

---

# Data Debugging is NP-hard for Classifiers Trained with SGD

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Data debugging is to find a subset of the training data such that the model obtained  
2 by retraining on the subset has a better accuracy. A bunch of heuristic approaches  
3 are proposed, however, none of them are guaranteed to solve this problem effec-  
4 tively. This leaves an open issue whether there exists an efficient algorithm to find  
5 the subset such that the model obtained by retraining on it has a better accuracy.  
6 To answer this open question and provide theoretical basis for further study on  
7 developing better algorithms for data debugging, we investigate the computational  
8 complexity of the problem named DEBUGGABLE. Given a machine learning  
9 model  $\mathcal{M}$  obtained by training on dataset  $D$  and a test instance  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$  where  
10  $\mathcal{M}(\mathbf{x}_{\text{test}}) \neq y_{\text{test}}$ , DEBUGGABLE is to determine whether there exists a subset  $D'$  of  
11  $D$  such that the model  $\mathcal{M}'$  obtained by retraining on  $D'$  satisfies  $\mathcal{M}'(\mathbf{x}_{\text{test}}) = y_{\text{test}}$ .  
12 To cover a wide range of commonly used models, we take SGD-trained linear  
13 classifier as the model and derive the following main results. (1) If the loss function  
14 and the dimension of the model are not fixed, DEBUGGABLE is NP-complete  
15 regardless of the training order in which all the training samples are processed  
16 during SGD. (2) For hinge-like loss functions, a comprehensive analysis on the  
17 computational complexity of DEBUGGABLE is provided; (3) If the loss function is a  
18 linear function, DEBUGGABLE can be solved in linear time, that is, data debugging  
19 can be solved easily in this case. These results not only highlight the limitations of  
20 current approaches but also offer new insights into data debugging.

## 21 1 Introduction

22 Given a machine learning model, data debugging is to find a subset of the training data such that  
23 the model will have a better accuracy if retrained on that subset [1]. Data debugging serves as a  
24 popular method of both data cleaning and machine learning interpretation. In the context of data  
25 cleaning, data debugging (*a.k.a.* training data debugging [2] or data cleansing [1]) can be used  
26 to improve the quality of the training data by removing the flaws leading to mispredictions [3–5].  
27 When it comes to ML interpretation, data debugging locates the part of the training data responsible  
28 for unexpected predictions of an ML model. Therefore it is also studied as a training data-based  
29 (*a.k.a.* instance-based [6]) interpretation, which is crucial for helping system developers and ML  
30 practitioners to debug ML system by reporting the harmful part of training data [7].

31 To solve the data debugging problem, existing researches adopt a two-phase score-based heuristic  
32 approach [2]. In the first phase, a score representing the estimated impact on the model accuracy is  
33 assigned to each training sample in the training data. It is hoped that the harmful part of training  
34 data gets a lower score than the other part. In the second phase, training samples with lower scores  
35 are removed greedily and the model is retrained on the modified training data. The two phases are  
36 carried out iteratively until a well-trained model is obtained. Most of the related works focus on

37 developing algorithms to estimate the scores efficiently in the first phase [8–16], but rarely study the  
38 effectiveness of the entire two-phase approach.

39 Since it is computationally intractable to estimate the score for all possible subsets of the training  
40 data, it is often assumed that the score representing the impact of a subset is approximately equal  
41 to the sum of the scores of each individual training samples from the subset. However, Koh et. al.  
42 [10] showed this is not always the case. For a bunch of subsets sampled from the training data,  
43 they empirically studied the difference between the estimated impact and the actual impact of each  
44 subset by taking influence functions as the scoring method. The estimated impact is calculated by  
45 summing up the score by influence function of each training samples in the subset, and the actual  
46 impact is measured by the improvement of accuracy of the model retrained after removing the subset  
47 from training data. They found that the estimated impact tends to underestimate the actual impact.  
48 Removing a large number of training samples could result in a large deviation between estimated  
49 and actual impacts. Although an upper bound of the deviation under certain assumptions has been  
50 derived, it is still unknown whether the deviation can be reduced or eliminated efficiently.

51 The above deviation also poses challenges to the effectiveness of the entire approach. Suppose the  
52 influence function is adopted as the scoring method, the accuracy of the model is not guaranteed  
53 to improve due to the deviation reported in [10] if a large group of training samples are removed  
54 during each iteration. Moreover, there is no theoretical analysis for the effectiveness of the greedy  
55 approach in the second phase. Even if only one training sample is removed during each iteration  
56 of the two-phase approach, the accuracy of the model is still not guaranteed to be improved. The  
57 effectiveness of the entire two-phase approach is therefore not assured. This leaves the following  
58 open problem:

59 **Problem 1.1.** Is there an efficient algorithm to find the subset of the training data, such that the  
60 model obtained by retraining on it has a better accuracy?

61 The computational complexity results presented in this paper demonstrate that it is unlikely to solve  
62 the data debugging problem efficiently in polynomial time. To figure out its hardness, we study the  
63 problem DEBUGGABLE which is the decision version of data debugging when the test set consists of  
64 only one instance. Formally, DEBUGGABLE is defined as follows:

65 **Problem 1.2** (DEBUGGABLE). Given a classifier  $\mathcal{M}$ , its training data  $T$ , a test instance  $(\mathbf{x}, y)$ . Is  
66 there a  $T' \subseteq T$ , such that  $\mathcal{M}$  predicts  $y$  on  $\mathbf{x}$  if retrained on  $T'$ ?

67 Basically, we prove that DEBUGGABLE is NP-complete, which means data debugging is unlikely  
68 to be solved in polynomial time. This result answers the open question mentioned above directly,  
69 this is, the large deviation of estimated impacts [10] cannot be reduced or eliminated efficiently. This  
70 is because if the impact of a subset of the training data could be accurately estimated as the sum of  
71 the impact of each training sample in the subset, data debugging can be solved in polynomial time,  
72 which is impossible unless  $P=NP$ .

73 Although DEBUGGABLE is generally intractable, we still hope to develop efficient algorithms tailored  
74 to specific cases. Thus it is necessary to figure out the root cause of the hardness for DEBUGGABLE.  
75 Previous research are always conducted based on the belief that the complexity of data debugging is  
76 due to the chosen model architecture is complicated. However, we show that at least for models trained  
77 by stochastic gradient descent (SGD), the hardness stems from the hyper-parameter configuration  
78 selected for the SGD training, which was not yet aware of by previous work. To cover a wide range of  
79 commonly used machine learning models, we take linear classifiers as the model and show that even  
80 for linear classifiers, DEBUGGABLE is NP-hard as long as they are trained by SGD. Moreover, we  
81 provided a comprehensive analysis on hyper-parameter configurations that affect the computational  
82 complexity of DEBUGGABLE, including the loss function, the model dimension and the training  
83 order. Training order, *a.k.a.* training data order [17] or order of training samples [18], refers to the  
84 order in which each training sample is considered during the SGD. Detailed complexity results are  
85 shown in Table 1.

86 Our contribution can be concluded as follows:

- 87 • We studied the computational complexity of data debugging and showed that data debugging  
88 is NP-hard for linear classifiers in the general setting for *all possible training orders*.
- 89 • We studied the complexity of DEBUGGABLE when the loss is fixed as the hinge-like  
90 function. For 2 or higher dimension, DEBUGGABLE is NP-complete when the training order

Table 1: Computational complexity of the data debugging problem

Loss Function	Dimension	Training Order	Complexity
Not Fixed	Not Fixed	-	NP-hard
Hinge-like	$\geq 2$	Adversarially Chosen	NP-hard
Hinge-like, $\beta < 0$	1	Adversarially Chosen	NP-hard
Hinge-like, $\beta \geq 0$	1	-	Linear Time
Linear	-	-	Linear Time

91 is adversarially chosen; For one-dimensional cases, DEBUGGABLE can be NP-hard when  
 92 the interception  $\beta < 0$ , and is solvable in linear time when  $\beta \geq 0$ .

93 • We proved that DEBUGGABLE is solvable in linear time when the loss function is linear.

94 Moreover, we have a discussion on the implications of these complexity results for machine learning  
 95 interpretability and data quality, as well as limitations of score-based greedy methods. Our results  
 96 suggest the further study as follows. (1) It is better to characterize the training sample and find the  
 97 criterion which can be used to decide the existence of efficient algorithms; (2) Designing algorithms  
 98 with CSP-solver is a potential way to solve data debugging more efficiently than the brute-force one;  
 99 (3) Developing random algorithms is a potential way to solve data debugging successfully with high  
 100 probability.

## 101 1.1 Related Works

102 The solution of data debugging has applications in database query results reliability enhancement  
 103 [2, 19], training data cleaning [1] and machine learning interpretation [9, 8, 10, 20, 21]. Existing  
 104 works on data debugging mainly adopt a two-phase approach, which scores the training samples in the  
 105 first phase and greedily deletes training samples with lower scores in the second phase. Most of the  
 106 research focus on the first phase. There are mainly two ways of scoring adopted for data debugging in  
 107 practice. Leave-one-out (LOO) retraining is a widely studied way, which evaluates the contribution of  
 108 a training sample through the difference in the model’s accuracy trained without that training sample.  
 109 To avoid the cost of model retraining, Koh and Liang took influence functions as an approximation of  
 110 LOO [8]. After that, various extensions and improvements of the influence function based method  
 111 are proposed, such as Fisher kernel [9], influence function for group impacts [10], second-order  
 112 approximations [11] and scalable influence functions [12]. Another way is Shapley-based scoring,  
 113 where the impact of a training sample is measured by its average marginal contribution to all subsets  
 114 of the training data [13]. Since Shapley-base scoring suffers from expensive computational cost [22],  
 115 recent works focus on techniques that efficiently estimate the Shapley value, including Monte-Carlo  
 116 sampling [13], group testing [14, 15] and using proxy models such as  $k$ -NN [16, 3]. However,  
 117 those methods do not admit any theoretical guarantee on the effectiveness. This paper discusses the  
 118 limitations of the above methods and suggests some future directions on data debugging.

## 119 2 Preliminaries and Problem Definition

120 **Linear classifiers.** Formally, a (binary) linear classifier is a function  $\lambda_{\mathbf{w}} : \mathbb{R}^d \rightarrow \{-1, 1\}$ , where  $d$  is  
 121 called its *dimension* and  $\mathbf{w} \in \mathbb{R}^d$  its parameter. Without loss of generality, the bias term of a linear  
 122 classifier is set as zero in this paper. All vectors in this paper are assumed to be *column* vectors. For  
 123 an input  $\mathbf{x}$ , the value of  $\lambda_{\mathbf{w}}$  is defined as

$$\lambda_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} \geq 0 \\ -1 & \text{otherwise.} \end{cases}$$

124 We denote the class of linear models as  $\Lambda$ .

125 **Training data.** A *training sample* is a pair  $(\mathbf{x}, y)$  in which  $\mathbf{x} \in \mathbb{R}^d$  is the input and  $y \in \{-1, 1\}$  is  
 126 the label of  $\mathbf{x}$ . The *training data* is a multiset of training samples. We employ  $\mathbf{w} \xrightarrow{T} \mathbf{w}'$  to denote  
 127 that the parameter  $\mathbf{w}'$  is obtained by training the parameter  $\mathbf{w}$  on the training data  $T$ , and employ  
 128  $\mathbf{w} \xrightarrow{(\mathbf{x}, y)} \mathbf{w}'$  to denote that  $\mathbf{w}'$  is obtained by training  $\mathbf{w}$  on the training sample  $(\mathbf{x}, y)$ .

129 **Loss functions and learning rates.** Binary linear classifiers typically use unary functions on  $y\mathbf{w}^\top \mathbf{x}$   
 130 as their loss functions [23]. Therefore we only consider loss functions of the form  $\mathcal{L} : y\mathbf{w}^\top \mathbf{x} \mapsto \mathbb{R}$   
 131 for the rest of the paper.

132 The *linear* loss is in the form of

$$\mathcal{L}_{\text{lin}}(y\mathbf{w}^\top \mathbf{x}) = -\alpha(y\mathbf{w}^\top \mathbf{x} + \beta).$$

133 The *hinge-like* loss function is defined as the following form

$$\mathcal{L}_{\text{hinge}}(y\mathbf{w}^\top \mathbf{x}) = \begin{cases} -\alpha(y\mathbf{w}^\top \mathbf{x} + \beta), & y\mathbf{w}^\top \mathbf{x} < \beta \\ 0, & \text{otherwise.} \end{cases}$$

134 We call  $\beta$  as the *interception* of  $\mathcal{L}_{\text{hinge}}$ . We represent the learning rate of a model using a vector  
 135  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_d)$ , where  $\eta_i \geq 0$  and each parameter  $w_i$  can be updated with the corresponding  
 136 learning rate  $\eta_i$ .

137 **Stochastic gradient descent.** The stochastic gradient descent (SGD) method updates parameter  $\mathbf{w}$   
 138 from its initial value  $\mathbf{w}^{(0)}$  through several epochs. During each epoch, the SGD goes through the  
 139 entire set of training samples in some training order through several iterations. The training order is  
 140 defined as a sequence of training samples, in the form of  $(\mathbf{x}_1, y_1) \dots (\mathbf{x}_n, y_n)$ . For  $1 \leq i < j \leq n$ ,  
 141  $(\mathbf{x}_i, y_i)$  is considered before  $(\mathbf{x}_j, y_j)$  during the SGD. We use  $w_i$  to denote the  $i$ -th coordinate of  $\mathbf{w}$ .  
 142 We also use  $\mathbf{w}^{(e,k)}$  to denote the value of  $\mathbf{w}$  at the end of  $k$ -th iteration of epoch  $e$  and use  $\mathbf{w}^{(e)}$  to  
 143 denote the value of  $\mathbf{w}$  after the end of epoch  $e$ . Assuming  $(\mathbf{x}, y)$  to be the training sample considered  
 144 at iteration  $k$ , the stochastic gradient descent (SGD) method updates parameter  $w_i$  for each  $i$  by

$$w_i^{(e,k)} \leftarrow w_i^{(e,k-1)} - \eta_i \cdot \frac{\partial \mathcal{L}(y(\mathbf{w}^{(e,k-1)})^\top \mathbf{x})}{\partial w_i} \quad (1)$$

145 In other words, we have

$$\mathbf{w}^{(e,k)} \leftarrow \mathbf{w}^{(e,k-1)} - \boldsymbol{\eta} \otimes \nabla \mathcal{L}(y(\mathbf{w}^{(e,k-1)})^\top \mathbf{x})$$

146 where  $\boldsymbol{\eta} \otimes \nabla \mathcal{L} = (\eta_1 \frac{\partial \mathcal{L}}{\partial w_1}, \dots, \eta_d \frac{\partial \mathcal{L}}{\partial w_d})$  is the Hadamard product. We say a training sample  $\mathbf{x}$   
 147 is *activated* at iteration  $k$  during epoch  $e$  if  $\nabla \mathcal{L}(y(\mathbf{w}^{(e,k-1)})^\top \mathbf{x}) \neq 0$ . The SGD terminates at  
 148 the end of epoch  $e$  if  $\|\mathbf{w}^{(e-1)} - \mathbf{w}^{(e)}\| < \varepsilon$  for threshold  $\varepsilon$  or  $e$  reached some predetermined  
 149 value. We denote  $\mathbf{w}^* = \mathbf{w}^{(e)}$ . A linear classifier trained by SGD with the meta-parameters  
 150 mentioned above is denoted as  $\text{SGD}_\Delta(\mathcal{L}, \boldsymbol{\eta}, \varepsilon, T) = \lambda_{\mathbf{w}^*}$ . With a slight abuse of notation, we define  
 151  $\text{SGD}_\Delta(\mathcal{L}, \boldsymbol{\eta}, \varepsilon, T, \mathbf{x}) = \lambda_{\mathbf{w}^*}(\mathbf{x})$ . We also use  $\text{SGD}_\Delta(T, \mathbf{x})$  to avoid cluttering when the context is clear.

152 **Problem definition.** With the above definitions, DEBUGGABLE for SGD-trained linear classifiers  
 153 can be formalized as follows:

DEBUGGABLE-LIN

**Input:** Training data  $T$ , loss function  $\mathcal{L}$ , initial parameter  $\mathbf{w}^{(0)}$ , learning  
 rate  $\boldsymbol{\eta}$ , threshold  $\varepsilon$  and instance  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$ .

**Output:** “Yes”: if  $\exists \Delta \subseteq T$  such that  $\text{SGD}_\Delta(\mathcal{L}, \boldsymbol{\eta}, \varepsilon, T \setminus \Delta, \mathbf{x}_{\text{test}}) = y_{\text{test}}$ ;  
 “No”: otherwise.

155 We say  $\text{SGD}_\Delta(\mathcal{L}, \boldsymbol{\eta}, \varepsilon, T)$  is *debuggable* on  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$  if  $(\mathcal{L}, \mathbf{w}^{(0)}, \boldsymbol{\eta}, \varepsilon, T, \mathbf{x}_{\text{test}}, y_{\text{test}})$  is a yes-instance  
 156 of DEBUGGABLE-LIN, and not *debuggable* on  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$  otherwise.

### 157 3 Results for Unfixed Loss Functions

158 In this section, we prove the NP-hardness of DEBUGGABLE-LIN. Intuitively, DEBUGGABLE-LIN is  
 159 to determine whether there exists a subset  $T' \subseteq T$  where activated training samples within  $T'$  drive  
 160 the parameter  $\mathbf{w}$  toward the region defined by  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} > 0$ . The activation of training samples  
 161 depends on the complex interaction between the training data and the model.

162 **Theorem 3.1.** DEBUGGABLE-LIN is NP-hard for all training orders.

163 We only show the proof sketch and leave the details in the appendix.

164 *Proof Sketch.* We build a reduction from an NP-hard problem MONOTONE 1-IN-3 SAT [24]:

**MONOTONE 1-IN-3 SAT**

**Input:** A 3-CNF formula  $\varphi$  with no negation signs.

**Output:** “Yes”: if  $\varphi$  has a 1-in-3 assignment, under which each clause contains exactly one true literal;  
“No”: otherwise.

166 For example,  $\varphi_1 = (x_1 \vee x_2 \vee x_3) \wedge (x_2 \vee x_3 \vee x_4)$  is a yes-instance because  $(x_1, x_2, x_3, x_4) =$   
167  $(\text{T}, \text{F}, \text{F}, \text{T})$  is an 1-in-3 assignment;  $\varphi_2 = (x_1 \vee x_2 \vee x_3) \wedge (x_2 \vee x_3 \vee x_4) \wedge (x_1 \vee x_2 \vee x_4) \wedge (x_1 \vee x_3 \vee x_4)$   
168 is a no-instance.

169 Given a 3-CNF formula  $\varphi$ , our goal is to construct a configuration of the training process, such that  
170 the resulting model outputs the correct answer if and only if its training data  $T'$  encodes an 1-in-3  
171 assignment  $\nu$  of  $\varphi$ . This can be done by carefully designing the encoding so that for each  $x_i \in \varphi$ ,  
172  $\nu(x_i) = \text{TRUE}$  if and only if  $\mathbf{t}_{x_i} \in T'$ . Finally, we can construct some  $T$  with  $T \supseteq T' \cup \{\mathbf{t}_{x_i} | x_i \in \varphi\}$ ,  
173 such that some classifier trained on  $T$  is a yes-instance of DEBUGGABLE-LIN if and only if  $\varphi$  is a  
174 yes-instance of MONOTONE 1-IN-3 SAT, thereby finishing our proof.

175 **The reduction.** Suppose  $\varphi$  has  $m$  clauses and  $n$  variables, let  $N = n + 2m + 1$ . We set the dimension  
176 of the linear classifier to  $N$ .

177 The input. Each coordinate of the input is named as

$$\mathbf{x} = (x_{c_1}, \dots, x_{c_m}, x_{x_1}, \dots, x_{x_n}, x_{b_1}, \dots, x_{b_m}, x_{\text{dummy}})^\top$$

178 We also use  $x_i$  to denote the  $i$ -th coordinate of  $\mathbf{x}$ .

179 The parameters. Each coordinate of the parameter is named as

$$\mathbf{w} = (w_{c_1}, \dots, w_{c_m}, w_{x_1}, \dots, w_{x_n}, w_{b_1}, \dots, w_{b_m}, w_{\text{dummy}})^\top$$

180 We also use  $w_i$  to denote the  $i$ -th coordinate of  $\mathbf{w}$ . Each  $w_{x_j}$  represents the truth value of variable  $x_j$ ,  
181 where 1 represents TRUE and -1 represents FALSE. Similarly, each  $w_{c_j}$  represents the truth value of  
182 clause  $c_j$  based on the value of its variables.  $w_{b_j}$  and  $w_{\text{dummy}}$  are used for convenience of proof.

183 The initial value of the parameter is set to

$$\mathbf{w}^{(0)} = \left( \overbrace{\frac{1}{2}, \dots, \frac{1}{2}}^m, \overbrace{-1, \dots, -1}^n, \overbrace{-1, \dots, -1}^m, 1 \right)^\top$$

184 Loss function. We denote  $U(x_0, \delta) := \{x | x_0 - \delta < x < x_0 + \delta\}$  as the  $\delta$ -neighborhood of  $x_0$  and  
185 define  $U(\pm x_0, \delta) = U(x_0, \delta) \cup U(-x_0, \delta)$ . We define the *local ramp function* as

$$r_{x_0, \delta}(x) = \begin{cases} 0 & , x \leq x_0 - \delta; \\ x - x_0 + \delta & , x \in U(x_0, \delta); \\ 2\delta & , x \geq x_0 + \delta. \end{cases}$$

186 The loss function is defined as

$$\mathcal{L} = -\frac{12N}{5} r_{-5, 0.01}(\mathbf{y}\mathbf{w}^\top \mathbf{x}) - r_{-\frac{1}{2}, 0.26}(\mathbf{y}\mathbf{w}^\top \mathbf{x}) - \frac{1}{1000N} \sum_{x_0 \in \{\pm 1, \pm 3\}} r_{x_0, 0.01}(\mathbf{y}\mathbf{w}^\top \mathbf{x}).$$

187  $\mathcal{L}$  is monotonically decreasing with derivatives

$$\frac{\partial \mathcal{L}}{\partial w_i} = \begin{cases} -\frac{12N}{5} \cdot yx_i & , \mathbf{y}\mathbf{w}^\top \mathbf{x} \in U(-5, 0.01); \\ -yx_i & , \mathbf{y}\mathbf{w}^\top \mathbf{x} \in U(-\frac{1}{2}, 0.26); \\ -\frac{1}{1000N} yx_i & , \mathbf{y}\mathbf{w}^\top \mathbf{x} \in \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U(x_0, 0.01); \\ 0 & , \text{otherwise.} \end{cases} \quad (2)$$

Table 2: Training data for  $\text{var}(i)$ 

$x_{x_i}$	$y$
5	1

Table 3: Training data for  $\text{clause}(i, i_1, i_2, i_3)$ 

$x_{c_i}$	$x_{x_{i_1}}$	$x_{x_{i_2}}$	$x_{x_{i_3}}$	$x_{b_i}$	$y$
1	1	1	1	$\frac{1}{2}$	1

188 Learning rate. The learning rate for SGD is set to be

$$\boldsymbol{\eta} = (\overbrace{5, \dots, 5}^m, \overbrace{\frac{1}{6N}, \dots, \frac{1}{6N}}^n, \overbrace{2000N, \dots, 2000N}^m, 1)^\top.$$

189 Training data. We define two gadgets,  $\text{var}(i)$  and  $\text{clause}(i, i_1, i_2, i_3)$ , as illustrated in Table 2 and  
 190 3. All the unspecified coordinates are set to zero. We use  $T_0$  to denote the training data.  $\text{var}(i)$   
 191 is contained in  $T_0$  if and only if  $x_i \in \varphi$ , and  $\text{clause}(i, i_1, i_2, i_3)$  is contained in  $T_0$  if and only if  
 192  $c_i = (x_{i_1} \vee x_{i_2} \vee x_{i_3}) \in \varphi$ .

193 Threshold and instance. The threshold  $\varepsilon$  can be any fixed value in  $\mathbb{R}_+$ . The instance is defined as  
 194  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$ , where  $y_{\text{test}} = 1$  and

$$\mathbf{x}_{\text{test}} = (\overbrace{1, \dots, 1}^m, \overbrace{0, \dots, 0}^{n+m}, \frac{-11m + 5}{2})^\top.$$

195 The following reduction works for all possible training orders. Intuitively, during the training process,  
 196 each  $\text{var}(i)$  in the training data will set  $w_{x_i}$  to around 1 (that is, mark  $x_i$  as TRUE) in the first epoch,  
 197 and each  $\text{clause}(i, i_1, i_2, i_3)$  will set  $w_{c_i}$  to near  $\frac{11}{2}$  in the second epoch, if and only if exactly one  
 198 of  $w_{x_{i_1}}, w_{x_{i_2}}, w_{x_{i_3}}$  is near 1 and the others near  $-1$  (that is, mark  $c_i$  as satisfied if exactly one of  
 199 its literals is TRUE and the others FALSE). The training process terminates at the end of the second  
 200 epoch.  $\square$

## 201 4 Results for Fixed Loss Functions

202 We have proved the NP-hardness for DEBUGGABLE-LIN when the loss function is not fixed. In  
 203 this section, we study the complexity when the loss function is fixed as linear and hinge-like  
 204 functions. Assuming that SGD terminates after only one epoch with a fixed order, we will show  
 205 that DEBUGGABLE-LIN is solvable in linear time for linear loss. For hinge-like loss functions,  
 206 DEBUGGABLE-LIN can be solved in linear time only when the dimension  $d = 1$  and the interception  
 207  $\beta \geq 0$ . For the rest cases, DEBUGGABLE-LIN becomes NP-hard.

### 208 4.1 The Easy Case

209 We start with the linear loss function  $\mathcal{L} = -\alpha(y\mathbf{w}^\top \mathbf{x} + \beta)$ , with which all the training data are  
 210 activated and  $\mathbf{w}^* = \mathbf{w}^*(T) = \mathbf{w}^{(0)} + \sum_{(\mathbf{x}, y) \in T} \alpha y \boldsymbol{\eta} \otimes \mathbf{x}$ . Since  $y_{\text{test}} \in \{-1, 1\}$ , DEBUGGABLE-LIN  
 211 is equivalent to deciding whether

$$\max_{T' \subseteq T} \{y_{\text{test}}(\mathbf{w}^*(T'))^\top \mathbf{x}_{\text{test}}\} > 0.$$

212 A training sample  $(\mathbf{x}, y)$  is “good” if  $y_{\text{test}}(\alpha y \boldsymbol{\eta} \otimes \mathbf{x})^\top \mathbf{x}_{\text{test}} > 0$  and “bad” otherwise. The *good*  
 213 *training-sample assessment* (GTA) algorithm, as shown in Algorithm 1, deals with this situation by  
 214 greedily picking all “good” training samples.

215 Denoting  $T^*$  as the set of all good data in  $T$ , it follows that

$$\begin{aligned} y_{\text{test}}(\mathbf{w}^*(T^*))^\top \mathbf{x}_{\text{test}} &= y_{\text{test}}(\mathbf{w}^{(0)})^\top \mathbf{x}_{\text{test}} + \sum_{(\mathbf{x}, y) \in T^*} y_{\text{test}}(\alpha y \boldsymbol{\eta} \otimes \mathbf{x})^\top \mathbf{x}_{\text{test}} \\ &\geq y_{\text{test}}(\mathbf{w}^{(0)})^\top \mathbf{x}_{\text{test}} + \sum_{(\mathbf{x}, y) \in T'} y_{\text{test}}(\alpha y \boldsymbol{\eta} \otimes \mathbf{x})^\top \mathbf{x}_{\text{test}} \end{aligned}$$

216 for all  $T' \subseteq T$ . Hence  $\max_{T' \subseteq T} \{y_{\text{test}}(\mathbf{w}^*(T'))^\top \mathbf{x}_{\text{test}}\} = y_{\text{test}}(\mathbf{w}^*(T^*))^\top \mathbf{x}_{\text{test}}$  and DEBUGGABLE-  
 217 LIN can be solved by GTA in linear time. The following theorem is straightforward.

218 **Theorem 4.1.** DEBUGGABLE-LIN is linear time solvable for linear loss functions.

---

**Algorithm 1:** Good Training-sample Assessment (GTA)

---

**Input:** Training data  $T$ , loss function  $\mathcal{L}$ , initial parameter  $\mathbf{w}^{(0)}$ , learning rate  $\eta$ , threshold  $\varepsilon$  and test instance  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$ .

**Output:** TRUE, iff  $\text{SGD}_\Lambda(\mathcal{L}, \eta, \varepsilon, T)$  is debuggable on  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$ .

```

1  $\mathbf{w} \leftarrow \mathbf{w}^{(0)}$ ;
2 for  $(\mathbf{x}, y) \in T$  do
3   if  $y_{\text{test}}(\alpha y \eta \otimes \mathbf{x})^\top \mathbf{x}_{\text{test}} > 0$  then
4      $\mathbf{w} \leftarrow \mathbf{w} + \alpha y \eta \otimes \mathbf{x}$ ;
5   end
6 end
7 if  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} \geq 0$  then
8   return TRUE;
9 end
10 return FALSE;

```

---

220 GTA is still effective for one-dimensional classifiers trained with hinge-like losses when  $\beta \geq 0$ .

221 **Theorem 4.2.** DEBUGGABLE-LIN is linear time solvable for hinge-like loss functions, when  $d = 1$   
 222 and  $\beta \geq 0$ .

223 *Proof.* It suffices to prove that if  $\exists T' \subseteq T$  such that  $\text{SGD}_\Lambda(T', x_{\text{test}}) = y_{\text{test}}$ ,  $\text{SGD}_\Lambda(T^*, x_{\text{test}}) = y_{\text{test}}$ .

224 a) Suppose all the data in  $T^*$  are activated, we have

$$\begin{aligned}
 y_{\text{test}} w^*(T^*) x_{\text{test}} &= y_{\text{test}} w^{(0)} x_{\text{test}} + \sum_{(x,y) \in T^*} y_{\text{test}} \alpha y \eta x x_{\text{test}} \\
 &\geq y_{\text{test}} w^{(0)} x_{\text{test}} + \sum_{(x,y) \in T' \cap T^*} y_{\text{test}} \alpha y \eta x x_{\text{test}} + \sum_{(x,y) \in T' \setminus T^*} y_{\text{test}} \alpha y \eta x x_{\text{test}} \\
 &= y_{\text{test}} w^*(T') x_{\text{test}} \geq 0
 \end{aligned}$$

225 b) Suppose  $(x, y) \in T^*$  is the first inactivated data during the training phase, and  $w$  is the current  
 226 parameter, we have  $ywx > \beta$ . Since  $\alpha \eta \cdot (xy) \cdot (x_{\text{test}} y_{\text{test}}) \geq 0$ , we have  $(x_{\text{test}} y_{\text{test}}) \cdot w \geq 0$ . Let  $T''$  be  
 227 the set of training data appeared before  $(x, y)$ , we have  $y_{\text{test}} w^*(T^*) x_{\text{test}} \geq y_{\text{test}} w^*(T'') x_{\text{test}} \geq 0$ .  $\square$

## 228 4.2 The Hard Case

229 The gradient of training data may not always be activated and could be affected by the training order.  
 230 When the training order is adversarially chosen, the following theorem shows that DEBUGGABLE-LIN  
 231 is NP-hard for all  $d \geq 2$  and  $\beta \in \mathbb{R}$ .

232 **Theorem 4.3.** If the training order is adversarially chosen and  $d \geq 2$ , DEBUGGABLE-LIN is NP-hard  
 233 for each hinge-like loss function at every constant learning rate.

234 *Proof sketch.* Since the result can be easily extended for all  $d > 2$  by padding the other  $d - 2$   
 235 dimensions with zeros, we only prove for the case of  $d = 2$ . We assume  $\beta \geq -1$  and leave the  
 236  $\beta < -1$  case to the appendix. To avoid cluttering, we further assume  $\eta = \mathbf{1}$  and  $\alpha = 1$ . The proof  
 237 can be easily generalized by appropriately re-scaling the constructed vectors.

238 We build a reduction from the subset sum problem, which is well-known to be NP-hard:

SUBSET SUM

**Input:** A set of positive integer  $S$ , and a positive integer  $t$ .

**Output:** “Yes”: if  $\exists S' \subseteq S$  such that  $\sum_{a \in S'} a = t$ ;  
 “No”: otherwise.

239

240 Suppose  $n = |S|$ ,  $m = \max_{a \in S} \{a\}$ ,  $\gamma = \max\{\beta, 1\}$  and  $S = \{a_1, a_2, \dots, a_n\}$ . We further assume  
 241  $n > 1$ . Let the training data be

$$T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} \cup \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$$

242 where  $\mathbf{x}_i y_i = (\frac{\sqrt{\gamma}}{n+1}, 3\sqrt{\gamma} a_i)$  for all  $1 \leq i \leq n$ ,  $\mathbf{x}_c y_c = ((18n^2 m^2 - 2)\sqrt{\gamma}, -3t\sqrt{\gamma})$ ,  $\mathbf{x}_b y_b =$   
 243  $(\sqrt{\gamma}, -\sqrt{\gamma})$ ,  $\mathbf{x}_a y_a = (\sqrt{\gamma}, \sqrt{\gamma})$ . Let  $\mathbf{w}^{(0)} = (-18n^2 m^2 \sqrt{\gamma}, 0)$ . Let the test instance  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$   
 244 satisfy  $\mathbf{x}_{\text{test}} y_{\text{test}} = (1, 0)$ .

245 Let the training order be  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)$ .

246 For each  $1 \leq i < n$ , suppose  $\mathbf{w}^{(0)} \xrightarrow{T \cap \{(\mathbf{x}_i, y_i) \mid 1 \leq j \leq i\}}$   $\mathbf{w}_i$ , we have

$$\begin{aligned} y_{i+1} \mathbf{w}_i^\top \mathbf{x}_{i+1} &\leq \frac{\sqrt{\gamma}}{n+1} (-18n^2 m^2 \sqrt{\gamma} + \frac{\sqrt{\gamma}^i}{n+1}) + 3\sqrt{\gamma} a_{i+1} \sum_{j=1}^i 3\sqrt{\gamma} a_j \\ &\leq \gamma \left( -\frac{n-1}{n+1} \cdot 9nm^2 + \frac{n}{(n+1)^2} \right) < -1 \leq \beta \end{aligned}$$

247 This means all the  $T \setminus \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$  can be activated. Thus the resulting parameter  
 248 trained by  $T \setminus \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$  is

$$\mathbf{w}_c = \mathbf{w}^{(0)} + \sum_{i=1}^n \mathbf{x}_i y_i = \left( -18n^2 m^2 \sqrt{\gamma} + \frac{\sqrt{\gamma} |T^*|}{n+1}, 3\sqrt{\gamma} \sum_{i=1}^n a_i \right).$$

249 It now suffices to prove that for all  $S' \subseteq S$ ,  $\sum_{a \in S'} a = t$  if and only if  $\exists T' \subseteq T$  such that  
 250  $\mathbf{w} : \mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}$  satisfies  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} > 0$ .

251 If: Suppose  $\exists S' \subseteq S$  such that  $\sum_{a \in S'} a = t$ , we prove that  $\exists T' \subseteq T$  such that  $y_{\text{test}} (\mathbf{w}^*)^\top \mathbf{x}_{\text{test}} > 0$   
 252 for  $\mathbf{w}^*$  satisfying  $\mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}^*$ .

253 Let  $T^* = \{(\mathbf{x}_i, y_i) \mid a_i \in S'\}$ ,  $T' = T^* \cup \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$ . We have

$$\mathbf{w}_c = (-18n^2 m^2 \sqrt{\gamma} + \frac{\sqrt{\gamma} |T^*|}{n+1}, 3\sqrt{\gamma} \sum_{a_i \in S'} a_i) = (-18n^2 m^2 \sqrt{\gamma} + \frac{\sqrt{\gamma} |T^*|}{n+1}, 3\sqrt{\gamma} t).$$

254 And therefore  $y_c \mathbf{w}_c^\top \mathbf{x}_c = \gamma \left( (-18n^2 m^2 + \frac{|T^*|}{n+1})(18n^2 m^2 - 2) - 9t^2 \right) < -1 \leq \beta$ , so

$$\mathbf{w}_c \xrightarrow{(\mathbf{x}_c, y_c)} \mathbf{w}_b = \mathbf{w}_c + \mathbf{x}_c y_c = (\sqrt{\gamma} (\frac{|T^*|}{n+1} - 2), 0).$$

255 Note that  $y_b \mathbf{w}_b^\top \mathbf{x}_b = \gamma (\frac{|T^*|}{n+1} - 2) < -1 \leq \beta$ , we have

$$\mathbf{w}_b \xrightarrow{(\mathbf{x}_b, y_b)} \mathbf{w}_a = \mathbf{w}_b + \mathbf{x}_a y_a = (\sqrt{\gamma} (\frac{|T^*|}{n+1} - 1), -\sqrt{\gamma})$$

256 Note also that  $y_a \mathbf{w}_a^\top \mathbf{x}_a = \gamma (\frac{|T^*|}{n+1} - 2) < -1 \leq \beta$ , we have

$$\mathbf{w}_a \xrightarrow{(\mathbf{x}_a, y_a)} \mathbf{w}^* = \mathbf{w}_a + \mathbf{x}_a y_a = (\frac{|T^*| \sqrt{\gamma}}{n+1}, 0)$$

257 Therefore,  $y_{\text{test}} (\mathbf{w}^*)^\top \mathbf{x}_{\text{test}} = \frac{|T^*| \sqrt{\gamma}}{n+1} > 0$ .

258 Only if: For each  $T' \subseteq T$ , let  $T^* = T' \setminus \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$ . If  $y_{\text{test}} (\mathbf{w}^*)^\top \mathbf{x}_{\text{test}} > 0$  for  
 259  $\mathbf{w}^*$  satisfying  $\mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}^*$ , we prove that  $\exists S' \subseteq S$  such that  $\sum_{a \in S'} a = t$ . We first show that for  
 260 each  $T' \subseteq T$ , if  $\mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}$  satisfying  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} > 0$ , we have  $\forall k \in \{a, b, c\}$ ,  $(\mathbf{x}_k, y_k) \in$   
 261  $T'$ ,  $y_k \mathbf{w}_k^\top \mathbf{x}_k < \gamma$ , where  $\mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}_c \xrightarrow{(\mathbf{x}_c, y_c)} \mathbf{w}_b \xrightarrow{(\mathbf{x}_b, y_b)} \mathbf{w}_a$ . Otherwise, suppose  $\exists k \in \{a, b, c\}$   
 262 such that  $(\mathbf{x}_k, y_k) \notin T'$  or  $y_k \mathbf{w}_k^\top \mathbf{x}_k \geq \gamma$ , we have

$$y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} \leq \sqrt{\gamma} (\frac{|T^*|}{n+1} - 1) < 0$$



263 which contradicts to the fact that  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} \geq 0$ .

264 Let  $S' = \{a_i | (\mathbf{x}_i, y_i) \in T^*\}$  and  $t' = \sum_{a_i \in S'} a_i$ , it suffices to prove  $t' = t$ . Notice that

$$\begin{aligned} \mathbf{w}^{(0)} \xrightarrow{T^* \cap \{(\mathbf{x}_i, y_i) | 1 \leq j \leq i\}} \mathbf{w}_c &= (\sqrt{\gamma}(-18n^2m^2 + \frac{|T^*|}{n+1}), 3\sqrt{\gamma} \sum_{a_i \in S'} a_i) \\ &= (\sqrt{\gamma}(-18n^2m^2 + \frac{|T^*|}{n+1}), 3\sqrt{\gamma}t') \end{aligned}$$

265 Hence  $y_c \mathbf{w}_c^\top \mathbf{x}_c = \gamma(-18n^2m^2 + \frac{|T^*|}{n+1})(18n^2m^2 - 2) - 9\gamma t t' < -1 \leq \beta$ , thus

$$\mathbf{w}_c \xrightarrow{(\mathbf{x}_c, y_c)} \mathbf{w}_b = \mathbf{w}_c + \mathbf{x}_c y_c = (\sqrt{\gamma}(\frac{|T^*|}{n+1} - 2), 3\sqrt{\gamma}(t' - t))$$

266 (1) If  $t' \leq t - 1$ , we have  $y_b \mathbf{w}_b^\top \mathbf{x}_b = \gamma(\frac{|T^*|}{n+1} - 2 + 3(t - t')) > \gamma \geq \beta$ , a contradiction.

267 (2) If  $t' \geq t + 1$ , we have  $y_a \mathbf{w}_a^\top \mathbf{x}_a = \gamma(\frac{|T^*|}{n+1} - 2 + 3(t' - t)) > \gamma \geq \beta$ , another contradiction.

268 Therefore  $t' = t$ , and this completes the proof.  $\square$

269 Moreover, DEBUGGABLE-LIN is NP-hard even when  $d = 1$  and  $\beta < 0$ .

270 **Theorem 4.4.** If the training order is adversarially chosen and  $d = 1$ , DEBUGGABLE-LIN remains  
271 NP-hard for *each* hinge-like loss function with  $\beta < 0$  at *every* constant learning rate.

272 **Remarks.** The training order in this section can be arbitrary as long as the last three training  
273 samples are  $(\mathbf{x}_c, y_c)$ ,  $(\mathbf{x}_b, y_b)$ ,  $(\mathbf{x}_a, y_a)$ , respectively. All the training samples are “good” since for  
274 each  $(\mathbf{x}, y) \in T$  we have  $\mathbf{x}^\top \mathbf{x}_{\text{test}} y y_{\text{test}} > 0$ . This implies that DEBUGGABLE-LIN is NP-hard even if  
275 all the training data are “good” training samples, and exemplifies why the GTA algorithm fails for  
276 higher dimensions.

## 277 5 Discussion and Conclusion

278 In this paper, we provided a comprehensive analysis on the complexity of DEBUGGABLE. We focus  
279 on the linear classifier that is trained using SGD, as it is a key component in the majority of popular  
280 models.

281 Since DEBUGGABLE is a special case of data debugging, the above results proved the intractability  
282 of data debugging and therefore gives a negative answer to Problem 1.1 declared in the introduction.  
283 The complexity results also demonstrated that it is not accurate to estimate the impact of subset of  
284 training data by summing up the score of each training samples in the subset, *as long as the scores*  
285 *can be calculated in polynomial time.*

286 In Section 4, a training sample is said to be “good” if it can help the resulting model to predict  
287 correctly on the test instance. That is, it can increase  $y_{\text{test}}(\mathbf{w}^*)^\top \mathbf{x}_{\text{test}}$ . However, in our proof we  
288 showed that DEBUGGABLE remains NP-hard even if all training samples are “good”. This suggests  
289 that the quality of a training sample does not depend only on some properties of itself but also on  
290 the interaction between the rest of the training data, which should be taken into consideration when  
291 developing data cleaning approaches.

292 Moreover, the NP-hardness of DEBUGGABLE implies that, it is in general intractable to figure out the  
293 causality between even the prediction of a linear classifier and its training data. This may be seem  
294 surprising since linear classifiers have long been considered “inherently interpretable”. As warned  
295 in [25], *a method being “inherently interpretable” needs to be verified before it can be trusted*, the  
296 concept of interpretability must be *rigorously defined*, or at least its boundaries specified.

297 Our results suggests the following directions for future research. Firstly, characterizing the training  
298 sample may be helpful in designing efficient algorithms for data debugging; Secondly, designing  
299 algorithms using CSP-solver is a potential way to solve data debugging more efficiently than the brute-  
300 force algorithms; Finally, developing random algorithms is a potential way to solve data debugging  
301 successfully with high probability.

## References

- 302  
303 [1] Satoshi Hara, Atsushi Nitanda, and Takanori Maehara. *Data Cleansing for Models Trained with SGD*.  
304 Curran Associates Inc., Red Hook, NY, USA, 2019.
- 305 [2] Weiyuan Wu, Lampros Flokas, Eugene Wu, and Jiannan Wang. Complaint-driven training data debugging  
306 for query 2.0. pages 1317–1334, 06 2020. doi: 10.1145/3318464.3389696.
- 307 [3] Bojan Karlaš, David Dao, Matteo Interlandi, Bo Li, Sebastian Schelter, Wentao Wu, and Ce Zhang. Data  
308 debugging with shapley importance over end-to-end machine learning pipelines, 2022.
- 309 [4] Felix Neutatz, Binger Chen, Ziawasch Abedjan, and Eugene Wu. From cleaning before ml to cleaning for  
310 ml. *IEEE Data Eng. Bull.*, 44:24–41, 2021. URL [https://api.semanticscholar.org/CorpusID:  
311 237542697](https://api.semanticscholar.org/CorpusID:237542697).
- 312 [5] Peng Li, Xi Rao, Jennifer Blase, Yue Zhang, Xu Chu, and Ce Zhang. Cleanml: A study for evaluating the  
313 impact of data cleaning on ml classification tasks. In *2021 IEEE 37th International Conference on Data  
314 Engineering (ICDE)*, pages 13–24, 2021. doi: 10.1109/ICDE51399.2021.00009.
- 315 [6] Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger Grosse. If influence functions are  
316 the answer, then what is the question? In *Proceedings of the 36th International Conference on Neural  
317 Information Processing Systems, NIPS '22*, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN  
318 9781713871088.
- 319 [7] Romila Pradhan, Jiongli Zhu, Boris Glavic, and Babak Salimi. Interpretable data-based explanations for  
320 fairness debugging. In *Proceedings of the 2022 International Conference on Management of Data,  
321 SIGMOD '22*, page 247–261, New York, NY, USA, 2022. Association for Computing Machinery.  
322 ISBN 9781450392495. doi: 10.1145/3514221.3517886. URL [https://doi.org/10.1145/3514221.  
323 3517886](https://doi.org/10.1145/3514221.3517886).
- 324 [8] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In  
325 *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, page  
326 1885–1894. JMLR.org, 2017.
- 327 [9] Rajiv Khanna, Been Kim, Joydeep Ghosh, and Oluwasanmi Koyejo. Interpreting black box predictions  
328 using fisher kernels. In *International Conference on Artificial Intelligence and Statistics*, 2018. URL  
329 <https://api.semanticscholar.org/CorpusID:53085397>.
- 330 [10] Pang Wei Koh, Kai-Siang Ang, Hubert Hua Kian Teo, and Percy Liang. On the accuracy of influence  
331 functions for measuring group effects. In *Neural Information Processing Systems*, 2019. URL [https://  
332 api.semanticscholar.org/CorpusID:173188850](https://api.semanticscholar.org/CorpusID:173188850).
- 333 [11] Samyadeep Basu, Xuchen You, and Soheil Feizi. On second-order group influence functions for black-  
334 box predictions. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20*.  
335 JMLR.org, 2020.
- 336 [12] Han Guo, Nazneen Rajani, Peter Hase, Mohit Bansal, and Caiming Xiong. FastIF: Scalable influence  
337 functions for efficient model interpretation and debugging. In Marie-Francine Moens, Xuanjing Huang,  
338 Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods  
339 in Natural Language Processing*, pages 10333–10350, Online and Punta Cana, Dominican Republic,  
340 November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.808.  
341 URL <https://aclanthology.org/2021.emnlp-main.808>.
- 342 [13] Amirata Ghorbani and James Y. Zou. Data shapley: Equitable valuation of data for machine learning.  
343 *ArXiv*, abs/1904.02868, 2019. URL <https://api.semanticscholar.org/CorpusID:102350503>.
- 344 [14] R. Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nicholas Hynes, Nezihe Merve Gürel, Bo Li,  
345 Ce Zhang, Dawn Xiaodong Song, and Costas J. Spanos. Towards efficient data valuation based on the  
346 shapley value. *ArXiv*, abs/1902.10275, 2019. URL [https://api.semanticscholar.org/CorpusID:  
347 67855573](https://api.semanticscholar.org/CorpusID:67855573).
- 348 [15] Ruoxi Jia, Fan Wu, Xuehui Sun, Jiachen Xu, David Dao, Bhavya Kailkhura, Ce Zhang, Bo Li, and Dawn  
349 Song. Scalability vs. utility: Do we have to sacrifice one for the other in data importance quantification?  
350 In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8235–8243,  
351 2021. doi: 10.1109/CVPR46437.2021.00814.
- 352 [16] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nezihe Merve Gürel, Bo Li, Ce Zhang, Costas  
353 Spanos, and Dawn Song. Efficient task-specific data valuation for nearest neighbor algorithms. *Proc.  
354 VLDB Endow.*, 12(11):1610–1623, jul 2019. ISSN 2150-8097. doi: 10.14778/3342263.3342637. URL  
355 <https://doi.org/10.14778/3342263.3342637>.

- 356 [17] Jeremy Mange. Effect of training data order for machine learning. In *2019 International Conference*  
357 *on Computational Science and Computational Intelligence (CSCI)*, pages 406–407, 2019. doi: 10.1109/  
358 CSCI49370.2019.00078.
- 359 [18] Ernie Chang, Hui-Syuan Yeh, and Vera Demberg. Does the order of training samples matter? improving  
360 neural data-to-text generation with curriculum learning. *ArXiv*, abs/2102.03554, 2021. URL <https://api.semanticscholar.org/CorpusID:231846815>.  
361
- 362 [19] YeJia Liu, Weiyuan Wu, Lampros Flokas, Jiannan Wang, and Eugene Wu. Enabling sql-based training  
363 data debugging for federated learning. *Proceedings of the VLDB Endowment*, 15:388–400, 02 2022. doi:  
364 10.14778/3494124.3494125.
- 365 [20] Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard Zemel. Understanding  
366 the origins of bias in word embeddings, 2019.
- 367 [21] Hao Wang, Berk Ustun, and Flavio P. Calmon. Repairing without retraining: Avoiding disparate impact  
368 with counterfactual distributions, 2019.
- 369 [22] Xiaotie Deng and Christos H. Papadimitriou. On the complexity of cooperative solution concepts. *Math.*  
370 *Oper. Res.*, 19:257–266, 1994. URL <https://api.semanticscholar.org/CorpusID:12946448>.
- 371 [23] Qi Wang, Yue Ma, Kun Zhao, and Yingjie Tian. A comprehensive survey of loss functions in machine  
372 learning. *Annals of Data Science*, 9, 04 2022. doi: 10.1007/s40745-020-00253-5.
- 373 [24] Erik D. Demaine, William Gasarch, and Mohammad Hajiaghayi. *Computational Intractability: A Guide to*  
374 *Algorithmic Lower Bounds*. MIT Press, 2024.
- 375 [25] Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable nlp systems: How should we define and  
376 evaluate faithfulness? In *Annual Meeting of the Association for Computational Linguistics*, 2020. URL  
377 <https://api.semanticscholar.org/CorpusID:215416110>.
- 378 [26] Victor Parque. Tackling the subset sum problem with fixed size using an integer representation scheme.  
379 In *2021 IEEE Congress on Evolutionary Computation (CEC)*, pages 1447–1453, 2021. doi: 10.1109/  
380 CEC45853.2021.9504889.

### 381 A Detailed Proofs for Section 3

382 **Notations.** Given some orderings  $\{o_e\}$  of training data, where  $o_e^t$  as the order of  $\mathbf{t}$  in epoch  $e$ . We  
 383 use  $w_{x_i}^{(e,l)}$  to denote the value of  $w_{x_i}$  after the  $l$ -th iteration in epoch  $e$ . We also denote  $\mathbf{x}_t$  and  $y_t$  as  
 384 the feature and the label of training data  $\mathbf{t}$ , respectively. We denote  $\mathbf{t}^{(e,l)}$  as the training sample being  
 385 considered during epoch  $e$ , iteration  $l$ .

386 **Lemma A.1.** Suppose  $T \subseteq T_0$  is the training data and let  $T_{l,r}^e = \{\mathbf{t}^{(e,l)}, \mathbf{t}^{(e,l+1)}, \dots, \mathbf{t}^{(e,r)}\}$   
 387 be the set of consecutive training samples considered during epoch  $e$  from iteration  $l$  to  $r$ . For  
 388  $1 \leq l \leq r \leq |T|$ , if  $\text{clause}(\gamma, i_1, i_2, i_3) \notin T_{l,r}^e$ , then  $w_{c_\gamma}^{(e,l-1)} = w_{c_\gamma}^{(e,r)}$ .

389 *Proof.* For each  $\mathbf{t} \in T_{l,r}^e$ , we have  $(x_t)_{c_\gamma} = 0$ . Therefore

$$\left| \frac{\partial \mathcal{L}}{\partial c_\gamma} \Big|_{\mathbf{t}} \right| \leq \max \left\{ \left| -\frac{12N}{5} y x_{c_\gamma} \right|, | - y x_{c_\gamma} |, \left| -\frac{1}{1000N} y x_{c_\gamma} \right|, 0 \right\} = 0$$

390 Hence  $\frac{\partial \mathcal{L}}{\partial c_\gamma} \Big|_{\mathbf{t}} = 0$ , and

$$w_{c_\gamma}^{(e,r)} = w_{c_\gamma}^{(e,l-1)} - \eta_{c_\gamma} \sum_{\mathbf{t} \in T_{l,r}^e} \frac{\partial \mathcal{L}}{\partial c_\gamma} \Big|_{\mathbf{t}} = w_{c_\gamma}^{(e,l-1)}$$

391 Similarly,  $(x_t)_{b_\gamma} = 0$ , and

$$\left| \frac{\partial \mathcal{L}}{\partial b_\gamma} \Big|_{\mathbf{t}} \right| \leq \max \left\{ \left| -\frac{12N}{5} y x_{b_\gamma} \right|, | - y x_{b_\gamma} |, \left| -\frac{1}{1000N} y x_{b_\gamma} \right|, 0 \right\} = 0$$

392 Hence  $\frac{\partial \mathcal{L}}{\partial b_\gamma} \Big|_{\mathbf{t}} = 0$ , and

$$w_{b_\gamma}^{(e,r)} = w_{b_\gamma}^{(e,l-1)} - \eta_{b_\gamma} \sum_{\mathbf{t} \in T_{l,r}^e} \frac{\partial \mathcal{L}}{\partial b_\gamma} \Big|_{\mathbf{t}} = w_{b_\gamma}^{(e,l-1)}$$

393 □

394 **Lemma A.2.** Suppose  $T \subseteq T_0$  is the training data and  $T_l := \{\mathbf{t}^{(1,1)}, \dots, \mathbf{t}^{(1,l)}\}$ .  $\forall 1 \leq i \leq n, 1 \leq$   
 395  $l \leq |T|$ ,  $w_{x_i}^{(1,l)} \in U(1, \frac{l+1}{6000N^2})$  if  $\text{var}(i) \in T_l$ ; Otherwise  $w_{x_i}^{(1,l)} \in U(-1, \frac{l+1}{6000N^2})$ .

396 *Proof.* We prove this lemma by induction.

397 **Basic Case:** Note that for all  $1 \leq i \leq n$ ,  $w_{x_i}^{(0)} = -1$ , and for all  $1 \leq \gamma \leq m$ ,  $w_{c_\gamma}^{(0)} = 1/2$ ,  $w_{b_\gamma}^{(0)} = -1$ .

398 We denote  $\mathbf{t} = \mathbf{t}^{(1,1)}$  to avoid cluttering. For any fixed  $i$ :

399 (1) If  $\mathbf{t} = \text{var}(i)$ . We have  $y_t(\mathbf{w}^{(0)})^\top \mathbf{x}'_t = 5w_{x_i}^{(0)} = -5$ , hence

$$\frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} = -\frac{12N}{5} y_t (x_t)_i = -12N$$

400 and

$$w_{x_i}^{(1,1)} = w_{x_i}^{(0)} - \eta_{x_i} \frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} = -1 - \frac{1}{6N} \left( -\frac{12N}{5} \right) = 1 \in U(1, \frac{2}{6000N^2})$$

401 (2) If  $\mathbf{t} = \text{clause}(\gamma, i, i', i'')$ . We have

$$y_t(\mathbf{w}^{(0)})^\top \mathbf{x}'_t = w_{x_i}^{(0)} + w_{x_{i'}}^{(0)} + w_{x_{i''}}^{(0)} + w_{c_\gamma}^{(0)} + \frac{1}{2} w_{b_\gamma}^{(0)} = -3$$

402 hence

$$\frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} = -\frac{1}{1000N} y_t (x_t)_{x_i} = -\frac{1}{1000N}$$

403 and

$$\begin{aligned} w_{x_i}^{(1,1)} &= w_{x_i}^{(0)} - \eta_{x_i} \left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = -1 - \frac{1}{6N} \left( -\frac{1}{1000N} \right) \\ &= -1 + \frac{1}{6000N^2} \in U\left(-1, \frac{2}{6000N^2}\right) \end{aligned}$$

404 (3) Otherwise,  $w_{x_i}$  will not be updated. Therefore  $w_{x_i}^{(1,1)} = w_{x_i}^{(0)} = -1 \in U\left(-1, \frac{2}{6000N^2}\right)$ .

405 Hence this lemma is true for  $l = 1$ .

406 **Induction Step:** Suppose the lemma is true for  $l < |T|$ . We prove that this lemma remains true for  
407  $l + 1$ . We denote  $\mathbf{t} = \mathbf{t}^{(1,l+1)}$  to avoid cluttering. This makes sense since  $l + 1 \leq |T|$  and thus  $\mathbf{t} \in T$ .  
408 For any fixed  $i$ :

409 (1) If  $\mathbf{t} = \text{var}(i)$ , then  $\text{var}(i) \notin T_l$  because there are at most one  $\text{var}(i)$  in  $T$  for each  $i$ .

410 Therefore  $w_{x_i}^{(1,l)} \in U\left(-1, \frac{l+1}{6000N^2}\right)$ . We have  $y_{\mathbf{t}}(\mathbf{w}^{(1,l)})^\top \mathbf{x}'_{\mathbf{t}} = 5w_{x_i}^{(1,l)} \in U(-5, 0.01)$ , and  
411  $\left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = -\frac{12N}{5} y_{\mathbf{t}}(x_{\mathbf{t}})_i = -12N$ . Hence

$$\begin{aligned} w_{x_i}^{(1,l+1)} &= w_{x_i}^{(1,l)} - \eta_{x_i} \left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = w_{x_i}^{(1,l)} - \frac{1}{6N} \left( -\frac{12N}{5} \right) \\ &= w_{x_i}^{(1,l)} + 2 \in U\left(1, \frac{l+2}{6000N^2}\right) \end{aligned}$$

412 (2) If  $\mathbf{t} = \text{clause}(\gamma, i, i', i'')$ . In this case,  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin T_{1,l}^1$  and by Lemma A.1 we have

413  $w_{c_\gamma}^{(1,l)} = w_{c_\gamma}^{(0)}$ ,  $w_{b_\gamma}^{(1,l)} = w_{b_\gamma}^{(0)}$ . From the induction hypothesis we have

$$w_{x_i}^{(1,l)}, w_{x_{i'}}^{(1,l)}, w_{x_{i''}}^{(1,l)} \in U\left(\pm 1, \frac{l+1}{6000N^2}\right)$$

414 and thus

$$\begin{aligned} y_{\mathbf{t}}(\mathbf{w}^{(1,l)})^\top \mathbf{x}'_{\mathbf{t}} &= w_{x_i}^{(1,l)} + w_{x_{i'}}^{(1,l)} + w_{x_{i''}}^{(1,l)} + w_{c_\gamma}^{(1,l)} + \frac{1}{2} w_{b_\gamma}^{(1,l)} \\ &= w_{x_i}^{(1,l)} + w_{x_{i'}}^{(1,l)} + w_{x_{i''}}^{(1,l)} \\ &\in \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U\left(x_0, \frac{3(l+1)}{6000N^2}\right) \subseteq \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U(x_0, 0.01) \end{aligned}$$

415 We have  $\left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = -\frac{1}{1000N}$  and  $w_{x_i}^{(1,l+1)} = w_{x_i}^{(1,l)} - \eta_{x_i} \left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = w_{x_i}^{(1,l)} + \frac{1}{6000N^2}$ . Consider the  
416 following cases:

417 • If  $\text{var}(i) \in T_l$ , then  $\text{var}(i) \in T_{l+1}$  and  $w_{x_i}^{(1,l)} \in U\left(1, \frac{l+1}{6000N^2}\right)$ . Therefore  $w_{x_i}^{(1,l+1)} \in$   
418  $U\left(1, \frac{l+2}{6000N^2}\right)$ .

419 • If  $\text{var}(i) \notin T_l$ , then  $\text{var}(i) \notin T_{l+1}$  and  $w_{x_i}^{(1,l)} \in U\left(-1, \frac{l+1}{6000N^2}\right)$ . Therefore  $w_{x_i}^{(1,l+1)} \in$   
420  $U\left(-1, \frac{l+2}{6000N^2}\right)$ .

421 (3) Otherwise,  $w_{x_i}$  will not be updated, and  $w_{x_i}^{(1,l+1)} = w_{x_i}^{(1,l)}$ . If  $\text{var}(i) \in T_l$  then  $\text{var}(i) \in T_{l+1}$  and  
422  $w_{x_i}^{(1,l+1)} \in U\left(1, \frac{l+2}{6000N^2}\right)$ ; Otherwise  $\text{var}(i) \notin T_{l+1}$  and  $w_{x_i}^{(1,l+1)} \in U\left(-1, \frac{l+2}{6000N^2}\right)$ .

423 Hence if the lemma is true for  $l < |T|$ , it is also true for  $l + 1$ . Therefore, the lemma is true for all  
424  $1 \leq l \leq |T|$ .  $\square$

425 **Corollary A.1.** Suppose  $T \subseteq T_0$  is the training data.  $\forall 1 \leq i \leq n, 1 \leq l \leq |T|$ , if  $\text{var}(i) \in T$ , then  
426  $w_{x_i}^{(1)} \in U\left(1, \frac{1}{6000N}\right)$ . Otherwise  $w_{x_i}^{(1)} \in U\left(-1, \frac{1}{6000N}\right)$ .

427 *Proof.* Note that  $w_{x_i}^{(1)} = w_{x_i}^{(1,|T|)}$  and  $N = 2m + n + 1$ . By Lemma A.2, if  $\text{var}(i) \in T$  we have

$$w_{x_i}^{(1,|T|)} \in U\left(1, \frac{|T|+1}{6000N^2}\right) \subseteq U\left(1, \frac{m+n+1}{6000N^2}\right) \subseteq U\left(1, \frac{1}{6000N}\right)$$

428 If  $\text{var}(i) \notin T$ , we have

$$w_{x_i}^{(1,|T|)} \in U\left(-1, \frac{|T|+1}{6000N^2}\right) \subseteq U\left(-1, \frac{m+n+1}{6000N^2}\right) \subseteq U\left(-1, \frac{1}{6000N}\right)$$

429

□

430 **Lemma A.3.** Suppose  $T \subseteq T_0$  is the training data.  $\forall 1 \leq \gamma \leq m$ , if  $\exists 1 \leq i_1, i_2, i_3 \leq n$  such that  
431  $\text{clause}(\gamma, i_1, i_2, i_3) \in T$ , then  $w_{b_\gamma}^{(1)} = 0, w_{c_\gamma}^{(1)} = \frac{1}{2} + \frac{1}{200N}$ ; Otherwise,  $w_{b_\gamma}^{(1)} = -1, w_{c_\gamma}^{(1)} = \frac{1}{2}$ .

432 *Proof.* (1) If such  $\mathbf{t}_\gamma = \text{clause}(\gamma, i_1, i_2, i_3)$  exists in  $T$ , by Lemma A.2 we have

$$w_{x_{i_1}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{x_{i_2}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{x_{i_3}}^{(1, o_{\mathbf{t}_\gamma}^1)} \in \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U\left(x_0, \frac{3(o_{\mathbf{t}_\gamma}^1 + 1)}{6000N^2}\right) \subseteq \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U(x_0, 0.01)$$

433 By Lemma A.1 we have  $w_{c_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} = w_{c_\gamma}^{(0)}$  and  $w_{b_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} = w_{b_\gamma}^{(0)}$  because  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin$

434  $T_{1, o_{\mathbf{t}_\gamma}^1 - 1}^1$ . Hence

$$\begin{aligned} y_{\mathbf{t}_\gamma} (\mathbf{w}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)})^\top \mathbf{x}'_{\mathbf{t}_\gamma} &= w_{x_{i_1}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{x_{i_2}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{x_{i_3}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{c_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} + \frac{1}{2} w_{b_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} \\ &= w_{x_{i_1}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{x_{i_2}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{x_{i_3}}^{(1, o_{\mathbf{t}_\gamma}^1)} + w_{c_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} \\ &\in \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U(x_0, 0.01) \end{aligned}$$

435 We have  $\frac{\partial \mathcal{L}}{\partial w_{c_\gamma}} \Big|_{\mathbf{t}_\gamma} = -\frac{1}{1000N}$ , and

$$w_{c_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1)} = w_{c_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} - \eta_{c_\gamma} \frac{\partial \mathcal{L}}{\partial w_{c_\gamma}} \Big|_{\mathbf{t}_\gamma} = \frac{1}{2} + 5 \times \frac{1}{1000N} = \frac{1}{2} + \frac{1}{200N}$$

436 Similarly,  $\frac{\partial \mathcal{L}}{\partial w_{b_\gamma}} \Big|_{\mathbf{t}_\gamma} = -\frac{1}{2000N}$  and

$$w_{b_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1)} = w_{b_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1 - 1)} - \eta_{b_\gamma} \frac{\partial \mathcal{L}}{\partial w_{b_\gamma}} \Big|_{\mathbf{t}_\gamma} = -1 - 2000N \times \left(-\frac{1}{2000N}\right) = 0$$

437 Note also that  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin T_{o_{\mathbf{t}_\gamma}^1, |T|}^1$ , by Lemma A.1 we have

438  $w_{c_\gamma}^{(1)} = w_{c_\gamma}^{(1, |T|)} = w_{c_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1)} = \frac{1}{2} + \frac{1}{200N}$  and  $w_{b_\gamma}^{(1)} = w_{b_\gamma}^{(1, |T|)} = w_{b_\gamma}^{(1, o_{\mathbf{t}_\gamma}^1)} = 0$ .

439 (2) If such  $\mathbf{t}_\gamma = \text{clause}(\gamma, i_1, i_2, i_3)$  does not exist in  $T$ , by Lemma A.1 we have  $w_{c_\gamma}^{(1)} = w_{c_\gamma}^{(0)} = \frac{1}{2}$   
440 and  $w_{b_\gamma}^{(1)} = w_{b_\gamma}^{(0)} = -1$ . □

441 **Lemma A.4.** Suppose  $T \subseteq T_0$  and  $C_l$  be the number of  $\text{clause}()$  in  $T_{1, l}^2$ .  $\forall 1 \leq i \leq n, 1 \leq l \leq |T|$ ,  
442  $w_{x_i}^{(2, l)} \in U\left(1, \frac{C_l + 1/2}{6N}\right)$  if  $\text{var}(i) \in T$ ; Otherwise  $w_{x_i}^{(2, l)} \in U\left(-1, \frac{C_l + 1/2}{6N}\right)$ .

443 *Proof.* Similar to the proof of A.2, we prove this lemma by induction.

444 **Basic Case:** Note that for all  $1 \leq i \leq n$ ,  $w_{x_i}^{(1)} = U\left(\pm 1, \frac{1}{6000N}\right)$ , and for all  $1 \leq \gamma \leq m$ ,  $w_{c_\gamma}^{(1)} \in$   
445  $\left\{\frac{1}{2}, \frac{1}{2} + \frac{1}{200N}\right\}$ ,  $w_{b_\gamma}^{(1)} \in \{-1, 0\}$ . We denote  $\mathbf{t} = \mathbf{t}^{(2, 1)}$  to avoid cluttering. For any fixed  $i$ :

446 (1) If  $\mathbf{t} = \text{var}(i)$ ,  $C_1 = 0$ . By Corollary A.1,  $w_{x_i}^{(1)} = U\left(1, \frac{1}{6000N}\right)$ . We have

$$y_{\mathbf{t}} (\mathbf{w}^{(1)})^\top \mathbf{x}'_{\mathbf{t}} = 5w_{x_i}^{(1)} \in U\left(5, \frac{1}{1200N}\right)$$

447 hence  $\left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = 0$ , and

$$w_{x_i}^{(2,1)} = w_{x_i}^{(1)} \in U\left(1, \frac{1}{6N}\right) = U\left(1, \frac{C_l + 1/2}{6N}\right)$$

448 (2) If  $\mathbf{t} = \text{clause}(\gamma, i, i', i'')$ ,  $C_1 = 1$ . By Lemma A.3, we have  $w_{c_\gamma}^{(1)} = \frac{1}{2} + \frac{1}{200N}$  and  $w_{b_\gamma}^{(1)} = 0$ .

449 Therefore,

$$\begin{aligned} y_{\mathbf{t}}(\mathbf{w}^{(1)})^\top \mathbf{x}'_{\mathbf{t}} &= w_{x_i}^{(1)} + w_{x_{i'}}^{(1)} + w_{x_{i''}}^{(1)} + w_{c_\gamma}^{(1)} + \frac{1}{2}w_{b_\gamma}^{(1)} \\ &= w_{x_i}^{(1)} + w_{x_{i'}}^{(1)} + w_{x_{i''}}^{(1)} + \frac{1}{2} - \frac{1}{200N} \\ &\in \bigcup_{x_0 \in \{\frac{1}{2} \pm 1, \frac{1}{2} \pm 3\}} U(x_0, 0.01) \end{aligned}$$

450 hence  $\left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} \in \{0, -y_{x_i}\} = \{-1, 0\}$ , and  $\eta_{x_i} \left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} \in \{-\frac{1}{6N}, 0\}$ .

451 By Corollary A.1, if  $\text{var}(i) \in T$ , we have

$$w_{x_i}^{(2,1)} = w_{x_i}^{(1)} - \eta_{x_i} \left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} \in U\left(1, \frac{3/2}{6N}\right) = U\left(1, \frac{C_l + 1/2}{6N}\right)$$

452 If  $\text{var}(i) \notin T$ , we have

$$w_{x_i}^{(2,1)} = w_{x_i}^{(1)} - \eta_{x_i} \left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} \in U\left(-1, \frac{3/2}{6N}\right) = U\left(-1, \frac{C_l + 1/2}{6N}\right)$$

453 (3) Otherwise,  $w_{x_i}$  will not be updated and  $C_1 \leq 1$ . Therefore if  $\text{var}(i) \in T$ ,

$$w_{x_i}^{(2,1)} = w_{x_i}^{(1)} \in U\left(1, \frac{3/2}{6N}\right) \subseteq U\left(1, \frac{C_l + 1/2}{6N}\right)$$

454 If  $\text{var}(i) \notin T$ ,

$$w_{x_i}^{(2,1)} = w_{x_i}^{(1)} \in U\left(-1, \frac{3/2}{6N}\right) \subseteq U\left(-1, \frac{C_l + 1/2}{6N}\right)$$

455 Hence this lemma is true for  $l = 1$ .

456 **Induction Step:** Suppose the lemma is true for  $l < |T|$ . We prove that this lemma remains true for  
457  $l + 1$ . We denote  $\mathbf{t} = \mathbf{t}^{(2,l+1)}$  to avoid cluttering. This makes sense since  $l + 1 \leq |T|$  and thus  $\mathbf{t} \in T$ .  
458 For any fixed  $i$ :

459 (1) If  $\mathbf{t} = \text{var}(i)$ ,  $C_{l+1} = C_l$ . By Corollary A.1,  $w_{x_i}^{(2,l)} \in U\left(1, \frac{C_l + 1/2}{6N}\right)$ .

460 We have  $y_{\mathbf{t}}(\mathbf{w}^{(2,l)})^\top \mathbf{x}'_{\mathbf{t}} = 5w_{x_i}^{(2,l)} \in U(5, 1/6)$  and  $\left. \frac{\partial \mathcal{L}}{\partial w_{x_i}} \right|_{\mathbf{t}} = 0$ . Hence  $w_{x_i}^{(2,l+1)} = w_{x_i}^{(2,l)} \in$   
461  $U\left(1, \frac{C_l + 1/2}{6N}\right)$ .

462 (2) If  $\mathbf{t} = \text{clause}(\gamma, i, i', i'')$ ,  $C_{l+1} = C_l + 1$ . In this case,  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin T_{1,l}^2$  and by Lemma  
463 A.1 and Lemma A.3 we have  $w_{c_\gamma}^{(2,l)} = w_{c_\gamma}^{(1)} = \frac{1}{2} + \frac{1}{200N}$ ,  $w_{b_\gamma}^{(2,l)} = w_{b_\gamma}^{(1)} = 0$ . From the induction  
464 hypothesis we have  $w_{x_i}^{(2,l)}, w_{x_{i'}}^{(2,l)}, w_{x_{i''}}^{(2,l)} \in U\left(\pm 1, \frac{C_l + 1/2}{6N}\right)$ . Noting that

$$\frac{C_l + 1/2}{6N} \leq \frac{m + 1/2}{6N} = \frac{m + 1/2}{(n + 2(m + 1/2))} \leq \frac{1}{12}$$

465 we have

$$\begin{aligned} y_{\mathbf{t}}(\mathbf{w}^{(2,l)})^\top \mathbf{x}'_{\mathbf{t}} &= w_{x_i}^{(2,l)} + w_{x_{i'}}^{(2,l)} + w_{x_{i''}}^{(2,l)} + w_{c_\gamma}^{(2,l)} + \frac{1}{2}w_{b_\gamma}^{(2,l)} \\ &= w_{x_i}^{(2,l)} + w_{x_{i'}}^{(2,l)} + w_{x_{i''}}^{(2,l)} + \frac{1}{2} + \frac{1}{200N} \\ &\in \bigcup_{x_0 \in \{\frac{1}{2} \pm 1, \frac{1}{2} \pm 3\}} U\left(x_0, \frac{3(C_l + 1/2)}{6N} + \frac{1}{200N}\right) \\ &\subseteq \bigcup_{x_0 \in \{\frac{1}{2} \pm 1, \frac{1}{2} \pm 3\}} U(x_0, 0.26) \end{aligned}$$

466 And thus  $\frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} \in \{0, -yx_{x_i}\} = \{-1, 0\}$ , and  $\eta_{x_i} \frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} \in \{-\frac{1}{6N}, 0\}$ .

467 By Corollary A.1, if  $\text{var}(i) \in T$ ,  $w_{x_i}^{(2,l+1)} = w_{x_i}^{(l)} - \eta_{x_i} \frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} \in U(1, \frac{C_l+3/2}{6N}) = U(1, \frac{C_{l+1}+1/2}{6N})$ ;

468 if  $\text{var}(i) \notin T$ ,  $w_{x_i}^{(2,l+1)} = w_{x_i}^{(l)} - \eta_{x_i} \frac{\partial \mathcal{L}}{\partial w_{x_i}} \Big|_{\mathbf{t}} \in U(-1, \frac{C_l+3/2}{6N}) = U(-1, \frac{C_{l+1}+1/2}{6N})$ .

469 (3) Otherwise,  $w_{x_i}$  will not be updated. We have  $C_{l+1} \leq C_l + 1$   $w_{x_i}^{(2,l+1)} = w_{x_i}^{(2,l)}$ . If  $\text{var}(i) \in T$   
470 then  $w_{x_i}^{(2,l+1)} \in U(1, \frac{C_{l+1}+1/2}{6N})$ ; If  $\text{var}(i) \notin T$  then  $w_{x_i}^{(2,l+1)} \in U(-1, \frac{C_{l+1}+1/2}{6N})$ .

471 Hence if the lemma is true for  $l < |T|$ , it is also true for  $l + 1$ . Therefore, the lemma is true for all  
472  $1 \leq l \leq |T|$ .  $\square$

473 **Corollary A.2.** Suppose  $T \subseteq T_0$  is the training data.  $\forall 1 \leq i \leq n$ , if  $\text{var}(i) \in T$ , then  $w_{x_i}^{(2)} \in$   
474  $U(1, 0.1)$ . Otherwise  $w_{x_i}^{(2)} \in U(-1, 0.1)$ .

475 *Proof.* Note that  $w_{x_i}^{(2)} = w_{x_i}^{(2,|T|)}$  and  $C_{|T|} \leq m$ . By Lemma A.4, if  $\text{var}(i) \in T$  we have

$$w_{x_i}^{(2,|T|)} \in U(1, \frac{C_{|T|} + 1/2}{6N}) \subseteq U(1, \frac{m + 1/2}{6N}) \subseteq U(1, \frac{1}{12}) \subseteq U(1, 0.1)$$

476 If  $\text{var}(i) \notin T$ , we have

$$w_{x_i}^{(1,|T|)} \in U(-1, \frac{C_{|T|} + 1/2}{6N}) \subseteq U(-1, \frac{m + 1/2}{6N}) \subseteq U(-1, \frac{1}{12}) \subseteq U(-1, 0.1)$$

477  $\square$

478 **Lemma A.5.** Suppose  $T \subseteq T_0$  is the training data.  $\forall 1 \leq i \leq m$ , if  $\exists 1 \leq i_1, i_2, i_3 \leq n$  such that  
479  $\text{clause}(i, i_1, i_2, i_3) \in T$ , then

480 1.  $w_{b_j}^{(2)} = 1000N$ ;

481 2.  $w_{c_j}^{(2)} = \frac{11}{2} + \frac{1}{200N}$  if exactly one of  $\text{var}(i_1)$ ,  $\text{var}(i_2)$ ,  $\text{var}(i_3)$  is in  $T$ . Otherwise  $w_{c_j}^{(2)} =$   
482  $\frac{1}{2} + \frac{1}{200N}$ .

483 Otherwise,  $w_{b_i}^{(2)} = -1$ ,  $w_{c_i}^{(2)} = \frac{1}{2}$ .

484 *Proof.* (1) If such  $\mathbf{t}_\gamma = \text{clause}(\gamma, i_1, i_2, i_3)$  exists in  $T$ , by Lemma A.4 we have

$$w_{x_{i_1}}^{(2, o_{\mathbf{t}_\gamma}^1)}, w_{x_{i_2}}^{(2, o_{\mathbf{t}_\gamma}^1)}, w_{x_{i_3}}^{(2, o_{\mathbf{t}_\gamma}^1)} \in U(\pm 1, \frac{m + 1/2}{6N}) \subseteq U(\pm 1, \frac{1}{12N})$$

485 By Lemma A.1 we have  $w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} = w_{c_\gamma}^{(1)} = \frac{1}{2} + \frac{1}{200N}$  and  $w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} = w_{b_\gamma}^{(1)} = 0$  because  
486  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin T_{1, o_{\mathbf{t}_\gamma}^1 - 1}^1$ . Consider the following two cases:

487 (a) If exactly one of  $\text{var}(i_1)$ ,  $\text{var}(i_2)$ ,  $\text{var}(i_3)$  is in  $T$ , by Corollary A.2 we have

$$\begin{aligned} y_{\mathbf{t}_\gamma} (\mathbf{w}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)})^\top \mathbf{x}'_{\mathbf{t}_\gamma} &= w_{x_{i_1}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_2}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_3}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + \frac{1}{2} w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} \\ &= w_{x_{i_1}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_2}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_3}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + \frac{1}{2} + \frac{1}{200N} \\ &\in U(-\frac{1}{2}, \frac{3}{12N} + \frac{1}{200N}) \subseteq U(-\frac{1}{2}, 0.26) \end{aligned}$$

488 Hence  $\frac{\partial \mathcal{L}}{\partial w_{c_\gamma}} \Big|_{\mathbf{t}_\gamma} = -1$ , and

$$w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} = w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} - \eta_{c_\gamma} \frac{\partial \mathcal{L}}{\partial w_{c_\gamma}} \Big|_{\mathbf{t}_\gamma} = \frac{1}{2} + \frac{1}{200N} + 5 = \frac{11}{2} + \frac{1}{200N}$$



489 Similarly,

$$w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} = w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} - \eta_{b_\gamma} \left. \frac{\partial \mathcal{L}}{\partial w_{b_\gamma}} \right|_{\mathbf{t}_\gamma} = 1000N$$

490 Note also that  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin T_{\alpha_{\mathbf{t}_\gamma}, |T|}^1$ , by Lemma A.1 we have  $w_{c_\gamma}^{(2)} = w_{c_\gamma}^{(2, |T|)} = w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} =$   
 491  $\frac{11}{2} - \frac{1}{200N}$  and  $w_{b_\gamma}^{(2)} = w_{b_\gamma}^{(2, |T|)} = w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} = 1000N$ .

492 (b) Otherwise, we have

$$\begin{aligned} y_{\mathbf{t}_\gamma}(\mathbf{w}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)})^\top \mathbf{x}'_{\mathbf{t}_\gamma} &= w_{x_{i_1}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_2}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_3}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + \frac{1}{2} w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} \\ &= w_{x_{i_1}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_2}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + w_{x_{i_3}}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} + \frac{1}{2} + \frac{1}{200N} \\ &\in \bigcup_{x_0 \in \{-\frac{1}{2}, \frac{1}{2}, \frac{5}{2}\}} U(x_0, \frac{3}{12N} + \frac{1}{200N}) \subseteq \bigcup_{x_0 \in \{-\frac{1}{2}, \frac{1}{2}, \frac{5}{2}\}} U(x_0, 0.26) \end{aligned}$$

493 Hence  $\left. \frac{\partial \mathcal{L}}{\partial w_{c_\gamma}} \right|_{\mathbf{t}_\gamma} = \left. \frac{\partial \mathcal{L}}{\partial w_{b_\gamma}} \right|_{\mathbf{t}_\gamma} = 0$ , so  $w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} = w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} = \frac{1}{2} + \frac{1}{200N}$ ,  $w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} = w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1 - 1)} = 0$ .

494 Note also that  $\text{clause}(\gamma, \cdot, \cdot, \cdot) \notin T_{\alpha_{\mathbf{t}_\gamma}, |T|}^1$ , by Lemma A.1 we have  $w_{c_\gamma}^{(2)} = w_{c_\gamma}^{(2, |T|)} = w_{c_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} =$   
 495  $\frac{1}{2} + \frac{1}{200N}$  and  $w_{b_\gamma}^{(2)} = w_{b_\gamma}^{(2, |T|)} = w_{b_\gamma}^{(2, o_{\mathbf{t}_\gamma}^1)} = 0$ .

496 (2) If such  $\mathbf{t}_\gamma = \text{clause}(\gamma, i_1, i_2, i_3)$  does not exist in  $T$ , by Lemma A.1 and Lemma A.3 we have  
 497  $w_{c_\gamma}^{(2)} = w_{c_\gamma}^{(1)} = \frac{1}{2}$  and  $w_{b_\gamma}^{(2)} = w_{b_\gamma}^{(1)} = -1$ .  $\square$

498 Moreover,  $\mathbf{w}$  reaches its fixpoint at the end of the second epoch and will no longer be updated.

499 **Lemma A.6.**  $\mathbf{w}^{(2)} = \mathbf{w}^{(3)}$ .

500 *Proof.* Suppose  $\mathbf{w}^{(2)} \neq \mathbf{w}^{(3)}$ , then there exists  $1 \leq i \leq N$  such that  $w_i^{(2)} \neq w_i^{(3)}$ , and there  
 501 are some training sample  $\mathbf{t}$  in the training data such that  $\left. \frac{\partial \mathcal{L}}{\partial w_i^{(2)}} \right|_{\mathbf{t}} \neq 0$ . Let  $\mathbf{t} = (\mathbf{x}_t, y_t)$  and  
 502  $\mathbb{I} = U(-5, 0.01) \cup U(-\frac{1}{2}, 0.26) \cup \left( \bigcup_{x_0 \in \{\pm 1, \pm 3\}} U(x_0, 0.01) \right)$ . By (2) we have  $y_t(\mathbf{w}^{(2)})^\top \mathbf{x}_t' \in \mathbb{I}$ .  
 503 At least one of the following is true:

- 504 1.  $\exists 1 \leq i \leq n, \mathbf{t} = \text{var}(i)$ . According to lemma A.2,  $y_t(\mathbf{w}^{(2)})^\top \mathbf{x}_t' = y w_{x_i}^{(2)} x_i \in$   
 505  $U(5, 0.5) \subseteq \mathbb{R} \setminus \mathbb{I}$ , contradicting to  $y_t(\mathbf{w}^{(2)})^\top \mathbf{x}_t' \in \mathbb{I}$ .
- 506 2.  $\exists 1 \leq i \leq m$  and  $1 \leq i_1, i_2, i_3 \leq n$ , such that  $\mathbf{t} = \text{clause}(i, i_1, i_2, i_3)$ . According to  
 507 lemma A.5, we have

$$\begin{aligned} y_t(\mathbf{w}^{(2)})^\top \mathbf{x}_t' &= w_{b_i}^{(2)} + w_{c_i}^{(2)} + w_{x_{i_1}}^{(2)} + w_{x_{i_2}}^{(2)} + w_{x_{i_3}}^{(2)} \\ &\geq 1000N + \frac{1}{2} + \frac{1}{200N} + 3 \times (-1 - 0.1) \\ &\geq 1000 - 3.3 \geq 996 \end{aligned}$$

508 We have  $y_t(w^{(2)})^\top \mathbf{x}_t' \notin \mathbb{I}$ , another contradiction.

509 Therefore  $\mathbf{w}^{(2)} = \mathbf{w}^{(3)}$ ,  $\mathbf{w}$  reaches its fixpoint at the end of the second epoch. In other words,  
 510  $\mathbf{w}^* = \mathbf{w}^{(2)}$ .  $\square$

511 We are now ready to give a rigorous proof of theorem 3.1.

512 *Proof of theorem 3.1.* It only suffices to prove the correctness of the reduction in section 3.

513 **If.** Suppose  $\varphi \in \text{MONOTONE 1-IN-3 SAT}$ , then there is a truth assignment  $\nu(\cdot)$  that assigns exactly  
 514 one variable in each clause of  $\varphi$  is true. Let  $\Delta = \{\text{var}(i) | \nu(x_i) = \text{FALSE}\}$ . Let  $\mathbf{w}'$  be the parameter  
 515 of  $\text{SGD}_\Lambda(T_0 \setminus \Delta)$ . By Lemma A.5,  $(w')_{c_\gamma} = \frac{11}{2} + \frac{1}{200N}$  for all  $1 \leq \gamma \leq m$ , hence

$$(\mathbf{w}')^\top \mathbf{x}_{\text{test}} = \sum_{\gamma=1}^m w'_{c_\gamma} \geq \frac{11m}{2} + \frac{-11m+5}{2} = \frac{5}{2} > 0$$

516 and  $\lambda_{\mathbf{w}'}(\mathbf{x}_{\text{test}}) = 1$ , thus  $\text{SGD}_\Lambda(T_0)$  is thus debuggable.

517 **Only if.** Suppose  $\text{SGD}_\Lambda(T_0)$  is debuggable, there will be a  $\Delta$  such that  $\text{SGD}_\Lambda(T_0, \mathbf{x}_{\text{test}}) = y_{\text{test}}$ . We  
 518 denote  $\mathbf{w}'$  as the parameter trained by SGD on  $T_0 \setminus \Delta$ . We have  $\lambda_{\mathbf{w}'}(\mathbf{x}_{\text{test}}) = 1$  and  $(\mathbf{w}')^\top \mathbf{x}_{\text{test}} \geq 0$ .  
 519 By Lemma A.5,  $w'_{c_\gamma} = \{\frac{1}{2} + \frac{1}{200N}, \frac{11}{2} + \frac{1}{200N}\}$ . Suppose  $w_{c^*} = \frac{1}{2} + \frac{1}{200N}$ , then

$$\begin{aligned} (\mathbf{w}')^\top \mathbf{x}_{\text{test}} &= w_{c^*} + \sum_{c_\gamma \neq c^*} w_{c_\gamma} \\ &\leq \frac{11}{2}(m-1) + \frac{1}{2} + \frac{m}{200N} - \frac{11m}{2} + \frac{5}{2} \\ &= -\frac{5}{2} + \frac{m}{200N} \\ &\leq -\frac{5}{2} + \frac{1}{200} = -2.495 < 0 \end{aligned}$$

520 leading to a contradiction.

521 As a consequence,  $w'_{c_\gamma} = \frac{11}{2} + \frac{1}{200N}$  for all  $1 \leq \gamma \leq m$ . By Lemma A.5, exactly one of  
 522  $\text{var}(i_1), \text{var}(i_2), \text{var}(i_3)$  is in  $T_0 \setminus \Delta$  for each  $c_\gamma = (x_{i_1} \vee x_{i_2} \vee x_{i_3})$ . Consider a truth assignment  $\nu$   
 523 that maps every  $x_i$  to FALSE where  $\text{var}(i) \in \Delta$ , and maps the rest to TRUE. Then  $\nu$  assigns exactly  
 524 one variable true in each  $c_\gamma = (x_{i_1} \vee x_{i_2} \vee x_{i_3})$  if and only if exactly one of  $\text{var}(i_1), \text{var}(i_2), \text{var}(i_3)$   
 525 is in  $T_0 \setminus \Delta$ . Hence  $\nu$  is a truth assignment that assigns true to exactly one variable in each clause of  
 526  $\varphi$ , and thus  $\varphi$  is a yes-instance of MONOTONE 1-IN-3 SAT.  $\square$

## 527 B Detailed Proofs for Section 4

### 528 B.1 Proof of Theorem 4.4

529 *Proof.* We build a reduction from the SUBSET SUM problem with a fixed size, which is NP-hard as a  
 530 particular case of the class of knapsack problems [26]. Formally, it is defined as:

SUBSET SUM with a fixed size  
**Input:** A set of positive integer  $S$ , and two positive integers  $t, k$ .  
**Output:** “Yes”: if  $\exists S' \subseteq S$  of size  $k$  such that  $\sum_{a \in S'} a = t$ ;  
 “No”: otherwise.

531 The ordered training data  $T$  is constructed as

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \cup \{(x_a, y_a)\}$$

533 where  $x_i y_i = \frac{2}{3} + \frac{a_i}{3 \sum_{a \in S} a}$  for all  $1 \leq i \leq n$  and  $x_a y_a = 1 + \frac{1}{6 \sum_{a \in S} a}$ . Let  $\eta = 1, \alpha = 1, \beta = -1$ ,  
 534  $w^{(0)} = -1 - \frac{2}{3}k - \frac{t}{3 \sum_{a \in S} a}$  and let the test instance  $(x_{\text{test}}, y_{\text{test}})$  satisfy  $x_{\text{test}} y_{\text{test}} = 1$ . It now suffices  
 535 to prove that  $\exists S' \subseteq S$  such that  $|S'| = k$  and  $\sum_{a \in S'} a = t$  if and only if  $\exists T' \subseteq T$  such that  
 536  $w : w^{(0)} \xrightarrow{T'} w$  satisfies  $y_{\text{test}} w x_{\text{test}} > 0$ .

537 **If:** Suppose  $\exists S' \subseteq S$  such that  $|S'| = k$  and  $\sum_{a \in S'} a = t$ . Let  $T^* = \{(x_i, y_i) | a_i \in S'\}$ , we prove  
 538 that  $y_{\text{test}} w^* x_{\text{test}} > 0$  for  $w^*$  satisfying  $w^{(0)} \xrightarrow{T'=T^* \cup \{(x_a, y_a)\}} w^*$ .

Since

$$\begin{aligned} w^{(0)} + \sum_{a_i \in S'} x_i y_i &= -1 - \frac{2}{3}k - \frac{t}{3 \sum_{a \in S} a} + \sum_{a_i \in S'} \left( \frac{2}{3} + \frac{a_i}{3 \sum_{a \in S} a} \right) \\ &= -1 - \frac{2}{3}k - \frac{t}{3 \sum_{a \in S} a} + \sum_{a_i \in S'} \frac{2}{3} + \frac{\sum_{a \in S'} a}{3 \sum_{a \in S} a} = -1 \end{aligned}$$

and  $\forall 1 \leq i \leq n, x_i y_i > \frac{2}{3}$ , for each  $1 \leq i < n$ , suppose  $w^{(0)} \xrightarrow{T^* \cap \{(x_j, y_j) | 1 \leq j \leq i\}}$   $w_i$ , we have

$$w_i x_{i+1} y_{i+1} < \left( w^{(0)} + \sum_{a_j \in S'} x_j y_j - \frac{2}{3} \right) \cdot \frac{2}{3} < -\frac{10}{9} < \beta.$$

539 That is, each training sample in  $T^*$  is activated. Then for  $w^{(0)} \xrightarrow{T^*} w_a$ , we have  $w_a = -1$ . Then,  
 540 since  $y_a w_a x_a = -(1 + \frac{1}{6 \sum_{a \in S} a}) < \beta$  and  $w_a \xrightarrow{(x_a, y_a)} w^*$  we have  $w^* = w_a + x_a y_a = \frac{1}{6 \sum_{a \in S} a}$ .  
 541 Therefore,  $y_{\text{test}} w^* x_{\text{test}} = \frac{1}{6 \sum_{a \in S} a} > 0$ .

542 Only if: For each  $T' \subseteq T$ , let  $T^* = T' \setminus \{(x_a, y_a)\}$  and  $c(T^*)$  be the set of training samples in  
 543  $T^*$  that are activated. If  $y_{\text{test}} w^* x_{\text{test}} \geq 0$  for  $w^*$  satisfying  $w^{(0)} \xrightarrow{T'} w^*$ , we prove that the set  
 544  $S' = \{a_i | (x_i, y_i) \in c(T^*)\}$  satisfies  $|S'| = k$  and  $\sum_{a \in S'} a = t$ .

545 We first show that  $y_{\text{test}} w_a x_{\text{test}} < 0$  for  $w^{(0)} \xrightarrow{c(T^*)} w_a$ . Otherwise, suppose  $y_{\text{test}} w_a x_{\text{test}} \geq 0$  we  
 546 have  $w_a \geq 0$ . Let  $(x, y)$  be the last training sample of  $c(T')$ , since  $\frac{2}{3} < xy \leq 1$ , we have  
 547  $w' \geq w_a - xy \geq -1$  for  $w' \xrightarrow{(x, y)} w_a$ . Thus  $y w' x \geq \beta$ , which contradicts to the definition of  $c(T^*)$ .

We next show that  $|S'| = k$ . Suppose  $|S'| \leq k - 1$ , we have

$$\begin{aligned} w_a &= w^{(0)} + \sum_{(x_i, y_i) \in c(T^*)} x_i y_i = -1 - \frac{2}{3}k - \frac{t}{3 \sum_{a \in S} a} + \sum_{a_i \in S'} \frac{2}{3} + \frac{\sum_{a \in S'} a}{3 \sum_{a \in S} a} \\ &< -1 - \frac{2}{3}k + \frac{2}{3}(k-1) + \frac{1}{3} = -\frac{4}{3} \end{aligned}$$

548 Thus  $w^* \leq w_a + x_a y_a < -\frac{4}{3} + (1 + \frac{1}{6 \sum_{a \in S} a}) < 0$  and then  $y_{\text{test}} w^* x_{\text{test}} < 0$ , which contradicts to  
 549 the fact that  $y_{\text{test}} w^* x_{\text{test}} \geq 0$ . Therefore  $|S'| \geq k$ .

Suppose  $|S'| \geq k + 1$ , we have

$$w_a = w^{(0)} + \sum_{(x_i, y_i) \in c(T^*)} x_i y_i \geq -1 - \frac{2}{3}k - \frac{1}{3} + \frac{2}{3}(k+1) = -\frac{2}{3}$$

550 Then  $y_a w_a x_a \geq (-\frac{2}{3}) \cdot (1 + \frac{1}{6 \sum_{a \in S} a}) \geq -\frac{7}{9} \geq \beta$ , that is,  $(x_a, y_a)$  is not activated and  $w^* = w_a$ .  
 551 Then since  $y_{\text{test}} w_a x_{\text{test}} < 0$ , we have  $y_{\text{test}} w^* x_{\text{test}} = y_{\text{test}} w_a x_{\text{test}} < 0$ , which contradicts to the fact that  
 552  $y_{\text{test}} w^* x_{\text{test}} \geq 0$ . Therefore  $|S'| = k$ .

It remains to prove that  $\sum_{a \in S'} a = t$ . Otherwise, suppose  $\sum_{a \in S'} a \leq t - 1$ , we have

$$\begin{aligned} w_a &= w^{(0)} + \sum_{(x_i, y_i) \in c(T^*)} x_i y_i \leq -1 - \frac{2}{3}k - \frac{t}{3 \sum_{a \in S} a} + \frac{2}{3}k + \frac{t-1}{3 \sum_{a \in S} a} \\ &= -1 - \frac{1}{3 \sum_{a \in S} a} \end{aligned}$$

553 Thus  $y_{\text{test}} w^* x_{\text{test}} \leq y_{\text{test}} (w_a + x_a y_a) x_{\text{test}} \leq -\frac{1}{6 \sum_{a \in S} a} < 0$ , which contradicts to the fact that  
 554  $y_{\text{test}} w^* x_{\text{test}} \geq 0$ . Therefore  $\sum_{a \in S'} a \geq t$ .

Suppose  $\sum_{a \in S'} a \geq t + 1$  we have

$$\begin{aligned} w_a &= w^{(0)} + \sum_{(x_i, y_i) \in c(T^*)} x_i y_i \geq -1 - \frac{2}{3}k - \frac{t}{3 \sum_{a \in S} a} + \frac{2}{3}k + \frac{t+1}{3 \sum_{a \in S} a} \\ &= -1 + \frac{1}{3 \sum_{a \in S} a} \end{aligned}$$

Thus

$$\begin{aligned} y_a w_a x_a &\geq \left(-1 + \frac{1}{3 \sum_{a \in S} a}\right) \cdot \left(1 + \frac{1}{6 \sum_{a \in S} a}\right) \\ &\geq -1 + \frac{1}{6 \sum_{a \in S} a} + \frac{1}{18 (\sum_{a \in S} a)^2} \geq \beta. \end{aligned}$$

555 That is,  $(x_a, y_a)$  is not activated and  $w^* = w_a$ . Then since  $y_{\text{test}} w_a x_{\text{test}} < 0$ , we have  $y_{\text{test}} w^* x_{\text{test}} =$   
556  $y_{\text{test}} w_a x_{\text{test}} < 0$ , which contradicts to the fact that  $y_{\text{test}} w^* x_{\text{test}} \geq 0$ . Therefore  $\sum_{a \in S'} a = t$ .  $\square$

## 557 B.2 Proof of Theorem 4.3 for $\beta < -1$

558 *Proof.* To avoid cluttering, we still assume  $\eta = 1$  and  $\alpha = 1$ . The proof can be generalized by  
559 appropriately re-scaling the constructed vectors.

560 Let  $M = -\beta(n+2) + 9\beta nm^2(n+1) + 3$ . Suppose  $n = |S| > 1$ ,  $m = \max_{a \in S} \{a\}$  and  
561  $S = \{a_1, a_2, \dots, a_n\}$ . We further assume  $n > 1$ . Let the ordered set of training samples be

$$T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} \cup \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$$

562 where  $\mathbf{x}_i y_i = (\frac{1}{n+1}, -3\beta a_i)$  for all  $1 \leq i \leq n$ ,  $\mathbf{x}_c y_c = (M + \frac{3}{2}\beta - 1, \beta(3t - \frac{1}{2}))$ ,  $\mathbf{x}_b y_b =$   
563  $(1, -1)$ ,  $\mathbf{x}_a y_a = (-\frac{3}{2}\beta, -\frac{3}{2}\beta)$ . Let  $\mathbf{w}^{(0)} = (-M, 0)$ . Let the test instance  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$  satisfy  
564  $\mathbf{x}_{\text{test}} y_{\text{test}} = (1, 0)$ .

565 For each  $1 \leq i < n$ , suppose  $\mathbf{w}^{(0)} \xrightarrow{T \cap \{(\mathbf{x}_i, y_i) | 1 \leq j \leq i\}}$   $\mathbf{w}_i$ , we have

$$\begin{aligned} y_{i+1} \mathbf{w}_i^\top \mathbf{x}_{i+1} &\leq -M \cdot \frac{1}{n+1} + \frac{i}{(n+1)^2} + 9\beta^2 a_{i+1} \sum_{j=1}^i a_j \\ &\leq -M \cdot \frac{1}{n+1} + \frac{n}{(n+1)^2} + 9\beta^2 nm^2 < \beta \end{aligned}$$

566 This means all the  $(\mathbf{x}_i, y_i) \in T \setminus \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$  can be activated and thus the resulting  
567 parameter trained by  $T \setminus \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$  is

$$\mathbf{w}_c = \mathbf{w}^{(0)} + \sum_{i=1}^n \mathbf{x}_i y_i = \left(-M + \frac{|T^*|}{n+1}, -3\beta \sum_{i=1}^n a_i\right)$$

568 It now suffices to prove that for all  $S' \subseteq S$ ,  $\sum_{a \in S'} a = t$  if and only if  $\exists T' \subseteq T$  such that  
569  $\mathbf{w} : \mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}$  such that  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} > 0$ .

570 If: Suppose  $\exists S' \subseteq S$  such that  $\sum_{a \in S'} a = t$ , we prove that  $\exists T' \subseteq T$  such that  $y_{\text{test}} (\mathbf{w}^*)^\top \mathbf{x}_{\text{test}} > 0$   
571 for  $\mathbf{w}^*$  satisfying  $\mathbf{w}^{(0)} \xrightarrow{T^*} \mathbf{w}^*$ .

572 Let  $T^* = \{(\mathbf{x}_i, y_i) | a_i \in S'\}$ ,  $T' = T^* \cup \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$ . We have

$$\mathbf{w}_c = \left(-M + \frac{|T^*|}{n+1}, -3\beta \sum_{a_i \in S'} a_i\right) = \left(-M + \frac{|T^*|}{n+1}, -3\beta t\right)$$

573 And  $y_c \mathbf{w}_c^\top \mathbf{x}_c = (-M + \frac{|T^*|}{n+1})(M + \frac{3}{2}\beta - 1) - 3t\beta^2(3t - \frac{1}{2}) < \beta$ , so

$$\mathbf{w}_c \xrightarrow{(\mathbf{x}_c, y_c)} \mathbf{w}_b = \mathbf{w}_c + \mathbf{x}_c y_c = \left(\frac{|T^*|}{n+1} + \frac{3}{2}\beta - 1, -\frac{1}{2}\beta\right)$$

574 Note that  $\beta < -1$ , we have  $y_b \mathbf{w}_b^\top \mathbf{x}_b = \frac{|T^*|}{n+1} + 2\beta < (\beta + \frac{|T^*|}{n+1}) + \beta < \beta$ , and

$$\mathbf{w}_b \xrightarrow{(\mathbf{x}_b, y_b)} \mathbf{w}_a = \mathbf{w}_b + \mathbf{x}_b y_b = \left(\frac{|T^*|}{n+1} + \frac{3}{2}\beta, -\frac{1}{2}\beta - 1\right)$$

575 Note also that  $y_a \mathbf{w}_a^\top \mathbf{x}_a = \frac{3}{2}(-\beta)(\frac{|T^*|}{n+1} - 1 + \beta) < \beta$ , we have

$$\mathbf{w}_a \xrightarrow{(\mathbf{x}_a, y_a)} \mathbf{w}^* = \mathbf{w}_a + \mathbf{x}_a y_a = \left(\frac{|T^*|}{n+1}, -2\beta - 1\right)$$

576 Therefore,  $y_{\text{test}}(\mathbf{w}^*)^\top \mathbf{x}_{\text{test}} = \frac{|T^*|}{n+1} \geq 0$ .

577 Only if: For each  $T' \subseteq T$ , let  $T^* = T' \setminus \{(\mathbf{x}_c, y_c), (\mathbf{x}_b, y_b), (\mathbf{x}_a, y_a)\}$ , if  $y_{\text{test}}(\mathbf{w}^*)^\top \mathbf{x}_{\text{test}}$  for  $\mathbf{w}^*$   
578 satisfying  $\mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w}^*$ , we prove that  $\exists S' \subseteq S$  such that  $\sum_{a \in S'} a = t$ . We first show that for  
579 each  $T' \subseteq T$ , if  $\mathbf{w}(\mathbf{w}^{(0)} \xrightarrow{T'} \mathbf{w})$  satisfying  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} \geq 0$ , we have  $\forall k \in \{a, b, c\}, (\mathbf{x}_k, y_k) \in$   
580  $T', y_k \mathbf{w}_k^\top \mathbf{x}_k < \beta$ , where  $\mathbf{w}^{(0)} \xrightarrow{T^*} \mathbf{w}_c \xrightarrow{(\mathbf{x}_c, y_c)} \mathbf{w}_b \xrightarrow{(\mathbf{x}_b, y_b)} \mathbf{w}_a$ . Otherwise, suppose  $\exists k \in \{a, b, c\}$   
581 such that  $(\mathbf{x}_k, y_k) \notin T'$  or  $y_k \mathbf{w}_k^\top \mathbf{x}_k \geq \beta$ , we have

$$\begin{aligned} y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} &\leq -M + \frac{|T^*|}{n+1} + M + \frac{3}{2}\beta - 1 + 1 - \frac{3}{2}\beta - \min \left\{ 1, M + \frac{3}{2}\beta - 1, -\frac{3}{2}\beta \right\} \\ &= \frac{|T^*|}{n+1} - 1 < 0 \end{aligned}$$

582 which contradicts to the fact that  $y_{\text{test}} \mathbf{w}^\top \mathbf{x}_{\text{test}} \geq 0$ .

583 Let  $S' = \{a_i | (\mathbf{x}_i, y_i) \in T^*\}$  and  $t' = \sum_{a_i \in S'} a_i$ , it suffices to prove  $t' = t$ . Notice that

$$\begin{aligned} \mathbf{w}^{(0)} \xrightarrow{T^*} \mathbf{w}_c &= \left(-M + \frac{|T^*|}{n+1}, -3\beta \sum_{a_i \in S'} a_i\right) \\ &= \left(-M + \frac{|T^*|}{n+1}, -3\beta t'\right) \end{aligned}$$

584 Hence  $y_c \mathbf{w}_c^\top \mathbf{x}_c = \left(-M + \frac{|T^*|}{n+1}\right)(M + \frac{3}{2}\beta - 1) - 3t'\beta^2(3t - \frac{1}{2}) < \beta$ , thus

$$\mathbf{w}_c \xrightarrow{(\mathbf{x}_c, y_c)} \mathbf{w}_b = \mathbf{w}_c + \mathbf{x}_c y_c = \left(\frac{|T^*|}{n+1} + \frac{3}{2}\beta - 1, -3\beta(t' - t) - \frac{1}{2}\beta\right)$$

585 (1) If  $t' \leq t - 1$ , we have

$$\begin{aligned} y_b \mathbf{w}_b^\top \mathbf{x}_b &= \frac{|T^*|}{n+1} - 1 + 2\beta + 3\beta(t' - t) \\ &\geq \frac{|T^*|}{n+1} - (1 + \beta) > 0 > \beta \end{aligned}$$

586 a contradiction. Hence  $\mathbf{w}_a = \mathbf{w}_b \xrightarrow{(\mathbf{x}_b, y_b)} \mathbf{w}_a = \left(\frac{|T^*|}{n+1} + \frac{3}{2}\beta, -3\beta(t' - t) - \frac{1}{2}\beta - 1\right)$ .

587 (2) If  $t' \geq t + 1$ , we have

$$\begin{aligned} y_a \mathbf{w}_a^\top \mathbf{x}_a &= -\frac{3\beta}{2} \left( \frac{|T^*|}{n+1} - 1 + \beta - 3\beta(t' - t) \right) \\ &\geq -\frac{3\beta}{2} \left( \frac{|T^*|}{n+1} - 1 - 2\beta \right) \\ &> -\frac{3\beta}{2} \left( \frac{|T^*|}{n+1} + 1 \right) > 0 > \beta \end{aligned}$$

588 another contradiction. Therefore  $t' = t$ , and this completes the proof.

589 □

## 590 C Limitations

591 It is important to emphasize that the complexity results in section 4 requires the training order to  
592 be adversarially chosen. The complexity of DEBUGGABLE for randomly chosen training order is  
593 unclear and needs to be figured out in the future research.

## 594 **NeurIPS Paper Checklist**

### 595 **1. Claims**

596 Question: Do the main claims made in the abstract and introduction accurately reflect the  
597 paper's contributions and scope?

598 Answer: [\[Yes\]](#)

599 Justification: The main results are discussed in section 3 and section 4.

600 Guidelines:

- 601 • The answer NA means that the abstract and introduction do not include the claims  
602 made in the paper.
- 603 • The abstract and/or introduction should clearly state the claims made, including the  
604 contributions made in the paper and important assumptions and limitations. A No or  
605 NA answer to this question will not be perceived well by the reviewers.
- 606 • The claims made should match theoretical and experimental results, and reflect how  
607 much the results can be expected to generalize to other settings.
- 608 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
609 are not attained by the paper.

### 610 **2. Limitations**

611 Question: Does the paper discuss the limitations of the work performed by the authors?

612 Answer: [\[Yes\]](#)

613 Justification: See section C in the appendix.

614 Guidelines:

- 615 • The answer NA means that the paper has no limitation while the answer No means that  
616 the paper has limitations, but those are not discussed in the paper.
- 617 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 618 • The paper should point out any strong assumptions and how robust the results are to  
619 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
620 model well-specification, asymptotic approximations only holding locally). The authors  
621 should reflect on how these assumptions might be violated in practice and what the  
622 implications would be.
- 623 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
624 only tested on a few datasets or with a few runs. In general, empirical results often  
625 depend on implicit assumptions, which should be articulated.
- 626 • The authors should reflect on the factors that influence the performance of the approach.  
627 For example, a facial recognition algorithm may perform poorly when image resolution  
628 is low or images are taken in low lighting. Or a speech-to-text system might not be  
629 used reliably to provide closed captions for online lectures because it fails to handle  
630 technical jargon.
- 631 • The authors should discuss the computational efficiency of the proposed algorithms  
632 and how they scale with dataset size.
- 633 • If applicable, the authors should discuss possible limitations of their approach to  
634 address problems of privacy and fairness.
- 635 • While the authors might fear that complete honesty about limitations might be used by  
636 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
637 limitations that aren't acknowledged in the paper. The authors should use their best  
638 judgment and recognize that individual actions in favor of transparency play an impor-  
639 tant role in developing norms that preserve the integrity of the community. Reviewers  
640 will be specifically instructed to not penalize honesty concerning limitations.

### 641 **3. Theory Assumptions and Proofs**

642 Question: For each theoretical result, does the paper provide the full set of assumptions and  
643 a complete (and correct) proof?

644 Answer: [\[Yes\]](#)

645 Justification: The proof of theorem 3.1 is available in section A; The proof of theorem 4.1  
646 and theorem 4.2 are available in section 4; The proof of theorem 4.3 is available in section 4  
647 and section B; The proof of theorem 4.4 is available in section B.

648 Guidelines:

- 649 • The answer NA means that the paper does not include theoretical results.
- 650 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
651 referenced.
- 652 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 653 • The proofs can either appear in the main paper or the supplemental material, but if  
654 they appear in the supplemental material, the authors are encouraged to provide a short  
655 proof sketch to provide intuition.
- 656 • Inversely, any informal proof provided in the core of the paper should be complemented  
657 by formal proofs provided in appendix or supplemental material.
- 658 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 659 4. Experimental Result Reproducibility

660 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
661 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
662 of the paper (regardless of whether the code and data are provided or not)?

663 Answer: [NA]

664 Justification: This paper does not include experiments.

665 Guidelines:

- 666 • The answer NA means that the paper does not include experiments.
- 667 • If the paper includes experiments, a No answer to this question will not be perceived  
668 well by the reviewers: Making the paper reproducible is important, regardless of  
669 whether the code and data are provided or not.
- 670 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
671 to make their results reproducible or verifiable.
- 672 • Depending on the contribution, reproducibility can be accomplished in various ways.  
673 For example, if the contribution is a novel architecture, describing the architecture fully  
674 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
675 be necessary to either make it possible for others to replicate the model with the same  
676 dataset, or provide access to the model. In general, releasing code and data is often  
677 one good way to accomplish this, but reproducibility can also be provided via detailed  
678 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
679 of a large language model), releasing of a model checkpoint, or other means that are  
680 appropriate to the research performed.
- 681 • While NeurIPS does not require releasing code, the conference does require all submis-  
682 sions to provide some reasonable avenue for reproducibility, which may depend on the  
683 nature of the contribution. For example
  - 684 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
685 to reproduce that algorithm.
  - 686 (b) If the contribution is primarily a new model architecture, the paper should describe  
687 the architecture clearly and fully.
  - 688 (c) If the contribution is a new model (e.g., a large language model), then there should  
689 either be a way to access this model for reproducing the results or a way to reproduce  
690 the model (e.g., with an open-source dataset or instructions for how to construct  
691 the dataset).
  - 692 (d) We recognize that reproducibility may be tricky in some cases, in which case  
693 authors are welcome to describe the particular way they provide for reproducibility.  
694 In the case of closed-source models, it may be that access to the model is limited in  
695 some way (e.g., to registered users), but it should be possible for other researchers  
696 to have some path to reproducing or verifying the results.

#### 697 5. Open access to data and code

698 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
699 tions to faithfully reproduce the main experimental results, as described in supplemental  
700 material?

701 Answer: [NA]

702 Justification: This paper does not include experiments requiring code.

703 Guidelines:

- 704 • The answer NA means that paper does not include experiments requiring code.
- 705 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
706 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 707 • While we encourage the release of code and data, we understand that this might not be  
708 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
709 including code, unless this is central to the contribution (e.g., for a new open-source  
710 benchmark).
- 711 • The instructions should contain the exact command and environment needed to run to  
712 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
713 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 714 • The authors should provide instructions on data access and preparation, including how  
715 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 716 • The authors should provide scripts to reproduce all experimental results for the new  
717 proposed method and baselines. If only a subset of experiments are reproducible, they  
718 should state which ones are omitted from the script and why.
- 719 • At submission time, to preserve anonymity, the authors should release anonymized  
720 versions (if applicable).
- 721 • Providing as much information as possible in supplemental material (appended to the  
722 paper) is recommended, but including URLs to data and code is permitted.

## 723 6. Experimental Setting/Details

724 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
725 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
726 results?

727 Answer: [NA]

728 Justification: This paper does not include experiments.

729 Guidelines:

- 730 • The answer NA means that the paper does not include experiments.
- 731 • The experimental setting should be presented in the core of the paper to a level of detail  
732 that is necessary to appreciate the results and make sense of them.
- 733 • The full details can be provided either with the code, in appendix, or as supplemental  
734 material.

## 735 7. Experiment Statistical Significance

736 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
737 information about the statistical significance of the experiments?

738 Answer: [NA]

739 Justification: This paper does not include experiments.

740 Guidelines:

- 741 • The answer NA means that the paper does not include experiments.
- 742 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
743 dence intervals, or statistical significance tests, at least for the experiments that support  
744 the main claims of the paper.
- 745 • The factors of variability that the error bars are capturing should be clearly stated (for  
746 example, train/test split, initialization, random drawing of some parameter, or overall  
747 run with given experimental conditions).
- 748 • The method for calculating the error bars should be explained (closed form formula,  
749 call to a library function, bootstrap, etc.)



- 750 • The assumptions made should be given (e.g., Normally distributed errors).
- 751 • It should be clear whether the error bar is the standard deviation or the standard error
- 752 of the mean.
- 753 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 754 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 755 of Normality of errors is not verified.
- 756 • For asymmetric distributions, the authors should be careful not to show in tables or
- 757 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 758 error rates).
- 759 • If error bars are reported in tables or plots, The authors should explain in the text how
- 760 they were calculated and reference the corresponding figures or tables in the text.

## 761 8. Experiments Compute Resources

762 Question: For each experiment, does the paper provide sufficient information on the com-  
 763 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
 764 the experiments?

765 Answer: [NA]

766 Justification: This paper does not include experiments.

767 Guidelines:

- 768 • The answer NA means that the paper does not include experiments.
- 769 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 770 or cloud provider, including relevant memory and storage.
- 771 • The paper should provide the amount of compute required for each of the individual
- 772 experimental runs as well as estimate the total compute.
- 773 • The paper should disclose whether the full research project required more compute
- 774 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 775 didn't make it into the paper).

## 776 9. Code Of Ethics

777 Question: Does the research conducted in the paper conform, in every respect, with the  
 778 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

779 Answer: [Yes]

780 Justification: They authors have made sure that the research conducted in the paper conform  
 781 with the NeurIPS Code of Ethics.

782 Guidelines:

- 783 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 784 • If the authors answer No, they should explain the special circumstances that require a
- 785 deviation from the Code of Ethics.
- 786 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
 787 eration due to laws or regulations in their jurisdiction).

## 788 10. Broader Impacts

789 Question: Does the paper discuss both potential positive societal impacts and negative  
 790 societal impacts of the work performed?

791 Answer: [NA]

792 Justification: The impacts are discussed in section 5

793 Guidelines:

- 794 • The answer NA means that there is no societal impact of the work performed.
- 795 • If the authors answer NA or No, they should explain why their work has no societal  
 796 impact or why the paper does not address societal impact.
- 797 • Examples of negative societal impacts include potential malicious or unintended uses  
 798 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
 799 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
 800 groups), privacy considerations, and security considerations.

- 801 • The conference expects that many papers will be foundational research and not tied  
802 to particular applications, let alone deployments. However, if there is a direct path to  
803 any negative applications, the authors should point it out. For example, it is legitimate  
804 to point out that an improvement in the quality of generative models could be used to  
805 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
806 that a generic algorithm for optimizing neural networks could enable people to train  
807 models that generate Deepfakes faster.
- 808 • The authors should consider possible harms that could arise when the technology is  
809 being used as intended and functioning correctly, harms that could arise when the  
810 technology is being used as intended but gives incorrect results, and harms following  
811 from (intentional or unintentional) misuse of the technology.
- 812 • If there are negative societal impacts, the authors could also discuss possible mitigation  
813 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
814 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
815 feedback over time, improving the efficiency and accessibility of ML).

## 816 11. Safeguards

817 Question: Does the paper describe safeguards that have been put in place for responsible  
818 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
819 image generators, or scraped datasets)?

820 Answer: [NA]

821 Justification: This paper only provides theoretical results and poses no such risks.

822 Guidelines:

- 823 • The answer NA means that the paper poses no such risks.
- 824 • Released models that have a high risk for misuse or dual-use should be released with  
825 necessary safeguards to allow for controlled use of the model, for example by requiring  
826 that users adhere to usage guidelines or restrictions to access the model or implementing  
827 safety filters.
- 828 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
829 should describe how they avoided releasing unsafe images.
- 830 • We recognize that providing effective safeguards is challenging, and many papers do  
831 not require this, but we encourage authors to take this into account and make a best  
832 faith effort.

## 833 12. Licenses for existing assets

834 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
835 the paper, properly credited and are the license and terms of use explicitly mentioned and  
836 properly respected?

837 Answer: [NA]

838 Justification: This paper does not use existing assets.

839 Guidelines:

- 840 • The answer NA means that the paper does not use existing assets.
- 841 • The authors should cite the original paper that produced the code package or dataset.
- 842 • The authors should state which version of the asset is used and, if possible, include a  
843 URL.
- 844 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 845 • For scraped data from a particular source (e.g., website), the copyright and terms of  
846 service of that source should be provided.
- 847 • If assets are released, the license, copyright information, and terms of use in the  
848 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
849 has curated licenses for some datasets. Their licensing guide can help determine the  
850 license of a dataset.
- 851 • For existing datasets that are re-packaged, both the original license and the license of  
852 the derived asset (if it has changed) should be provided.

853 • If this information is not available online, the authors are encouraged to reach out to  
854 the asset’s creators.

### 855 13. New Assets

856 Question: Are new assets introduced in the paper well documented and is the documentation  
857 provided alongside the assets?

858 Answer: [NA]

859 Justification: This paper does not release new assets.

860 Guidelines:

- 861 • The answer NA means that the paper does not release new assets.
- 862 • Researchers should communicate the details of the dataset/code/model as part of their  
863 submissions via structured templates. This includes details about training, license,  
864 limitations, etc.
- 865 • The paper should discuss whether and how consent was obtained from people whose  
866 asset is used.
- 867 • At submission time, remember to anonymize your assets (if applicable). You can either  
868 create an anonymized URL or include an anonymized zip file.

### 869 14. Crowdsourcing and Research with Human Subjects

870 Question: For crowdsourcing experiments and research with human subjects, does the paper  
871 include the full text of instructions given to participants and screenshots, if applicable, as  
872 well as details about compensation (if any)?

873 Answer: [NA]

874 Justification: This paper does not involve crowdsourcing nor research with human subjects.

875 Guidelines:

- 876 • The answer NA means that the paper does not involve crowdsourcing nor research with  
877 human subjects.
- 878 • Including this information in the supplemental material is fine, but if the main contribu-  
879 tion of the paper involves human subjects, then as much detail as possible should be  
880 included in the main paper.
- 881 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
882 or other labor should be paid at least the minimum wage in the country of the data  
883 collector.

### 884 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 885 Subjects

886 Question: Does the paper describe potential risks incurred by study participants, whether  
887 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
888 approvals (or an equivalent approval/review based on the requirements of your country or  
889 institution) were obtained?

890 Answer: [NA]

891 Justification: This paper does not involve crowdsourcing nor research with human subjects.

892 Guidelines:

- 893 • The answer NA means that the paper does not involve crowdsourcing nor research with  
894 human subjects.
- 895 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
896 may be required for any human subjects research. If you obtained IRB approval, you  
897 should clearly state this in the paper.
- 898 • We recognize that the procedures for this may vary significantly between institutions  
899 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
900 guidelines for their institution.
- 901 • For initial submissions, do not include any information that would break anonymity (if  
902 applicable), such as the institution conducting the review.