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

Chengliang Chai Beijing Institute of Technology Jiabin Liu Beijing Institute of Technology Nan Tang QCRI, Qatar

Ju Fan Renmin University of China

Dongjing Miao Harbin Institute of Technology Jiayi Wang Tsinghua University Yuyu Luo Tsinghua University (Presenter) Guoliang Li Tsinghua University













# Data-effective ML

- Data cleaning: obtain high-quality data
  - Inconsistency
  - Outlier

(F

- Miss Labels
- Missing Values



-						
Name	Gender	Depart.	Age	Working Years		
Nancy	F	Sales	24	2		
Amy	М	Finance	42	10		
Lei	М	Market	35	2		
Lily	F	Market	36	8		
James	М	Sales	38	2		
Tom	М	Finance	36	6		



**High-quality Data** 



Name	Gender	Depart.	Age	Working Years	
Nancy	F	Sales	24	2	
Amy		Finance	42	10	
Lei	М		35		
Lily	F	Market		8	
James	М	Sales	38	2	
Tom	М		36	6	

Low-quality Data (Missing Value)

# Data-efficient ML

- > Algorithm:
  - Stochastic Gradient Decent (SGD)
- Data:

S

Coreset: Given a large train set D, Coreset is a core subset of D, denoted by C(D), which is selected to represent D such that M(C(D))≈M(D).



**High-quality Data** 

## Our Goal

#### Our goal is to achieve both *data-efficienct* and *data-effective* ML !

#### ML Model (Acc. : 96%)

10 min for training (12X)

Name	Gender	Depart.	Age	Working Years
Nancy	F	Sales	24	2
Amy		Finance	42	10
Lei	М		35	
Lily	F	Market		8
James	М	Sales	38	2
Tom	М		36	6

**Coreset Selection Over Incomplete Data** 

Name	Gender	Depart.	Age	Working Years
Nancy	F	Sales	24	2
Lily	F	Market	36	8
Tom	М	Finance	36	6

High-quality Data (Coreset)

Low-quality Data (Missing Value)

## Introduction to coreset selection

- How to select a coreset (a small subset) that has the same performance theoretically with the full?
- Gradient Descent

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

Gradient-based coreset selection

s.t.  $|C| \leq K$ 

$$C^{*} = \underset{C \subseteq D_{c}, w_{j} \geq 0}{\operatorname{arg min}} \max_{\theta \in \vartheta} \|\underbrace{\sum_{i=1}^{n} \nabla f_{i}(\theta)}_{\text{full gradient}} - \underbrace{\sum_{j=1}^{|C|} w_{j} \nabla f_{\gamma(j)}(\theta)}_{\text{coreset gradient}} \|,$$

$$\underbrace{full \operatorname{gradient}}_{\text{gradient approximation error}} \|$$

(2)

## Example of coreset selection

Dataset *D*<sub>c</sub>



 $\phi(i) = j$  indicates the mapping from  $D_c$  to C

 $w_i$  is the weight of each tuple in the coreset.

## Gradient-based coreset selection

> Minimize the gradient approximation error  $\|\sum_{i=1}^{n} \nabla f_i(\theta) - \sum_{i=1}^{|C|} w_j \nabla f_{\gamma(j)}(\theta)\|$ 

Triangle equation (Upper bound minimization of error)

$$\left\|\sum_{i=1}^{n} \nabla f_{i}(\theta) - \sum_{j=1}^{|C|} w_{j} \nabla f_{\gamma(j)}(\theta)\right\| \leq \sum_{i=1}^{n} \left\|\nabla f_{i}(\theta) - \nabla f_{\gamma(\phi_{\theta}(i))}(\theta)\right\|$$

> Assign every tuple to the tuple in C with most gradient similarity

$$\left\|\sum_{i=1}^{n} \nabla f_{i}(\theta) - \sum_{j=1}^{|C|} w_{j} \nabla f_{\gamma(j)}(\theta)\right\| \leq \sum_{i=1}^{n} \min_{c_{j} \in C} \left\|\nabla f_{i}(\theta) - \nabla f_{\gamma(j)}(\theta)\right\|$$

For convex problem, the bound of gradiet difference of a pair of tuples.

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

> Problem definition  $C^* = \underset{C \subseteq D_c}{\operatorname{arg\,min}} \sum_{i=1}^n \underset{c_j \in C}{\min} \|\mathbf{x}_i - \mathbf{x}_{\gamma(j)}\|$  **Greedy Solution** 

$$C^* = \underset{C \subseteq D_c}{\operatorname{arg\,min}} \sum_{i=1}^n \min_{c_j \in C} \|\mathbf{x}_i - \mathbf{x}_{\gamma(j)}\|$$

Greedy solution with approximate ratio (sub-modular):



Full train set  $D_c$  (size=n)

Next: Add which row in **D**<sub>c</sub> to **C** ?

**Greedy Solution**  $C^* = \underset{C \subseteq D_c}{\operatorname{arg\,min}} \sum_{i=1}^n \underset{c_j \in C}{\min} \|\mathbf{x}_i - \mathbf{x}_{\gamma(j)}\|$ 

Next: Add which row in **D**<sub>c</sub> to **C** ?



### **Coreset Selection over Incomplete Data**



- 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))

# Intuitive baselines (sequential)



Solution	Accuracy	Human Cost	Machine Cost
$(1) \mathbf{C}(\mathbf{H}(D))$	High	High	Low
$(2) \mathbf{C}(\mathbf{A}(D))$	Low	None	Low
$(3) \operatorname{H}(\operatorname{C}(D))$	Low	Low	Low
$(4) \mathbf{A}(\mathbf{C}(D))$	Low	None	Low

# Challenge

- Computing a good coreset from dirty data is to *accurately estimate the* ground truth of each missing value, which has multiple possible repairs.
- > The *combinations of all possible repairs* constitute a huge search space.

# Key idea of our Solutions (GoodCore)

- > Model the combinations of possible repairs as *possible worlds* of the original dirty data
- Selecting an expected optimal coreset that can represent the possible worlds of D by gradient approximation without training in advance

#### $D \rightarrow G(D)$ : Very Time-consuming



## An Example: Expected Optimal Coreset

Problem Definition:

$$C^* = \underset{C \subseteq D}{\arg\min} E[C], \text{ s.t. } |C| \le K$$

$$\mathbf{E}[C] = \sum_{k=1}^{|I_W|} p_k S_k$$



NP-hard with sub-modular property: can be solved by a greedy algorithm

## **Expectation Computation**

Brute-force method: prohibitively expensive to compute

$$E[C] = \sum_{k=1}^{|I_W|} p_k(\sum_{i=1}^n \min_{c_j \in C_k} s_{ij}).$$

Tuple-based expectation computation:





# A Three-loop Framework

3rd loop

First loop: Add tuples one-by-one

Second loop: Sample a bacth of tuples, and select the one with the largest utility to add.



Third loop: Compute the utility of each tuple in the second loop, considering all the possible worlds.

# Imputation-in-the-loop Strategies

- One tuple each iteration:
  - Impute (Auto/Human) the tuple with the highest <u>Algorithm 3: Batch algorithm of GoodCore</u> utility at each iteration.
     Imput: D, K, h, batch size b. Output: A coreset C, weight W.
  - Time complexity  $O(KhnL^2)_{1}$
- One batch each iteration with human-in-the-loop:
  - One tuple per iteration by humans is too clostly.
  - Impute a batch of tuples after several iterations.

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

## Experiments

Dataset	D	М	# 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

Table 1. Statistics of datasets

## Effectiveness



# Efficiency



Fig. 9. Efficiency of different methods.

### Human cost

Datasets	$G(D, \mathcal{O}^H)$	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	-

#### Table 2. Human cost of different methods

# Thanks!