# PrIU: A Provenance-Based Approach for Incrementally Updating Regression Models

Yinjun Wu University of Pennsylvania wuyinjun@seas.upenn.edu Val Tannen University of Pennsylvania val@cis.upenn.edu Susan B. Davidson University of Pennsylvania susan@cis.upenn.edu

### **ABSTRACT**

The ubiquitous use of machine learning algorithms brings new challenges to traditional database problems such as incremental view update. Much effort is being put in better understanding and debugging machine learning models, as well as in identifying and repairing errors in training datasets. Our focus is on how to assist these activities when they have to retrain the machine learning model after removing problematic training samples in cleaning or selecting different subsets of training data for interpretability. This paper presents an efficient provenance-based approach, PrIU, and its optimized version, PrIU-opt, for incrementally updating model parameters without sacrificing prediction accuracy. We prove the correctness and convergence of the incrementally updated model parameters, and validate it experimentally. Experimental results show that up to two orders of magnitude speed-ups can be achieved by PrIU-opt compared to simply retraining the model from scratch, yet obtaining highly similar models.

### **CCS CONCEPTS**

• Information systems  $\rightarrow$  Data cleaning,Incremental maintenance,Data provenance; • Mathematics of computing  $\rightarrow$  Exploratory data analysis; • Theory of computation  $\rightarrow$  Convex optimization.

### **KEYWORDS**

Data provenance, machine learning, deletion propagation

### **ACM Reference Format:**

Yinjun Wu, Val Tannen, and Susan B. Davidson. 2020. PrIU: A Provenance-Based Approach for Incrementally Updating Regression Models. In 2020 ACM SIGMOD International Conference on

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

SIGMOD '20, June 14–19, 2020, Portland, OR, USA © 2020 Association for Computing Machinery. ACM ISBN 978-1-4503-6735-6/20/06...\$15.00 https://doi.org/10.1145/3318464.3380571

Management of Data (SIGMOD'20), June 14–19, 2020, Portland, OR, USA. ACM, New York, NY, USA, 28 pages. https://doi.org/10.1145/3318464.3380571

# 1 INTRODUCTION

In database terminology, this paper is about *efficient incremental view updates*, specifically about using provenance annotations to propagate the effect of deletions from the input data to the output. However, the views that we consider are *regression models* (linear and binomial/multinomial logistic regression) and the input data consists of the samples used to train these models.

The need for incremental techniques to efficiently update regression models arises in several contexts, for example data cleaning and interpretability. Data cleaning has been extensively studied by the database community [8, 12, 17, 46], and is typically an iterative and interactive process, allowing data analysts to alternate between analysis and cleaning tasks, as well as to interact with other parties such as IT staff and data curators [33]. Machine learning techniques are particularly sensitive to dirty data in training datasets, since it can result in erroneous models and counter-intuitive predictions for test datasets [8]. A number of techniques have therefore recently been proposed for detecting and repairing dirty data in machine learning, e.g., [25, 32]. The work presented in this paper can be incorporated into these data cleaning pipelines by assuming that dirty data in the training set has already been detected, and addresses the next step by providing a solution for incrementally updating the machine learning model after the dirty data is removed.

Interpretability is also a major concern in machine learning (see, for example, the general discussions in [15, 38], the extensive human subjects experiments in [45], as well the many references in these papers). The problem is being studied from several different perspectives (see Sec. 2). The data-driven approaches of [15, 35] discover factors of interpretability by performing repeated retraining of models using multiple different subsets of a training dataset to understand the relationship between samples with certain feature characteristics and the model behavior. Such repeated retraining also occurs in model debugging [25, 28, 34] and deletion diagnostics [9].

In this respect, our work shares goals with [30], which develops an *influence function* to approximately quantify the influence of a *single* training sample on the model parameters and prediction results; this can also be used for estimating the model parameter change after the removal of one training sample. However, extending the influence function approach to multiple training samples significantly weakens prediction accuracy. *In contrast, our techniques are not only efficient but significantly more accurate.* 

**Connection to Provenance.** Note that the problem of incrementally updating the model after removing a subset of the training samples can be seen as a question of *data provenance* [4, 7, 20]. Data provenance tracks the dependencies between input and output data; in particular, the *provenance semiring framework* [20, 21] has been used for applying incremental updates (specifically deletions) to views.

In the semiring framework, input data is annotated with provenance tokens which are carried through the operators performed on the data (e.g. select, project, join, union). Output data is then annotated with *provenance polynomials* expressed in terms of the provenance tokens. When an input tuple is deleted, the effect on the output can be efficiently calculated by essentially "zeroing out" its token in the provenance polynomial. Recently, the framework has been extended to include basic linear algebra operations: matrix addition and multiplication [52]. In this extension, the provenance polynomials play the role of scalars and multiplication with scalars plays the role of annotating matrices and vectors with provenance.

As an example, suppose that p, q, r, s are provenance tokens that annotate samples in a training dataset. Our methods will show that vectors of interest (such as the vector of model parameters) can be expressed with provenanceannotated expressions such as:

$$\mathbf{w} = (p^2q * \mathbf{u}) + (qr^4 * \mathbf{v}) + (ps * \mathbf{z})$$

Here, u, v, z are numerical vectors signifying contributions to the answer w and they are annotated (algebraic operation \*) with  $p^2q$ ,  $qr^4$ , ps which are provenance polynomials to be read as follows: the provenance  $p^2q$  represents the use of both data items labeled p and q and, in fact, the first item is used twice. Now suppose the data item annotated with r is deleted while those annotated p, q, s are retained. We can express the updated value of w under this deletion by setting r to the "provenance 0 polynomial", denoted  $0_{prov}$ which signifies absence, and p, q, s to the "provenance 1 polynomial", denoted 1prov, which signifies "neutral" presence, no need to track further. The algebraic properties of provenance polynomials and of their annotation of matrices/vectors ensure what one would expect, e.g,  $0_{prov} \cdot r^4 = 0_{prov}$  as well as  $0_{\mathrm{prov}} * \mathbf{v} = \mathbf{0}$  (the all-zero vector) and  $1_{\mathrm{prov}} * \mathbf{z} = \mathbf{z}$ . It follows that under this deletion  $\mathbf{w} = \mathbf{u} + \mathbf{z}$ .

**Approach.** In this paper, we use the extension of the semiring framework to matrix operations to track the provenance of input samples through the training of logistic regression and linear regression models using gradient descent and its variants. In each iteration of the training phase, a gradient-based "update rule" updates the model parameters, which can be annotated with provenance polynomials. For logistic regression, we can achieve this via piecewise linear interpolation over the non-linear components in the gradient update rule.

In addition to enabling provenance tracking, the linearization of the gradient update rule allows us to separate the contributions of the training samples from the contributions of the model parameters from the previous iteration. As a result, the effect of deleting training samples on the gradient update rule can be obtained by "zeroing out" the provenance tokens corresponding to those samples.

Challenges. Reasoning over provenance to enable incremental updates introduces significant overhead in the gradient descent calculation. To speed up incremental updates over model parameters for dense datasets, we use several optimizations in our implementation, PrIU: First, between iterations during the training phase over the full training dataset, we cache intermediate results (some matrix expression) that capture only the contribution of the training samples. These are annotated with provenance. Then during the model update phase, the propagation of the deletion of a subset of samples comes down to a subtraction of the "zeroedout" contributions of the removed samples. Second, we apply singular value decomposition (SVD) over the intermediate results to reduce their dimensions. An optimized version of PrIU, PrIU-opt, is also designed for further optimizations over datasets with small feature sets using incremental updates to eigenvalues. (For logistic regression, it is used by terminating provenance tracking early when provenance expressions stabilize. See Section 5 for more details). But the optimizations above cannot work for sparse datasets, for which we use only the linearization of the update rule for logistic regression.

As we shall see, PrIU and PrIU-opt can lead to speed-ups of up to 2 orders of magnitude when compared to a baseline of retraining the model from the updated input data; however, for sparse datasets the speedup is only 10%. While the practical impact of this speed-up may be small for an engineer who only deletes one subset of training samples, especially if retraining takes only a few minutes, the impact is much greater for an engineer who repeatedly removes multiple different subsets of training samples, e.g. when exploring factors of interpretability. In this case, even one order of magnitude speed-up reduces exploration from several hours to a few minutes.

**Contributions** of this paper include:

- (1) A theoretical framework which enables data provenance to be tracked and used for fast incremental model updates when subsets of training samples are removed. The framework extends the approach in [19, 20, 52] to linear regression and (binary and multinomial) logistic regression models.
- (2) Analytical results showing the convergence and accuracy of the updated model parameters for logistic regression, which are approximately computed by applying piecewise linear interpolation over the non-linear operations in the model parameter update rules.
- (3) Efficient provenance-based algorithms, PrIU and PrIUopt, which achieve fast model updates after removing subsets of training samples.
- (4) Extensive experiments showing the effectiveness and accuracy of PrIU and PrIU-opt in incrementally updating the linear regression and logistic regression models compared to the straightforward approach of retraining from scratch, as well as compared to implementing an extension of the influence function in [30].
- (5) Enabling work on interpretability that seeks to understand the effect of removing *subsets* of the training data, rather than just of a single training sample.

The remainder of the paper is organized as follows. In Section 2, we describe related work in incremental model maintenance, data provenance, data cleaning, and machine learning model interpretability. Section 3 reviews the basic concepts of linear regression and logistic regression. The theoretical development of how to use provenance in the update rules of linear regression and logistic regression is presented in Section 4, and its implementation provided in Section 5. Experimental results comparing our approach to other solutions are presented in Section 6. We conclude in Section 7.

To our knowledge, this is the first work to use provenance for the purpose of incrementally updating machine learning model parameters.

### 2 RELATED WORK

Incremental model maintenance. There have been several proposals for materializing machine learning models for future reuse. [13, 22, 40] target the problem of efficiently updating the model as the training data changes, which focus primarily on linear regression and Naive Bayes models, and use closed-form solutions (rather than iterative algorithms, e.g., gradient-based approaches) of the model parameters to determine incremental updates in light of additions and deletions of training samples while [23] deals with how to merge prematerialized models to construct new models based on user requests. In addition, [22] also deals with incremental updates of the model parameters based on the Mixture Weight

Methods (a variant of gradient descent) for logistic regression. The method, however, puts additional training samples into another batch and averages the pre-computed parameters derived from other batches (over the original data) with the parameters computed over the additional batch. This cannot be used for incremental deletions which is our focus in this paper.

The basic ideas of [13, 22, 40] on how to incrementally update linear regression models are somewhat similar. Due to the existence of the matrix inverse operations in the closed-form solution for linear regression, only the intermediate results built with linear operations are maintained as views. They are updated when insertion or deletion happens in the input training data. After that, matrix inversion is used to compute the final updated model parameters. In contrast, our approach proceeds directly to a gradient descent-based linear regression. As we shall see, our experiments show that our approach is more efficient than the closed-form update.

Data provenance. Data provenance captures where data comes from and how it is processed. Within the database community, various approaches have been proposed to track provenance through queries, e.g. where and why provenance [4], and semiring provenance [2, 20]. Provenance is used to identify the source of errors in computational processes, such as workflows [1] and network diagnostics [53]. It is also used to support efficient incremental updates through database queries and schema mappings [18, 19, 26] and workflow computation [16]. Provenance support for linear algebra operations in the context of machine learning tasks has also been recently studied [52]. This work was mentioned in [5] as a first step in using data provenance for interpretability of machine learning models.

Data cleaning. The goal of data cleaning is to detect and fix errors in data, and is a crucial step in preparing data for data analytics/machine learning tasks [14, 34]. However, if erroneous/dirty data is detected after the model has been trained, the machine learning algorithm must be rerun to obtain the updated model parameters. This repetitive training can cause significant delays when large volumes of data are processed. One approach is to start each training phase by setting the initial model parameters to the ones generated by the previous training phase over the dirty data [34]. Our contribution is orthogonal to this approach, and updates the machine learning model parameters directly by reasoning over provenance rather than retraining from scratch.

Interpreting and understanding ML models. Fully understanding the behavior of ML models, especially deep neural network models, is difficult due to their complexity. Moreover, there are different perspectives on what we should understand. For example, one approach separates model components into "shape" functions, one for each feature, in generalized additive models, in particular for linear and

logistic regression [6, 39]. Closest to our perspective is the idea of *influence function* [30] (similar problem is also mentioned in [44]), which originates from *deletion diagnostics* in statistics [9]. [30] estimates the effect of removing a *single* training sample on the already obtained model, *without* retraining the model. The influence function uses the Taylor expansion of the derivative of a customized objective function for the model parameter. The calculation (and thus the approximation of model parameter change) is only based on lower-order terms in Taylor expansion.

This can be seen as a method for incremental model update for just one sample deletion. In fact, we have observed that the method could be extended to deleting an arbitrary number of samples, which led us to compare it experimentally to our approach. The results (see Section 6) show that this approach leads to very inaccurate results when multiple training samples are deleted.

### 3 PRELIMINARIES

We give an overview of linear and logistic regression along with the gradient-based method for learning model parameters. Assume a training dataset ( $\mathbf{X}$ ,  $\mathbf{Y}$ ), where  $\mathbf{X}$  is an  $n \times m$  matrix representing the feature matrix while  $\mathbf{Y}$  is an  $n \times 1$  vector representing the labels, i.e.:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \end{bmatrix}^T \mathbf{Y} = \begin{bmatrix} y_1, y_2, \dots, y_n \end{bmatrix}^T$$
 (1)

For both linear and logistic regression we only focus on a common case: L2-regularization. The objective functions of linear regression, binary logistic regression and multinomial logistic regression with L2-regularization are presented in Equations 2-4 respectively  $^1$ 

$$h(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \frac{\lambda}{2} ||\mathbf{w}||_2^2$$
 (2)

$$h(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} \ln(1 + \exp\{-y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i\}) + \frac{\lambda}{2} ||\mathbf{w}||_2^2$$
 (3)

$$h(\mathbf{w}) = \frac{1}{n} \sum_{k=1}^{q} \sum_{y_i = k} (\ln(\sum_{j=1}^{q} e^{\mathbf{w}_j^{\mathsf{T}} \mathbf{x}_i}) - \mathbf{w}_k^{\mathsf{T}} \mathbf{x}_i) + \frac{\lambda}{2} ||\mathbf{w}||_2^2$$

$$\mathbf{w} = vec([\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q])$$
(4)

where  $\mathbf{w}$  is the vector of model parameters and  $\lambda$  is the *regularization rate*. For simplicity, we denote  $\mathbf{w} = vec([\mathbf{w}_1, \mathbf{w}_2, \ldots, \mathbf{w}_q])$  for multinomial logistic regression where q represents the number of possible classes. Typical learning methods for computing  $\mathbf{w}$  are to apply gradient descent (GD) or its variant, stochastic gradient descent (SGD) or mini-batch stochastic gradient method (mb-SGD) [47] to minimize the objective function  $h(\mathbf{w})$  iteratively. GD, SGD and mb-SGD are the same in nature since mb-SGD can be regarded as

a generalization of GD and SGD. They are therefore called Gradient-based method (GBM) for short, and hereafter we will only take mb-SGD as an example. Considering the similarities between binary logistic regression and multinomial logistic regression and the complexity of the computation related to the latter one, we will only present the formulas related to binary logistic regression below. All the theorems that hold for binary logistic regression can be also proven to be true for multinomial logistic regression.

At each iteration, mb-SGD updates the  $\mathbf{w}^{(t)}$  by using the average gradient of  $h(\mathbf{w})$  over a randomly selected minibatch from the training dataset. Specifically, for linear regression and logistic regression, the rule for updating  $\mathbf{w}^{(t)}$  under mb-SGD is presented below (Equations 5 and 6 respectively):

$$\mathbf{w}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \frac{2\eta_t}{B} \sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i (\mathbf{x}_i^T \mathbf{w}^{(t)} - y_i) \quad (5)$$

$$\mathbf{w}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}^{(t)} + \frac{\eta_t}{B} \sum_{i \in \mathcal{D}^{(t)}} y_i \mathbf{x}_i \left(1 - \frac{1}{1 + \exp\{-y_i \mathbf{w}^{(t)T} \mathbf{x}_i\}}\right)$$
(6)

where  $\eta_t$  is called the *learning rate* and  $\mathcal{B}^{(t)}$  represents a mini-batch of B training samples. For SGD,  $\mathcal{B}^{(t)}$  includes only one sample (B = 1), while for GD,  $\mathcal{B}^{(t)}$  includes all the training samples (B = n).

### 4 ITERATION MODELS

In this section, we first discuss the annotation with provenance of the gradient-bases update rules in our approach. Next, we discuss for the non-linear operations in logistic regression the linearization that makes our provenance annotation framework usable. Finally, we give a rigorous theoretical analysis of the convergence of the iterative process with provenance-annotated update rules for both linear and logistic regression model and the similarity to the expected results after linearization for logistic regression models.

# 4.1 Provenance annotations for matrices

In the semiring framework [2, 20, 21] one begins by annotating input data with elements of a set T of provenance tokens. These annotations are then propagated through query operators as they combine according to two operations: "+" that records alternative use of information, as in relational union or projection, and "·", that records joint use of information, as in relational join. With these, the annotations become provenance polynomials whose indeterminates are tokens and with coefficients in  $\mathbb{N}$ . For example, the monomial  $p^2q$  is the provenance of a result for which the data item annotated p was used twice together with the item annotated q used once. We denote the set of polynomials by  $\mathbb{N}[T]$ .

 $<sup>^1\</sup>mathrm{Here}$  we assume that the two possible labels in binary logistic regression are 1 and -1.

In the extension of the framework to matrix algebra [52], annotation formally becomes a *multiplication of vectors with scalars* as in linear algebra. The role of scalars is played by provenance polynomials and the role of vectors, of course, is played by matrices (generalizing their row vectors and the transposes of these).

Matrices annotated with provenance polynomials form a nice algebraic structure that extends matrix multiplication and addition. We denote multiplication with scalars by "\*" writing p\*A for the matrix A annotated with the provenance polynomial p. For space reasons we cannot repeat here the technical development in [52], however, we mention a crucial algebraic property of annotated matrix multiplication, which also illustrates combining provenance in joint use:

$$(\mathfrak{p}_1 * \mathbf{A}_1)(\mathfrak{p}_2 * \mathbf{A}_2) = (\mathfrak{p}_1 \cdot \mathfrak{p}_2) * (\mathbf{A}_1 \mathbf{A}_2)$$

We apply this framework to tracking input training samples through GBM's in which the update involves only matrix multiplication and addition. Let the training dataset be  $(\mathbf{X}, \mathbf{Y})$  where  $\mathbf{X}$  is an  $n \times m$  feature matrix and  $\mathbf{Y}$  is an  $n \times 1$  column vector of sample labels. For  $i = 1, \ldots, n$ , we annotate every sample  $(\mathbf{x}_i, y_i)$   $(\mathbf{x}_i$  and  $[y_i]$  are the i'th rows in  $\mathbf{X}$  respectively  $\mathbf{Y}$ ) with a distinct provenance token  $p_i$ . Next, we decompose  $\mathbf{X}$  and  $\mathbf{Y}$  as algebraic expressions in terms of  $p_1 * \mathbf{x}_1, \ldots, p_n * \mathbf{x}_n, p_1 * [y_1], \ldots, p_n * [y_n]$  and some matrices made up of the reals 0 and 1. These "helper" matrices are annotated with the provenance polynomial  $1_{\text{prov}} \in \mathbb{N}[T]$  (has only a term of degree zero which is the natural number 1) meaning "always available, no need to track". We illustrate with the provenance-annotated  $\mathbf{X}$  when n = 2:

$$\mathbf{X} = (1_{\text{prov}} * \begin{bmatrix} 1 \\ 0 \end{bmatrix})(p_1 * \mathbf{x}_1) + (1_{\text{prov}} * \begin{bmatrix} 0 \\ 1 \end{bmatrix})(p_2 * \mathbf{x}_2) =$$

$$= (p_1 * \begin{bmatrix} \mathbf{x}_1 \\ 0 \dots 0 \end{bmatrix}) + (p_2 * \begin{bmatrix} 0 \dots 0 \\ \mathbf{x}_2 \end{bmatrix})$$

When  $\mathbf{X}$  is transposed, a similar decomposition applies in terms of the annotated column vectors  $p_i * \mathbf{x}_i^T$ . We also note that the algebra of annotated matrices follows the same laws as the usual matrix algebra. Consequently, we can perform in the algebra of provenance-annotated matrices the calculations involved in the gradient-based update rules. For illustration, a calculation involving  $\mathbf{X}$  that without provenance takes the form  $\sum_{i=1}^n \alpha_i \mathbf{x}_i \mathbf{x}_i^T$  (where  $\alpha_i$  are some real numbers) becomes with provenance annotations

$$\sum_{i=1}^{n} (1_{\text{prov}} * [\alpha_i])(p_i * \mathbf{x}_i)(p_i * \mathbf{x}_i^T) = \sum_{i=1}^{n} p_i^2 * (\alpha_i \mathbf{x}_i \mathbf{x}_i^T).$$

And here is the provenance-annotated expression for the update rule of linear regression (i.e. Equation 5):

$$\mathcal{W}^{(t+1)} \leftarrow [(1 - \eta_t \lambda)(1_k * \mathbf{I}) - \frac{2\eta_t}{\mathcal{P}^{(t)}} \sum_{i \in \mathcal{B}^{(t)}} p_i^2 * \mathbf{x}_i \mathbf{x}_i^T] \mathcal{W}^{(t)} + \frac{2\eta_t}{\mathcal{P}^{(t)}} \sum_{i \in \mathcal{B}^{(t)}} p_i^2 * \mathbf{x}_i y_i$$
(7)

where  $\mathbf{W}^{(t)}$  represents the provenance-annotated expression for the vector  $\mathbf{w}^{(t)}$  of model parameters while  $\mathcal{P}^{(t)}$  represents a provenance-annotated expression for the number of samples in the min-batch  $\mathcal{B}^{(t)}$ , for example, following the approach to aggregation in [2],  $\mathcal{P}^{(t)} = \sum_{i \in \mathcal{B}^{(t)}} p_i * 1$ .

In the semiring framework there is no division operation so we used fractions with denominator  $\mathcal{P}^{(t)}$  in Equation 7 only for notational purposes. As we shall see immediately below, in incremental update  $\mathcal{P}^{(t)}$  can be replaced with an integer.

As with the other applications of the semiring framework, deletion propagation is done by "zeroing-out" the deleted samples. That is, if sample i is deleted we set the corresponding provenance token  $p_i = 0_{\text{prov}} \in \mathbb{N}[T]$  (has only a term of degree zero which is the natural number 0). The challenge, as detailed in the following section is how to do this efficiently throughout the gradient descent.

For the samples that remain we obtain (after we stop the iterations) a provenance-annotated expression that can be put in the form  $W = \sum \mathfrak{m}_k * \mathbf{u}_k$  where  $\mathfrak{m}_k$  is a *monomial* in the provenance tokens and each  $\mathbf{u}_k$  is a vector of contributions to the model parameters. To get the updated vector of model parameters we set each remaining provenance token to  $1_{\text{prov}}$  obtaining  $\mathbf{w}^{\text{upd}} = \sum \mathbf{u}_k$ . And, as promised, we notice that when all the provenance tokens are set to  $0_{\text{prov}}$  or  $1_{\text{prov}}$  the provenance expression  $\mathcal{P}^{(t)}$  comes down to an integer. Denoting this integer by  $B_U^{(t)}$  and denoting the set of the indexes of the removed training samples by  $\mathcal{R}$ , the provenance-annotated update rule for  $W^{(t+1)}$  becomes:

$$\mathcal{W}_{U}^{(t+1)} \leftarrow \left[ (1 - \eta_{t} \lambda) (1_{\text{prov}} * \mathbf{I}) - \frac{2\eta_{t}}{B_{U}^{(t)}} \sum_{\substack{i \in \mathcal{B}^{(t)} \\ i \notin \mathcal{R}}} p_{i}^{2} * \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right] \mathcal{W}_{U}^{(t)} + \frac{2\eta_{t}}{B_{U}^{(t)}} \sum_{\substack{i \in \mathcal{B}^{(t)} \\ i \notin \mathcal{R}}} p_{i}^{2} * \mathbf{x}_{i} y_{i}$$
(8)

## 4.2 Linearization for logistic regression

The model in [52] supports tracking provenance through matrix addition and multiplication. In order to apply it to GBM for logistic expression, we linearize, using *piecewise linear interpolation*, the non-linear operations in the corresponding update rules, i.e. Equation 6.

In Equation 6, the non-linear operations can be abstracted as  $f(x) = 1 - \frac{1}{1+e^{-x}}$ , where the value of the product  $y_i \mathbf{w}^{(t)T} \mathbf{x}_i$  is assigned to the variable x in Equation 6. Then f(x) can be approximated by applying 1-D piecewise linear interpolation [31]. So for each  $\mathbf{x}_i$  and  $\mathbf{w}^{(t)}$ ,  $f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)$  can be approximated by  $s(y_i \mathbf{w}^{(t)T} \mathbf{x}_i) = a^{i,(t)} y_i \mathbf{w}^{(t)T} \mathbf{x}_i + b^{i,(t)}$ , where

 $a^{i,(t)}$  and  $b^{i,(t)}$  are the linear coefficients produced by the linearizations, which depends on which sub-interval (defined by piecewise linear interpolation) the value of  $y_i \mathbf{w}^{(t)T} \mathbf{x}_i$  locates and thus should be varied between different  $\mathbf{x}_i$  and different  $\mathbf{w}^{(t)T}$  (see the associated superscript).

Throughout the paper, we will consider the case in which the variable x in f(x) is defined within an interval [-a, a] (a = 20) that is equally partitioned into  $10^6$  sub-intervals; for x outside [-a, a], we assume that s(x) is a constant since when |x| > a, the value of f(x) is very close to its bound (0 or 1). We will show that the length of each sub-interval influences the approximation rate.

In terms of multinomial logistic regression, the non-linear operations in its update rule is the softmax function, which is a vector-valued function and thus requires piecewise linear interpolation in multiple dimensions, which can be achieved by using the interpolation method proposed in [51].

After the interpolation step over the update rules for binary logistic regression, Equation 6 is approximated as:

$$\mathbf{w}_{L}^{(t+1)} \approx \left[ (1 - \eta_{t} \lambda) \mathbf{I} + \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} a^{i,(t)} \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right] \mathbf{w}_{L}^{(t)}$$

$$+ \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} b^{i,(t)} y_{i} \mathbf{x}_{i}$$
(9)

in which  $\mathbf{w}_L^{(t)}$  represents the model parameter after linearization at  $t^{\text{th}}$  iteration. By annotating each training sample  $\mathbf{x}_i$  with provenance token  $p_i$  and by taking the similar derivation of Equation 8, after the removal of the subset of training samples the provenance expression becomes:

$$\mathcal{W}_{LU}^{(t+1)} \leftarrow \left[ (1 - \eta_t \lambda)(1_{\text{prov}} * \mathbf{I}) + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathcal{R}} p_i^2 * (a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T) \right] \mathcal{W}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathcal{R}} p_i^2 * (b^{i,(t)} y_i \mathbf{x}_i)$$

$$(10)$$

By setting all the  $p_i$  in Equation 10 as  $1_{prov}$ , we can get the update rule for the updated model parameter  $\mathbf{w}_{LU}^{(t)}$ , i.e.:

$$\mathbf{w}_{LU}^{(t+1)} \approx \left[ (1 - \eta_t \lambda) \mathbf{I} + \frac{\eta_t}{B_U^{(t)}} \sum_{\substack{i \in \mathscr{B}^{(t)} \\ , i \notin \mathcal{R}}} a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T \right] \mathbf{w}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{\substack{i \in \mathscr{B}^{(t)} \\ , i \notin \mathcal{R}}} b^{i,(t)} y_i \mathbf{x}_i$$
(11)

# 4.3 Convergence analysis for provenance-annotated iterations

One concern in using GBM is whether the model parameters ultimately converge. This has been extensively studied in the machine learning community [3, 29, 36, 48, 49]. In [3], convergence conditions have been provided for GD and SGD over strong convex objective functions. Those convergence conditions can exactly fit linear regression and logistic

regression with L2-regularization because their objective functions are strong convex.

A similar concern occurs when GBM is coupled with provenance, i.e. whether the provenance expression  $\mathcal{W}_{U}^{(t)}$  in Equation 8 and  $\mathcal{W}_{LU}^{(t)}$  in Equation 10 converge in the case when the original model parameter  $\mathbf{w}^{(t)}$  converges. We propose the following definition for the convergence of provenance-annotated expressions.

*Definition 1.* **Convergence of provenance-annotated expressions.** The expression  $W^{(t)} = \sum_i \mathfrak{p}_i^{(t)} * \mathbf{u}_i^{(t)}$  converges when  $t \to \infty$  iff every matrix  $\mathbf{u}_i^{(t)}$  converges when  $t \to \infty$ .

As mentioned before, we hope that the convergence of  $W_U^{(t)}$  and  $W_{LU}^{(t)}$  can be achieved when  $\mathbf{w}^{(t)}$  can converge. The convergence conditions of  $\mathbf{w}^{(t)}$  are presented below:

LEMMA 1. Convergence conditions for general mb-SGD. [3] Given an objective function  $h(\mathbf{w})$ , which is L-Lipschitz continuous and  $\lambda$ -strong convex once the learning rate  $\eta_t$  satisfies: 1)  $\eta_t < \frac{1}{L}$ ; 2)  $\eta_t$  is a constant across all the iterations (denoted by  $\eta$ ), then  $\mathbf{w}^{(t)}$  converges when mb-SGD is used.

Unfortunately, our theoretical analysis shows that there is no convergence guarantee for  $\mathcal{W}_{U}^{(t)}$  and  $\mathcal{W}_{LU}^{(t)}$  under the convergence conditions from Lemma 1, i.e.:

Theorem 2.  $W_U^{(t)}$  in Equation 8 and  $W_{LU}^{(t)}$  in Equation 10 need not converge under the conditions in Lemma 1.  $^2$ 

However,  $W_U^{(t)}$  in Equation 8 and  $W_{LU}^{(t)}$  in Equation 10 converge under the conditions in Lemma 1 with one more assumption about the provenance expression, i.e.:

THEOREM 3. The expectation of  $W_U^{(t)}$  in Equation 8 and of  $W_{LU}^{(t)}$  in Equation 10, converge when  $t \to \infty$  if we also assume that provenance polynomial multiplication is idempotent.

Intuitively speaking, the assumption of multiplication idempotence for provenance polynomials means that we do not track *multiple joint uses of the same data sample*, which is not problematic for deletion propagation.

# 4.4 Accuracy analysis for linearized logistic regression

The next question is whether the approximated model parameters after linearization of Equation 6 (i.e.  $\mathbf{w}_L^{(t)}$  in Equation 9) is close enough to the real model parameters from Equation 6. By following the approximation property of piecewise linear interpolation, we can prove that the distance between  $\mathbf{w}^{(t)}$  and  $\mathbf{w}_L^{(t)}$  is very small.

 $<sup>^2</sup>$ Due to space limitations the proofs of the theorems are omitted. They will appear in the full version of the paper.

Theorem 4.  $||E(\mathbf{w}^{(t)} - \mathbf{w}_L^{(t)})||_2$  is bounded by  $O((\Delta x)^2)$  where  $\Delta x$  is an arbitrarily small value representing the length of the longest sub-interval used in piecewise linear interpolations.

Furthermore, in terms of the updated model parameters for logistic regression, we also need to guarantee that the updated parameters  $\mathbf{w}_{LU}^{(t)}$  are close to the real updated model parameters without linearization (denoted by  $\mathbf{w}_{RU}$ ), i.e.:

$$\mathbf{w}_{RU}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}_{RU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{\substack{i \in \mathscr{B}^{(t)} \\ i \notin \mathcal{R}}} y_i \mathbf{x}_i f(y_i \mathbf{w}_{RU}^{(t)} \mathbf{x}_i)$$
(12)

Recall that  $f(x) = 1 - \frac{1}{1 + e^{-x}}$ . Note that the linear coefficients  $a^{i,(t)}$  and  $b^{i,(t)}$  in Equation 11 are actually derived in the training phase where all samples exist (rather than in the model update phase), which implies that a larger difference between  $\mathbf{w}_{LU}^{(t)}$  and  $\mathbf{w}_{RU}^{(t)}$  should be expected. Surprisingly, we can prove that the distance between  $\mathbf{w}_{LU}^{(t)}$  and  $\mathbf{w}_{RU}^{(t)}$  is still small enough.

THEOREM 5.  $||E(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$ , where  $\Delta n$  is the number of the removed samples and  $\Delta x$  is defined in Theorem 4.

### 5 IMPLEMENTATION

We now discuss how the ideas in the previous section are implemented in PrIU and PrIU-opt, for both linear and logistic regression. Along the way, time and space complexity analyses, as well as theorems that justify our approximation strategies used in PrIU and PrIU-opt, are provided.

### 5.1 PrIU: Linear regression

In Equation 8, by setting all the  $p_i$  as  $1_{\text{prov}}$ , the expression  $\sum_{i \in \mathscr{B}^{(t)}, i \notin \mathscr{R}} p_i^2 * \mathbf{x}_i \mathbf{x}_i^T$  becomes  $\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T - \sum_{i \in \mathscr{B}^{(t)}, i \in \mathscr{R}} \mathbf{x}_i \mathbf{x}_i^T$ , in which the first term,  $\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T$ , can be regarded as provenance information and thus cached as an intermediate result for each mini-batch during the training phase for the original model parameter  $\mathbf{w}^{(t)}$ . Thus we only need to compute the latter term during the incremental update phase.  $\sum_{i \in \mathscr{B}^{(t)}, i \notin \mathscr{R}} \mathbf{x}_i y_i$  can be computed in a similar way. In the end, Equation 8 is then rewritten as follows for the purpose of incremental updates:

$$\mathbf{w}_{U}^{(t+1)} \leftarrow \left[ (1 - \eta_{t} \lambda) \mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}} \sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right]$$

$$- \sum_{\substack{i \in \mathscr{B}^{(t)} \\ i \in \mathscr{O}}} \mathbf{x}_{i} \mathbf{x}_{i}^{T} \mathbf{w}_{U}^{(t)} + \frac{2\eta_{t}}{B_{U}^{(t)}} (\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i} y_{i} - \sum_{i \in \mathscr{B}^{(t)}, i \in \mathscr{R}} \mathbf{x}_{i} y_{i})$$

$$(13)$$

Note that  $\sum_{i \in \mathscr{B}^{(t)}, i \in \mathscr{R}} \mathbf{x}_i \mathbf{x}_i^T$  can be rewritten into matrix form, i.e.  $\Delta \mathbf{X}_{\mathscr{B}^{(t)}}^T \Delta \mathbf{X}_{\mathscr{B}^{(t)}}$  where  $\Delta \mathbf{X}_{\mathscr{B}^{(t)}}$  is a matrix consisting of the removed samples in the mini-batch  $\mathscr{B}^{(t)}$ . The associativity

property of matrix multiplication can also be used to avoid expensive matrix-matrix multiplications (i.e.  $\Delta \mathbf{X}_{\mathscr{B}^{(t)}}^T \Delta \mathbf{X}_{\mathscr{B}^{(t)}})$  by conducting more efficient matrix-vector multiplications instead (e.g. computing  $\Delta \mathbf{X}_{\mathscr{B}^{(t)}} \mathbf{w}_U^{(t+1)}$  first and then multiplying the result by  $\Delta \mathbf{X}_{\mathscr{B}^{(t)}}^T$ ).

Suppose that  $\Delta B$  samples are removed from each minibatch on average, then the time complexity of updating the model parameters in each iteration using Equation 13 will be  $O(\Delta Bm + m^2)$  (recall that the dimension of  $\mathbf{X}$  is  $n \times m$ ). In contrast, the time complexity for retraining from scratch (i.e. not caching  $\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T$  in Equation 13) will be  $O((B - \Delta B)m)$ . Of course, performance predictions based on asymptotic complexity give only very rough guidance, and we conduct experiments for realistic assessments. Still, the bounds above suggest, for example, that for small  $\Delta B$  and  $m \ll B$  incremental deletions with PrIU work better than retraining (and our experiments verify this, see Section 6).

Typically, however, a smaller mini-batch size B is used. To deal with the case in which m > B, we notice that the rank of the intermediate result,  $\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T$  should be no more than B, thus smaller than m when B is smaller than m. This motivates us to reduce the dimension of the intermediate results using SVD, i.e.  $\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T = \mathbf{U}^{(t)} \mathbf{S}^{(t)} \mathbf{V}^{T,(t)}$ , where  $\mathbf{S}^{(t)}$  is a diagonal matrix whose diagonal elements represent the singular values, while  $\mathbf{U}^{(t)}$  and  $\mathbf{V}^{(t)}$  are the left and right singular vectors. Suppose after SVD, we only keep the r largest singular values and the corresponding singular vectors where  $r \ll B$ , then  $\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T$  is approximated by  $\mathbf{U}_{1...r}^{(t)}, \mathbf{S}_{1...r}^{(t)} \mathbf{V}_{1...r}^{(t)T}$  ( $\mathbf{U}_{1...r}^{(t)}, \mathbf{V}_{1...r}^{(t)}$  represents the submatrix composed of the first r columns and  $\mathbf{S}_{1...r}^{(t)}$  is a diagonal matrix composed of the first r eigenvalues in  $\mathbf{S}^{(t)}$ ). Thus Equation 13 is rewritten as:

$$\mathbf{w}_{U}^{(t+1)} \leftarrow \left[ (1 - \eta_{t} \lambda) \mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}} (\mathbf{U}_{1...r}^{(t)} \mathbf{S}_{1...r}^{(t)} \mathbf{V}_{1...r}^{(t)T} \right]$$

$$- \Delta \mathbf{X}_{\mathscr{B}^{(t)}}^{T} \Delta \mathbf{X}_{\mathscr{B}^{(t)}}) \mathbf{w}_{U}^{(t)} + \frac{2\eta_{t}}{B_{U}^{(t)}} (\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i} y_{i} - \sum_{\substack{i \in \mathscr{B}^{(t)} \\ i \in \mathscr{R}}} \mathbf{x}_{i} y_{i})$$

$$(14)$$

Here we can cache the results of  $\mathbf{U}_{1...r}^{(t)}\mathbf{S}_{1...r}^{(t)}$  (denoted by  $\mathbf{P}_{1...r}^{(t)}$ ) and  $\mathbf{V}_{1...r}^{(t)}$  for efficient updates, both of which have dimensions  $m \times r$ .

**Time complexity.** The time complexity to update the model parameters using this approach is  $O(rm + \Delta Bm)$  for each iteration since the computation time is dominated by the matrix-vector computation, e.g. the multiplications of  $\mathbf{V}_{1..r}^{T,(t)}$  and  $\mathbf{w}_{U}^{(t)}$ . This is more efficient than retraining from scratch, which has time complexity  $O((B - \Delta B)m)$ . So the total complexity for PrIU is  $O(\tau rm + \tau \Delta B)$ , where  $\tau$  is the number of iterations in the training phase.

**Space complexity.** Using this approximation, at each iteration we only need to cache  $\mathbf{P}_{1..r}^{(t)}$  and  $\mathbf{V}_{1..r}^{(t)}$ , which require space O(rm). So the total space complexity will be  $O(\tau rm)$  for  $\tau$  iterations.

Theorem 6. Approximation ratio Under the convergence conditions for  $\mathbf{w}^{(t)}$ ,  $||\mathbf{w}^{(t)}||$  should be bounded by some constant C. Suppose  $\frac{||\mathbf{U}_{1..r}^{(t)}\mathbf{S}_{1..r}^{(t)}\mathbf{V}_{1..r}^{T,(t)}||_2}{||\mathbf{U}^{(t)}\mathbf{S}^{(t)}\mathbf{V}^{T,(t)}||_2} \geq 1 - \epsilon$  where  $\epsilon$  is a small value, then the change of model parameters caused by the approximation will be bounded by  $O(\epsilon)$ .

This shows that with proper choice of r in the SVD approximation, the updated model parameters computed by PrIU or PrIU-opt should be still very close to the expected result. So in our implementations, r is chosen based on  $\epsilon$  (say 0.01) such that the inequality in Theorem 6 is satisfied.

# 5.2 PrIU-opt: Optimizations for linear regression

When the feature space is small, additional optimizations can be used for linear regression. Note that according to [29], the model parameters derived by both SGD and mb-SGD will end up with statistically the same results as GD. This means that the update rule in Equation 5 and Equation 13 could be approximated by its alternative using GD, i.e.:

$$\mathbf{w}^{(t+1)} \leftarrow ((1 - \eta_t \lambda)\mathbf{I} - \frac{2\eta_t}{n}\mathbf{X}^T\mathbf{X})\mathbf{w}^{(t)} + \frac{2\eta_t}{n}\mathbf{X}^T\mathbf{Y}$$
(15)
$$\mathbf{w}_U^{(t+1)} \leftarrow ((1 - \eta_t \lambda)\mathbf{I} - \frac{2\eta_t}{n - \Delta n}(\mathbf{X}^T\mathbf{X} - \Delta \mathbf{X}^T\Delta \mathbf{X}))\mathbf{w}_U^{(t)}$$

$$+ \frac{2\eta_t}{n - \Delta n}(\mathbf{X}^T\mathbf{Y} - \Delta \mathbf{X}^T\Delta \mathbf{Y})$$
(16)

in which  $(\Delta \mathbf{X}, \Delta \mathbf{Y})$  represent the removed samples while  $\Delta n$  represents the number of those samples. Let  $\mathbf{M}$  and  $\mathbf{N}$  denote  $\mathbf{X}^T\mathbf{X}$  and  $\mathbf{X}^T\mathbf{Y}$  respectively. Then eigenvalue decomposition can be applied over  $\mathbf{M}$ , i.e.  $\mathbf{M} = \mathbf{Q}$  diag  $(\{c_i\}_{i=1}^n)$   $\mathbf{Q}^{-1}$  (where  $c_i$  represents the eigenvalues of  $\mathbf{M}$ ). This is then plugged into Equation 15 and computed recursively, which results in the following formula:

$$\mathbf{w}^{(t+1)} = \mathbf{Q} \ diag \ (\{\Pi_{j=1}^{t} (1 - \eta_{j}\lambda - \frac{2\eta_{j}}{n} c_{i})\}_{i=1}^{n}) \ \mathbf{Q}^{-1} \mathbf{w}^{(0)}$$

$$+ \mathbf{Q} diag (\sum_{l=1}^{t-1} \eta_{l} \{\Pi_{j=l+1}^{t} (1 - \eta_{j}\lambda - \frac{2\eta_{j}}{n} c_{i})\}_{i=1}^{n}) \ \mathbf{Q}^{-1} \frac{2\mathbf{N}}{n}$$

$$(17)$$

This indicates that once the eigenvalues and eigenvectors of each  ${\bf M}$  are given, we can derive  ${\bf w}^{(t)}$  by simply computing the product,  $\Pi_{j=1}^t(1-\eta_j\lambda-\frac{2\eta_j}{n}c_i)$ , and the sum of the product,  $\sum_{l=1}^{t-1}\eta_l\Pi_{j=l+1}^t(1-\eta_j\lambda-\frac{2\eta_j}{n}c_i)$ , on diagonal entries. The overhead of this is only  $O(\tau m)$  (recall that  $\tau$  represents the total iteration number), and thus we avoid the repetitive matrix multiplication operations through the for-loops. Also, observe that  ${\bf M}'={\bf X}^T{\bf X}-\Delta{\bf X}^T\Delta{\bf X}$  can be regarded as a small

change over  $\mathbf{M}$  when  $\Delta n$  is small. Thus we can use the results on incremental updates over eigenvalues in [41], i.e. when the difference between the eigenvectors of  $\mathbf{M}'$  and that of  $\mathbf{M}$  is negligible then the eigenvalues of  $\mathbf{M}'$  are estimated as:

$$\mathbf{Q}^{-1}\mathbf{M}'\mathbf{Q} = diag(\{c_i'\}_{i=1}^n)$$
 (18)

Here,  $c'_i$  represents the approximated  $i^{th}$  eigenvalue of  $\mathbf{M}'$ . It indicates that we can apply eigenvalue decomposition over  $\mathbf{M}$  offline before the model incremental update phase and use Equation 18 to get the updated eigenvalues online.

**Time complexity.** The time complexity for updating the model parameters is dominated by the computation of  $c_i'$ , which is followed by the computation over each  $c_i'$  as Equation 17 does. These have time complexities  $O(\min\{\Delta n, m\}m^2)$  and  $O(\tau m)$ , respectively. So the total time complexity is  $O(\min\{\Delta n, m\}m^2) + O(\tau m)$ , which can be more efficient than the closed-form solution (see experiments in Section 6).

**Space complexity.** The method avoids caching the provenance information at each iteration, and only requires caching Q,  $Q^{-1}$  and all the eigenvalues  $c_i$ , which takes space  $O(m^2)$ 

THEOREM 7. (Approximation ratio) The approximation of PrIU-opt over the model parameters is bounded by  $O(||\Delta X^T \Delta X||)$ 

This shows that with small number of removed samples, the approximation ratio should be very small.

## 5.3 PrIU: Logistic regression

As the first step of the implementation of PrIU for logistic regression, non-linear operations are linearized using piecewise linear interpolation. Then, based on the analysis in Section 4, given the ids of the samples to be removed in dense datasets,  $\mathcal{R}$ , Equation 11 is rewritten as follows:

$$\mathbf{w}_{LU}^{(t+1)} \leftarrow \left[ (1 - \eta_t \lambda) \mathbf{I} + \frac{\eta_t}{B_U^{(t)}} (\mathbf{C}^{(t)} - \Delta \mathbf{C}^{(t)}) \right] \mathbf{w}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} (\mathbf{D}^{(t)} - \Delta \mathbf{D}^{(t)})$$

$$(19)$$

where  $\mathbf{C}^{(t)}$ ,  $\mathbf{D}^{(t)}$ ,  $\Delta \mathbf{C}^{(t)}$ ,  $\Delta \mathbf{D}^{(t)}$  are:

$$\mathbf{C}^{(t)} = \sum_{i \in \mathcal{B}^{(t)}} a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T, \Delta \mathbf{C}^{(t)} = \sum_{i \in \mathcal{R}, i \in \mathcal{B}^{(t)}} a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T$$

$$\mathbf{D}^{(t)} = \sum_{i \in \mathcal{B}^{(t)}} b^{i,(t)} y_i \mathbf{x}_i, \Delta \mathbf{D}^{(t)} = \sum_{i \in \mathcal{R}, i \in \mathcal{B}^{(t)}} b^{i,(t)} y_i \mathbf{x}_i$$

Similar to linear regression, the intermediate results  $\mathbf{C}^{(t)}$  and  $\mathbf{D}^{(t)}$  are cached and the dimension of  $\mathbf{C}^{(t)}$  can be reduced by using SVD before the model update phase, which can happen offline. Suppose after SVD,  $\mathbf{C}^{(t)} \approx \mathbf{P}_{1..r}^{(t)} \mathbf{V}_{1..r}^{T,(t)}$ , in which  $\mathbf{P}_{1..r}^{(t)}$  and  $\mathbf{V}_{1..r}^{(t)}$  are two matrices with dimension  $m \times r$ .

In the end, Equation 19 is modified as below for incremental model updates:

$$\mathbf{w}_{LU}^{(t+1)} \leftarrow \left[ (1 - \eta_t \lambda) \mathbf{I} + \frac{\eta_t}{B_U^{(t)}} (\mathbf{P}_{1..r}^{(t)} \mathbf{V}_{1..r}^{T,(t)} - \Delta \mathbf{C}^{(t)}) \right] \mathbf{w}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} (\mathbf{D}^{(t)} - \Delta \mathbf{D}^{(t)})$$
(20)

Time complexity. To apply Equation 20 in the model update phase, the computation of  $\mathbf{P}_{1...r}^{(t)}\mathbf{V}_{1...r}^{T,(t)}\mathbf{w}_{LU}^{(t)}$  and  $\Delta\mathbf{C}^{(t)}\mathbf{w}_{LU}^{(t)}$  become the major overhead, which have time complexity O(rm) and  $O(\Delta Bm)$ , respectively. Suppose there are  $\tau$  iterations in total, then the total time complexity is  $O(\tau(rm+\Delta Bm))$ . In comparison, the time complexity of retraining from scratch is  $O(\tau((B-\Delta B)m+C_{non}m))$ , where  $C_{non}$  represents the overhead of the non-linear operations. When  $r\ll B$  and  $\Delta B\ll B$ , we can therefore expect PrIU to be more efficient than retraining from scratch.

**Space complexity analysis** Through this approximation, we need to cache  $\mathbf{P}_{1...r}^{(t)}$  and  $\mathbf{V}_{1...r}^{(t)}$  at each iteration, which requires  $O(\tau rm)$  space in total. Plus,  $O(n\lceil \frac{\tau B}{n} \rceil)$  extra space is necessary to cache the linear coefficients. So the total space complexity will be  $O(\tau rm) + O(n\lceil \frac{\tau B}{n} \rceil)$ .

THEOREM 8. (Approximation ratio) Similar to Theorem 6, the deviation caused by the SVD approximation will be bounded by  $O(\epsilon)$ , given the ratio  $\frac{||P_{1..r}^{(t)}V_{1..r}^{T,(t)}||_2}{||P^{(t)}V_{1..t}^{T,(t)}||_2} \geq 1-\epsilon$ . So using Theorem 5,  $||E(\boldsymbol{w}_{LU}^{(t)}-\boldsymbol{w}_{RU}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n}\Delta x)+O((\frac{\Delta n}{n})^2)+O((\Delta x)^2)+O(\epsilon)$ .

This indicates that  $\mathbf{w}_{LU}^{(t)}$  should be very close to  $\mathbf{w}_{RU}^{(t)}$  (similar to the discussion after Theorem 6).

**Discussion** Notice that for sparse datasets with large feature space, we can utilize the efficient sparse matrix operations by retraining from scratch. Also note that the intermediate result  $\mathbf{C}^{(t)}$  will be a sparse matrix for such datasets. However, after SVD, there is no guarantee that  $\mathbf{P}^{(t)}$  and  $\mathbf{V}^{(t)}$  are sparse matrices. Therefore, for sparse training datasets, we will simply use the linearized update rule, i.e. Equation 11 directly, without considering the strategies above.

# 5.4 PrIU-opt: Optimizations for logistic regression

Again, when the feature space is small additional optimizations are possible. In particular, we observe that for each sample i the change in the coefficients  $a^{i,(t)}$  and  $b^{i,(t)}$  from one iteration to the next becomes smaller and smaller as  $\mathbf{w}^{(t)}$  converges. This suggests that we can stop capturing new provenance information at some earlier iteration, call it  $t_s$ , and continue with the same provenance until convergence. Suppose that for each sample i we approximate  $a^{i,(t)}$ ,  $b^{i,(t)}$  by  $a^{i,*}$  and  $b^{i,*}$  after the iteration  $t_s$ . Therefore the matrices  $\mathbf{C}^{(t)}$ ,  $\mathbf{D}^{(t)}$ ,  $\Delta \mathbf{C}^{(t)}$  and  $\Delta \mathbf{D}^{(t)}$  will be approximated using the

coefficients  $a^{i,*}$  and  $b^{i,*}$  and will remain the same for all iterations  $t \ge t_s$  allowing us to avoid their recomputation. In the experiments, we found that a rule of thumb that takes  $t_s$  to be 70% of the total number of iterations works well.

This has the same form as for linear regression, motivating us to use the same techniques from PrIU-opt for linear regression, i.e. conducting eigenvalue decomposition over  $\mathbf{C}^{(t)}$ , followed by incrementally updating the eigenvalues given the changes  $\Delta \mathbf{C}^{(t)}$ , thus avoiding recomputations after the iteration  $t_s$ .

**Time complexity.** Before and after the iteration  $t_s$ , the total time complexity is  $O(t_s(rm+\Delta Bm))$  and  $O(\min\{\Delta n,m\}m^2)+O((\tau-t_s)m)$  (see the time complexity analysis in Section 5.2) respectively. Thus the total time complexity is  $O(t_s(rm+\Delta Bm))+O(\min\{\Delta n,m\}m^2)+O((\tau-t_s)m)$ .

**Space complexity.** After the iteration  $t_s$ , we only need to keep the eigenvectors of  $\mathbb{C}^{(t)}$ , which requires  $O(m^2)$  space. Including the space overhead for the first  $t_s$  iterations, the total space complexity is  $O(m^2) + O(t_s rm) + O(n \lceil \frac{t_s B}{n} \rceil)$ .

Theorem 9. (Approximation ratio) Suppose that after the iteration  $t_s$  the gradient of the objective function is smaller than  $\delta$ , then the approximations of PrIU-opt can lead to deviations of the model parameters bounded by  $O((\tau - t_s)\delta) + O(||\Delta X^T \Delta X||)$ . By combining the analysis in Theorem 5,  $||E(\boldsymbol{w}_{LU}^{(t)} - \boldsymbol{w}_{RU}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2) + O((\tau - t_s)\delta) + O(||\Delta X^T \Delta X||)$ 

This thus indicates that  $\mathbf{w}_{LU}^{(t)}$  should be very close to  $\mathbf{w}_{RU}^{(t)}$ . **Discussion.** Our current framework handles linear and logistic models with L2 regularization. Our solutions cannot handle L1 regularization since in this case the gradient of the objective function is not continuous, thus invalidating some of the error bound analysis above. How to handle L1 regularization will be our future work.

### **6 EXPERIMENTS**

### 6.1 Experimental setup

**Platform.** We conduct extensive experiments in Python 3.6 and use PyTorch 1.3.0 [42] for the experiments for dense datasets and scipy 1.3.1 [27] for the experiments for sparse datasets. All experiments were conducted on a Linux server with an Intel(R) Xeon(R) CPU E5-2630 v4 @ 2.20GHz and 64GB of main memory.

**Datasets.** Six datasets were used in our experiments: (1) the UCI SGEMM GPU dataset<sup>3</sup>; (2) the UCI Covtype dataset <sup>4</sup>; (3) the UCI HIGGS dataset <sup>5</sup>; (4) the RCV1 dataset <sup>6</sup> (5)

 $<sup>^3</sup> https://archive.ics.uci.edu/ml/datasets/SGEMM+GPU+kernel+performance$ 

<sup>&</sup>lt;sup>4</sup>https://archive.ics.uci.edu/ml/datasets/covertype

 $<sup>^5</sup> https://archive.ics.uci.edu/ml/datasets/HIGGS \\$ 

<sup>&</sup>lt;sup>6</sup>simplified version from https://scikit-learn.org/0.18/datasets/rcv1.html

the Kaggle ECG Heartbeat Categorization Dataset<sup>7</sup>; (6) the CIFAR-10 dataset <sup>8</sup>, which are referenced as SGEMM, Cov, HIGGS, RCV1, Heartbeat and cifar10 hereafter.

SGEMM has continuous label values, therefore we use it in experiments with *linear regression* while the rest of them have values that are appropriate for classification. Each dataset is partitioned into *training* (90% of the samples) and *validation* (10% of the samples) datasets, the latter used for measuring the accuracy of models trained from the former.

The characteristics of these datasets are listed in Table 1, which indicates that RCV1 and cifar10 have extremely large feature space (over 30k model parameters) while other datasets have much fewer parameters (Heartbeat has around 1000 while others have less than 500).

### 6.2 Experiment design

We conduct two sets of experiments, the first of which aims to evaluate the performance of PrIU and PrIU-opt with respect to the deletion of *one subset* of the training samples. We do this over different types of datasets (dense VS sparse, large feature space VS small) with varied configurations (how many samples to be removed, mini-batch size, iteration numbers etc.), and compare against retraining from scratch. The second set of experiments simulate the scenario where users *repetitively* remove different subsets of training samples.

In the first set of experiments we simulate the cleaning scenario. To specify the samples to be removed from the training datasets, we introduce dirty samples, which are a selected subset of samples from the original dataset  $\mathcal T$  that are modified to incorrect values by rescaling. The resulting dataset is denoted  $\mathcal T_{\text{dirty}}$ , over which the initial model  $\mathcal M_{\text{init}}$  is constructed. The dirty samples are then removed in the model update phase. The goal is to compare the robustness of PrIU, PrIU-opt and the *influence function* [30] method when dirty data exists. In the experiment we vary the number of erroneous samples generated. The ratio between the erroneous samples and the original training dataset is called the **deletion rate**, and we give it values ranging from 0.0001 (i.e. 0.01%) to 0.2 (i.e. 20%).

In the second set of the experiments, we simulate the scenario in which users debug or interpret models by removing different subsets of samples, necessitating repeated incremental model update operations. We assume that the datasets are very large; to simulate this, we create three synthetic datasets  $\mathcal{T}_{\text{cat}}$  by concatenating 4 copies of HIGGS, 20 copies of Cov and 130 copies of Heartbeat such that the total number of training samples is around 40 million, 11 million and 11 million, respectively, which are denoted HIGGS (extended), Cov (extended) and Heartbeat (extended), respectively. In the

experiments, ten different subsets are removed and for each of them the deletion rate is about 0.1% of randomly picked samples out of the full training set. The hyperparameters for this set of experiments are listed in Table 2.

**Baseline.** For both  $\mathcal{T}_{dirty}$  and  $\mathcal{T}_{cat}$ , we simulate what users (presumably unaware of errors) would do, and train an initial model  $\mathcal{M}_{init}$  using the following standard method: Manually derive the formula for the gradient of the objective function and then program explicitly the GBM iterations. The erroneous or chosen samples are then removed from  $\mathcal{T}_{dirty}$  or  $\mathcal{T}_{cat}$ . For linear regression (except for GBM), we also compare PrIU and PrIU-opt against close-form formula solutions for incremental updates [13, 22, 23, 40], denoted by Closed-form.

**Incrementality.** To update the model  $\mathcal{M}_{\text{init}}$ , the straightforward solution is to retrain from the scratch by using the same standard method as before but exclude the removed samples from each mini-batch. We denote this solution by BaseL. In contrast, our approach uses PrIU or PrIU-opt to incrementally update the model. The time taken by BaseL, PrIU or PrIU-opt to produce the updated model is reported in the experiments as the *update time*, and is compared over the two solutions: retraining with BaseL vs. incremental update with PrIU or PrIU-opt.

Note that our PrIU/ PrIU-opt approach uses provenance information collected from the whole training dataset. This phase is *offline* for the PrIU/ PrIU-opt algorithms and is *not* included in their reported running times. In practice, for the first set of experiments (cleaning of erroneous samples) provenance collection is done during the training of  $\mathcal{M}_{\text{init}}$  from  $\mathcal{T}_{\text{dirty}}$ . For the second set of experiments (repeated deletions of subsets for debugging or interpretability) provenance collection is done during an initial training of  $\mathcal{M}_{\text{init}}$  from the entire dataset  $\mathcal{T}_{\text{cat}}$ , which only needs to be done *once* even if many deletions of subsets are performed subsequently.

Since PrIU-opt is the optimized version for datasets with small feature space we only record the update time of PrIU over RCV1 and cifar10, which have very large feature spaces.

**Accuracy.** We compare the quality of the updated model obtained by BaseL and PrIU/PrIU-opt. The goal is to show that the improvement in update time is not achieved at the expense of accuracy. For experiments with linear regression, we use the *mean squared error (MSE)* over the validation datasets as a measure for accuracy. A lower MSE corresponds to higher accuracy over the validation set. For experiments with binary or multinomial logistic regression, we use the updated model to classify the samples in the validation datasets and report their *validation accuracy*.

**Model comparison.** We also compare the updated models *structurally* by comparing the vector of updated model parameters obtained via PrIU/PrIU-opt against the ones by using BaseL. This is done in two different ways: 1) Using *distance*, that is, the *L2-norm* of the difference between the

<sup>&</sup>lt;sup>7</sup>https://www.kaggle.com/shayanfazeli/heartbeat

<sup>&</sup>lt;sup>8</sup>https://www.cs.toronto.edu/~kriz/cifar.html

two vectors, for both linear and logistic regression, and 2) Using *similarity*, that is, the *cosine* of the angle between the two vectors. The latter is only done for logistic regression since the angle is only relevant for classification techniques. For both linear regression and logistic regression, we also record the changes of the signs and magnitude of individual coordinate of the updated model parameters by PrIU and PrIU-opt compared to the ones obtained by BaseL.

Comparison with influence function. As indicated in Section 2, the *influence function* method in [30] can be extended to handle the removal of multiple training samples by us (details omitted). We denote the resulting method INFL and compare it against PrIU/PrIU-opt in the experiments. We predicted and verified experimentally that this approach produces models with poor validation accuracy since the derivation of INFL relies on the approximation of the Taylor expansion, which can be inaccurate. We also notice that the Taylor expansion used in INFL involves the computation of the Hessian matrix, which is very expensive for datasets with extremely large feature space. So we did not run INFL over RCV1 and cifar10 in the experiments; the comparison between PrIU/PrIU-opt and INFL over other datasets is enough to show the benefits of our approaches.

Effect of the hyperparameters and feature space size As discussed in Section 5, the performance of PrIU and PrIUopt is influenced by the mini-batch size, the number of iterations and the size of the feature space. To explore the effect of the first two parameters for logistic regression, three different combinations of mini-batch size and number of iterations are used over Cov, denoted Cov (small), Cov (large 1) and Cov (large 2) (see Table 2). Since the datasets used for logistic regression have different feature space sizes, the performance difference with respect to feature space size is also compared. Since there is only one dataset for linear regression, SGEMM, we extend this dataset by adding 1500 random features for each sample to determine the effect of feature space size. The extended version of SGEMM is denoted SGEMM (extended) (see Table 2). Other hyperparameters used in the experiments are shown in Table 2. Note that since erroneous samples exist in the training datasets for the first set of experiments, some values of the learning rate need to be very small to make sure that the convergence can be reached.

In the experiments, we answer the following questions:

- (Q1) Do the optimizations used in PrIU-opt compared to PrIU lead to a significant improvement in update time without sacrificing accuracy when the number of features in the training set is small?
- (Q2) Do PrIU and PrIU-opt afford significant gains in efficiency compared to BaseL?
- **(Q3)** Are the efficiency gains provided by PrIU and PrIU-opt achieved without sacrificing the accuracy of the updated model?

- (Q4) Can we experimentally validate the theoretical analysis in Sections 4.4 and 5, i.e. that the updated model derived through the approximations in PrIU and PrIU-opt is very close to the one obtained by BaseL?
- **(Q5)** Does the influence function approach, INFL, provide a competitive alternative to PrIU and PrIU-opt?
- (Q6) Can we experimentally show the effect of the hyperparameters, such as mini-batch size and iteration numbers over the performance gains of PrIU and PrIU-opt?
- (Q7) Can we experimentally show the effect of the feature space size (i.e. the number of model parameters, which equals to the feature number times the number of classes for multi-nomial logistic regression)?
- **(Q8)** What is the memory overhead of PrIU and PrIU-opt for caching the provenance information?

### 6.3 Experimental results

We report the results of our experiments in this subsection.

(Q1) We compare the update time of PrIU and PrIU-opt for linear regression using SGEMM (extended) in Figure 1b. The results show that the update time of PrIU-opt is significantly better than that of PrIU except when the deletion rate is approaching 20%. We also see from Table 4 that PrIU-opt and BaseL yield models that have exactly the same validation accuracy. Therefore, although PrIU-opt uses additional

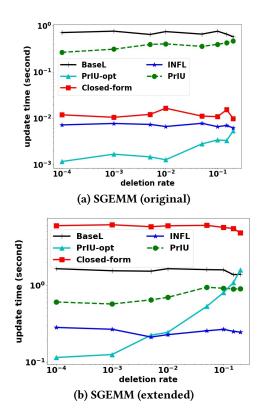


Figure 1: Update time using linear regression

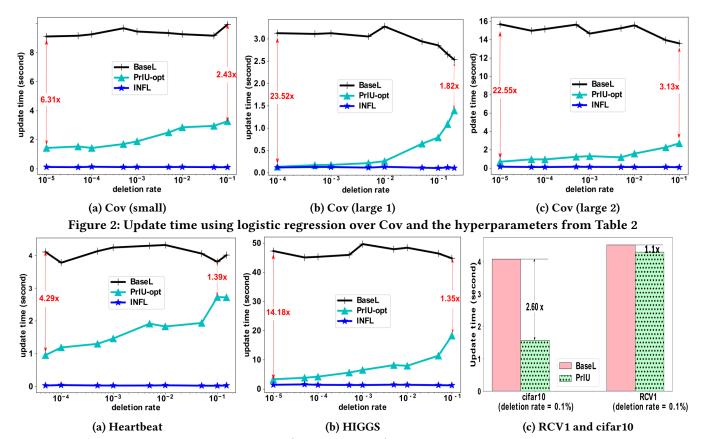


Figure 3: Update time using logistic regression

Table 1: Summary of datasets

name	# features	# classes	# samples
SGEMM	18		241,600
Cov	54	7	581,012
HIGGS	28	2	11,000,000
RCV1	47,236	2	23,149
Heartbeat	188	7	87,553
cifar10	3072	10	50,000

approximations for optimization, they do not hurt the predictive power of the updated models. This shows that the optimization strategies in Sections 5.2 and 5.4 are worth the design and implementation effort. Consequently, we will only compare PrIU-opt against other approaches except for cifar10 and RCV1 which have extremely large feature spaces.

(Q2) Figures 1a-1b compare the update time in BaseL and PrIU-opt using linear regression (ignore the INFL lines for the moment), while Figures 2-3 show the same results for logistic regression for single model update operation. Observe that for both linear and logistic regression, when the deletion rate is small (<0.01), PrIU-opt can achieve significant speed-up compared to BaseL: up to two orders of

Table 2: Summary of hyperparameters used in the experiments

	::	# of			
name	mini-		other hyper-		
	batch	iterations	parameters		
	size		$(\eta, \lambda)$		
SGEMM (original)	200	2000	$(5 \times 10^{-3}, 0.1)$		
SGEMM (extended)	200	2000	$(5 \times 10^{-3}, 0.1)$		
Cov (small)	200	10000	$(1 \times 10^{-4}, 0.001)$		
Cov (large 1)	10000	500	$(1 \times 10^{-4}, 0.001)$		
Cov (large 2)	10000	3000	$(1 \times 10^{-4}, 0.001)$		
HIGGS	2000	20000	$(1 \times 10^{-5}, 0.01)$		
Cov (extended)	1000	40000	$(1 \times 10^{-4}, 0.001)$		
HIGGS	2000	20000	$(1 \times 10^{-5}, 0.01)$		
HIGGS (extended)	2000	60000	$(1 \times 10^{-5}, 0.01)$		
Heartbeat	500	5000	$(1 \times 10^{-5}, 0.1)$		
Heartbeat	eartbeat		(110=5.0.1)		
(extended)	500	40000	$(1 \times 10^{-5}, 0.1)$		
RCV1	500	3000	$(1 \times 10^{-6}, 0.5)$		
cifar10	500	1000	(0.001, 0.1)		

magnitude for linear regression and up to around 23x for logistic regression (for Cov (large 1) and Cov (large 2) with low deletion rate). Even when the feature spaces are extremely

Dataset	BaseL	PrIU	PrIU-opt
Cov (small)	0.71	4.30	4.34
Cov (large 1)	0.87	4.02	3.49
Cov (large 2)	1.34	21.0	17.4
HIGGS	5.09	8.40	8.40
SGEMM (original)	2.43	2.45	2.48
SGEMM (extended)	4.94	6.66	5.74
Heartbeat	0.46	6.01	5.69
RCV1	0.28	0.3	-
cifar10	0.79	26.59	-

large, with deletion rate 0.1%, there is around a 2.6x speed-up for dense datasets (cifar10 in Figure 3c) and only 10% for sparse datasets (RCV1 in Figure 3c), respectively (similar speed-ups were observed for other small deletion rates). The former shows the effectiveness of the optimization strategies in PrIU over dense datasets with a large feature space while the latter is due to the fact that the optimization strategies for dense datasets were not applied over the sparse ones. Notice that for linear regression, PrIU-opt is always faster than Closed-form. Figure 4 shows the results of repetitive model updates; PrIU-opt achieves an order of magnitude speed-up for HIGGS (extended).

(Q3) Table 4 (validation accuracy for PrIU and PrIU-opt column) compares the quality of the models obtained by PrIU/PrIU-opt with that of the models obtained by BaseL. For these results we chose the highest deletion rate in the experiments, i.e. 20%. For all the experiments, the validation accuracy (MSE in the case of linear regression) of the updated models obtained by PrIU and PrIU-opt *match exactly* the accuracy of the ones obtained by BaseL. Combined with the answer to Q2, we can conclude that *PrIU-opt speeds up the model update time by up to two orders of magnitude without sacrificing any validation accuracy*.

**(Q4)** We investigate why PrIU-opt has the same validation accuracy as BaseL by measuring the distance and similarity between the updated models computed by PrIU-opt and BaseL. The results are presented in Table 4 (again, ignore the columns for INFL). The results indicate that the updated model parameters computed by PrIU-opt are very close to the ones obtained by BaseL since the cosine similarity is almost 1 (see the "similarity" column) while the L2-dist is very small (see the "distance" column). An even finer-grained analysis, comparing the signs and magnitude of each coordinate in the model parameters updated by PrIU-opt and BaseL shows that there is no sign flipping and only negligible magnitude changes for PrIU-opt compared to BaseL when the deletion rate is small. Even with a large deletion rate of 20% in HIGGS, only 2 out of 58 coordinates flip their signs with small magnitude change.

(Q5) The model update time of INFL is also included in Figures 2 and 3. Note that it can be up to one order of magnitude better than PrIU-opt, which is expected since using INFL to update the model parameters does not require an iterative computation. However, there is a significant drop in validation accuracy of the updated model derived by INFL compared to BaseL and PrIU-opt (see Table 4), which is due to the significantly higher L2-dist (see the "distance" column) and lower cosine similarity (see the "similarity" column) of its updated model compared to the model derived by BaseL. We conclude that PrIU and PrIU-opt produce much better models than INFL yet can still achieve comparable speed-ups.

(Q6) Effect of mini-batch size. The effect of mini-batch size is seen by comparing Cov (large 1) and Cov (small). One observation is that with larger mini-batch size, the maximal speed-up of PrIU-opt is around 23x, while with the smaller mini-batch size it is only about 6x, see Figures 2a and 2b This confirms the analysis in Section 5. In the second set of experiments, we used a small mini-batch size for Cov (1000) and Heartbeat (500), resulting in only 4.62x and 3.2x speed-ups by PrIU-opt, respectively (see Figure 4).

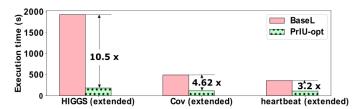


Figure 4: The execution time of repetitively removing 10 different subsets

**Effect of number of iterations.** A comparison of Cov (large 1) and Cov (large 2), which have the same mini-batch size but a different number of iterations, can be found in Figures 2b and 2c. We observe that no matter how many iterations the program runs for, at the same deletion rate PrIU-opt achieves a similar speed-up against BaseL. For example, we have up to around 23x speed-up for small deletion rates and smaller speed-up for higher deletion rates (note the difference in y-axis scale between Figures 2b and 2c). However, increasing the number of iterations increases the amount of provenance information cached for PrIU-opt, thus requiring more memory. As Table 3 indicates, since there are 6x iterations for Cov (large 2) compared to Cov (large 1), roughly 6x memory is needed, confirming the analysis in Section 5. However for Cov, with a large mini-batch size and 500 iterations, convergence is achieved and we do not observe a difference in validation accuracy between Cov (large 1) and Cov (large 2). Note that according to [10, 37, 43], the theoretical optimal number of passes for logistic regression

Table 4: Accuracy and similarity comparison between
PrIU-opt and INFL with deletion rate 0.2

Dataset	Validation accuracy		distance		similarity	
	BaseL =	INFL	PrIU-	INFL	PrIU-	INFL
	PrIU-opt		opt		opt	
Cov (small)	48.76%	36.93%	0.184	1.287	0.992	0.624
Cov (large 1)	48.76%	37.99%	0.0016	1.047	1.0	0.738
Cov (large 2)	48.76%	46.38%	0.0003	1.430	1.0	0.471
HIGGS	52.99%	47.99%	0.0004	0.006	0.979	-0.040
Heartbeat	82.78%	74.34%	0.0016	0.583	1.00	0.143
SGEMM	0.001	0.002	0.027	0.140	-	-
(origin)						
SGEMM	0.001	0.002	0.029	0.141	-	-
(extended)						

using mb-SGD (one pass equals to the total number of iterations divided by the number of iterations used for going through the full training set) is quite small. However, for Cov (large 2) the number of passes over the full training set is quite large  $(3000/(581012/10000) \approx 60)$ . Such a high memory usage should therefore not arise in practice.

(Q7) In terms of the update time for experiments over datasets with a comparable mini-batch size but with different feature space sizes (Heartbeat VS HIGGS), we notice that a larger number of model parameters leads to poorer performance by PrIU-opt (compare Figures 3a and 3b). This is also validated through a second set of experiments in which HIGGS (extended) achieves significant speed-up compared to Heartbeat (extended) (see Figure 4). This confirms the analysis in Section 5, where we show how the asymptotic execution time of PrIU and PrIU-opt depends on the number of the model parameters.

(Q8) Table 3 shows that in most cases, both PrIU and PrIU-opt only consume no more than 5x memory compared to BaseL (ignore the number for Cov (large 2) since, as discussed earlier, it is a rare case in practice). However, with a large number of model parameters (like cifar10 and Heartbeat) there is over 10x memory consumption for PrIU and PrIU-opt. How to decrease the memory usage for dense datasets with large feature space is left for future work.

**Discussion.** Extensive experiments using linear regression and logistic regression over the datasets above show the feasibility of our approach. PrIU and PrIU-opt can achieve up to two orders of magnitude speed-up for incrementally updating model parameters compared to the baseline, especially for large datasets with a small feature space. This is done without sacrificing the correctness of the results

(measured by similarity to the updated model parameters by BaseL) and the prediction performance. The experiments also show that the optimizations used in PrIU-opt give significant performance gains compared to PrIU with only a small loss of accuracy. We observe that INFL is not a good solution because of the poor quality of models produced when more than one sample is removed.

**Limitations.** Our experiments also show the limitations of our solutions. They concern the memory footprint when the feature space or the number of iterations is large (anticipated by several analyses in Section 5) and the marginal speed-up for large sparse datasets (See Section 5.3). We shall endeavor to approach these limitations in future work.

### 7 CONCLUSIONS

In this paper, we build a connection between data provenance and incremental machine learning model updates, which is useful in many machine learning and data science applications. Building on an extension of the provenance semiring framework [21] to include basic linear algebra operations [52], we capture provenance in the training phase of linear regression and (binary and multinomial) logistic regression and address non-linear operations in logistic regression using piecewise linear interpolation. We prove that linearization does not harm convergence of the updated parameters and similarity to the expected results. Based on these theoretical results, we construct solutions, PrIU and PrIU-opt, which are optimized to reduce the time and space overhead. The benefits of our solutions are experimentally verified through extensive evaluations over various datasets. Looking forward, we believe that these solutions for simpler machine learning models are likely to extend to generalized additive models [24] and they also pave the way toward solutions for more complicated machine learning models such as deep neural networks.

### ACKNOWLEDGMENTS

This material is based upon work that is in part supported by the Defense Advanced Research Projects Agency (DARPA) under Contract No. HR001117C0047. Partial support was provided by NSF Awards 1547360 and 1733794. Tannen's work at the National University of Singapore was supported in part by the Kwan Im Thong Hood Cho Temple/Avalokiteśvara.

### **REFERENCES**

- Yael Amsterdamer, Susan B Davidson, Daniel Deutch, Tova Milo, Julia Stoyanovich, and Val Tannen. 2011. Putting lipstick on pig: Enabling database-style workflow provenance. *Proceedings of the VLDB Endow*ment 5, 4 (2011), 346–357.
- [2] Yael Amsterdamer, Daniel Deutch, and Val Tannen. 2011. Provenance for aggregate queries. In Proceedings of the thirtieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems. ACM, 153–164.
- [3] Léon Bottou, Frank E Curtis, and Jorge Nocedal. 2018. Optimization methods for large-scale machine learning. *Siam Review* 60, 2 (2018), 223–311.
- [4] Peter Buneman, Sanjeev Khanna, and Tan Wang-Chiew. 2001. Why and where: A characterization of data provenance. In *International* conference on database theory. Springer, 316–330.
- [5] Peter Buneman and Wang-Chiew Tan. 2018. Data Provenance: What next? ACM SIGMOD Record 47, 3 (2018), 5–16.
- [6] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015. Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015. 1721– 1730
- [7] James Cheney, Laura Chiticariu, and Wang Chiew Tan. 2009. Provenance in Databases: Why, How, and Where. Foundations and Trends in Databases 1, 4 (2009), 379–474.
- [8] Xu Chu, Ihab F Ilyas, Sanjay Krishnan, and Jiannan Wang. 2016. Data cleaning: Overview and emerging challenges. In Proceedings of the 2016 International Conference on Management of Data. ACM, 2201–2206.
- [9] R Dennis Cook. 1977. Detection of influential observation in linear regression. *Technometrics* 19, 1 (1977), 15–18.
- [10] John Darzentas. 1984. Problem complexity and method efficiency in optimization. *Journal of the Operational Research Society* 35, 5 (1984), 455–455.
- [11] Shagnik Das. [n.d.]. A brief note on estimates of binomial coefficients.
- [12] Tamraparni Dasu and Theodore Johnson. 2003. Exploratory data mining and data cleaning. Vol. 479. John Wiley & Sons.
- [13] Amol Deshpande and Samuel Madden. 2006. MauveDB: supporting model-based user views in database systems. In Proceedings of the 2006 ACM SIGMOD international conference on Management of data. ACM, 73–84.
- [14] Mohamad Dolatshah, Mathew Teoh, Jiannan Wang, and Jian Pei. 2018. Cleaning crowdsourced labels using oracles for statistical classification. Proceedings of the VLDB Endowment 12, 4 (2018), 376–389.
- [15] Finale Doshi-Velez and Been Kim. 2017. A roadmap for a rigorous science of interpretability. arXiv preprint arXiv:1702.08608 150 (2017).
- [16] Tommy Ellkvist, David Koop, Erik W Anderson, Juliana Freire, and Cláudio Silva. 2008. Using provenance to support real-time collaborative design of workflows. In *International Provenance and Annotation* Workshop. Springer, 266–279.
- [17] Wenfei Fan and Floris Geerts. 2012. Foundations of Data Quality Management. Morgan & Claypool Publishers.
- [18] Todd J Green, Grigoris Karvounarakis, Zachary G Ives, and Val Tannen. 2007. Update exchange with mappings and provenance. In *Proceedings* of the 33rd international conference on Very large data bases. VLDB Endowment, 675–686.
- [19] Todd J Green, Grigoris Karvounarakis, Zachary G Ives, and Val Tannen. 2010. Provenance in ORCHESTRA. (2010).
- [20] Todd J Green, Grigoris Karvounarakis, and Val Tannen. 2007. Provenance semirings. In Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems. ACM, 31–40.

- [21] Todd J. Green and Val Tannen. 2017. The Semiring Framework for Database Provenance. In Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2017, Chicago, IL, USA, May 14-19, 2017. 93-99.
- [22] Priyank Gupta, Nick Koudas, Europa Shang, Ryan Johnson, and Calisto Zuzarte. 2015. Processing analytical workloads incrementally. arXiv preprint arXiv:1509.05066 (2015).
- [23] Sona Hasani, Saravanan Thirumuruganathan, Abolfazl Asudeh, Nick Koudas, and Gautam Das. 2018. Efficient construction of approximate ad-hoc ML models through materialization and reuse. *Proceedings of the VLDB Endowment* 11, 11 (2018), 1468–1481.
- [24] Trevor Hastie and Robert Tibshirani. 1986. Generalized additive models. *Statist. Sci.* 1, 3 (1986), 297–318.
- [25] Alireza Heidari, Joshua McGrath, Ihab F Ilyas, and Theodoros Rekatsinas. 2019. HoloDetect: Few-Shot Learning for Error Detection. arXiv preprint arXiv:1904.02285 (2019).
- [26] Zachary G Ives, Todd J Green, Grigoris Karvounarakis, Nicholas E Taylor, Val Tannen, Partha Pratim Talukdar, Marie Jacob, and Fernando Pereira. 2008. The ORCHESTRA collaborative data sharing system. ACM Sigmod Record 37, 3 (2008), 26–32.
- [27] Eric Jones, Travis Oliphant, Pearu Peterson, et al. 2001–. SciPy: Open source scientific tools for Python. http://www.scipy.org/
- [28] Daniel Kang, Deepti Raghavan, Peter Bailis, and Matei Zaharia. 2018. Model Assertions for Debugging Machine Learning. https://www-cs.stanford.edu/~matei/papers/2018/mlsys\_model\_assertions.pdf Preprint.
- [29] Hamed Karimi, Julie Nutini, and Mark Schmidt. 2016. Linear convergence of gradient and proximal-gradient methods under the polyak-lojasiewicz condition. In Joint European Conference on Machine Learning and Knowledge Discovery in Databases. Springer, 795–811.
- [30] Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 1885–1894.
- [31] Rainer Kress. 1998. *Interpolation*. Springer New York, New York, NY, 151–188. https://doi.org/10.1007/978-1-4612-0599-9\_8
- [32] Sanjay Krishnan, Michael J Franklin, Ken Goldberg, and Eugene Wu. 2017. Boostclean: Automated error detection and repair for machine learning. arXiv preprint arXiv:1711.01299 (2017).
- [33] Sanjay Krishnan, Daniel Haas, Michael J Franklin, and Eugene Wu. 2016. Towards reliable interactive data cleaning: a user survey and recommendations. In *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*. ACM, 9.
- [34] Sanjay Krishnan, Jiannan Wang, Eugene Wu, Michael J Franklin, and Ken Goldberg. 2016. ActiveClean: interactive data cleaning for statistical modeling. *Proceedings of the VLDB Endowment* 9, 12 (2016), 948–959.
- [35] Sanjay Krishnan and Eugene Wu. 2017. Palm: Machine learning explanations for iterative debugging. In Proceedings of the 2nd Workshop on Human-In-the-Loop Data Analytics. ACM, 4.
- [36] Raunak Kumar and Mark Schmidt. 2017. Convergence rate of expectation-maximization. In 10th NIPS Workshop on Optimization for Machine Learning.
- [37] Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. 2012. Efficient backprop. In Neural networks: Tricks of the trade. Springer, 9–48.
- [38] Zachary Lipton. 2016. The Mythos of Model Interpretability. CoRR abs/1606.03490 (2016).
- [39] Yin Lou, Rich Caruana, and Johannes Gehrke. 2012. Intelligible models for classification and regression. In The 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '12, Beijing, China, August 12-16, 2012. 150–158.

- [40] Milos Nikolic, Mohammed Elseidy, and Christoph Koch. 2014. LIN-VIEW: incremental view maintenance for complex analytical queries. In *International Conference on Management of Data*, SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014. 253–264.
- [41] Huazhong Ning, Wei Xu, Yun Chi, Yihong Gong, and Thomas S Huang. 2010. Incremental spectral clustering by efficiently updating the eigensystem. *Pattern Recognition* 43, 1 (2010), 113–127.
- [42] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In NIPS-W.
- [43] Boris T Polyak and Anatoli B Juditsky. 1992. Acceleration of stochastic approximation by averaging. SIAM Journal on Control and Optimization 30, 4 (1992), 838–855.
- [44] Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. 2017. Data management challenges in production machine learning. In Proceedings of the 2017 ACM International Conference on Management of Data. ACM, 1723–1726.
- [45] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach. 2018. Manipulating and measuring model interpretability. arXiv preprint arXiv:1802.07810 (2018).

- [46] Erhard Rahm and Hong Hai Do. 2000. Data cleaning: Problems and current approaches. IEEE Data Eng. Bull. 23, 4 (2000), 3–13.
- [47] Herbert Robbins and Sutton Monro. 1951. A stochastic approximation method. The annals of mathematical statistics (1951), 400–407.
- [48] Mark Schmidt. 2014. Convergence rate of stochastic gradient with constant step size. (2014).
- [49] Jennifer She and Mark Schmidt. 2017. Linear convergence and support vector identifiation of sequential minimal optimization. In 10th NIPS Workshop on Optimization for Machine Learning. 5.
- [50] W Sun and Lars Erik Sjöberg. 2001. Convergence and optimal truncation of binomial expansions used in isostatic compensations and terrain corrections. *Journal of Geodesy* 74, 9 (2001), 627–636.
- [51] Alan Weiser and Sergio E Zarantonello. 1988. A note on piecewise linear and multilinear table interpolation in many dimensions. *Math. Comp.* 50, 181 (1988), 189–196.
- [52] Zhepeng Yan, Val Tannen, and Zachary G Ives. 2016. Fine-grained Provenance for Linear Algebra Operators.. In *TaPP*.
- [53] Wenchao Zhou, Micah Sherr, Tao Tao, Xiaozhou Li, Boon Thau Loo, and Yun Mao. 2010. Efficient querying and maintenance of network provenance at internet-scale. In Proceedings of the 2010 ACM SIGMOD International Conference on Management of data. ACM, 615–626.

### A APPENDIX

#### A.1 Notations

A.1.1 Notations for objective functions, gradients and update rule. The objective functions for linear regression, binary logistic regression and multinomial logistic regression are shown as below (They are Equation (2)-(4) in the paper):

$$h(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \frac{\lambda}{2} ||\mathbf{w}||_2^2$$
 (21)

$$h(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} \ln(1 + \exp\{-y_i \mathbf{w}^{\mathsf{T}} \mathbf{x}_i\}) + \frac{\lambda}{2} ||\mathbf{w}||_2^2$$
 (22)

$$h(\mathbf{w}) = \frac{1}{n} \sum_{k=1}^{q} \sum_{y_i = k} \left( \ln\left(\sum_{j=1}^{q} e^{\mathbf{w}_j^{\mathsf{T}} \mathbf{x}_i}\right) - \mathbf{w}_k^{\mathsf{T}} \mathbf{x}_i \right) + \frac{\lambda}{2} ||vec([\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q])||_2^2$$

$$\mathbf{w} = vec([\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q])$$
(23)

Note that  $h(\mathbf{w})$  can be rewritten as  $h(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} h_i(\mathbf{w}) + \frac{\lambda}{2} ||\mathbf{w}^{(t)}||$ . For example for Equation (3),  $h_i(\mathbf{w}) = (y_i - \mathbf{x}_i^T \mathbf{w})^2$ For linear regression and logistic regression, the rule for updating  $\mathbf{w}^{(t)}$  under mb-SGD is presented below (They are Equation (5)-(6) in the paper and we use  $\nabla^{(t)} h(\mathbf{w}^{(t)})$  to denote the average gradients evaluated over the mini-batch at the  $t_{th}$  iteration):

$$\mathbf{w}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \frac{2\eta_t}{B} \sum_{i \in \mathcal{B}^{(t)}} \mathbf{x}_i (\mathbf{x}_i^T \mathbf{w}^{(t)} - y_i)$$

$$= (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \frac{\eta_t}{B} \sum_{i \in \mathcal{B}^{(t)}} \nabla h_i (\mathbf{w}^{(t)}) = (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \eta_t \nabla^{(t)} h(\mathbf{w}^{(t)})$$
(24)

$$\mathbf{w}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}^{(t)} + \frac{\eta_t}{B} \sum_{i \in \mathscr{B}^{(t)}} y_i \mathbf{x}_i (1 - \frac{1}{1 + \exp\{-y_i \mathbf{w}^{(t)T} \mathbf{x}_i\}})$$

$$= (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \frac{\eta_t}{B} \sum_{i \in \mathscr{B}^{(t)}} \nabla h_i(\mathbf{w}^{(t)}) = (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \eta_t \nabla^{(t)} h(\mathbf{w}^{(t)})$$
(25)

Note that in Equation (25), the non-linear part can be abstracted as  $f(x) = 1 - \frac{1}{1 + e^{-x}}$ . So this formula can be also represented as:

$$\mathbf{w}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}^{(t)} + \frac{\eta_t}{B} \sum_{i \in \mathscr{B}^{(t)}} y_i \mathbf{x}_i (1 - \frac{1}{1 + \exp\{-y_i \mathbf{w}^{(t)T} \mathbf{x}_i\}})$$

$$= (1 - \eta_t \lambda) \mathbf{w}^{(t)} + \frac{\eta_t}{B} \sum_{i \in \mathscr{B}^{(t)}} y_i \mathbf{x}_i f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)$$

$$= (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \frac{\eta_t}{B} \sum_{i \in \mathscr{B}^{(t)}} \nabla h_i(\mathbf{w}^{(t)}) = (1 - \eta_t \lambda) \mathbf{w}^{(t)} - \eta_t \nabla^{(t)} h(\mathbf{w}^{(t)})$$

$$(26)$$

So

$$\nabla h^{(t)}(\mathbf{w}) = -\frac{1}{B} \sum_{i \in \mathcal{B}^{(t)}} y_i \mathbf{x}_i f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)$$

Also we can explicitly evaluate  $\nabla^2 h^{(t)}(\mathbf{w})$  as:

$$\nabla^2 h^{(t)}(\mathbf{w}^{(t)}) = -\frac{1}{B} \sum_{\mathscr{D}^{(t)}} \mathbf{x}_i \mathbf{x}_i^T f'(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)$$
(27)

in which  $-\sum_{i\in\mathscr{B}^{(t)}}\mathbf{x}_i\mathbf{x}_i^Tf'(y_i\mathbf{w}^{(t)T}\mathbf{x}_i)$  should be a semi-definite matrix since  $f(x)=1-\frac{1}{1+exp\{-x\}}$  is a monotonically decreasing function and thus f'(x) should be negative for any x.

A.1.2 Notations for the linearized update rule. After the interpolation step over the update rules for binary logistic regression, Equation (25) can be approximated as (It is Equation (9) in the paper):

$$\mathbf{w}_{L}^{(t+1)} \approx \left[ (1 - \eta_{t}\lambda)\mathbf{I} + \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} a^{i,(t)} \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right] \mathbf{w}_{L}^{(t)} + \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} b^{i,(t)} y_{i} \mathbf{x}_{i}$$

$$(28)$$

which can be also represented as:

$$\mathbf{w}_{L}^{(t+1)} \approx \left[ (1 - \eta_{t} \lambda) \mathbf{I} + \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} a^{i,(t)} \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right] \mathbf{w}_{L}^{(t)} + \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} b^{i,(t)} y_{i} \mathbf{x}_{i}$$

$$= (1 - \eta_{t} \lambda) \mathbf{w}_{L}^{(t)} + \frac{\eta_{t}}{B} \sum_{i \in \mathcal{B}^{(t)}} y_{i} \mathbf{x}_{i} s(y_{i} \mathbf{w}_{L}^{(t)T} \mathbf{x}_{i})$$

$$(29)$$

in which  $s(x) = a^{i,(t)}x + b^{i,(t)}$ .

Suppose after removing certain subset (the number of those samples is  $\Delta n$  and the corresponding indices are  $\mathcal{R}$ ), Equation (29) becomes (It is Equation (11) in the paper):

$$\mathbf{w}_{LU}^{(t+1)} \approx \left[ (1 - \eta_t \lambda) \mathbf{I} + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathcal{R}} a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T \right] \mathbf{w}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathcal{R}} b^{i,(t)} y_i \mathbf{x}_i$$

$$= (1 - \eta_t \lambda) \mathbf{w}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathcal{R}} y_i \mathbf{x}_i s(y_i \mathbf{w}_{LU}^{(t)T} \mathbf{x}_i)$$
(30)

For the linearized version of the update rule of logistic regression in Equation (29) and the update rule in Equation (30), we represent  $\nabla T^{(t)}(\mathbf{w}_L^{(t)})$  and  $\nabla R^{(t)}(\mathbf{w}_{LU}^{(t)})$  as:

$$\nabla T_i^{(t)}(\mathbf{w}_L^{(t)}) = -y_i \mathbf{x}_i s(y_i \mathbf{w}_L^{(t)T} \mathbf{x}_i) = (-a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_L^{(t)} - b^{i,(t)} y_i \mathbf{x}_i$$
(31)

$$\nabla T^{(t)}(\mathbf{w}_L^{(t)}) = \frac{1}{B} \sum_{i \in \mathcal{R}^{(t)}} \nabla T_i^{(t)}(\mathbf{w}_L^{(t)})$$
(32)

$$\nabla R_i^{(t)}(\mathbf{w}_{LU}^{(t)}) = -y_i \mathbf{x}_i s(y_i \mathbf{w}_{LU}^{(t)T} \mathbf{x}_i) = -a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T \mathbf{w}_{LU}^{(t)} - b^{i,(t)} y_i \mathbf{x}_i$$
(33)

$$\nabla R^{(t)}(\mathbf{w}_{LU}^{(t)}) = \frac{1}{B_U^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} \nabla R_i^{(t)}(\mathbf{w}_{LU}^{(t)})$$
(34)

where  $\nabla T^{(t)}(\mathbf{w}_L^{(t)})$  and  $\nabla R^{(t)}(\mathbf{w}_{LU}^{(t)})$  can be considered as pseudo-derivative in Equation (29). So Equation (29) and Equation (30) can be rewritten as:

$$\mathbf{w}_{L}^{(t+1)} = (1 - \eta_{t}\lambda)\mathbf{w}_{L}^{(t)} + \frac{\eta_{t}}{B_{U}^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} y_{i}\mathbf{x}_{i}s(y_{i}\mathbf{w}_{L}^{(t)T}\mathbf{x}_{i}) = (1 - \eta_{t}\lambda)\mathbf{w}_{L}^{(t)} - \eta_{t}\nabla T^{(t)}(\mathbf{w}_{L}^{(t)})$$

$$(35)$$

$$\mathbf{w}_{LU}^{(t+1)} = (1 - \eta_t \lambda) \mathbf{w}_{LU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathcal{D}^{(t)}, i \notin \mathcal{R}} y_i \mathbf{x}_i s(y_i \mathbf{w}_{LU}^{(t)T} \mathbf{x}_i) = (1 - \eta_t \lambda) \mathbf{w}_{LU}^{(t)} - \eta_t \nabla R^{(t)}(\mathbf{w}_{LU}^{(t)})$$
(36)

In contrast, by computing the model parameter from the scratch for logistic regression after removing the same set of training samples, the update rule is (It is Equation (12) in the paper):

$$\mathbf{w}_{RU}^{(t+1)} \leftarrow (1 - \eta_t \lambda) \mathbf{w}_{RU}^{(t)} + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathscr{D}^{(t)} \ i \notin \mathscr{R}} y_i \mathbf{x}_i f(y_i \mathbf{w}_{RU}^{(t)} \mathbf{x}_i) = (1 - \eta_t \lambda) \mathbf{w}_{RU}^{(t)} - \eta_t \nabla^{(t)} g(\mathbf{w}_{RU}^{(t)})$$

$$(37)$$

which aims at minimizing the following objective function:

$$g(\mathbf{w}) = \frac{1}{n - \Delta n} \sum_{i \in \mathcal{R}} h_i(\mathbf{w}) + \frac{\lambda}{2} ||\mathbf{w}||_2^2$$
(38)

in which  $\mathcal{R}$  represents the ids of the samples that are removed and  $\Delta n$  represents the number of the removed samples and

$$\nabla^{(t)}g(\mathbf{w}_{RU}^{(t)}) = \frac{1}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathcal{R}} \nabla h_i(\mathbf{w}_{RU}^{(t)}).$$

Similarly after removing certain subset, the update rule for linear regression model is (It is Equation (13) in the paper):

$$\mathbf{w}_{U}^{(t+1)} \leftarrow \left[ (1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}} \sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i}\mathbf{x}_{i}^{T} \right] \\
- \sum_{i \in \mathscr{B}^{(t)}, i \in \mathcal{R}} \mathbf{x}_{i}\mathbf{x}_{i}^{T} \mathbf{w}_{U}^{(t)} + \frac{2\eta_{t}}{B_{U}^{(t)}} (\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i}y_{i} - \sum_{i \in \mathscr{B}^{(t)}, i \in \mathcal{R}} \mathbf{x}_{i}y_{i})$$
(39)

The provenance expression for the model parameters of linear regression model and logistic regression model after removing subset of training samples are (They are Equation (8) and Equation (10) in the paper):

$$\mathcal{W}_{U}^{(t+1)} \leftarrow \left[ (1 - \eta_{t} \lambda) (1_{\text{prov}} * \mathbf{I}) - \frac{2\eta_{t}}{B_{U}^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} p_{i}^{2} * \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right] \mathcal{W}_{U}^{(t)} + \frac{2\eta_{t}}{B_{U}^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} p_{i}^{2} * \mathbf{x}_{i} y_{i}$$

$$(40)$$

$$\mathcal{W}_{LU}^{(t+1)} \leftarrow \left[ (1 - \eta_t \lambda)(1_{\text{prov}} * \mathbf{I}) + \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} p_i^2 * (a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T) \right] \mathcal{W}_{LU}^{(t)}$$

$$+ \frac{\eta_t}{B_U^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} p_i^2 * (b^{i,(t)} y_i \mathbf{x}_i)$$

$$(41)$$

### A.2 proof preliminary

There are some useful properties related to matrix theory, matrix norm, real analysis and SGD convergence, which will be used in the follow-up proof.

Lemma 2 (SGD convergence, [3]). (Full version of Lemma (1) in the paper) Suppose that the stochastic gradient estimates are correlated with the true gradient, and bounded in the following way. There exist two scalars  $J_1 \geq J_2 > 0$  such that for arbitrary  $\mathcal{B}_t$ , the following two inequalities hold:

$$\nabla h(\mathbf{w}_t)^T \mathbb{E} \frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla h_i(\mathbf{w}_t) \ge J_2 \|\nabla h(\mathbf{w}_t)\|^2$$
(42)

$$\|\mathbb{E}\frac{1}{B_t}\sum_{i\in\mathscr{B}_t}\nabla h_i(\mathbf{w}_t)\| \leq J_1\|\nabla h(\mathbf{w}_t)\|$$
(43)

Also, for two scalars  $J_3, J_4 \ge 0$  we have:

$$Var\left(\frac{1}{B_t}\sum_{i\in\mathcal{B}_t}\nabla h_i\left(\mathbf{w}_t\right)\right) \leq J_3 + J_4 \|\nabla h\left(\mathbf{w}_t\right)\|^2$$
(44)

By combining equations (42)-(44), the following inequality holds:

$$\mathbb{E}\|\frac{1}{B_t}\sum_{i\in\mathscr{B}_t}\nabla F_i(\mathbf{w}_t)\|^2 \le J_3 + J_5\|\nabla F(\mathbf{w}_t)\|^2$$
(45)

where  $J_5 = J_4 + J_1^2 \ge J_2^2 \ge 0$ .

Then stochastic gradient descent with fixed step size  $\eta_t = \eta \leq \frac{J_2}{L_{15}}$  has the convergence rate:

$$\mathbb{E}\left[h\left(\mathbf{w}_{t}\right) - h\left(\mathbf{w}^{*}\right)\right] \leq \frac{\eta L J_{3}}{2\mu J_{2}} + (1 - \eta \mu J_{2})^{t-1} \left(h\left(\mathbf{w}_{1}\right) - h\left(\mathbf{w}^{*}\right) - \frac{\eta L J_{3}}{2\mu J_{2}}\right) \to \frac{\eta L J_{3}}{2\mu J_{2}}$$
(46)

If the gradient estimates are unbiased, then  $\mathbb{E}\frac{1}{B_t}\sum_{i\in\mathscr{B}_t}\nabla h_i\left(\mathbf{w}_t\right)=\frac{1}{n}\sum_{i=1}^n\nabla h_i\left(\mathbf{w}_t\right)=\nabla h\left(\mathbf{w}_t\right)$  and thus  $J_1=J_2=1$ . Moreover,  $J_3\sim 1/B$ , where B is the minibatch size, because  $J_2$  is the variance of the stochastic gradient.

So the convergence condition for fixed step size becomes  $\eta_t = \eta \le \frac{1}{LJ_5}$ , in which  $J_5 = J_4 + J_1^2 = J_4 + 1 \ge 1$ . So  $\eta_t = \eta \le \frac{1}{LJ_5} \le \frac{1}{L}$  suffices to ensure convergence.

So in what follows, we will simply consider the case where the learning rate is a constant across all the iterations as Lemma 2 indicates.

LEMMA 3. For a matrix A, its L2-norm equals to its largest singular value and the maximal eigenvalue of matrix  $A^T A$ , i.e.:  $||A||_2 = \sigma_{max}(A) = \sqrt{C_{max}(A^T A)}$ .

where  $\sigma_{max}$  and  $C_{max}$  represents the largest singular value and the largest eigenvalue of certain matrix.

If A is a semi-definite matrix, its eigenvalue is the same as its singular value, then the equation above can be rewritten as:  $||A||_2 = \sigma_{max}(A) = C_{max}(A)$ .

LEMMA 4. If an  $n \times n$  matrix A is a real symmetric matrix, then we can find n mutually orthogonal eigenvectors for A.

LEMMA 5. Given an iteration formula  $\mathbf{u}^{(t+1)} = A\mathbf{u}^{(t)} + \mathbf{b}$  where  $\mathbf{A}$  is a matrix while  $\mathbf{u}^{(t)}$  is a vector to be derived iteratively, if  $\mathbf{I} - \mathbf{A}$  is invertible, then the following statements are equivalent:

- (1)  $\mathbf{u}^{(t)}$  will get converged
- (2)  $||\mathbf{B}||_p < 1$  for some matrix norm  $||||_p$

LEMMA 6. **Cauchy schwarz inequality** For any two matrix **A** and **B**, their norm should satisfy the Cauchy schwarz inequality, i.e.:  $||AB||_X \le ||A||_X ||B||_X$  where  $||\cdot||_X$  represents any matrix norm

LEMMA 7. Weyl's inequality For any three Hermitian matrices, M, N, P satisfying M = N + P, the eigenvalues of M is:  $\mu_1 \ge \mu_2 \ge \mu_3 \cdots \ge \mu_n$ ;

the eigenvalues of N is:  $v_1 \ge v_2 \ge v_3 \cdots \ge v_n$ ;

and the eigenvalues of **P** is:  $\rho_1 \ge \rho_2 \ge \rho_3 \cdots \ge \rho_n$ ;

the following inequalities hold:  $v_i + \rho_n \le \mu_i \le v_i + \rho_1$ 

The following lemma requires the definition of Lipschitz-continuity and Strong-convexity, which are provided below:

**Lipschitz-continuous** A function f(x) is Lipschitz-continuous (L-continuous) if there exists a constant L such that the following inequality is satisfied for all x, y:

$$|f(y) - f(x)| \le L||y - x||_2^2 \tag{47}$$

Another form of Equation (47) is:

$$f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||y - x||_2^2$$
(48)

**Strong convexity** A function f(x) is  $\lambda$ -strong convexity iff there exists a constant  $\lambda$  such that the following inequality is satisfied for all x, y:

$$f(y) \ge f(x) + \langle \nabla f(x), y - x \rangle + \frac{\lambda}{2} ||y - x||^2$$
(49)

Other equivalent forms of Equation (49) are:

$$(\nabla f(x) - \nabla f(y))(x - y) \ge \lambda ||x - y||_2^2 \tag{50}$$

$$\nabla^2 f(x) \ge \lambda \tag{51}$$

Then there is a useful lemma about  $\lambda$ -strong convexity, i.e.

LEMMA 8. a function f(x) is a strong convex function iff  $f(x) - \frac{\lambda}{2}||x||_2^2$  is a convex function,

LEMMA 9. **Piecewise linear interpolation** In Piecewise linear interpolation [31], we assume that the function to be approximated is a continuous function f(x) where  $x \in [a,b]$ . Piecewise linear interpolation starts by picking up a series of breaking points,  $x_i$  such that  $a < x_1 < x_2 < \cdots < x_p < b$  and then constructs a linear interpolant s(x) over each interval  $[x_{j-1}, x_j]$  as follows:

$$s(x) = \frac{x - x_{j-1}}{x_j - x_{j-1}} f(x_j) + \frac{x_j - x}{x_j - x_{j-1}} f(x_{j-1})$$
  
=  $a_j x + b_j, x \in [x_{j-1}, x_j)$  (52)

The following property holds on how close the value of s(x) is compared to the original function f(x):

$$|f(x) - s(x)| \le \frac{1}{8} (\Delta x)^2 \max_{a \le x \le b} |f''(x)| = O((\Delta x)^2)$$

$$|f'(x) - s'(x)| \le \frac{1}{2} (\Delta x) \max_{a \le x \le b} |f''(x)| = O((\Delta x))$$
(53)

Lemma 10. Expectation of the number of the removed samples Because of the randomness from mb-SGD, the  $\Delta n$  removed samples can be viewed as uniformly distributed within all n training samples, which can be considered as a 0-1 Bernoulli distribution with probability  $\frac{\Delta n}{n}$ . In other words, we can define a random variable  $\mathbf{S}_i$  for each sample, which is 1 with probability  $\frac{\Delta n}{n}$  and 0 with probability  $1-\frac{r}{n}$ . So within a single mini-batch  $\mathcal{B}^t$ , we can have

$$\mathbb{E}(\sum_{i \in \mathscr{D}^t} S_i) = \mathbb{E}(\Delta B_t) = B \frac{r}{n}$$

and

$$Var(\sum_{i \in \mathscr{L}^t} S_i) = B\frac{r}{n}(1 - \frac{r}{n})$$

So in terms of the random variable  $\frac{\Delta B_t}{B}$ , its expectation and variance will be

$$\mathbb{E}(\frac{\Delta B_t}{B}) = \frac{r}{n} \tag{54}$$

and

$$Var(\frac{\Delta B_t}{B}) = \frac{r}{Bn}(1 - \frac{r}{n}) \tag{55}$$

In mb-SGD, a typical assumption is used for the convergence analysis of the model parameter  $\mathbf{w}^{(t)}$  in Equation (24)-(25) and the update rules for other general models, i.e.:

Lemma 11. For any randomly selected sample  $i_j$  in some batch, the expectation of its gradient should be the same as the gradient over the all the samples, i.e.:

 $E(\nabla h_{i_i}(\mathbf{w})) = \nabla h(\mathbf{w})$ 

which also implies that the following equality holds for mb-SGD:

 $E(\nabla(\frac{1}{B}\sum_{i\in\mathscr{B}^{(t)}}h_i(\mathbf{w}))) = \nabla h(\mathbf{w})$ 

where E is the expectation value with respect to the sampling over the entire training samples.

In what follows, our analysis is based on the following assumptions:

Assumption 1. every  $h_i(\mathbf{w})$  (i = 1, 2, ..., n) is L-Lipschitz continuous. Since  $h_i(\mathbf{w})$  has L2-norm regularization term, then we also know that  $h_i(\mathbf{w})$  is  $\lambda$ -strong convex.

Assumption 2. each  $\nabla h_i(\mathbf{w}^{(t)})$  is bounded by some constant  $c_1$  for each  $\mathbf{w}^{(t)}$ .

Assumption 3. The function f'(\*) is  $c_2$ —Lipschitz continuous, which means that the following inequality holds:

$$|f'(x) - f'(y)| \le c_2 ||x - y||$$

# A.3 Main results and proofs

THEOREM 10. (It is Theorem 2 in the paper)  $W_U^{(t)}$  in Equation (40) and  $W_{LU}^{(t)}$  in Equation (41) need not converge under the conditions in Lemma 2.

PROOF. Let us take linear regression as an example. Note that we can explicitly evaluate the second order derivative of  $h(\mathbf{w})$  for linear regression, i.e.  $\nabla^2 h(\mathbf{w})$ , then according to Assumption 1,  $\nabla^2 h(\mathbf{w})$  should satisfy the following inequality:

$$\lambda \le ||\nabla^2 h(\mathbf{w})||_2 = ||\frac{2}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T + \lambda \mathbf{I}||_2 \le L$$
(56)

In order to prove Theorem 10, we need to show that there exists a case where  $\mathcal{W}_{U}^{(t)}$  cannot converge under the conditions in Lemma 2. This is achieved by considering gradient descent (GD) without excluding any original training samples, i.e.

 $\{p_{i_1}, p_{i_2}, \dots, p_{i_z}\} = \{1, 2, \dots, n\}$ , every  $B_U^{(t)} = n$  in Equation (40) and every  $\mathcal{B}^{(t)}$  includes all n samples in Equations (24) and 40. We can then apply the update rule in Equations (24) and (40) recursively, which ends up with:

$$\mathbf{w}^{(t+1)} = ((1 - \eta \lambda)\mathbf{I} - \frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T})^{t+1} \mathbf{w}^{(0)}$$

$$+ (\sum_{i=1}^{t} ((1 - \eta \lambda)\mathbf{I} - \frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T})^{j}) \frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} y_{i}$$

$$(57)$$

$$\mathcal{W}_{U}^{(t)} = ((1 - \eta \lambda) 1_{\text{prov}} * \mathbf{I} - \frac{2\eta}{n} \sum_{i=1}^{n} p_{i}^{2} * \mathbf{x}_{i} \mathbf{x}_{i}^{T})^{t} \mathcal{W}_{U}^{(0)} 
+ (\sum_{i=1}^{t} ((1 - \eta \lambda) 1_{\text{prov}} * \mathbf{I} - \frac{2\eta}{n} \sum_{i=1}^{n} p_{i}^{2} * \mathbf{x}_{i} \mathbf{x}_{i}^{T})^{j}) \frac{2\eta}{n} \sum_{i=1}^{n} p_{i}^{2} * \mathbf{x}_{i} y_{i}$$
(58)

According to Assumption 1, the following inequality should be satisfied:

$$||\eta \nabla^2 h(\mathbf{w})||_2 = ||\frac{2\eta}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T + \lambda \eta \mathbf{I}||_2 \le \eta L \le 1$$

$$(59)$$

which implies that

$$||(1 - \eta \lambda)\mathbf{I} - \frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T}||_{2} \le 1$$

$$(60)$$

Also since every  $\mathbf{x}_i \mathbf{x}_i^T$  is a semi-positive definite matrix, by using Lemma 3, Equation (59) also implies that:

$$1 \ge ||\frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T} + \lambda \eta \mathbf{I}||_{2} = C_{max} (\frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T} + \lambda \eta \mathbf{I})$$

$$\ge C_{max} (\frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T}) = ||\frac{2\eta}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{T}||_{2}$$
(61)

Then by applying Equation (7), Equation (61) also leads to:

$$\left|\left|\frac{2\eta}{n}\mathbf{x}_{i}\mathbf{x}_{i}^{T}\right|\right|_{2} = C_{max}\left(\frac{2\eta}{n}\mathbf{x}_{i}\mathbf{x}_{i}^{T}\right) < C_{max}\left(\frac{2\eta}{n}\sum_{i=1}^{n}\mathbf{x}_{i}\mathbf{x}_{i}^{T}\right) \le 1$$
(62)

Then we can expand the first term in the right-hand side of Equation (58), the tensor product with provenance monomial  $p_i^t$  should be  $p_i^t*\binom{t}{\frac{t}{2}}(1-\eta\lambda)^{\frac{t}{2}}(-\frac{2\eta}{n}\mathbf{x}_i\mathbf{x}_i^T)^{\frac{t}{2}}$ . According to the convergence conditions in Lemma 2,  $\eta<\frac{1}{L}$  and thus  $||\binom{t}{\frac{t}{2}}(1-\eta\lambda)^{\frac{t}{2}}(-\frac{2\eta}{n}\mathbf{x}_i\mathbf{x}_i^T)^{\frac{t}{2}}||_2 \geq \binom{t}{\frac{t}{2}}||(-\frac{\eta}{n}\mathbf{x}_i\mathbf{x}_i^T)^{\frac{t}{2}}(\frac{2(L-\lambda)}{L})^{\frac{t}{2}}||_2$ . According to [11, 50], when  $t\to\infty$ ,  $\binom{t}{\frac{t}{2}}$  should be very close to  $2^t$  and thus when  $||-\frac{\eta}{n}\mathbf{x}_i\mathbf{x}_i^T||_2 \geq \frac{L}{2(L-\lambda)}$  (note that  $||-\frac{2\eta}{n}\mathbf{x}_i\mathbf{x}_i^T||_2$  can be any value between 0 and 1 according to Equation (62)),  $\binom{t}{\frac{t}{2}}||(-\frac{\eta}{n}\mathbf{x}_i\mathbf{x}_i^T)^{\frac{t}{2}}||_2\to\infty$ , which means that the tensor product with provenance monomial  $p_i^t$  cannot converge and thus  $\mathbf{W}_U^{(t)}$  cannot converge.

Theorem 11. (It is Theorem 3 in the paper) The expectation of  $W_U^{(t)}$  in Equation (40) and of  $W_{LU}^{(t)}$  in Equation (41), converge when  $t \to \infty$  if we also assume that provenance polynomial multiplication is idempotent.

Convergence proof for linear regression We simply need to consider whether  $W^{(t)}(\{p_{i_1}, p_{i_2}, \dots, p_{i_z}\})$  converges (suppose there are  $\Delta n$  provenance tokens in total that are set as 0, which corresponds to the deletion of  $\Delta n$  samples), which equals to the update rule in Equation (39) and leads to a new objective function without the  $\Delta n$  removed samples (denoted by  $(\Delta X, \Delta Y)$ ), i.e:

In what follows, we will only prove the convergence of binary logistic regression, i.e. the convergence of  $W_{LU}^{(t)}$ , which is the same as proving the convergence of  $\mathbf{w}_{LU}^{(t)}$ . The convergence of linear regression and multi-nomial logistic regression can be proven in the similar way.

According to Lemma 9, the following equation holds:

$$\|\nabla T_i^{(t)}(\mathbf{w}^{(t)}) - \nabla h_i^{(t)}(\mathbf{w}^{(t)})\| = \|y_i \mathbf{x}_i (s(y_i \mathbf{w}^{(t)T} \mathbf{x}_i) - f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i))\| = O((\Delta x)^2)$$
(63)

$$\|\nabla^{2} T_{i}^{(t)}(\mathbf{w}^{(t)}) - \nabla^{2} h_{i}^{(t)}(\mathbf{w}^{(t)})\| = \|-a^{i,(t)} \mathbf{x}_{i} \mathbf{x}_{i}^{T} + f'(y_{i} \mathbf{w}^{(t)T} \mathbf{x}_{i}) \mathbf{x}_{i} \mathbf{x}_{i}^{T}\| = O((\Delta x))$$
(64)

for all  $\mathbf{w}^{(t)}$  (rather than  $\mathbf{w}_{RU}^{(t)}$  or  $\mathbf{w}_{LU}^{(t)}$  since  $a^{i,(t)}$  and  $b^{i,(t)}$  are evaluated when  $\mathbf{w}$  is  $\mathbf{w}^{(t)}$ ). Since  $\lambda \leq \|\nabla^2 h^{(t)}(\mathbf{w}^{(t)})\| \leq L$ , then  $\lambda - O(\Delta x) \leq \|\nabla^2 T^{(t)}(\mathbf{w}^{(t)})\| \leq L + O(\Delta x)$ . Then by bringing in the definition of  $\nabla^2 T^{(t)}(\mathbf{w}^{(t)})$ , we know that:

$$\lambda - O(\Delta x) \le \|\nabla^2 T^{(t)}(\mathbf{w}^{(t)})\| = \|-a^{i,(t)}\mathbf{x}_i\mathbf{x}_i^T + \lambda \mathbf{I}\| \le L + O(\Delta x)$$

$$\tag{65}$$

By using the fact that  $\|*\|_2 = C_{max}(*)$ , the following formula also holds:

$$\lambda - O(\Delta x) \le C_{max}(-a^{i,(t)}\mathbf{x}_i\mathbf{x}_i^T + \lambda \mathbf{I}) \le L + O(\Delta x)$$
(66)

But note that since every  $a^{i,(t)}$  is a negative value, then  $-a^{i,(t)}\mathbf{x}_i\mathbf{x}_i^T$  is a semi-positive definite matrix and thus:

$$\|-a^{i,(t)}\mathbf{x}_{i}\mathbf{x}_{i}^{T}+\lambda\mathbf{I}\| \geq \lambda \tag{67}$$

Then we bound  $\mathbf{I} - \eta(-a^{i,(t)}\mathbf{x}_i\mathbf{x}_i^T + \lambda \mathbf{I})$  as:

$$\|\mathbf{I} - (-\eta a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T + \eta \lambda \mathbf{I})\| \le 1 - \eta \lambda \tag{68}$$

Similarly we know that the following inequality holds for any  $\mathbf{w}$ :

$$\|\mathbf{I} - \eta \nabla^2 h^{(t)}(\mathbf{w})\| = \|(1 - \eta \lambda)\mathbf{I} - \eta [-\mathbf{x}_i \mathbf{x}_i^T f'(y_i \mathbf{w}^T \mathbf{x}_i)]\| \le 1 - \eta \lambda \tag{69}$$

Then the convergence of  $\mathbf{w}_{II}$  can be derived below (in the case of constant  $\eta_t$  according to Lemma 2):

$$||\mathbf{w}_{LU}^{(t+1)} - \mathbf{w}_{LU}^{*}||_{2}$$

$$= ||\mathbf{w}_{LU}^{(t)} - \eta \nabla^{(t)} R^{(t)} (\mathbf{w}_{LU}^{(t)}) - \mathbf{w}_{LU}^{*} + \eta \nabla^{(t)} R^{(t)} (\mathbf{w}_{LU}^{*}) - \eta \nabla^{(t)} R^{(t)} (\mathbf{w}_{LU}^{*})||_{2}$$
(70)

Then by using the fact that if  $\mathbf{w}_{LU}^{(t)}$  converges,  $\|\nabla^{(t)}R^{(t)}(\mathbf{w}_{LU}^*)\| \leq C$  for some constant value C for all t. By using the triangle inequality and the definition of  $\nabla R^{(t)}(*)$ , the formula above is bounded as:

$$\leq \|\{\mathbf{I} - \frac{1}{B_U^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} \eta[-a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T + \lambda \mathbf{I}]\} (\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{LU}^*) \| + \eta C$$

$$\leq \|\mathbf{I} - \frac{1}{B_U^{(t)}} \sum_{i \in \mathscr{B}^{(t)}, i \notin \mathscr{R}} \eta[-a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T + \lambda \mathbf{I}] \|\|(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{LU}^*)\| + \eta C$$

By using the result from Equation (68), we get:

$$\leq (1 - \eta \lambda) \| \| (\mathbf{w}_{II}^{(t)} - \mathbf{w}_{II}^*) \| + \eta C$$

By applying the formula above recursively, we get:

$$||\mathbf{w}_{LU}^{(t+1)} - \mathbf{w}_{LU}^*||_2 \le \frac{C}{\lambda}$$

This finishes the proof.

THEOREM 12. (This is Theorem 4 in the paper)  $||E(\mathbf{w}^{(t)} - \mathbf{w}_{L}^{(t)})||_{2}$  is bounded by  $O((\Delta x)^{2})$  where  $\Delta x$  is an arbitrarily small value representing the length of the longest sub-interval used in piecewise linear interpolations.

PROOF. By subtracting Equation (25) by Equation (35) and taking the matrix norm, we get:

$$\begin{split} &||\mathbb{E}(\mathbf{w}^{(t+1)} - \mathbf{w}_{L}^{(t+1)})||_{2} \\ &= \mathbb{E}(||\mathbf{w}^{(t)} - \eta \nabla h^{(t)}(\mathbf{w}^{(t)}) - (\mathbf{w}_{L}^{(t)} - \eta \nabla T^{(t)}(\mathbf{w}_{L}^{(t)}))||_{2}) \\ &= \mathbb{E}(||\mathbf{w}^{(t)} - \mathbf{w}_{L}^{(t)} - \eta[\nabla T^{(t)}(\mathbf{w}^{(t)}) - \nabla T^{(t)}(\mathbf{w}_{L}^{(t)})] - \eta[\nabla h^{(t)}(\mathbf{w}^{(t)}) - \nabla T^{(t)}(\mathbf{w}^{(t)})]||_{2}) \\ &= \mathbb{E}(||[\mathbf{I} - (-\frac{\eta}{B} \sum_{i \in \mathscr{B}^{(t)}} a^{i,(t)} \mathbf{x}_{i} \mathbf{x}_{i}^{T} + \eta \lambda \mathbf{I})](\mathbf{w}^{(t)} - \mathbf{w}_{L}^{(t)}) - \eta[\nabla h^{(t)}(\mathbf{w}^{(t)}) - \nabla T^{(t)}(\mathbf{w}^{(t)})]||_{2}) \end{split}$$

Then by using triangle inequality and the bound from Equation (68), the formula above is further bounded as:

$$\leq E(\|[\mathbf{I} - (-\frac{\eta}{B} \sum_{i \in \mathscr{B}^{(t)}} a^{i,(t)} \mathbf{x}_i \mathbf{x}_i^T + \eta \lambda \mathbf{I})] \|\|(\mathbf{w}^{(t)} - \mathbf{w}_L^{(t)})\|_2 + \eta \|[\nabla h^{(t)}(\mathbf{w}^{(t)}) - \nabla T^{(t)}(\mathbf{w}^{(t)})]\|_2) \\
\leq (1 - \eta \lambda) \|(\mathbf{w}^{(t)} - \mathbf{w}_L^{(t)})\|_2 + \eta \|[\nabla h^{(t)}(\mathbf{w}^{(t)}) - \nabla T^{(t)}(\mathbf{w}^{(t)})]\|_2$$

Then by using the result from Equation (63), the formula above is rewritten as:

$$= (1 - \eta \lambda) \|(\mathbf{w}^{(t)} - \mathbf{w}_{I}^{(t)})\|_{2} + \eta O((\Delta x)^{2})$$

Then by applying the formula above recursively, we have:

$$= (1 - \eta \lambda)^{t} \|(\mathbf{w}^{(0)} - \mathbf{w}_{L}^{(0)})\|_{2} + \frac{1 - (1 - \eta \lambda)^{t}}{\eta \lambda} \eta O((\Delta x)^{2})$$

Since  $\mathbf{w}^{(0)} = \mathbf{w}_L^{(0)}$  and  $\eta \leq \frac{1}{L}$ , then the formula above is bounded as:

$$\leq \frac{1}{\lambda}O((\Delta x)^2) = O((\Delta x)^2)$$

According to Assumption 2 and Equation (63) and by using the triangle inequality and the Theorem above, we have:

$$\|\nabla T_{i}^{(t)}(\mathbf{w}_{L}^{(t)})\| = \|\nabla T_{i}^{(t)}(\mathbf{w}_{L}^{(t)}) - \nabla T_{i}^{(t)}(\mathbf{w}^{(t)}) + \nabla T_{i}^{(t)}(\mathbf{w}^{(t)}) - \nabla h_{i}^{(t)}(\mathbf{w}^{(t)}) + \nabla h_{i}^{(t)}(\mathbf{w}^{(t)})\|$$

$$\leq \|\nabla T_{i}^{(t)}(\mathbf{w}_{L}^{(t)}) - \nabla T_{i}^{(t)}(\mathbf{w}^{(t)})\| + \|\nabla T_{i}^{(t)}(\mathbf{w}^{(t)}) - \nabla h_{i}^{(t)}(\mathbf{w}^{(t)})\| + \|\nabla h_{i}^{(t)}(\mathbf{w}^{(t)})\|$$

$$\leq \|-a^{i,(t)}\mathbf{x}_{i}\mathbf{x}_{i}^{T} + \lambda \mathbf{I}\|O((\Delta x)^{2}) + O((\Delta x)^{2}) + c_{1} := c_{2}$$

$$(71)$$

Theorem 13.  $||\mathbb{E}(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n})$ .

Proof. By using the definition of  $\mathbf{w}_{RU}^{(t)}$  and  $\mathbf{w}^{(t)}$ , i.e. Equation (37) and Equation (25), we have:

$$\begin{split} & \mathbb{E}(\|\mathbf{w}_{RU}^{(t+1)} - \mathbf{w}^{(t+1)}\|) \\ &= \mathbb{E}(\|\mathbf{w}_{RU}^{(t)} - \eta(\lambda \mathbf{w}_{RU}^{(t)} + \frac{1}{B} \sum_{i \in \mathcal{B}^{(t)}} y_i \mathbf{x}_i f(y_i \mathbf{w}_{RU}^{(t)T} \mathbf{x}_i)) - [\mathbf{w}^{(t)} - \eta \lambda \mathbf{w}^{(t)} - \eta \frac{1}{B_U^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} y_i \mathbf{x}_i f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)]\|) \\ &= \mathbb{E}(\|(1 - \eta \lambda)(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)}) + \frac{\eta}{B} \sum_{i \in \mathcal{B}^{(t)}} y_i \mathbf{x}_i [f(y_i \mathbf{w}_{RU}^{(t)T} \mathbf{x}_i) - f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)] + \frac{\eta}{B} \sum_{i \in \mathcal{B}^{(t)}} y_i \mathbf{x}_i f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i) \\ &- \frac{\eta}{B_U^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} y_i \mathbf{x}_i f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)\|) \end{split}$$

Then by using the cauchy mean value theorem over  $f(y_i \mathbf{w}_{RU}^{(t)T} \mathbf{x}_i) - f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)$ , the formula above is bounded as:

$$\leq \mathbb{E}(\|[(1-\eta\lambda)\mathbf{I} + \frac{\eta}{B}\sum_{i\in\mathscr{B}^{(t)}}\mathbf{x}_{i}\mathbf{x}_{i}^{T}f'(p)](\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\|)$$

$$+ \mathbb{E}(\|\frac{\eta}{B}\sum_{i\in\mathscr{B}^{(t)}}y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - \frac{\eta}{B_{U}^{(t)}}\sum_{i\in\mathscr{B}^{(t)},i\notin\mathscr{R}}y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i})\|)$$

$$= \mathbb{E}(\|[(1-\eta\lambda)\mathbf{I} + \frac{\eta}{B}\sum_{i\in\mathscr{B}^{(t)}}\mathbf{x}_{i}\mathbf{x}_{i}^{T}f'(p)](\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\|) + \mathbb{E}(\|\frac{\eta}{B}\sum_{i\in\mathscr{B}^{(t)}}\nabla h_{i}(\mathbf{w}^{(t)}) - \frac{\eta}{B_{U}^{(t)}}\sum_{i\in\mathscr{B}^{(t)},i\notin\mathscr{R}}\nabla h_{i}(\mathbf{w}^{(t)})\|)$$

By rewriting the formula above and using the upper bound on  $\|\nabla h_i(\mathbf{w}^{(t)})\|$ , we get:

$$\begin{split} &= \mathbb{E}(\|[(1-\eta\lambda)\mathbf{I} + \frac{\eta}{B}\sum_{i\in\mathcal{B}^{(t)}}\mathbf{x}_{i}\mathbf{x}_{i}^{T}f'(p)](\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\|) + \mathbb{E}(\|\frac{\eta}{B}\sum_{i\in\mathcal{B}^{(t)},} \nabla h_{i}(\mathbf{w}^{(t)}) + (\frac{\eta}{B} - \frac{\eta}{B_{U}^{(t)}})\sum_{i\in\mathcal{B}^{(t)},} \nabla h_{i}(\mathbf{w}^{(t)})\|) \\ &\leq \mathbb{E}(\|[(1-\eta\lambda)\mathbf{I} + \frac{\eta}{B}\sum_{i\in\mathcal{B}^{(t)}}\mathbf{x}_{i}\mathbf{x}_{i}^{T}f'(p)](\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\|) + \mathbb{E}(\frac{\eta}{B}\sum_{i\in\mathcal{B}^{(t)},} c_{1} + (\frac{\eta}{B} - \frac{\eta}{B_{U}^{(t)}})\sum_{i\in\mathcal{B}^{(t)},i\notin\mathcal{R}} c_{1}) \\ &= \mathbb{E}(\|[(1-\eta\lambda)\mathbf{I} + \frac{\eta}{B}\sum_{i\in\mathcal{B}^{(t)}}\mathbf{x}_{i}\mathbf{x}_{i}^{T}f'(p)](\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\|) + \mathbb{E}(\frac{2\eta\Delta B^{(t)}}{B}c_{1}) \end{split}$$

Then by using the result from Equation (69) and Equation (54), the formula above is bounded as:

$$\leq (1 - \eta \lambda) \|\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)}\| + 2\eta c_1 \frac{\Delta n}{n}$$

By applying the formula recursively, we get:

$$\leq 2 \frac{1}{c_1 \lambda} \frac{\Delta n}{n} = O(\frac{\Delta n}{n})$$

Theorem 14. (It is Theorem 5 in the paper)  $||E(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$ , where  $\Delta n$  is the number of the removed samples and  $\Delta x$  is defined in Theorem 12

Proof. By using the definition of  $\mathbf{w}_{LU}$  and  $\mathbf{w}_{RU}$  and subtracting the former one from the latter one, we have:

$$\begin{split} & \mathbb{E}(\|\mathbf{w}_{LU}^{(t+1)} - \mathbf{w}_{RU}^{(t+1)}\|_{2}) \\ & = \mathbb{E}(\|\mathbf{w}_{LU}^{(t)} - \lambda \eta \mathbf{w}_{LU}^{(t)} - \eta \nabla R^{(t)}(\mathbf{w}_{LU}^{(t)}) - (\mathbf{w}_{RU}^{(t)} - \lambda \eta \mathbf{w}_{RU}^{(t)} - \eta \nabla g^{(t)}(\mathbf{w}_{RU}^{(t)}))\|_{2}) \\ & = \mathbb{E}(\|\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)} - \eta (\nabla R^{(t)}(\mathbf{w}_{LU}^{(t)}) - \nabla R^{(t)}(\mathbf{w}_{RU}^{(t)})) - \lambda \eta (\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)}) - \eta (\nabla R^{(t)}(\mathbf{w}_{RU}^{(t)}) - \nabla g^{(t)}(\mathbf{w}_{RU}^{(t)}))\|_{2}) \\ & \leq \mathbb{E}(\|[\mathbf{I} - \eta \frac{1}{B_{U}^{(t)}} \sum_{i \in \mathcal{B}^{(t)}, i \notin \mathcal{R}} (-a^{i,(t)} \mathbf{x}_{i} \mathbf{x}_{i}^{T} + \lambda \mathbf{I})](\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})\| + \eta \|\nabla R^{(t)}(\mathbf{w}_{RU}^{(t)}) - \nabla g^{(t)}(\mathbf{w}_{RU}^{(t)})\|) \end{split}$$

Then by using the result from Equation (68), the formula above is bounded as:

$$\leq (1 - \eta \lambda) \|\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)}\| + \eta \mathbb{E}(\|\nabla R^{(t)}(\mathbf{w}_{RU}^{(t)}) - \nabla g^{(t)}(\mathbf{w}_{RU}^{(t)})\|) 
= (1 - \eta \lambda) \|\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)}\| + \eta \mathbb{E}(\|\frac{1}{B_{U}^{(t)}} \sum_{\substack{i \in \mathcal{B}^{(t)} \\ i \notin \mathcal{R}}} y_{i} \mathbf{x}_{i} [s(y_{i} \mathbf{w}_{RU}^{(t)T} \mathbf{x}_{i}) - f(y_{i} \mathbf{w}_{RU}^{(t)T} \mathbf{x}_{i})]\|)$$
(72)

in which we bound  $y_i \mathbf{x}_i s(y_i \mathbf{w}_{RU}^{(t)T} \mathbf{x}_i) - y_i \mathbf{x}_i f(y_i \mathbf{w}_{RU}^{(t)T} \mathbf{x}_i)$  as below:

$$\begin{aligned} &\|y_{i}\mathbf{x}_{i}s(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i})\| \\ &= \|y_{i}\mathbf{x}_{i}s(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}s(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) + y_{i}\mathbf{x}_{i}s(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) + y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) + y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) + y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i})\| + \|y_{i}\mathbf{x}_{i}s(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i})\| \\ &= \|a^{i,(t)}\mathbf{x}_{i}\mathbf{x}_{i}^{T}(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)}) + y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - y_{i}\mathbf{x}_{i}f(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i})\| + O((\Delta x)^{2}) \end{aligned}$$

The last step uses the result from Equation (63). Then by using the Cauchy mean value theorem on  $f(y_i \mathbf{w}^{(t)T} \mathbf{x}_i) - f(y_i \mathbf{w}_{RU}^{(t)T} \mathbf{x}_i)$ , we know that:

$$= \|a^{i,(t)}\mathbf{x}_{i}\mathbf{x}_{i}^{T}(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)}) + y_{i}\mathbf{x}_{i}[\int_{0}^{1} f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i} + x(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i} - y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}))dx](y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i} - y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i})\| + O((\Delta x)^{2})$$

$$\leq \|[a^{i,(t)} - \int_{0}^{1} f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i} + x(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i} - y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}))dx]\mathbf{x}_{i}\mathbf{x}_{i}^{T}\|\|(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\| + O((\Delta x)^{2})$$

Then by adding and subtracting  $f'(y_i \mathbf{w}^{(t)T} \mathbf{x}_i)$  in the first term and using the fact from Equation (64) and Assumption 3, the formula above is bounded as:

$$= \|[a^{i,(t)} - f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) + f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - \int_{0}^{1} f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i} + x(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i} - y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}))dx]\mathbf{x}_{i}\mathbf{x}_{i}^{T}\|\|(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\| + O((\Delta x)^{2})$$

$$\leq [\|(a^{i,(t)} - f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}))\mathbf{x}_{i}\mathbf{x}_{i}^{T}\| + \|\mathbf{x}_{i}\mathbf{x}_{i}^{T}\int_{0}^{1} [f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}) - f'(y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i} + x(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i} - y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i}))]dx\|]\|(\mathbf{w}_{RU}^{(t)T} - \mathbf{w}^{(t)T})\|$$

$$+ O((\Delta x)^{2})$$

$$\leq O((\Delta x))\|(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\| + \|\mathbf{x}_{i}\mathbf{x}_{i}^{T}\int_{0}^{1} c_{2}x(y_{i}\mathbf{w}_{RU}^{(t)T}\mathbf{x}_{i} - y_{i}\mathbf{w}^{(t)T}\mathbf{x}_{i})dx\|\|(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\| + O((\Delta x)^{2})$$

$$\leq O((\Delta x))\|(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\| + \frac{c_{2}}{2}\|\mathbf{x}_{i}\mathbf{x}_{i}^{T}y_{i}\mathbf{x}_{i}^{T}\|\|(\mathbf{w}_{RU}^{(t)} - \mathbf{w}^{(t)})\|^{2} + O((\Delta x)^{2})$$

Then by using the result from Theorem 13, the formula above is bounded as:

$$\leq O((\Delta x))O(\frac{\Delta n}{n}) + \frac{c_2}{2} \|\mathbf{x}_i \mathbf{x}_i^T y_i \mathbf{x}_i^T \| (O(\frac{\Delta n}{n}))^2 + O((\Delta x)^2) = O(\frac{\Delta n}{n} \Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$$

which is then plugged into Equation (72), we have:

$$\mathbb{E}(\|\mathbf{w}_{LU}^{(t+1)} - \mathbf{w}_{RU}^{(t+1)}\|_{2})$$

$$\leq (1 - \eta\lambda)\|\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)}\| + \eta[O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^{2}) + O((\Delta x)^{2})]$$

which is then used recursively. Then we have:

$$\leq \frac{1}{\lambda} \left[ O(\frac{\Delta n}{n} \Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2) \right] = O(\frac{\Delta n}{n} \Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$$

Theorem 15. Approximation ratio (It is Theorem 6 in the paper) Under the convergence conditions for  $\mathbf{w}^{(t)}$ ,  $||\mathbf{w}^{(t)}||$  should be bounded by some constant C. Suppose  $\frac{||\mathbf{U}_{1,...}^{(t)}\mathbf{S}_{1,...}^{(t)}\mathbf{V}_{1,...}^{T,(t)}||_2}{||\mathbf{U}^{(t)}\mathbf{S}^{(t)}\mathbf{V}^{T,(t)}||_2} \geq 1 - \epsilon$  where  $\epsilon$  is a small value, then the change of model parameters caused by the approximation will be bounded by  $O(\epsilon)$ .

PROOF. The approximate update rule for linear regression by using SVD after removing subsets of training samples is:

$$\mathbf{w}_{U}^{(t+1)'} \leftarrow \left[ (1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}}(\mathbf{U}_{1..r}^{(t)}\mathbf{S}_{1..r}^{(t)}\mathbf{V}_{1..r}^{T,(t)} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}})\right]\mathbf{w}_{U}^{(t)'} + \frac{2\eta_{t}}{B_{U}^{(t)}}(\sum_{i\in\mathscr{B}^{(t)}}\mathbf{x}_{i}y_{i} - \sum_{i\in\mathscr{B}^{(t)}, i\in\mathscr{R}}\mathbf{x}_{i}y_{i})$$

$$(73)$$

By comparing Equation (39) against this approximate update rule, the only difference is  $\mathbf{U}_{1..r}^{(t)}\mathbf{S}_{1..r}^{(t)}\mathbf{V}_{1..r}^{T,(t)}$  in Equation (73) and  $\sum_{i\in\mathscr{B}^{(t)}}\mathbf{x}_i\mathbf{x}_i^T$  in Equation (39). Then according to the condition in this theorem,  $||\sum_{i\in\mathscr{B}^{(t)}}\mathbf{x}_i\mathbf{x}_i^T - \mathbf{U}_{1..r}^{(t)}\mathbf{S}_{1..r}^{(t)}\mathbf{V}_{1..r}^{T,(t)}||_2 = ||\mathbf{U}^{(t)}\mathbf{S}^{(t)}\mathbf{V}^{T,(t)} - \mathbf{U}_{1..r}^{(t)}\mathbf{S}_{1..r}^{(t)}\mathbf{V}_{1..r}^{T,(t)}||_2 = O(\epsilon)$ . So by subtracting Equation (39) by Equation (73), the result becomes:

$$||\mathbf{w}_{U}^{(t+1)'} - \mathbf{w}_{U}^{(t+1)}||_{2} \leftarrow ||[(1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}}(\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i}\mathbf{x}_{i}^{T} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}})](\mathbf{w}_{U}^{(t)'} - \mathbf{w}_{U}^{(t)})$$

$$+ (\frac{2\eta_{t}}{B_{U}^{(t)}}(\mathbf{U}_{1...r}^{(t)}\mathbf{S}_{1...r}^{T}\mathbf{V}_{1...r}^{T,(t)} - \sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i}\mathbf{x}_{i}^{T}))\mathbf{w}_{U}^{(t)}||_{2}$$

$$\leq ||[(1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}}(\sum_{i \in \mathscr{B}^{(t)}} \mathbf{x}_{i}\mathbf{x}_{i}^{T} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}})]||_{2}||(\mathbf{w}_{U}^{(t)'} - \mathbf{w}_{U}^{(t)})||_{2} + O(\epsilon)$$

$$(74)$$

Then we evaluate the expectation between both sides, which ends up with:

$$\mathbb{E}(||\mathbf{w}_{U}^{(t+1)'} - \mathbf{w}_{U}^{(t+1)}||_{2})$$

$$= ||[(1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{n - \Lambda n}(\mathbf{X}^{T}\mathbf{X} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}})]||_{2}||(\mathbf{w}_{U}^{(t)'} - \mathbf{w}_{U}^{(t)})||_{2} + O(\epsilon)$$
(75)

in which  $\lambda \mathbf{I} + \frac{2}{n-\Delta n} (\mathbf{X}^T \mathbf{X} - \Delta \mathbf{X}^T_{\mathscr{B}^{(t)}} \Delta \mathbf{X}_{\mathscr{B}^{(t)}})$  equals to the hessian matrix for the object function over the remaining samples after removing subsets of samples, which should be still  $\lambda$ -strong convex and L-lipschitz continuous. It indicates that:

$$1 - \eta_t \lambda \ge ||(1 - \lambda \eta_t)\mathbf{I} - \frac{2\eta_t}{n - \Lambda n}(\mathbf{X}^T \mathbf{X} - \Delta \mathbf{X}_{\mathscr{B}^{(t)}}^T \Delta \mathbf{X}_{\mathscr{B}^{(t)}})||_2 \ge 1 - \eta_t L$$
(76)

Since according to Lemma 2,  $\eta_t L < 1$ , then:

$$1 > 1 - \eta_t \lambda \ge ||(1 - \lambda \eta_t)\mathbf{I} - \frac{2\eta_t}{n - \Delta n}(\mathbf{X}^T \mathbf{X} - \Delta \mathbf{X}_{\mathscr{B}^{(t)}}^T \Delta \mathbf{X}_{\mathscr{B}^{(t)}})||_2 \ge 1 - \eta_t L > 0$$

$$(77)$$

So we can compute Equation (75) recursively and thus the following inequality holds:

$$E(||\mathbf{w}_{U}^{(t+1)'} - \mathbf{w}_{U}^{(t+1)}||_{2}) \le O(\epsilon)$$
(78)

Theorem 16. (Approximation ratio) (It is Theorem 7 in the paper) The approximation of PrIU-opt over the model parameters is bounded by  $O(||\Delta X^T \Delta X||)$ 

PROOF. The update rule through gradient descent for linear regression model is as below:

$$\mathbf{w}_{U}^{(t+1)} \leftarrow ((1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{n - \Delta n}(\mathbf{X}^{T}\mathbf{X} - \Delta\mathbf{X}^{T}\Delta\mathbf{X}))\mathbf{w}_{U}^{(t)} + \frac{2\eta_{t}}{n - \Delta n}(\mathbf{X}^{T}\mathbf{Y} - \Delta\mathbf{X}^{T}\Delta\mathbf{Y})$$

$$(79)$$

while the approximated update rule through the approximations by incremental updates over eigenvalue is:

$$\mathbf{w}_{U}^{(t+1)'} \leftarrow ((1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{n - \Delta n}\mathbf{Q}^{-1}diag[c'_{1}, c'_{2}, \dots, c'_{m}]\mathbf{Q})\mathbf{w}_{U}^{(t)'} + \frac{2\eta_{t}}{n - \Delta n}(\mathbf{X}^{T}\mathbf{Y} - \Delta\mathbf{X}^{T}\Delta\mathbf{Y})$$

$$(80)$$

According to [41], the difference between  $\mathbf{Q}^{-1}diag[c_1', c_2', \dots, c_m']\mathbf{Q}$  and  $\mathbf{X}^T\mathbf{X} - \Delta\mathbf{X}^T\Delta\mathbf{X}$  is bounded by  $O(\Delta\mathbf{X}^T\Delta\mathbf{X})$ . So by subtracting Equation (79) by Equation (80), the results become:

$$||\mathbf{w}_{U}^{(t+1)'} - \mathbf{w}_{U}^{(t+1)'}||_{2} \leftarrow ||[(1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}}(\mathbf{X}^{T}\mathbf{X} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}})](\mathbf{w}_{U}^{(t)'} - \mathbf{w}_{U}^{(t)})$$

$$+ (\frac{2\eta_{t}}{B_{U}^{(t)}}(\mathbf{Q}^{-1}diag[c'_{1}, c'_{2}, \dots, c'_{m}]\mathbf{Q} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}}))\mathbf{w}_{U}^{(t)}||_{2}$$

$$\leq ||[(1 - \eta_{t}\lambda)\mathbf{I} - \frac{2\eta_{t}}{B_{U}^{(t)}}(\sum_{i \in \mathscr{B}^{(t)}}\mathbf{x}_{i}\mathbf{x}_{i}^{T} - \Delta\mathbf{X}_{\mathscr{B}^{(t)}}^{T}\Delta\mathbf{X}_{\mathscr{B}^{(t)}})]||_{2}||(\mathbf{w}_{U}^{(t)'} - \mathbf{w}_{U}^{(t)})||_{2} + O(\Delta\mathbf{X}^{T}\Delta\mathbf{X})$$

$$(81)$$

Through the similar analysis to Theorem 15, the above formula is computed recursively, which ends up with:

$$||\mathbf{w}_{U}^{(t+1)'} - \mathbf{w}_{U}^{(t+1)}||_{2} \le O(\Delta \mathbf{X}^{T} \Delta \mathbf{X})$$
(82)

Theorem 17. (Approximation ratio) (It is Theorem 8 in the paper) Similar to Theorem 15, the deviation caused by the SVD approximation will be bounded by  $O(\epsilon)$ , given the ratio  $\frac{||P_{1...r}^{(t)}V_{1...r}^{T,(t)}||_2}{||P^{(t)}V^{T,(t)}||_2} \ge 1 - \epsilon$ . So using Theorem 14,  $||E(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2) + O(\epsilon)$ .

PROOF. Let's assume that the incremental updated model parameter without SVD approximation is  $\mathbf{w}_{LU0}^{(t)}$ . According to the results in Theorem 14, the  $||\mathbf{w}_{LU0}^{(t)} - \mathbf{w}_{RU}^{(t)}|| \le O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$ . After the SVD approximation, similar analysis to Theorem 15 can be done such that  $||\mathbf{w}_{LU0}^{(t)} - \mathbf{w}_{LU}^{(t)}||_2 \le O(\epsilon)$ . So  $||E(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})||_2 \le O(\epsilon) + O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$ 

Theorem 18. (Approximation ratio) (It is Theorem 9 in the paper) Suppose that after iteration  $t_s$  the gradient of the objective function is smaller than  $\delta$ , then the approximations of PrIU-opt can lead to deviations of the model parameters bounded by  $O((\tau - t_s)\delta) + O(||\Delta X^T \Delta X||)$ . By combining the analysis in Theorem 14,  $||E(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})||_2$  is bounded by  $O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2) + O((\tau - t_s)\delta) + O(||\Delta X^T \Delta X||)$ 

PROOF. Let's assume that after  $t_s^{th}$  iteration, the incremental updated model parameter without only linearization approximation is  $\mathbf{w}_{LU\,0}^{(t)}$ . Then  $||\mathbb{E}(\mathbf{w}_{LU\,0}^{(t)} - \mathbf{w}_{RU}^{(t)})||_2 \le O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^2) + O((\Delta x)^2)$  based on Theorem 14. Also let's assume that after  $t_s^{th}$  iteration, the incremental updated model parameter with both linearization approximation and SVD approximation is  $\mathbf{w}_{LU\,1}^{(t)}$ . Then:

$$||\mathbf{w}_{LU1}^{(t)} - \mathbf{w}_{LU0}^{(t)}||_{2} = ||\mathbf{w}_{LU}^{(t-1)}||_{1} - \nabla R^{(t-1)}(\mathbf{w}_{LU1}^{(t-1)}) - (\mathbf{w}_{LU0}^{(t)} - \nabla R^{(t-1)}(\mathbf{w}_{LU0}^{(t-1)}))||_{2}$$

$$\leq ||\mathbf{w}_{LU1}^{(t-1)} - \mathbf{w}_{LU0}^{(t-1)}||_{2} + 2\delta \leq ||\mathbf{w}_{LU1}^{(t_s)} - \mathbf{w}_{LU0}^{(t_s)}||_{2} + 2(t - t_s - 1)\delta = O((t - t_s)\delta)$$
(83)

Finally, by using Theorem 16,  $||\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{LU1}^{(t)}||_2 \le O(||\Delta \mathbf{X}^T \Delta \mathbf{X}||)$ . By combining those results together, we have:

$$||E(\mathbf{w}_{LU}^{(t)} - \mathbf{w}_{RU}^{(t)})||_{2} \le O(\frac{\Delta n}{n}\Delta x) + O((\frac{\Delta n}{n})^{2}) + O((\Delta x)^{2}) + O((\tau - t_{s})\delta) + O(||\Delta \mathbf{X}^{T}\Delta \mathbf{X}||)$$
(84)