_4, t_6\}$  from  $\{t_1, t_3, t_4, t_6\}$ . Then the second loop iterates the three tuples and computes the utility for each one (the third loop).

**The third loop (line 7)** will loop through all tuples in  $D$ , so as to compute the utility of tuple  $t$  used in the second loop. To be specific, the core part of the utility computation (*i.e.*, `ComputeUtility`) is to compute  $E[C] = \sum_{k=1}^{|I_W|} p_k S_k = \sum_{k=1}^{|I_W|} p_k (\sum_{i=1}^n \min_{c_j \in C_k} s_{ij})$ , from which we can see that it is inevitable to iterate the  $n$  tuples in  $D$ . However, the most challenging part is that we also have to enumerate a large number of possible worlds. We will illustrate how to solve this in details in Section 4.

**The imputation step (line 16).** After GoodCore selects the coresset  $C$  using the above 3 loops, we can leverage a human or automatic method to impute the tuples that are incomplete in  $C$ , which correspond to Case (5) and Case (6) in Section 1 respectively.

**Weights computation (lines 18-21).** It computes the weight of each tuple in  $C$ , which will be used to approximate the full gradient during training. For training, tuples in the coresset are randomly shuffled. Afterwards, suppose that in each step of the gradient decent, when we use  $c_j \in C$  to update the gradient, we compute the gradient  $(\nabla f_j)$  of  $c_j$  first, and then use  $w_j \nabla f_j$  to update the model parameters.  $w_j$  is the number of tuples in  $D$  assigned to  $c_j$ . The above steps repeat until the model converges.

**Optimization.** Unfortunately, the 3-loop computation of the strategy is rather expensive due to the large number of possible worlds (Section 4). To address this, we can integrate either human-in-the-loop or the automatic method into GoodCore framework (Section 5). It iteratively imputes one incomplete tuple or a mini-batch of incomplete tuples. Once the tuple(s) is (are) computed and added to the coresset within the first loop, the number of possible worlds can be significantly reduced, and so does the computational cost. In addition, how to select an appropriate coresset size  $K$  is an important problem, which will be discussed in Section 7.3.

## 4 GOODCORE ALGORITHM

In this section, we will illustrate GoodCore algorithm in details for solving Eq. 6, which is proven to be prohibitively expensive (Section 4.1). Then we focus on how to compute the expectation using possible worlds (Section 4.2) in the algorithm.

### 4.1 Problem Complexity

Let us first discuss the time complexity of finding the optimum of Eq. 6.

**THEOREM 1.** *The problem of expected optimal coreset selection over incomplete data is NP-hard.*

**PROOF.** Let us consider a special case that there is no missing value in  $D$ . Our problem becomes the typical coreset selection problem over complete data, which has been proven to be NP-hard by reduction from the Minimum Vertex Cover problem [23, 54, 55]. Hence, our problem is also NP-hard.  $\square$

**THEOREM 2.** *The problem of expected optimal coreset selection over incomplete data has the sub-modular property.*

**PROOF.** First, we regard  $E[C] = \sum_{k=1}^{|I_w|} p_k S_k$  as a utility function, where  $S_k = \sum_{i=1}^n \min_{c_j \in C_k} s_{ij}$ . In fact,  $S_k$  can be regarded as a function of the coreset score computation over complete data, which has already proven to have the sub-modular property [35, 54, 55]. Therefore, consider the property that a non-negative linear combination of sub-modular functions is also sub-modular [44]. To be specific, given any sub-modular function  $g_1, g_2, \dots, g_k$  and non-negative numbers  $\alpha_1, \alpha_2, \dots, \alpha_k$ . Then the function  $G$  defined by  $G = \sum_{i=1}^k \alpha_i g_i$  is sub-modular. Hence, we can conclude that our studied problem is a sub-modular problem because  $E[C] = \sum_{k=1}^{|I_w|} p_k S_k$ , where  $p_k > 0$ .  $\square$

**The greedy algorithm.** Given the sub-modular property, naturally, we can design a greedy algorithm with an approximate ratio. As shown in Algorithm 1, we greedily add one tuple to the coreset at each iteration. The added tuple should have the largest utility computed by  $E[t|C] = E[C] - E[C \cup \{t\}]$ . Hence, the key component is that given the original train data ( $D$ ) and a coreset ( $C$  or  $C \cup \{t\}$ ), how to compute the expectation of GA error ( $E[C]$  or  $E[C \cup \{t\}]$ ) of the coreset. However, it is non-trivial because of the large number of possible worlds. We will first introduce how to compute the probability  $p_k$ , and describe the expectation computation in Section 4.2. After  $K$  tuples are added, we can impute missing tuples in the coreset generated by GoodCore.

## 4.2 Expectation Computation

**Possible world probability.** To compute the expectation, it is inevitable to derive the probability of each possible world, which can be taken as a pre-processing step in our framework. To be specific, since tuples with missing values are always imputed independently [52], given a possible world  $W_k$ , the probability  $p_k$  can be computed by  $p_k = \prod_{t \in W_k, \mathbb{I}[t]=1} p_k^t$ , where  $p_k^t$  denotes the probability of the appearance of tuple  $t$  with  $\mathbb{I}[t] = 1$ . Besides, apparently  $p_k^t = 1$  when  $\mathbb{I}[t] = 0$ , so  $p_k = 1$  if there are only complete tuples. Therefore, our focus is on how to get the value of  $p_k^t$ , which can be solved by many approaches, like statistic methods and learning-based methods (see [52] for a survey). In this paper, we use the learning-based method [12] with a Python library [1] to generate the probability, which can be easily replaced by other libraries or domain-specific methods. During training, learning-based methods take  $D$  as input and learn a model  $\mathcal{M}$  to describe the joint data distribution. For inference, we have  $P(\mathcal{A}_i|\mathbf{x}, v_{mask}) = \mathcal{M}(\mathbf{x}, v_{mask}, \omega^*)$ , where the model takes as input the feature vector  $\mathbf{x}$  of  $t$ , the mask vector  $v_{mask}$  (indicating which attributes are missing) and the model parameter  $\omega^*$ , outputs the probability distribution of a missing attribute  $\mathcal{A}_i$ .

Suppose that  $t$  just has one missing attribute  $\mathcal{A}_i$ , and then  $v_{mask}$  is a one-hot vector with  $v_{mask}[i] = 0$ . Hence, we can directly obtain  $p_k^t$  from the distribution  $P(\mathcal{A}_i|\mathbf{x}, v_{mask})$ . For  $t$  with multiple missing attributes, we can also compute  $p_k^t$  using the chain rule. If  $t$  has two missing values of  $\mathcal{A}_i$  and  $\mathcal{A}_j$ , to compute  $p_k^t$ , we have to compute  $P(\mathcal{A}_i, \mathcal{A}_j|\mathbf{x}, v_{mask})$ , abbreviated as  $P(\mathcal{A}_i, \mathcal{A}_j) = P(\mathcal{A}_i)P(\mathcal{A}_j|\mathcal{A}_i)$ .  $P(\mathcal{A}_i)$  can be obtained by masking the  $i$ -th and  $j$ -th attribute in  $v_{mask}$ . Then, we only mask the  $j$ -th attribute and impute different values of  $\mathcal{A}_i$  to obtain  $P(\mathcal{A}_j|\mathcal{A}_i)$ .

**EXAMPLE 5.** In Figure 5(a), suppose that for the first possible world, we have to compute  $p_1 = p_1^2 \times p_1^3 \times p_1^4 \times p_1^6$ . For instance, to compute  $p_1^3$ , given the trained deep learning model, we feed

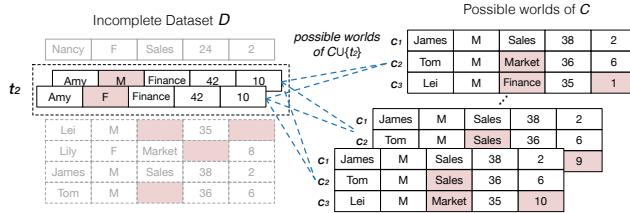


Fig. 7. Tuple-based expectation computation.

$\{\text{Lei}, \text{M}, \text{Mask}, 35, \text{Mask}\}$  and a one-hot vector  $\{1, 1, 0, 1, 0\}$  into the model and compute the probability distribution of this tuple, from which we can get  $p_1^3$ , i.e., the probability of  $\{\text{Lei}, \text{M}, \text{Sales}, 35, 1\}$ .

Compared with statistical approaches, deep learning-based methods use more powerful models with good learning capacity and consider the correlation between attributes. For practitioners, they can use any ad-hoc method to compute the probability.

**Brute-force expectation computation.** Recap that  $E[C] = \sum_{k=1}^{|\mathcal{I}_W|} p_k (\sum_{i=1}^n \min_{c_j \in C_k} s_{ij})$ . Intuitively, the brute-force method is to enumerate each possible world, compute the probability and finally get the expectation. However, there are a huge number of possible worlds, which makes the computation prohibitively expensive. Specifically, we assume the attribute number  $M$  and  $|\mathcal{A}_m|, m \in [1, M]$  are constants, so the number of possible worlds of each tuple is a constant, denoted by  $L$ . Suppose that the number of tuples with missing values is  $O(n)$ , so the number of possible worlds ( $|\mathcal{I}_W|$ ) is  $O(L^n)$ . Given a coresset  $C$ , the time complexity to compute  $E[C]$  is  $O(nL^n)$ , which is rather expensive.

**Tuple-based expectation computation.** To further elaborate, we can easily expand  $E[C]$  as follows:

$$\begin{aligned} E[C] = & p_1 (\underbrace{\min_{c_j \in C_1} s_{1j} + \min_{c_j \in C_1} s_{2j} + \cdots + \min_{c_j \in C_1} s_{nj}}_{\cdots}) \\ & + p_2 (\underbrace{\min_{c_j \in C_2} s_{1j} + \min_{c_j \in C_2} s_{2j} + \cdots + \min_{c_j \in C_2} s_{nj}}_{\cdots}) + \cdots \\ & + p_{|\mathcal{I}_W|} (\underbrace{\min_{c_j \in C_{|\mathcal{I}_W|}} s_{1j} + \min_{c_j \in C_{|\mathcal{I}_W|}} s_{2j} + \cdots + \min_{c_j \in C_{|\mathcal{I}_W|}} s_{nj}}_{\cdots}). \end{aligned}$$

We can see from the above equation that these underlined terms are only related to  $t_1 \in D$  as well as  $\{C_1, C_2, \dots, C_{|\mathcal{I}_W|}\}$ , i.e., the coresets corresponding to the  $|\mathcal{I}_W|$  possible worlds. However, as the coresset  $C$  is much smaller than the full data  $D$ , the number of possible worlds of  $C$  will be also much smaller than  $|\mathcal{I}_W|$ , and thus there will be many duplicates among  $\{C_1, C_2, \dots, C_{|\mathcal{I}_W|}\}$ . Therefore, many of these underlined terms have identical variable parts, i.e.,  $\min_{c_j \in C_k} s_{1j}$ , when they are associated with the same  $C_k$ . These terms are *like terms*. Combining these like terms (i.e.,  $\sum_{k=1}^{|\mathcal{I}_W|} p_k \min_{c_j \in C_k} s_{1j}$ ), we can get the expectation of  $\min_{c_j \in C} s_{1j}$ , denoted by  $E[\min_{c_j \in C} s_{1j}]$ .

In short, we can convert the expectation computation over the possible worlds of the entire training set  $D$  to the sum of expectation of each tuple in  $D$ , as follows:

$$E[C] = \sum_{k=1}^{|\mathcal{I}_W|} p_k (\sum_{i=1}^n \min_{c_j \in C_k} s_{ij}) = \sum_{i=1}^n E[\min_{c_j \in C} s_{ij}] \quad (7)$$

**Algorithm 2:** ComputeUtility (3rd-loop to compute  $E[t|C]$ )

---

**Input:** Incomplete train data  $D$ , current coreset  $C$ , a sampled tuple  $t$ .  
**Output:** The expectation  $E[t|C]$ .

```

1   $\hat{C} = C \cup \{t\}$ ;
2   $E[\hat{C}] = 0$ ;
3  for each tuple  $t_i \in D$  do
4    if  $\mathbb{I}[t] = 0$  and  $\mathbb{I}[t_i] = 0$  then
5       $E[\hat{C}] += \min_{c_j \in \hat{C}} s_{ij}$ ;
6    else
7      Get the possible worlds of  $\hat{C} \cup \{t_i\}$ ;
8      Compute  $E[\min_{c_j \in \hat{C}} s_{ij}]$  using these possible worlds and their probabilities;
9       $E[\hat{C}] += E[\min_{c_j \in \hat{C}} s_{ij}]$ ;
10    $E[t|C] = E[C] - E[\hat{C}]$ ;
11   return  $E[t|C]$ ;

```

---

EXAMPLE 6. Figure 7 shows how to compute  $E[\min_{c_j \in C} s_{2j}]$ . Instead of enumerating  $|\mathcal{I}_W|$  possible worlds by the brute-force method, we can enumerate a much smaller number of possible worlds of  $C \cup t_2$ , compute the corresponding probabilities and finally get the tuple expectation. Specifically, The left part of Figure 7 shows the possible worlds of the tuple, the right part shows the possible worlds of the coreset, and their combination is the possible worlds of  $C \cup t_2$ . Then, following Eq. 7, we can iterate the tuples in  $D$ , compute their expectations and sum them up to derive  $E[C]$ .

**Time complexity.** Since the coreset size is  $K$ , and the number of tuples with missing values in the coreset is  $O(K)$ , the time complexity of computing  $E[C]$  using tuple-based method is  $O(nL^K)$ , where  $K$  is much smaller than  $n$ , compared with the brute-force method. However, note that computing  $E[C]$  is just the third loop in the entire framework. Besides, the first two loops incrementally add  $K$  tuples into the coreset, and sample  $h$  tuples for tuple selection respectively. Hence, the overall time complexity of coreset selection over incomplete data is  $O(KhnL^K)$ , which is still expensive when  $K$  is not small enough. In the next section, we involve the imputation-in-the-loop strategies to achieve further improvement.

## 5 OPTIMIZED GOODCORE WITH IMPUTATION-IN-THE-LOOP

As discussed above, it is rather expensive to directly compute all the  $K$  tuples in the coreset. Hence, in this section, we propose to involve the imputation-in-the-loop mechanism that asks the human, *i.e.*, Case (7), or automatic method, *i.e.*, Case (8) to impute these missing values iteratively while they are generated by Algorithm 1.

The advantages of this optimization are two-fold. First, with more and more missing values being imputed, the number of possible worlds is greatly reduced, which reduces the machine cost a lot. Second, for human-in-the-loop imputation, it allows us to gradually impute the tuples accurately, and thus the coreset score computation can be more and more accurate, which produces a better coreset.

### 5.1 One Tuple Each Iteration

In fact, we can just slightly modify Algorithm 1 to achieve the imputation-in-the-loop strategy. To be specific, in the first loop, we will iteratively impute the tuple once an incomplete tuple  $t^*$  is computed by GoodCore, rather than conducting the imputation after  $K$  tuples are computed, as discussed in Section 4. To this end, we move the imputation step (lines 15-16 in Algorithm 1) inside

**Algorithm 3:** Batch algorithm of GoodCore

---

```

Input:  $D, K, h$ , batch size  $b$ .
Output: A coresset  $C$ , weight  $\mathbb{W}$ .
1  $C = \emptyset, cnt = 0;$ 
2 while  $|C| < K$  do
3   Sample  $h$  tuples as  $T_{sample} \subseteq D \setminus C$ 
4   for each tuple  $t \in T_{sample}$  do
5      $E[t|C] = \text{ComputeUtility}(t, C, D);$ 
6    $t^* = \arg \max_{t \in T_{sample}} E[t|C];$ 
7    $C = C \cup \{t^*\};$ 
8   if  $I[t^*] = 1$  then
9      $cnt++;$ 
10  if  $cnt = b$  then
11    Ask the human to impute the incomplete tuples;
12     $cnt = 0;$ 
13 Compute the weight  $\mathbb{W}.$ 
14 return  $C, \mathbb{W};$ 

```

---

the first loop of Algorithm 1, *i.e.*, imputing each selected  $t^*$  by a human or automatic method in each iteration after line 17.

Afterwards, we will add the next tuple into the coresset, so another loop starts and  $h$  tuples are sampled. In the following, we will expand the third loop, *i.e.*, the function `ComputeUtility` (line 7) of Algorithm 1 under this one tuple per iteration scenario.

As shown in Algorithm 2, at the beginning, we temporarily add the sampled tuple  $t$  to the current coresset, so as to compute the benefit of  $t$ , *i.e.*,  $E[t|C]$ . To this end, we have to first compute the expectation of GA error bound of  $\hat{C}$  (*i.e.*, computing  $E[\hat{C}]$  in the for-loop lines 3-12). And the expectation w.r.t.  $C$  (*i.e.*,  $E[C]$ ) has been computed in the last loop. Then we can compute  $E[t|C] = E[C] - E[\hat{C}]$  (line 10).

Specifically, to compute  $E[\hat{C}]$ , we will use the tuple-based expectation computation method proposed in Section 4.2. For each tuple  $t_i \in D$ , if  $t_i$  and  $t$  are both complete, we can directly compute  $\min_{c_j \in \hat{C}} s_{ij}$  because there is no incomplete data in  $\hat{C}$  (lines 4-5). Otherwise, we will enumerate the possible worlds of  $\hat{C} \cup \{t_i\}$ , compute their probabilities and compute  $E[\min_{c_j \in \hat{C}} s_{ij}]$  (lines 7-8). Note that since there are at most two tuples (*i.e.*,  $t_i$  and  $t$ ) have missing values, the number of possible worlds is small because other missing values in  $\hat{C}$  have been imputed by humans in previous iterations.

**Time complexity analysis.** As discussed above, using this human-in-the-loop strategy, the number of possible worlds to be considered is greatly reduced. For Algorithm 2, the time complexity is  $O(nL^2)$  because there are at most two incomplete tuples in  $\hat{C}$ . For the entire three loops framework, the time complexity is  $O(KhnL^2)$ , which is much lower than the solution without imputation in the loop.

However, if we utilize the human for imputation, the above method will incorporate many human iterations. In the following, we propose to ask human to impute a small batch of missing tuples in each iteration, so as to reduce the number of human iterations.

## 5.2 One Batch Each Iteration with Human-in-the-loop

In Section 5.1, one tuple per iteration by humans requires many human iterations. However, if we just incorporate a single human iteration like Section 4.2, it is infeasible to compute the tuples to

be imputed due to the large number of possible worlds. Therefore, in this subsection, we propose a trade-off solution that asks the human to impute a small batch of tuples per human iteration.

To be specific, as shown in Algorithm 3, compared with the one tuple per human iteration algorithm (*i.e.*, the modified Algorithm 1 at the beginning of Section 5.1), we additionally take the batch size  $b$  as input (when  $b = 1$ , Algorithm 3 is in fact the modified Algorithm 1). Algorithm 3 also incorporates 3 loops, but the main difference is that we do not instantly ask the human to impute the most beneficial tuple  $t^*$  among  $T_{sample}$ . Instead, we just add  $t^*$  into the coreset  $C$  (line 7). When there have been  $b$  incomplete tuples, we ask the human to impute these tuples together (line 10-12). Finally we compute the weight (line 13), same as Algorithm 1. Although this approach reduces the number of human iterations, it takes a longer time to compute  $E[t|C]$  (line 5) than Algorithm 1 because there are more incomplete tuples, which indicates more possible worlds. Specifically, the time complexity of computing  $E[t|C]$  is  $O(nL^b)$ , which is also expensive. Hence, we propose a heuristic method to accelerate this process as follows.

**Reducing the number of possible worlds.** A straightforward method of improving the efficiency is to reduce the number of possible worlds. To this end, intuitively, we should focus more on the possible world with a high probability, so these possible worlds with low probabilities can be pruned without sacrificing the accuracy of expectation computation much. Note that for each possible world, the probability is computed by the multiplication of the probabilities of incomplete tuples in the world because the tuples can be considered independent [52]. Therefore, we can remove the possible worlds of each tuple with low probabilities (*i.e.*, reducing  $L$ ), and thus the number of possible worlds of the entire coreset is greatly reduced. For example, we can keep top- $l$  (*e.g.*,  $l = 3$ ) possible worlds (*i.e.*, 3 different possible imputations of  $t$  with high probabilities) of a tuple  $t$ . Then for the batch of  $b$  incomplete tuples, the number of possible worlds is  $l^b$  and the complexity of computing  $E[t|C]$  is  $O(nl^b)$ , where both  $l$  and  $b$  are small enough. Overall, the time complexity is  $O(KhnL^b)$ . Besides, we can also apply this heuristic method to make the algorithm in Section 4 practical, which is evaluated in Section 7.4.

### 5.3 Convergence Rate Analysis

Convergence rate is often used to reflect the speed of finding the optimal parameters for the machine learning algorithm. With a higher convergence rate, we can take fewer epochs to make the model converge. To compute the convergence rate, we have to compute the distance between the parameter  $\theta$  and the optimal parameter  $\theta^*$  in the  $t$ -th and the  $(t + 1)$ -th epoch. Since  $f$  is a strongly convex function,  $\forall \theta, \theta'$  we have

$$f(\theta) - f(\theta') \geq \nabla f(\theta')(\theta - \theta') + \frac{\eta}{2} \|\theta - \theta'\|^2 \quad (8)$$

where  $\eta$  is a constant. We denote the stepsize as  $\zeta_t = \frac{\zeta_0}{k^\tau}$  for the  $t$ -th epoch, where  $\tau$  is a constant. After using gradient descent in each step, we have  $\|\theta^{t+1} - \theta^*\|^2 = \|\theta^t - \zeta_t \sum_{j=1}^{|C|} w_j \nabla f_{Y(j)}(\theta_{j-1}^t) - \theta^*\|^2$ . Then, following Eq. 8, we have

$$\begin{aligned} \|\theta^{t+1} - \theta^*\|^2 &\leq \|\theta^t - \theta^*\|^2 - 2\zeta_t \sum_{j=1}^{|C|} (f_j(\theta^t) - f_j(\theta^*)) \\ &\quad + 2\zeta_t \sum_{j=1}^{|C|} (f_j(\theta_{j-1}^t) - f_j(\theta^t)) + \zeta_t^2 \sum_{j=1}^{|C|} \|w_j \nabla f_j(\theta_{j-1}^t)\|^2 \end{aligned} \quad (9)$$

Recap that we select a coreset that minimizes  $E[C]$  through converting gradient difference to feature distance ( $s_{ij}$ ) computation. Obviously, given a dataset,  $s_{ij}$  can be bounded (suppose that  $s_{ij} \leq s_0$ ). Then we have  $E[\min_{c_j \in C} s_{ij}] = \sum_{k=1}^{|I_w|} p_k (\min_{c_j \in C_k} s_{ij}) \leq \sum_{k=1}^{|I_w|} p_k * s_0 = s_0$ , and

thus  $E[C] = \sum_{i=1}^n E[\min_{c_j \in C} s_{ij}] \leq n * s_0 = \kappa_1$ . Besides, we also have  $\max_{\theta \in \theta^*} \|\sum_{i=1}^n \nabla f_i(\theta) - \sum_{j=1}^{|C|} w_j \nabla f_{Y(j)}(\theta)\| \leq \sum_{i=1}^n \min_{c_j \in C} \|\nabla f_i(\theta) - \nabla f_{Y(j)}(\theta)\| \leq \sum_{i=1}^n \min_{c_j \in C} s_{ij} \leq \kappa_1$ . Following the definition of convex function, we have  $f_j(\theta^t) - f_j(\theta^*) \leq w_j \nabla f_j(\theta^*)(\theta^t - \theta^*) + \frac{\eta}{2} \|\theta^t - \theta^*\|^2$ . Based on the above things, we can apply Cauchy-Schwarz inequality [68] and derive

$$\begin{aligned} & -2\zeta_t \sum_{j=1}^{|C|} (f_j(\theta^t) - f_j(\theta^*)) \\ & \leq -\eta \zeta_t \|\theta^t - \theta^*\|^2 + 2\zeta_t \left\| \sum_{j=1}^{|C|} w_j \nabla f_j(\theta^*) \right\| \|\theta^t - \theta^*\| \\ & \leq -\eta \zeta_t \|\theta^t - \theta^*\|^2 + \frac{2\zeta_t^2 |C| \kappa_1 \kappa_2}{\eta} \end{aligned} \quad (10)$$

where  $\kappa_2$  can be regarded as the upper bound of  $\|\theta_t - \theta^*\|$ . Since  $f$  is convex, thus, for item  $f_j(\theta_{j-1}^t) - f_j(\theta^t)$ , we have  $f_j(\theta_{j-1}^t) - f_j(\theta^t) \leq \|w_j \nabla f_j(\theta^t)\| \zeta_t \sum_{i=1}^{j-1} \|w_i \nabla f_i(\theta_{i-1}^t)\|$ . In addition, we can assume that  $\max_{j \in \{1, \dots, |C|\}} \|\nabla f_j(\theta)\| \leq \kappa_3$ . Then, we have

$$\begin{aligned} & 2\zeta_t \sum_{j=1}^{|C|} (f_j(\theta_{j-1}^t) - f_j(\theta^t)) + \zeta_t^2 \sum_{j=1}^{|C|} \|w_j \nabla f_j(\theta_{j-1}^t)\|^2 \\ & \leq 2\zeta_t \sum_{j=1}^{|C|} \|w_j \nabla f_j(\theta^t)\| \zeta_t \sum_{i=1}^{j-1} \|w_i \nabla f_i(\theta_{i-1}^t)\| + \zeta_t^2 \sum_{j=1}^{|C|} \|w_j \nabla f_j(\theta_{j-1}^t)\|^2 \\ & \leq 2\zeta_t^2 (|C|^2 - |C|) w_{\max}^2 \kappa_3^2 + \zeta_t^2 |C| w_{\max}^2 \kappa_3^2 \end{aligned} \quad (11)$$

Thus, from Eq. 9 to Eq. 11, we can get

$$\|\theta^{t+1} - \theta^*\| \leq (1 - \eta \zeta_t) \|\theta^t - \theta^*\|^2 + \frac{2\zeta_t^2 |C| \kappa_1 \kappa_2}{\eta} + \zeta_t^2 |C|^2 w_{\max}^2 \kappa_3^2 \quad (12)$$

Finally, following Lemma 4 in [21], the convergence rate of Algorithm 1 is at the same rate of  $O(\frac{1}{\sqrt{k}})$  as the convergence rate on the entire dataset [59]. Therefore, theoretically, the selected coresnet can converge with the same number of epochs as training on the full data. In this way, since coresnet has a much smaller size than the full data, the efficiency can be much improved.

## 6 GROUP-BASED ACCELERATION

As discussed above, we can observe that in Section 5.1, even with the most efficient imputation-in-the-loop strategy, *i.e.*, one tuple in each iteration, the time complexity is  $O(KhnL^2)$ , where  $K$  is the size of the coresnet,  $h$  is the sample size,  $L$  is a small constant and  $n$  is the number of entire dataset. Therefore, obviously, the efficiency is dominated by  $n$ , which is still low when  $n$  is large, and thus it is necessary to further accelerate this process.

**Key observation.** Recap that in Figure 4, we can observe that given a tuple  $c$  in the coresnet, the tuples in the origin full train set  $D_c$  represented by  $c$  are likely to be closer to each other than other tuples not represented by  $c$ . Based on this observation, we propose to first cluster the full train set into groups, and then compute the coresnet based on these groups. This can achieve much acceleration because the number of groups is much smaller than  $n$ .

At the following, we will theoretically and empirically show the groups-based solution can accelerate the coresnet selection process without sacrificing the effectiveness much.

### 6.1 Group-based Solution Overview

As shown in Figure ??, one of the core parts of coresnet computation is to compute the tuple-tuple distance, *i.e.*,  $s_{ij}$ . For the group-based solution, we just need to consider the relationship between

**Algorithm 4:** Group-based Overall Solution (imputation-in-the-loop)

---

```

Input: Incomplete train data  $D$ , coresset size  $K$ , sample size  $h$ .
Output: A coresset  $C \subseteq D$ , weight  $\mathbb{W} = \{w_j\}, |C| = |\mathbb{W}| = K$ .
1  $C = \emptyset;$ 
2 Cluster  $D$  into groups  $\mathcal{G} = \{G_1, G_2, \dots, G_{|\mathcal{G}|}\}$ .
3 while  $|C| < K$  do
4   /*1st loop*/
5   Sample  $h$  tuples as  $T_{sample} \subseteq D \setminus C$ 
6   for each tuple  $t \in T_{sample}$  do
7     /*2nd loop*/
8      $\hat{C} = C \cup \{t\};$ 
9     for each group  $G_k \in \mathcal{G}$  do
10      /*3rd loop*/
11       $E[\hat{C}] += E[\min_{c_j \in \hat{C}} b_{jk} \times |G_k|]$ , where  $b_{jk}$  is the upper bound of
12       $\bar{s}_{Y(j)k} = \max_{v \in G_k} s_{Y(j)v}, s_{Y(j)v} = \|\mathbf{x}_v - \mathbf{x}_{Y(j)}\|, Y(j) \in [1, n];$ 
13       $E[t|C] = E[C] - E[\hat{C}];$ 
14       $t^* = \arg \max_{t \in T_{sample}} E[t|C];$ 
15      if  $I[t^*] = 1$  then
16        | Impute  $t$  by a human or automatic method.
17         $C = C \cup \{t^*\};$ 
18 for  $j = 1$  to  $|C|$  do
19   for  $i = 1$  to  $n$  do
20     if  $c_j = \arg \min_{c_{j'} \in C} \max_{\theta \in \mathcal{G}} \|\nabla f_i(\theta) - \nabla f_{Y(j')}(\theta)\|$  then
21       |  $w_j += 1;$ 
22 return  $C, \mathbb{W};$ 

```

---

tuples and these pre-computed groups, namely tuple-group distance, rather than the large amount of tuple-tuple distances. As we will discuss below, the computation of tuple-group distance does not need to iterate all tuples in the group, and thus the overall efficiency can be much improved.

At a high level, the overall process of group-based solution is shown in Algorithm 4. To be more explicit, we illustrate Algorithm 4 in comparison with Algorithm 1 in Section 3.2 without clustering (the modified parts are highlighted in blue fonts).

To be specific,

## 6.2 Group-based GA Error Bound

## 6.3 Algorithm

### 6.3.1 Pre-processing.

### 6.3.2 Computing the upper bound.

## 7 EXPERIMENT

In this section, we sufficiently compare our proposed methods with multiple baselines on real datasets to demonstrate our effectiveness and efficiency.

Table 1. Statistics of datasets

Dataset	$ D $	$M$	# Incomp. Tuples	Task
Nursery	10960	9	3218	Multi-Class.
HR	18287	12	5475	Binary Class.
Adult	32842	14	10752	Binary Class.
Credit	131,000	11	76813	Binary Class.
BikeShare	13300	15	4821	Regression
Air	437,200	18	128,372	Regression

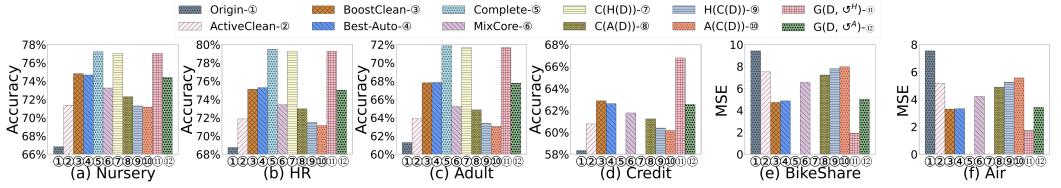


Fig. 8. Effectiveness of different methods.

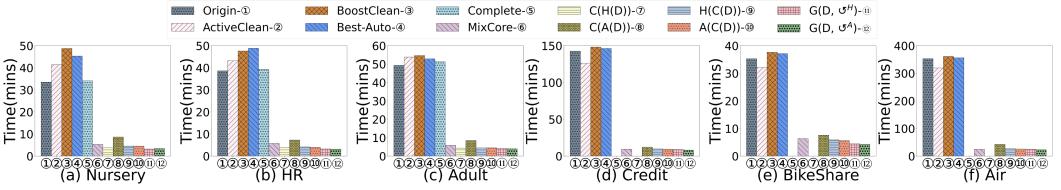


Fig. 9. Efficiency of different methods. Note that only machine cost (i.e., runtime of machine) is considered.

## 7.1 Experimental Settings

**Dataset.** We evaluate on 6 real-world datasets that are often used in the field of data imputation [33, 40, 43, 47], as shown in Table 1, where  $M$  denotes the number of attributes.

(1) Nursery [2] is a multi-classification task, which predicts “*the level of recommendation for whether a child goes to school*”. There are five different levels, *i.e.*, {not\_recom, priority, recommend, spec\_prior, very\_recom}. (2) HR [19] is a binary classification task of “*predicting whether an employee would change the job*”. (3) Adult [3] is a binary classification task that predicts “*if the annual revenue of a people is over 50000 dollars*”. (4) Credit [4] is a binary classification task that predicts “*whether the loan will be deferred based on a person’s economic situation*”. (5) BikeShare [5] is a regression task that predicts “*the number of bike sharing in a given time*”. (6) Air [6] is a regression task that predicts “*the air quality at a certain time*”.

For datasets (1)-(3), we follow existing works [34, 70, 72] to manually inject missing values until the rate of missing tuples is 30%, and we will vary the percentage of incomplete tuples in Section 7.6. Datasets (4)-(6) already contain missing values. For all datasets, we randomly split them for 80%/10%/10% as train/validation/test sets.

**Evaluation metrics.** We mainly evaluate the effectiveness and efficiency of GoodCore and baselines. For effectiveness, we use the *prediction accuracy* for the classification task and use the *mean square error* ( $MSE = \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}$ , where  $N$  denotes the size of test set) for the regression task.

For efficiency, we focus on the machine cost (*i.e.*, the runtime of machine) as well as the human cost (the number of tuples imputed by human for human-involved methods). For datasets (1)-(3), we have the ground truth of missing tuples, so we use them to simulate the human imputation.

For datasets (4)-(6), we leverage the expert to impute missing values in the coreset by looking at the top-5 values recommended by the automatic method as a reference. Note that we only involve humans when it is affordable. For baselines that require humans to impute a lot of missing tuples (*i.e.*, Complete and C(H(D)) as below), we will not apply them on datasets (4)-(6).

**Baselines.** We compare GoodCore with a variety of baselines.

(1) Origin refers to just training on  $D$ .

(2) ActiveClean [40] is an iterative data cleaning framework, which estimates the impact of tuples and prioritizes cleaning the tuples that much affect the model performance. In each iteration, it can ask the human to clean a sample subset of tuples. We set the sample size to 50, same as the paper.

(3) BoostClean [38] is an automatic data cleaning method that iteratively selects a cleaning method from several pre-defined algorithms, applies to the train dataset and updates the model. We use MICE [65], MISSForest [67], GAIN [72] as pre-defined algorithms.

(4) Best-Auto uses MICE [65], MISSForest [67], GAIN [72] to respectively impute the train set and selects the one that achieves the highest accuracy on the validation set.

(5) Complete is an ideal case that trains on the ground truth, *i.e.*,  $D_c$ . Note that only datasets (1)-(3) have the ground truth to evaluate this baseline. Datasets (4)-(6) do not have the ground truth and it is too expensive to ask the human to impute so many missing values.

(6) MixCore is a baseline that selects a coreset from all complete tuples, and then we randomly select some incomplete tuples to impute. We set the number of incomplete tuples to be imputed equal to that of other baselines for fair comparison. Finally we train with the tuples in the coreset plus the imputed ones.

(7) C(H(D)) first involves human to impute the dataset  $D$  and then selects a coreset. Similar to Complete, only datasets (1)-(3) can be evaluated on it because they have the ground truth. The coreset selection solution is the algorithm in [54], which is a greedy algorithm by modifying Algorithm 1 without considering the possible worlds.

(8) C(A(D)) first uses automatic data imputation methods to impute the dataset  $D$ , and then selects a coreset using the same method of baseline (7).

(9) H(C(D)) directly selects a coreset based on  $D$  and then asks human to impute the incomplete tuples of the coreset.

(10) A(C(D)) also directly selects a coreset from  $D$ , it then uses MICE [65] to impute the incomplete tuples in the coreset.

**Our solutions.** We compare GoodCore and its variants.

(11)  $G(D, \cup^H)$  uses GoodCore to select the coreset and iteratively asks human to impute incomplete tuples (one tuple per human iteration) during the coreset selection process.

(12)  $G(D, \cup^A)$  is similar to  $G(D, \cup^H)$ , but the automatic MICE method is used.

Besides, since the coreset of  $H(G(D))$  (*or*  $A(G(D))$ ) is too expensive to compute due to the large number of possible worlds, we do not directly compare with it. Instead, we will limit the number of possible worlds of each tuple to 3 as discussed in Section 5.2 and evaluate in Section 7.4.

**Hyper-parameter setting.** We use SVM and linear regression as the default downstream model for classification and regression tasks, respectively. We vary the downstream models in Section 7.6. For model training, we use SGD and k-inverse decay scheduling, *i.e.*,  $\alpha_k = \alpha_0 / (1 + bk)$  ( $\alpha_0$  and  $b$  are hyper-parameters to be tuned independently for different methods). The sample size  $h$  is set to 200 as default and we vary the size in Section 7.6. The number of training epochs is set as 20. We also impute the test data using the same method that is applied to the train data before testing.

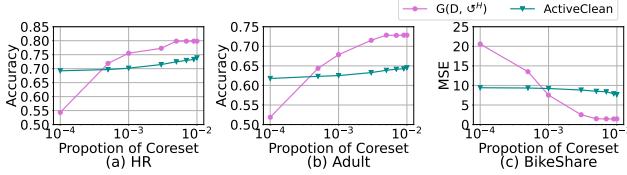


Fig. 10. Coreset size selection of GoodCore.

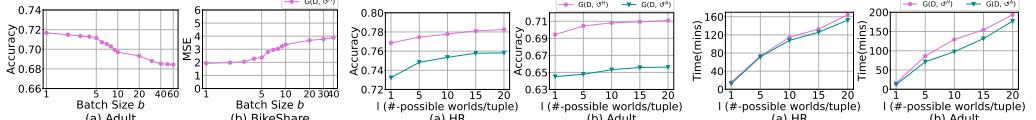
Fig. 11. GoodCore for batch algo- Fig. 12. Effectiveness when varying Fig. 13. Efficiency when varying  $l$ .  
rithm.  $l$ .

Table 2. Human cost of different methods

Dataset	$G(D, \cup^H)$	$H(C(D))$	$C(H(D))$
Nursery	37	22	3278
HR	44	32	5475
Adult	63	81	10752
Credit	52	67	-
BikeShare	38	25	-
Air	98	102	-

## 7.2 Overall Evaluation

In this part, we compare GoodCore solutions with baselines. We use  $\rho = \frac{K}{|D|}$  to denote the proportion of coreset to the entire train set. We set  $\rho = 0.005$  for datasets (1)-(4),  $\rho = 0.001$  and  $\rho = 0.0005$  for larger datasets (5) and (6) respectively. We will further conduct evaluation by varying the coreset size in Section 7.3.

**7.2.1 Evaluation of model accuracy.** The results are provided in Figure 8. To summarize, the result could be generally ranked as  $G(D, \cup^H)/C(H(D))/\text{Complete} > G(D, \cup^A)/\text{BoostClean}/\text{Best-Auto} > C(A(D)) > \text{MixCore} > \text{ActiveClean} > H(C(D))/A(C(D)) > \text{Origin}$ . Next, we explain the results.

In general, on all datasets, our method  $G(D, \cup^H)$ , Complete and  $C(H(D))$  perform the best. Complete and  $C(H(D))$  achieve a high accuracy because they ask the human to impute missing values accurately, but incur a high human cost. For example, Complete and  $C(H(D))$  achieve accuracy of 71.9% and 71.7% on Adult.  $G(D, \cup^H)$  is competitive with them because it selects a good coreset that can well represent the unknown ground truth via gradient approximation. In addition, we can observe that  $G(D, \cup^H)$  performs better than  $G(D, \cup^A)$  because human imputation is more accurate than automatic methods. For example, on Adult,  $G(D, \cup^H)$  has an accuracy of 71.7%, while  $G(D, \cup^A)$  and others are below 68%.  $G(D, \cup^A)$ , BoostClean and Best-Auto have competitive performance on accuracy. BoostClean and Best-Auto can have a not bad performance because they impute all tuples and train on the entire dataset, but they cannot achieve efficient training (see 7.2.2). But we can train on the much smaller coreset generated by  $G(D, \cup^A)$  with a good accuracy, because GoodCore considers the possible repairs to derive the coreset that can approximate the full gradient of the entire dataset. Given the same number of tuples to be imputed by human,  $G(D, \cup^H)$  also outperforms ActiveClean because we have theoretical guarantees on the gradient approximation. For other baselines,  $H(C(D))$  and  $A(C(D))$  do not perform well because they select

the coresets from an incomplete dataset.  $C(A(D))$  cannot achieve a good performance because the selected coresets can not well represent the complete entire dataset, as it does not consider possible repairs as our method. MixCore does not perform well (*e.g.*, 65.2% on Adult) because  $G(D, \cup^H)$  and  $G(D, \cup^A)$  select a better coresset considering the full data. For Origin, on Adult, the model has an accuracy of 61.3% because of the incomplete tuples.

**7.2.2 Evaluation of the efficiency.** We evaluate the efficiency of all methods, including the machine cost and human cost.

**Machine cost.** Machine cost is shown in Figure 9. The results could be ranked as  $H(C(D))/A(C(D))/G(D, \cup^H)/G(D, \cup^A)/C(H(D))/MixCore < C(A(D)) < Complete < Origin < ActiveClean < BoostClean/Best - Auto$ . We can observe that the first 5 methods in the ranking have low machine cost, mainly because they train based on the selected coressets and do not need iterative training.  $G(D, \cup^H)$  and  $G(D, \cup^A)$  are slightly slower because they need to iterate several possible repairs during the process of coresset selection. But  $G(D, \cup^H)$  is still more efficient than Origin, Complete, BoostClean and Best-Auto by more than one order of magnitude, because they need to train on the entire training data. Moreover, ActiveClean and BoostClean are not efficient either because they incorporate multiple training times, so as to estimate the gradient while data imputation. Best-Auto is slow because training multiple imputation models takes time.

**Human cost.** In terms of the human cost,  $C(H(D))$ ,  $H(C(D))$ ,  $G(D, \cup^H)$  and ActiveClean involve human. As shown in Table 2,  $C(H(D))$  is very expensive because it asks the human to impute all missing tuples. For example, on dataset Adult, 10752 tuples have to be imputed. We do not compare Credit, BikeShare and Air for  $C(H(D))$  because they do not have the ground truth. But  $H(C(D))$  and  $G(D, \cup^H)$  are cost-effective because human just needs to impute missing tuples in the much smaller coresset. For example, they only cost 81 and 63 tuples on dataset Adult respectively. ActiveClean asks the human to iteratively impute the data. Given the same number of tuples to impute, our method can achieve much higher accuracy. We will evaluate it in details in next subsection.

**Summary.** Based on the results, we have the following conclusions. (1) Our proposed methods  $G(D, \cup^H)$  and  $G(D, \cup^A)$  can achieve high model accuracy because the selected coressets can well represent the underlying ground truth by gradient approximation considering possible repairs. Meanwhile, they are practical because of the low machine cost. (2) Compared with  $C(H(D))$  that involves human to impute the entire dataset  $D$ , the human cost of  $G(D, \cup^H)$  is much lower, as observed in Table 2, *e.g.*, 37 vs. 3278 on the Nursery dataset. Thus, we can choose  $G(D, \cup^H)$  when we want to achieve high model accuracy and afford a certain human cost. (3) If we neither care very much about the accuracy nor consider to incur human cost, the much more efficient  $G(D, \cup^A)$  is a good choice.

### 7.3 Coreset Size Selection of GoodCore

Recap that GoodCore needs the user-specified coresset size as input. Thus, we discuss how to select an appropriate coresset size. We adopt a simple yet effective solution that starts from a coresset with a small size, train over it and evaluate via a validation set, enlarge the coresset and iteratively train until the performance does not improve much. To be specific, initially, we begin with  $\rho = 10^{-4}$ , and enlarge the coresset by 2 times iteratively. If the performance on validation set varies no more than 0.5% within three successive iterations, we will stop. Figure 10 shows the performance on dataset HR, Adult and BikeShare when varying the coresset size. We can see that the performance first improves rapidly, then remains stable just after several iterations. For example, on dataset Adult, when  $\rho = 5 \times 10^{-3}$ , the accuracy has improved to 72.85% on the validation set. Empirically, an ideal coresset size is between  $\rho = 10^{-3}$  to  $10^{-2}$ .

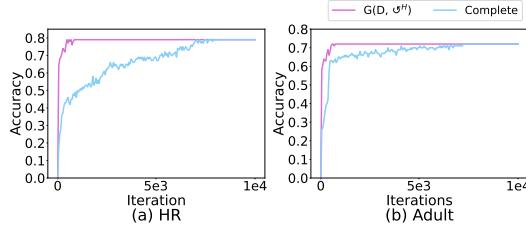


Fig. 14. Convergence of GoodCore.

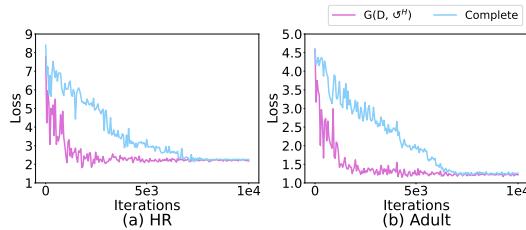


Fig. 15. Loss of GoodCore.

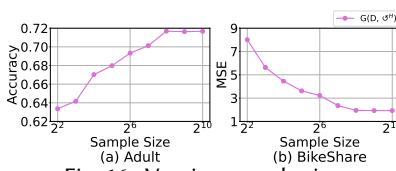


Fig. 16. Varying sample size.

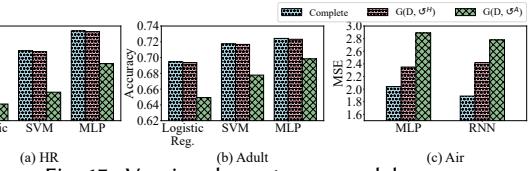


Fig. 17. Varying downstream model.

**Summary.** The results show that coresets size is not difficult to determine. If the user is willing to specify a coresset size like in Section 7.2 based on the empirical finding, we can directly compute a coresset without training. If she cannot, we can also get a good coresset with just several training iterations over small coresets, which is also efficient.

**Compare with ActiveClean.** Figure 10 also reports an interesting comparison with ActiveClean. Specifically, in ActiveClean, we use the coreset size  $K$  as the budget, i.e., number of tuples to be imputed by human in each active cleaning iteration. We can observe that at the beginning, when the coreset size is very small, ActiveClean is better because it trains with the entire dataset including the imputed tuples, while we train the model using only few tuples in the coreset. However, as with the increase of the coreset size, we can see that  $G(D, \cup^H)$  outperforms ActiveClean. This is because ActiveClean uses a heuristic method to estimate the impact of tuples to the overall gradient, which is not theoretically bounded (e.g., with gradient bounds like Coreset) and thus not accurate enough. For  $G(D, \cup^H)$ , it can achieve high accuracy with a proper coresset size, which is not large.

#### 7.4 Batch Algorithm of GoodCore

In Section 7.2,  $G(D, \cup^H)$  outperforms other baselines on accuracy, but requires many human iterations. In this part, we evaluate the batch algorithm of GoodCore by varying the batch size  $b$ , i.e., Algorithm 3 to reduce the number of iterations. Intuitively, the algorithm is  $G(D, \cup^H)$  when  $b = 1$ . Then we increase  $b$  until a single batch with a size  $b$  can contain all incomplete tuples in the coreset with size  $K$ , which is in fact the algorithm  $H(G(D))$ . Due to the large number of possible worlds, we adopt the heuristic method in Section 5.2 to set  $l = 3$  when  $b > 1$ .

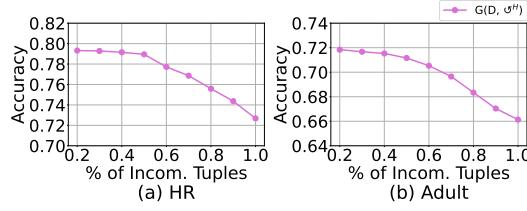


Fig. 18. Varying missing tuple rate.

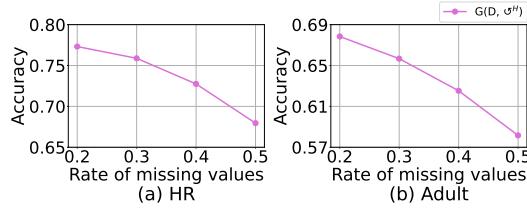


Fig. 19. Varying missing value rate.

Table 3. The number of possible worlds on different datasets

Method	Nursery	HR	Adult
$H(G(D))$	$10^{201}$	$10^{201}$	$10^{202}$
$A(G(D))$	$10^{201}$	$10^{201}$	$10^{202}$
$G(D, \cup^H)$	$10^3$	$10^3$	$10^4$
$G(D, \cup^A)$	$10^3$	$10^3$	$10^4$

In Figure 11, the  $x$ -axis denotes the batch size and the  $y$ -axis denotes the test performance on dataset Adult and BikeShare . We can see that when  $b$  is small (*i.e.*,  $b \leq 5$ ), the performance does not significantly decrease (*e.g.*, on Adult, the accuracy decrease from 71.7% to 71.2%). When  $b$  keeps increasing, the performance slightly decreases. Thus, GoodCore is not sensitive to the batch size  $b$  and the Algorithm 3 can reduce the number of human iterations without sacrificing much model performance.

In this part, we also vary the number of possible worlds by varying  $l$ , which is the number of possible world per tuple. The larger  $l$ , the larger number of possible worlds we have. The results are shown in Figures 12 and 13. In terms of the accuracy, we can see that with  $l$  increasing (fixing  $b = 10$ ), the accuracy increases first and then remains stable soon, but the time keeps increasing because more possible worlds indicate more computation. Hence, we do not need a large  $l$ .

When it comes to the number of possible worlds, we would like to clarify that we do not compare with  $H(G(D))$  and  $A(G(D))$  because the number of possible worlds of  $D$  is very large, which is infeasible to compute. We show the number in Table 3, where we also report the numbers of possible worlds of  $G(D, \cup^H)$  and  $G(D, \cup^A)$  in each iteration, which are practical to compute.

## 7.5 Convergence Evaluation

In Section 5.3, we have shown the convergence rate of GoodCore theoretically. In this part, we test the convergence of training over the coreset ( $G(D, \cup^H)$ ) and entire data (Complete) empirically. Figure 14 shows the test accuracy of two methods with the number of training iterations increasing. We can observe that on both datasets, training on the coreset converges much faster than training on the full data. For example, on dataset Adult, it takes  $\sim 40$  iterations for GoodCore to converge,

which is  $180\times$  faster than Complete. This is because GoodCore has the same convergence rate with training over the entire dataset as discussed in the theoretical result of Section 5.3, but the entire dataset (*e.g.*, Adult) is  $200\times$  (similar to  $180\times$ ) larger than the coresset ( $\rho = 0.005$ ). That is, GoodCore converges with the same number of epochs as training on the entire dataset. Since the size of coresset is much smaller, GoodCore is more efficient. Also, we can achieve competitive accuracy as training on full data by approximating the full gradient with a theoretical bound.

Furthermore, we report the loss change to reflect the relation between actual convergence rate and theoretical results. In Figure 15, on dataset Adult, the initial loss is 8.4. According to the theoretical convergence rate  $O(\frac{1}{\sqrt{k}})$  (this  $k$  denotes the  $k$ -th epoch), the loss should decrease to around 3.8 at the end of 5-th epoch ( $\approx 3200$ -th iteration). Actually, the actual loss decreases to 3.25 at that time, which is close to the theoretical value.

## 7.6 Sensitivity Analysis

**Varying the sample size.** In this part, we vary the sample size  $h$  and evaluate the performance. The experimental results are shown in Figure 16. We vary the sample size  $h$  from  $2^2$  to  $2^{10}$ . We can see that when  $h$  is too small, the performance is low. The reason is that GoodCore cannot precisely estimate the utilities of tuples when  $h$  is small. When the sample size increases, we can see that the performance improves rapidly and finally becomes stable, which indicates that GoodCore is not much sensitive to the sample size when  $h$  is not too small.

**Varying the downstream models.** Recap that GoodCore can be used on different convex models. Thus, in this part, we apply GoodCore on different convex models and evaluate the performance. We evaluate logistic regression and SVM for classification tasks. For regression tasks, we evaluate linear regression, ridge regression and SVM regression. We can see that in Figure 17 (a) and (b), on dataset Adult,  $G(D, \cup^H)$  achieved 71.7% accuracy for SVM, better than on logistic regression (69.4%). Although different downstream models may have different performance, GoodCore can improve the model performance for the specific downstream model. In order to show that GoodCore can be used for other types of models like neural networks, we compare with Multilayer Perceptron (MLP, fully connected networks of 2 hidden layers with 256 nodes for each layer), although GoodCore does not hold theoretical guarantee for this non-convex model. As shown in Figure 17, we can see that MLP achieves almost the same performance as the ground truth. This is because the coresset selected by GoodCore can also represent the full dataset. However, in Figure 17(c), on a large dataset Air (with metric MSE, the lower the better), neural network based methods (we also implement RNN, 2 hidden layers with 128 nodes for each layer) can have a better accuracy but the coresset cannot perfectly achieve the same performance. This may because this large dataset has more informative things to learn, and it is hard for the coresset-based solution to well represent the dataset without the theoretical guarantee.

**Varying the percentage of incomplete tuples.** In this part, we vary the rate of missing tuples and evaluate the performance, as shown in Figure 18. Note that the rate denotes the percentage of incomplete tuples rather than the cell values. Even if a tuple just has one missing attribute, it is regarded as incomplete. We vary the percentage from 20% to 100%. We can observe that the performance does not decrease much with the percentage increasing from 20% to 50%, which indicates that GoodCore is not very sensitive to the percentage of incomplete tuples in this range. After that, the accuracy decreases because there are more missing tuples.

Besides, we also vary the rate of missing cell values in Figure 19. In this scenario, for example, 50% missing values of a dataset indicates more number of missing cell values than the scenario of 50% missing tuples. Hence, we can see that the accuracy decreases more quickly than Figure 18.

## 8 RELATED WORK

**Task-agnostic incomplete data imputation.** Data imputation has been widely studied for years. Existing methods can be divided into two categories: statistic-based methods and learning-based methods. The former one always uses the statistic information [24, 48] (like mean, median or mode) to impute the missing values. Also, some methods compute the similarity of the incomplete tuples to the complete tuples and use the most similar one to impute the missing values [8, 32, 69]. Recently, to improve the imputation accuracy, many learning-based methods focus on how to use ML to learn the data distribution (*e.g.*, MissForest imputation [67], MICE [65], IIM [73]), and then use the trained model to predict the missing values. Besides traditional ML models, some deep learning models are also used for data imputation (*e.g.*, autoencoder [28, 49, 58], GANs [66, 72]).

**Coreset selection.** Huang et al. [31] studied how to compute and continuously update the coresets while training. But it is rather time-consuming because of the training process. To solve this problem, works [13, 14] selected the coresets without training in advance, but they can only be customized to particular model types respectively. Another line of works [37, 54, 56] focused on selecting the coresets to approximate the full gradient without training in advance for multiple model types, which is regarded as an optimization problem that can be solved by the three nested loops framework (see Section 2.2).

In short, none of them considers coreset selection over incomplete data. Different from them, we make the first attempt to select coresets over incomplete data.

**Data cleaning for ML.** Recently, there have been several works that clean the data to optimize the ML model. In contrast to the above discussion about task-agnostic incomplete data imputation, data cleaning for ML is task-aware, which triggers new technical challenges. SampleClean [39] focuses on cleaning selected samples, so as to answer SQL aggregate queries more efficiently, but it is not for any model. CPClean [33] proposes certain prediction to impute missing data for optimizing ML models. Different from us, it is customized to nearest neighbor classifiers rather than convex models solved by the gradient decent algorithm. BoostClean [38] regards data cleaning as a boosting problem that iteratively selects from a predefined set of cleaning algorithms, so as to continuously maximize the accuracy of a validation set with training iteratively. Closer to our work is ActiveClean [40], which progressively cleans the data tuples that are likely to much influence the model measured by the gradients. Different from us, given a budget  $K$ , we can select the coresets without training, but ActiveClean needs to train iteratively and label a set of validation dataset. We empirically show that our method outperforms ActiveClean on model accuracy and efficiency in Section 7.

**Data Preparation for AI.** Data preparation [16] can be utilized to improve the effectiveness of ML model, including data discovery [19, 45, 46], data cleaning [15, 29], data labeling [17, 18, 42], data cleaning [26, 27, 50–52, 71] and data exploration [62, 63].

## 9 CONCLUSION

In this paper, we propose the GoodCore framework to select a good coresset over the incomplete data, which achieves data-effective and data-efficient ML. We formulate it as an expected optimal coresset selection problem, which is NP-hard. Then we propose a greedy algorithm with an approximation ratio. We also propose to involve imputation-in-the-loop strategies into GoodCore to further improve the efficiency. We conduct experiments on real-world datasets to verify the effectiveness and efficiency of GoodCore.

## REFERENCES

- [1] 2022. <https://github.com/awslabs/datawig>.
- [2] 2022. <https://archive.ics.uci.edu/ml/datasets/nursery>.
- [3] 2022. <https://archive.ics.uci.edu/ml/datasets/adult>.
- [4] 2022. <https://www.kaggle.com/>.
- [5] 2022. <https://ride.capitalbikeshare.com/system-data>.
- [6] 2022. <https://auctus.vida-nyu.org/>.
- [7] Zeyuan Allen-Zhu, Yang Yuan, and Karthik Sridharan. 2016. Exploiting the structure: Stochastic gradient methods using raw clusters. *NeurIPS* 29 (2016).
- [8] Naomi S Altman. 1992. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician* 46, 3 (1992), 175–185.
- [9] Marcelo Arenas, Leopoldo E. Bertossi, and Jan Chomicki. 1999. Consistent Query Answers in Inconsistent Databases. In *PODS*. ACM Press, 68–79.
- [10] Leopoldo E. Bertossi. 2011. *Database Repairing and Consistent Query Answering*. Morgan & Claypool Publishers.
- [11] Leopoldo E. Bertossi. 2019. Database Repairs and Consistent Query Answering: Origins and Further Developments. In *PODS*. ACM, 48–58.
- [12] Felix Biessmann, Tammo Rukat, Phillip Schmidt, Prathik Naidu, Sebastian Schelter, Andrey Taptunov, Dustin Lange, and David Salinas. 2019. DataWig: Missing Value Imputation for Tables. *JMLR* 20, 175 (2019), 1–6.
- [13] Vladimir Braverman, Dan Feldman, and Harry Lang. 2016. New Frameworks for Offline and Streaming Coreset Constructions. *CoRR* abs/1612.00889 (2016).
- [14] Trevor Campbell and Tamara Broderick. 2018. Bayesian Coreset Construction via Greedy Iterative Geodesic Ascent. In *ICML 2018*, Vol. 80. PMLR, 697–705.
- [15] Chengliang Chai, Lei Cao, Guoliang Li, Jian Li, Yuyu Luo, and Samuel Madden. 2020. Human-in-the-loop outlier detection. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 19–33.
- [16] Chengliang Chai, Guoliang Li, Ju Fan, and Yuyu Luo. 2020. Crowdchart: Crowdsourced data extraction from visualization charts. *IEEE Transactions on Knowledge and Data Engineering* 33, 11 (2020), 3537–3549.
- [17] Chengliang Chai, Guoliang Li, Jian Li, Dong Deng, and Jianhua Feng. 2016. Cost-effective crowdsourced entity resolution: A partial-order approach. In *SIGMOD*. 969–984.
- [18] Chengliang Chai, Guoliang Li, Jian Li, Dong Deng, and Jianhua Feng. 2018. A partial-order-based framework for cost-effective crowdsourced entity resolution. *The VLDB Journal* 27 (2018), 745–770.
- [19] Chengliang Chai, Jiaxin Liu, Nan Tang, Guoliang Li, and Yuyu Luo. 2022. Selective data acquisition in the wild for model charging. *PVLDB* 15, 7 (2022), 1466–1478.
- [20] Chengliang Chai, Jiayi Wang, Yuyu Luo, Zeping Niu, and Guoliang Li. 2022. Data management for machine learning: A survey. *TKDE* (2022).
- [21] Kai Lai Chung. 1954. On a stochastic approximation method. *The Annals of Mathematical Statistics* (1954), 463–483.
- [22] Ting Deng, Wenfei Fan, and Floris Geerts. 2016. Capturing Missing Tuples and Missing Values. *ACM Trans. Database Syst.* 41, 2 (2016), 10:1–10:47.
- [23] Irit Dinur and Samuel Safra. 2005. On the hardness of approximating minimum vertex cover. *Annals of mathematics* (2005), 439–485.
- [24] Alireza Farhangfar, Lukasz A. Kurgan, and Witold Pedrycz. 2007. A Novel Framework for Imputation of Missing Values in Databases. *IEEE Trans. Syst. Man Cybern. Part A* 37, 5 (2007), 692–709.
- [25] Dan Feldman. 2020. Introduction to Core-sets: an Updated Survey. *CoRR* abs/2011.09384 (2020).
- [26] Yunjun Gao and Xiaoye Miao. 2018. Query processing over incomplete databases. *Synthesis Lectures on Data Management* 10, 2 (2018), 1–122.
- [27] Congcong Ge, Yunjun Gao, Xiaoye Miao, Bin Yao, and Haobo Wang. 2020. A hybrid data cleaning framework using markov logic networks. *TKDE* 34, 5 (2020), 2048–2062.
- [28] Lovedeep Gondara and Ke Wang. 2017. Multiple Imputation Using Deep Denoising Autoencoders. *CoRR* abs/1705.02737 (2017).
- [29] Shuang Hao, Chengliang Chai, Guoliang Li, Nan Tang, Ning Wang, and Xiang Yu. 2020. Outdated fact detection in knowledge bases. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. IEEE, 1890–1893.
- [30] Thomas Hofmann, Aurelien Lucchi, Simon Lacoste-Julien, and Brian McWilliams. 2015. Variance reduced stochastic gradient descent with neighbors. *Advances in Neural Information Processing Systems* 28 (2015).
- [31] Jiawei Huang, Ruomin Huang, Wenjie Liu, Nikolaos M. Freris, and Hu Ding. 2021. A Novel Sequential Coreset Method for Gradient Descent Algorithms. In *ICML 2021*, Vol. 139. PMLR, 4412–4422.
- [32] José M. Jerez, Ignacio Molina, Pedro J. García-Laencina, Emilio Alba, Nuria Ribelles, Miguel Martín, and Leonardo Franco. 2010. Missing data imputation using statistical and machine learning methods in a real breast cancer problem. *Artif. Intell. Medicine* 50, 2 (2010), 105–115.

- [33] Bojan Karlas, Peng Li, Renzhi Wu, Nezihe Merve Gürel, Xu Chu, Wentao Wu, and Ce Zhang. 2020. Nearest Neighbor Classifiers over Incomplete Information: From Certain Answers to Certain Predictions. *Proc. VLDB Endow.* 14, 3 (2020), 255–267.
- [34] Shahidul Islam Khan and Abu Sayed Md. Latiful Hoque. 2020. SICE: an improved missing data imputation technique. *J. Big Data* 7, 1 (2020), 37.
- [35] Krishnateja Killamsetty, S Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. 2021. Grad-match: Gradient matching based data subset selection for efficient deep model training. In *ICML*. 5464–5474.
- [36] KrishnaTeja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh K. Iyer. 2021. GLISTER: Generalization based Data Subset Selection for Efficient and Robust Learning. In *AAAI 2021*. AAAI Press, 8110–8118.
- [37] Katrin Kirchhoff and Jeff A. Bilmes. 2014. Submodularity for Data Selection in Machine Translation. In *EMNLP 2014*. ACL, 131–141.
- [38] Sanjay Krishnan, Michael J. Franklin, Ken Goldberg, and Eugene Wu. 2017. BoostClean: Automated Error Detection and Repair for Machine Learning. *CoRR* abs/1711.01299 (2017).
- [39] Sanjay Krishnan, Jannan Wang, Michael J. Franklin, Ken Goldberg, Tim Kraska, Tova Milo, and Eugene Wu. 2015. SampleClean: Fast and Reliable Analytics on Dirty Data. *IEEE Data Eng. Bull.* 38, 3 (2015), 59–75.
- [40] Sanjay Krishnan, Jannan Wang, Eugene Wu, Michael J. Franklin, and Ken Goldberg. 2016. ActiveClean: Interactive Data Cleaning For Statistical Modeling. *Proc. VLDB Endow.* 9, 12 (2016), 948–959.
- [41] Claude Lemaréchal. 2012. Cauchy and the gradient method. *Doc Math Extra* 251, 254 (2012), 10.
- [42] Guoliang Li, Chengliang Chai, Ju Fan, Xueping Weng, Jian Li, Yudian Zheng, Yuanbing Li, Xiang Yu, Xiaohang Zhang, and Haitao Yuan. 2018. CDB: A crowd-powered database system. *Proceedings of the VLDB Endowment* 11, 12 (2018), 1926–1929.
- [43] Peng Li, Xi Rao, Jennifer Blase, Yue Zhang, Xu Chu, and Ce Zhang. 2021. CleanML: A Study for Evaluating the Impact of Data Cleaning on ML Classification Tasks. In *ICDE*. 13–24.
- [44] Hui Lin and Jeff Bilmes. 2011. A class of submodular functions for document summarization. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*. 510–520.
- [45] Jiabin Liu, Chengliang Chai, Yuyu Luo, Yin Lou, Jianhua Feng, and Nan Tang. 2022. Feature augmentation with reinforcement learning. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 3360–3372.
- [46] Jiabin Liu, Fu Zhu, Chengliang Chai, Yuyu Luo, and Nan Tang. 2021. Automatic data acquisition for deep learning. *Proceedings of the VLDB Endowment* 14, 12 (2021), 2739–2742.
- [47] Tongyu Liu, Ju Fan, Yingqiang Luo, Nan Tang, Guoliang Li, and Xiaoyong Du. 2021. Adaptive data augmentation for supervised learning over missing data. *Proc. VLDB Endow.* 14, 7 (2021), 1202–1214.
- [48] Chris Mayfield, Jennifer Neville, and Sunil Prabhakar. 2010. ERACER: a database approach for statistical inference and data cleaning. In *SIGMOD*. ACM, 75–86.
- [49] John T McCoy, Steve Kroon, and Lidia Auret. 2018. Variational autoencoders for missing data imputation with application to a simulated milling circuit. *IFAC-PapersOnLine* 51, 21 (2018), 141–146.
- [50] Xiaoye Miao, Yunjun Gao, Su Guo, and Wanqi Liu. 2018. Incomplete data management: a survey. *Frontiers of Computer Science* 12 (2018), 4–25.
- [51] Xiaoye Miao, Yangyang Wu, Lu Chen, Yunjun Gao, Jun Wang, and Jianwei Yin. 2021. Efficient and effective data imputation with influence functions. *Proceedings of the VLDB Endowment* 15, 3 (2021), 624–632.
- [52] Xiaoye Miao, Yangyang Wu, Lu Chen, Yunjun Gao, and Jianwei Yin. 2022. An Experimental Survey of Missing Data Imputation Algorithms. *TKDE* (2022).
- [53] Baharan Mirzasoleiman, Ashwinkumar Badanidiyuru, Amin Karbasi, Jan Vondrák, and Andreas Krause. 2015. Lazier than lazy greedy. In *AAAI*, Vol. 29.
- [54] Baharan Mirzasoleiman, Jeff A. Bilmes, and Jure Leskovec. 2020. Coresets for Data-efficient Training of Machine Learning Models. In *ICML 2020*, Vol. 119. 6950–6960.
- [55] Baharan Mirzasoleiman, Kaidi Cao, and Jure Leskovec. 2020. Coresets for robust training of deep neural networks against noisy labels. *NeurIPS* 33 (2020), 11465–11477.
- [56] Baharan Mirzasoleiman, Kaidi Cao, and Jure Leskovec. 2020. Coresets for Robust Training of Deep Neural Networks against Noisy Labels. In *NeurIPS 2020*.
- [57] Alexander Munteanu and Chris Schwiegelshohn. 2018. Coresets-methods and history: A theoreticians design pattern for approximation and streaming algorithms. *KI-Künstliche Intelligenz* 32, 1 (2018), 37–53.
- [58] Alfredo Nazábal, Pablo M. Olmos, Zoubin Ghahramani, and Isabel Valera. 2020. Handling incomplete heterogeneous data using VAEs. *Pattern Recognit.* 107 (2020), 107501.
- [59] Angelia Nedić and Dimitri Bertsekas. 2001. Convergence rate of incremental subgradient algorithms. In *Stochastic optimization: algorithms and applications*. Springer, 223–264.
- [60] Felix Neutatz, Binger Chen, Ziawasch Abedjan, and Eugene Wu. 2021. From Cleaning before ML to Cleaning for ML. *IEEE Data Eng. Bull.* (2021).

- [61] Andrew Ng. 2021. MLOPs: From Model-centric to Data-centric AI.
- [62] Xuedi Qin, Chengliang Chai, Yuyu Luo, Nan Tang, and Guoliang Li. 2020. Interactively discovering and ranking desired tuples without writing sql queries. In *SIGMOD*. 2745–2748.
- [63] Xuedi Qin, Chengliang Chai, Yuyu Luo, Tianyu Zhao, Nan Tang, Guoliang Li, Jianhua Feng, Xiang Yu, and Mourad Ouzzani. 2021. Ranking desired tuples by database exploration. In *ICDE*. IEEE, 1973–1978.
- [64] LKPJ Rduesseeun and P Kaufman. 1987. Clustering by means of medoids. In *Proceedings of the statistical data analysis based on the L1 norm conference, neuchatel, switzerland*, Vol. 31.
- [65] Patrick Royston and Ian R White. 2011. Multiple imputation by chained equations (MICE): implementation in Stata. *Journal of statistical software* 45 (2011), 1–20.
- [66] Indro Spinelli, Simone Scardapane, and Aurelio Uncini. 2020. Missing data imputation with adversarially-trained graph convolutional networks. *Neural Networks* 129 (2020), 249–260.
- [67] Daniel J. Stekhoven and Peter Bühlmann. 2012. MissForest - non-parametric missing value imputation for mixed-type data. *Bioinform.* 28, 1 (2012), 112–118.
- [68] Gilbert Strang. 2006. *Linear algebra and its applications*. Belmont, CA: Thomson, Brooks/Cole.
- [69] Bhekisipho Twala, Michelle Cartwright, and Martin J. Shepperd. 2005. Comparison of various methods for handling incomplete data in software engineering databases. In *ISESE 2005*. 105–114.
- [70] Richard Wu, Aoqian Zhang, Ihab F. Ilyas, and Theodoros Rekatsinas. 2020. Attention-based Learning for Missing Data Imputation in HoloClean. In *MLSys 2020*. mlsys.org.
- [71] Yangyang Wu, Xiaoye Miao, Yuchen Peng, Lu Chen, Yunjun Gao, and Jianwei Yin. 2022. An Interactive Data Imputation System. In *DASFAA*. Springer, 495–499.
- [72] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. 2018. GAIN: Missing Data Imputation using Generative Adversarial Nets. In *ICML*, Vol. 80. 5675–5684.
- [73] Aoqian Zhang, Shaoxu Song, Yu Sun, and Jianmin Wang. 2019. Learning Individual Models for Imputation. In *ICDE*. IEEE, 160–171.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009