
Privacy-Preserving Logistic Regression Training with A Faster Gradient Variant

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Logistic regression training over encrypted data has been an attractive idea to
2 security concerns for years. In this paper, we propose a faster gradient variant
3 called quadratic gradient to implement logistic regression training in a ho-
4 momorphic encryption domain, the core of which can be seen as an extension
5 of the simplified fixed Hessian [5]. We enhance Nesterov’s accelerated gradient
6 (NAG) and Adaptive Gradient Algorithm (Adagrad) respectively with this gradient
7 variant and evaluate the enhanced algorithms on several datasets. Experimental
8 results show that the enhanced methods have a state-of-the-art performance in
9 convergence speed compared to the naive first-order gradient methods. We then
10 adopt the enhanced NAG method to implement homomorphic logistic regression
11 training and obtain a comparable result by only 3 iterations.

12 1 Introduction

13 Logistic regression (LR) is a widely used classification technology especially in medical risk assess-
14 ment due to its simplicity but powerful performance. Given a person’s healthcare data related to a
15 certain disease, we can train an LR model capable of telling whether or not this person is likely to
16 develop this disease. However, such personal health information is highly private to individuals. The
17 privacy concern, therefore, becomes a major obstacle for individuals to share their biomedical data.
18 The most secure solution is to encrypt the data into ciphertexts first by Homomorphic Encryption
19 (HE) and then securely outsource the ciphertexts to the cloud, without allowing the cloud to access
20 the data directly. iDASH is an annual competition that aims to call for implementing interesting cryp-
21 tographic schemes in a biological context. Since 2014, iDASH has included the theme of genomics
22 and biomedical privacy. The third track of the 2017 iDASH competition and the second track of
23 the 2018 iDASH competition were both to develop homomorphic-encryption-based solutions for
24 building an LR model over encrypted data.

25 Several studies on logistic regression models are based on homomorphic encryption. Aono et al. [2]
26 only used an additive HE scheme and left some of the challenging HE computations to a trusted
27 client. Kim et al. [14] discussed the problem of performing LR training in an encrypted environment.
28 They used the full batch gradient descent in the training process and the least square method to get the
29 approximation of the sigmoid function. In the iDASH 2017 competition, Bonte and Vercauteren [5],
30 Kim et al. [12], Chen et al. [6], and Crawford et al. [8] all investigated the same problem that Kim et
31 al. [14] studied. In the iDASH competition of 2018, Kim et al. [13] and Blatt et al. [3] further worked
32 on it for an efficient packing and semi-parallel algorithm. The papers most relevant to this work
33 are [5] and [12]. Bonte and Vercauteren [5] developed a practical algorithm called the simplified
34 fixed Hessian (SFH) method. Our study complements their work and adopts the ciphertext packing
35 technique proposed by Kim et al. [12] for efficient homomorphic computation.

36 Our specific contributions in this paper are as follows:

- 37 1. We propose a new gradient variant, quadratic gradient, which can unite the first-order
38 gradient method and the second-order Newton’s method as one.
- 39 2. We develop two enhanced gradient methods by equipping the original methods with
40 quadratic gradient. The resulting methods show a better performance in the con-
41 vergence speed.
- 42 3. We adopt the enhanced Nesterov’s accelerated gradient to implement privacy-preserving
43 logistical regression training, to our best knowledge, which seems to be the best candidate
44 without compromising much on computation and storage.

45 2 Preliminaries

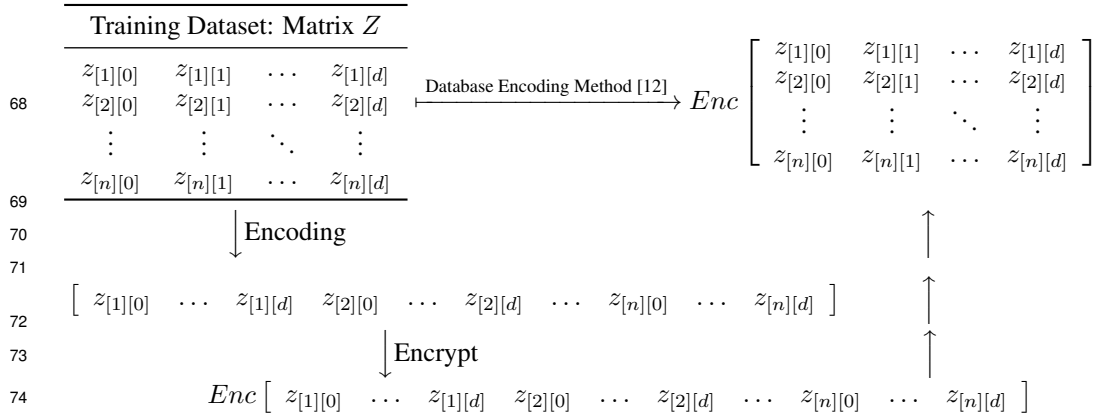
46 We adopt the square brackets “[]” to denote the index of a vector or matrix element in what follows.
47 For example, for a vector $\mathbf{v} \in \mathbb{R}^{(n)}$ and a matrix $M \in \mathbb{R}^{m \times n}$, $\mathbf{v}[i]$ or $\mathbf{v}_{[i]}$ means the i -th element of
48 vector \mathbf{v} and $M[i][j]$ or $M_{[i][j]}$ the j -th element in the i -th row of M .

49 2.1 Fully Homomorphic Encryption

50 Fully Homomorphic Encryption (FHE) is a type of cryptographic scheme that can be used to compute
51 an arbitrary number of additions and multiplications directly on the encrypted data. It was not until
52 2009 that Gentry constructed the first FHE scheme via a bootstrapping operation [9]. FHE schemes
53 themselves are computationally time-consuming; the choice of dataset encoding matters likewise
54 to the efficiency. In addition to these two limits, how to manage the magnitude of plaintext [11]
55 also contributes to the slowdown. Cheon et al. [7] proposed a method to construct an HE scheme
56 with a rescaling procedure which could eliminate this technical bottleneck effectively. We adopt
57 their open-source implementation HEAAN while implementing our homomorphic LR algorithms. It
58 is inevitable to pack a vector of multiple plaintexts into a single ciphertext for yielding a better
59 amortized time of homomorphic computation. HEAAN supports a parallel technique (aka SIMD) to
60 pack multiple numbers in a single polynomial by virtue of the Chinese Remainder Theorem and
61 provides rotation operation on plaintext slots. The underlying HE scheme in HEAAN is well described
62 in [12, 14, 10].

63 2.2 Database Encoding Method

64 Kim et al. [12] devised an efficient and promising database-encoding method by using SIMD technique,
65 which could make full use of the computation and storage resources. Suppose that a database has a
66 training dataset consisting of n samples with $(1 + d)$ covariates, they packed the training dataset Z
67 into a single ciphertext in a row-by-row manner:



75 Using this encoding scheme, we can manipulate the data matrix Z by performing HE operations on the
76 ciphertext $\text{Enc}[Z]$, with the help of only three HE operations - rotation, addition and multiplication.
77 For example, if we want the first column of $\text{Enc}[Z]$ alone and filter out the other columns, we can

78 design a constant matrix F consisting of ones in the first column and zeros in the rest columns and
 79 then multiply $Enc[Z]$ by $Enc[F]$, obtaining the resulting ciphertext $Enc[Z_p]$:

$$Enc[F] \otimes Enc[Z] = Enc[Z_p] \quad (\text{where “}\otimes\text{” means the component-wise HE multiplication})$$

$$= Enc \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix} \otimes Enc \begin{bmatrix} z_{[1][0]} & z_{[1][1]} & \dots & z_{[1][d]} \\ z_{[2][0]} & z_{[2][1]} & \dots & z_{[2][d]} \\ \vdots & \vdots & \ddots & \vdots \\ z_{[n][0]} & z_{[n][1]} & \dots & z_{[n][d]} \end{bmatrix} = Enc \begin{bmatrix} z_{[1][0]} & 0 & \dots & 0 \\ z_{[2][0]} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ z_{[n][0]} & 0 & \dots & 0 \end{bmatrix}.$$

80 Han et al. [10] introduced several operations to manipulate the ciphertexts, such as a procedure named
 81 “SumColVec” to compute the summation of the columns of a matrix. By dint of these basic operations,
 82 more complex calculations such as computing the gradients in logistic regression models are feasible.

83 2.3 Logistic Regression

84 Logistic regression is widely used in binary classification tasks to infer whether a binary-valued
 85 variable belongs to a certain class or not. LR can be generalized from linear regression [15] by
 86 mapping the whole real line ($\beta^T \mathbf{x}$) to $(0, 1)$ via the sigmoid function $\sigma(z) = 1/(1 + \exp(-z))$, where
 87 the vector $\beta \in \mathbb{R}^{(1+d)}$ is the main parameter of LR and the vector $\mathbf{x} = (1, x_1, \dots, x_d) \in \mathbb{R}^{(1+d)}$ the
 88 input covariate. Thus logistic regression can be formulated with the class label $y \in \{\pm 1\}$ as follows:

$$\Pr(y = +1 | \mathbf{x}, \beta) = \sigma(\beta^T \mathbf{x}) = \frac{1}{1 + e^{-\beta^T \mathbf{x}}},$$

$$\Pr(y = -1 | \mathbf{x}, \beta) = 1 - \sigma(\beta^T \mathbf{x}) = \frac{1}{1 + e^{+\beta^T \mathbf{x}}}.$$

89 LR sets a threshold (usually 0.5) and compares its output with it to decide the resulting class label.

90 The logistic regression problem can be transformed into an optimization problem that seeks a param-
 91 eter β to maximize $L(\beta) = \prod_{i=1}^n \Pr(y_i | \mathbf{x}_i, \beta)$ or its log-likelihood function $l(\beta)$ for convenience in
 92 the calculation:

$$l(\beta) = \ln L(\beta) = - \sum_{i=1}^n \ln(1 + e^{-y_i \beta^T \mathbf{x}_i}),$$

93 where n is the number of examples in the training dataset. LR does not have a closed form of
 94 maximizing $l(\beta)$ and two main methods are adopted to estimate the parameters of an LR model:
 95 (a) gradient descent method via the gradient; and (b) Newton’s method by the Hessian matrix. The
 96 gradient and Hessian of the log-likelihood function $l(\beta)$ are given by, respectively:

$$\nabla_{\beta} l(\beta) = \sum_i (1 - \sigma(y_i \beta^T \mathbf{x}_i)) y_i \mathbf{x}_i,$$

$$\nabla_{\beta}^2 l(\beta) = \sum_i (y_i \mathbf{x}_i) (\sigma(y_i \beta^T \mathbf{x}_i) - 1) \sigma(y_i \beta^T \mathbf{x}_i) (y_i \mathbf{x}_i)$$

$$= X^T S X$$

97 where S is a diagonal matrix with entries $S_{ii} = (\sigma(y_i \beta^T \mathbf{x}_i) - 1) \sigma(y_i \beta^T \mathbf{x}_i)$ and X the dataset.

98 The log-likelihood function $l(\beta)$ of LR has at most a unique global maximum [1], where its gradient
 99 is zero. Newton’s method is a second-order technique to numerically find the roots of a real-valued
 100 differentiable function, and thus can be used to solve the β in $\nabla_{\beta} l(\beta) = 0$ for LR.

101 3 Technical Details

It is quite time-consuming to compute the Hessian matrix and its inverse in Newton’s method for
 each iteration. One way to limit this downside is to replace the varying Hessian with a fixed matrix
 \bar{H} . This novel technique is called the fixed Hessian Newton’s method. Böhning and Lindsay [4] have
 shown that the convergence of Newton’s method is guaranteed as long as $\bar{H} \leq \nabla_{\beta}^2 l(\beta)$, where \bar{H} is

a symmetric negative-definite matrix independent of β and “ \leq ” denotes the Loewner ordering in the sense that the difference $\nabla_{\beta}^2 l(\beta) - \bar{H}$ is non-negative definite. With such a fixed Hessian matrix \bar{H} , the iteration for Newton’s method can be simplified to:

$$\beta_{t+1} = \beta_t - \bar{H}^{-1} \nabla_{\beta} l(\beta).$$

102 Böhning and Lindsay also suggest the fixed matrix $\bar{H} = -\frac{1}{4} X^T X$ is a good lower bound for the
 103 Hessian of the log-likelihood function $l(\beta)$ in LR.

104 3.1 the Simplified Fixed Hessian method

Bonte and Vercauteren [5] simplify this lower bound \bar{H} further due to the need for inverting the fixed Hessian in the encrypted domain. They replace the matrix \bar{H} with a diagonal matrix B whose diagonal elements are simply the sums of each row in \bar{H} . They also suggest a specific order of calculation to get B more efficiently. Their new approximation B of the fixed Hessian is:

$$B = \begin{bmatrix} \sum_{i=0}^d \bar{h}_{0i} & 0 & \dots & 0 \\ 0 & \sum_{i=0}^d \bar{h}_{1i} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sum_{i=0}^d \bar{h}_{di} \end{bmatrix},$$

105 where \bar{h}_{ki} is the element of \bar{H} . This diagonal matrix B is in a very simple form and can be obtained
 106 from \bar{H} without much difficulty. The inverse of B can be approximated in the encrypted form by
 107 means of computing the inverse of every diagonal element of B via the iterative of Newton’s method
 108 with an appropriate start value. Their simplified fixed Hessian method can be formulated as follows:

$$\begin{aligned} \beta_{t+1} &= \beta_t - B^{-1} \cdot \nabla_{\beta} l(\beta), \\ &= \beta_t - \begin{bmatrix} b_{00} & 0 & \dots & 0 \\ 0 & b_{11} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & b_{dd} \end{bmatrix} \cdot \begin{bmatrix} \nabla_0 \\ \nabla_1 \\ \vdots \\ \nabla_d \end{bmatrix} = \beta_t - \begin{bmatrix} b_{00} \cdot \nabla_0 \\ b_{11} \cdot \nabla_1 \\ \vdots \\ b_{dd} \cdot \nabla_d \end{bmatrix}, \end{aligned}$$

109 where b_{ii} is the reciprocal of $\sum_{i=0}^d \bar{h}_{0i}$ and ∇_i is the element of $\nabla_{\beta} l(\beta)$.

110 Consider a special situation: if b_{00}, \dots, b_{dd} are all the same value $-\eta$ with $\eta > 0$, the iterative formula
 111 of the SFH method can be given as:

$$\beta_{t+1} = \beta_t - (-\eta) \cdot \begin{bmatrix} \nabla_0 \\ \nabla_1 \\ \vdots \\ \nabla_d \end{bmatrix} = \beta_t + \eta \cdot \nabla_{\beta} l(\beta),$$

112 which is the same as the formula of the naive gradient *ascent* method. Such coincident is just what
 113 the idea behind this work comes from: there is some relation between the Hessian matrix and the
 114 learning rate of the gradient (descent) method. We consider $b_{ii} \cdot \nabla_i$ as a new enhanced gradient
 115 variant’s element and assign a new learning rate to it. As long as we ensure that this new learning
 116 rate decreases from a positive floating-point number greater than 1 (such as 2) to 1 in a bounded
 117 number of iteration steps, the fixed Hessian Newton’s method guarantees the algorithm will converge
 118 eventually.

119 The SFH method proposed by Bonte and Vercauteren [5] has two limitations: (a) in the construction
 120 of the simplified fixed Hessian matrix, all entries in the symmetric matrix \bar{H} need to be non-positive.
 121 For machine learning applications the datasets will be in advance normalized into the range [0,1],
 122 meeting the convergence condition of the SFH method. However, for other cases such as numerical
 123 optimization, it doesn’t always hold; and (b) the simplified fixed Hessian matrix B that Bonte and
 124 Vercauteren [5] constructed can still be singular, especially when the dataset is a high-dimensional
 125 sparse matrix, such as the datasets from the MNIST database. We complement their work by removing
 126 these limitations so as to generalize this simplified fixed Hessian to be invertible in any case and
 127 propose a faster gradient variant, which we term `quadratic gradient`.

128 3.2 Quadratic Gradient

129 Suppose that a differentiable scalar-valued function $F(\mathbf{x})$ has its gradient \mathbf{g} and Hessian matrix H ,
 130 with any matrix $\bar{H} \leq H$ in the Loewner ordering as follows:

$$\mathbf{g} = \begin{bmatrix} g_0 \\ g_1 \\ \vdots \\ g_d \end{bmatrix}, \quad H = \begin{bmatrix} \nabla_{00}^2 & \nabla_{01}^2 & \cdots & \nabla_{0d}^2 \\ \nabla_{10}^2 & \nabla_{11}^2 & \cdots & \nabla_{1d}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \nabla_{d0}^2 & \nabla_{d1}^2 & \cdots & \nabla_{dd}^2 \end{bmatrix}, \quad \bar{H} = \begin{bmatrix} \bar{h}_{00} & \bar{h}_{01} & \cdots & \bar{h}_{0d} \\ \bar{h}_{10} & \bar{h}_{11} & \cdots & \bar{h}_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ \bar{h}_{d0} & \bar{h}_{d1} & \cdots & \bar{h}_{dd} \end{bmatrix},$$

131 where $\nabla_{ij}^2 = \nabla_{ji}^2 = \frac{\partial^2 F}{\partial x_i \partial x_j}$. We construct a new Hessian matrix \tilde{B} as follows:

$$\tilde{B} = \begin{bmatrix} -\varepsilon - \sum_{i=0}^d |\bar{h}_{0i}| & 0 & \cdots & 0 \\ 0 & -\varepsilon - \sum_{i=0}^d |\bar{h}_{1i}| & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -\varepsilon - \sum_{i=0}^d |\bar{h}_{di}| \end{bmatrix},$$

132 where ε is a small positive constant to avoid division by zero (usually set to $1e-8$).

133 As long as \tilde{B} satisfies the convergence condition of the above fixed Hessian method, $\tilde{B} \leq H$, we
 134 can use this approximation \tilde{B} of the Hessian matrix as a lower bound. Since we already assume that
 135 $\bar{H} \leq H$, it will suffice to show that $\tilde{B} \leq \bar{H}$. We prove $\tilde{B} \leq \bar{H}$ in a similar way that [5] did.

136 **Lemma 1.** *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix, and let B be the diagonal matrix whose diagonal*
 137 *entries $B_{kk} = -\varepsilon - \sum_i |A_{ki}|$ for $k = 1, \dots, n$, then $B \leq A$.*

138 *Proof.* By definition of the Loewner ordering, we have to prove the difference matrix $C = A - B$
 139 is non-negative definite, which means that all the eigenvalues of C need to be non-negative. By
 140 construction of C we have that $C_{ij} = A_{ij} + \varepsilon + \sum_{k=1}^n |A_{ik}|$ for $i = j$ and $C_{ij} = A_{ij}$ for $i \neq j$.
 141 By means of Gerschgorin's circle theorem, we can bound every eigenvalue λ of C in the sense that
 142 $|\lambda - C_{ii}| \leq \sum_{i \neq j} |C_{ij}|$ for some index $i \in \{1, 2, \dots, n\}$. We conclude that $\lambda \geq A_{ii} + \varepsilon + |A_{ii}| \geq$
 143 $\varepsilon > 0$ for all eigenvalues λ and thus that $B \leq A$. \square

144 **Definition 3.1** (Quadratic Gradient). *Given such a \tilde{B} above, we define the quadratic gradient*
 145 *as $G = \tilde{B} \cdot \mathbf{g}$ with a new learning rate η , where \tilde{B} is a diagonal matrix with diagonal entries*
 146 *$\tilde{B}_{kk} = 1/|\tilde{B}_{kk}|$, and η should be always no less than 1 and decrease to 1 in a limited number*
 147 *of iteration steps. Note that G is still a column vector of the same size as the gradient \mathbf{g} . To*
 148 *maximize or minimize the function $F(\mathbf{x})$, we can use the iterative formulas: $\mathbf{x}_{k+1} = \mathbf{x}_k + \eta \cdot G$ or*
 149 *$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \cdot G$, just like the naive gradient.*

150 We point out here that \bar{H} could be the Hessian matrix H itself and \tilde{B} further optimized to: $\tilde{B}_{kk} =$
 151 $\bar{h}_{kk} + |\bar{h}_{kk}| + \varepsilon - \sum_{i=0}^d |\bar{h}_{ki}|$. In our experiments, we use $\bar{H} = -\frac{1}{4} X^T X$ to construct our \tilde{B} .

152 3.3 Two Enhanced Methods

153 Quadratic Gradient can be used to enhance NAG and Adagrad.

154 NAG is a different variant of the momentum method to give the momentum term much more
 155 prescience. The iterative formulas of the gradient ascent method for NAG are as follows:

$$V_{t+1} = \beta_t + \alpha_t \cdot \nabla J(\beta_t), \quad (3)$$

$$\beta_{t+1} = (1 - \gamma_t) \cdot V_{t+1} + \gamma_t \cdot V_t, \quad (4)$$

156 where V_{t+1} is the intermediate variable used for updating the final weight β_{t+1} and $\gamma_t \in (0, 1)$ is a
 157 smoothing parameter of moving average to evaluate the gradient at an approximate future position
 158 [12]. The enhanced NAG is to replace (3) with $V_{t+1} = \beta_t + (1 + \alpha_t) \cdot G$. Our enhanced NAG
 159 method is described in Algorithm 1.

160 Adagrad is a gradient-based algorithm suitable for dealing with sparse data. The updated operations
 161 of Adagrad and its quadratic-gradient version, for every parameter $\beta_{[i]}$ at each iteration step t , are as

Algorithm 1 The enhanced Nesterov's accelerated gradient method

Input: training dataset $X \in \mathbb{R}^{n \times (1+d)}$; training label $Y \in \mathbb{R}^{n \times 1}$; and the number κ of iterations;
Output: the parameter vector $V \in \mathbb{R}^{(1+d)}$

- 1: Set $\bar{H} \leftarrow -\frac{1}{4}X^T X$ $\triangleright \bar{H} \in \mathbb{R}^{(1+d) \times (1+d)}$
- 2: Set $V \leftarrow \mathbf{0}, W \leftarrow \mathbf{0}, \bar{B} \leftarrow \mathbf{0}$ $\triangleright V \in \mathbb{R}^{(1+d)}, W \in \mathbb{R}^{(1+d)}, \bar{B} \in \mathbb{R}^{(1+d) \times (1+d)}$
- 3: **for** $i := 0$ to d **do**
- 4: $\bar{B}[i][i] \leftarrow \varepsilon$ $\triangleright \varepsilon$ is a small positive constant such as $1e-8$
- 5: **for** $j := 0$ to d **do**
- 6: $\bar{B}[i][j] \leftarrow \bar{B}[i][i] + |\bar{H}[i][j]|$
- 7: **end for**
- 8: **end for**
- 9: Set $\alpha_0 \leftarrow 0.01, \alpha_1 \leftarrow 0.5 \times (1 + \sqrt{1 + 4 \times \alpha_0^2})$
- 10: **for** $count := 1$ to κ **do**
- 11: Set $Z \leftarrow \mathbf{0}$ $\triangleright Z \in \mathbb{R}^n$ is the inputs for sigmoid function
- 12: **for** $i := 1$ to n **do**
- 13: **for** $j := 0$ to d **do**
- 14: $Z[i] \leftarrow Z[i] + Y[i] \times V[j] \times X[i][j]$
- 15: **end for**
- 16: **end for**
- 17: Set $\sigma \leftarrow \mathbf{0}$ $\triangleright \sigma \in \mathbb{R}^n$ is to store the outputs of the sigmoid function
- 18: **for** $i := 1$ to n **do**
- 19: $\sigma[i] \leftarrow 1/(1 + \exp(-Z[i]))$
- 20: **end for**
- 21: Set $g \leftarrow \mathbf{0}$
- 22: **for** $j := 0$ to d **do**
- 23: **for** $i := 1$ to n **do**
- 24: $g[j] \leftarrow g[j] + (1 - \sigma[i]) \times Y[i] \times X[i][j]$
- 25: **end for**
- 26: **end for**
- 27: Set $G \leftarrow \mathbf{0}$
- 28: **for** $j := 0$ to d **do**
- 29: $G[j] \leftarrow \bar{B}[j][j] \times g[j]$
- 30: **end for**
- 31: Set $\eta \leftarrow (1 - \alpha_0)/\alpha_1, \gamma \leftarrow 1/(n \times count)$ $\triangleright n$ is the size of training data
- 32: **for** $j := 0$ to d **do**
- 33: $w_{temp} \leftarrow V[j] + (1 + \gamma) \times G[j]$
- 34: $V[j] \leftarrow (1 - \eta) \times w_{temp} + \eta \times W[j]$
- 35: $W[j] \leftarrow w_{temp}$
- 36: **end for**
- 37: $\alpha_0 \leftarrow \alpha_1, \alpha_1 \leftarrow 0.5 \times (1 + \sqrt{1 + 4 \times \alpha_0^2})$
- 38: **end for**
- 39: **return** V

162 follows, respectively:

$$\beta_{[i]}^{(t+1)} = \beta_{[i]}^{(t)} - \frac{\eta}{\varepsilon + \sqrt{\sum_{k=1}^t \mathbf{g}_{[i]}^{(k)} \cdot \mathbf{g}_{[i]}^{(k)}}} \cdot \mathbf{g}_{[i]}^{(t)},$$

$$\beta_{[i]}^{(t+1)} = \beta_{[i]}^{(t)} - \frac{1 + \eta}{\varepsilon + \sqrt{\sum_{k=1}^t G_{[i]}^{(k)} \cdot G_{[i]}^{(k)}}} \cdot G_{[i]}^{(t)}.$$

163 **Performance Evaluation** We evaluate the performance of various algorithms in the clear using the
 164 Python programming language on the same desktop computer with an Intel Core CPU G640 at
 165 1.60 GHz and 7.3 GB RAM. Since our focus is on how fast the algorithms converge in the training
 166 phase, the loss function, maximum likelihood estimation (MLE), is selected as the only indicator. We
 167 evaluate four algorithms, NAG, Adagrad, and their quadratic-gradient versions (denoted as Enhanced
 168 NAG and Enhanced Adagrad, respectively) on the datasets that Kim et al. [12] adopted: the iDASH
 169 genomic dataset (iDASH), the Myocardial Infarction dataset from Edinburgh (Edinburgh), Low Birth
 170 weight Study (lbw), Nhanes III (nhanes3), Prostate Cancer study (pcs), and Umaru Impact Study
 171 datasets (uis). The genomic dataset is provided by the third task in the iDASH competition of 2017,
 172 which consists of 1579 records. Each record has 103 binary genotypes and a binary phenotype
 173 indicating if the patient has cancer. The other five datasets all have a single binary dependent variable.
 174 Figures 1 and 2 show that except for the enhanced Adagrad method on the iDASH genomic dataset
 our enhanced methods all converge faster than their original ones in other cases.

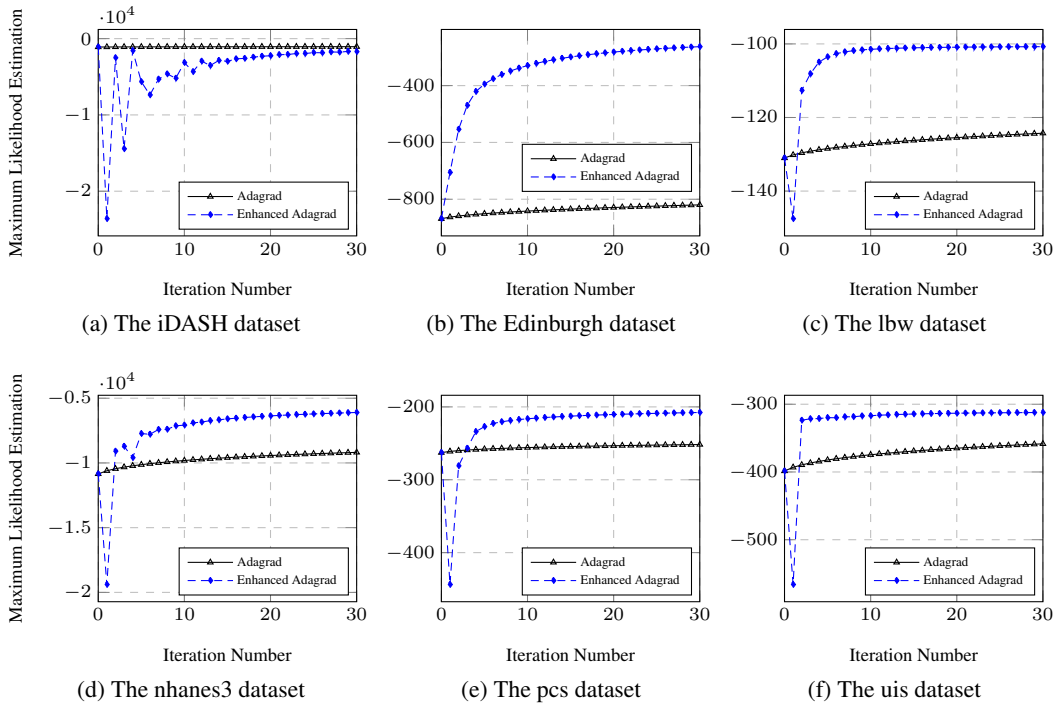


Figure 1: Training results in the clear for Adagrad and Enhanced Adagrad

175

176 In all the Python experiments, the time to calculate the \bar{B} in quadratic gradient G before running the
 177 iterations and the time to run each iteration for various algorithms are negligible (few seconds).

178 4 Privacy-preserving Logistic Regression Training

179 Adagrad method is not a practical solution for homomorphic LR due to its frequent inversion
 180 operations. It seems plausible that the enhanced NAG is probably the best choice for privacy-preserving
 181 LR training. We adopt the enhanced NAG method to implement privacy-preserving
 182 logistic regression training. The difficulty in applying the quadratic gradient is to invert the diagonal

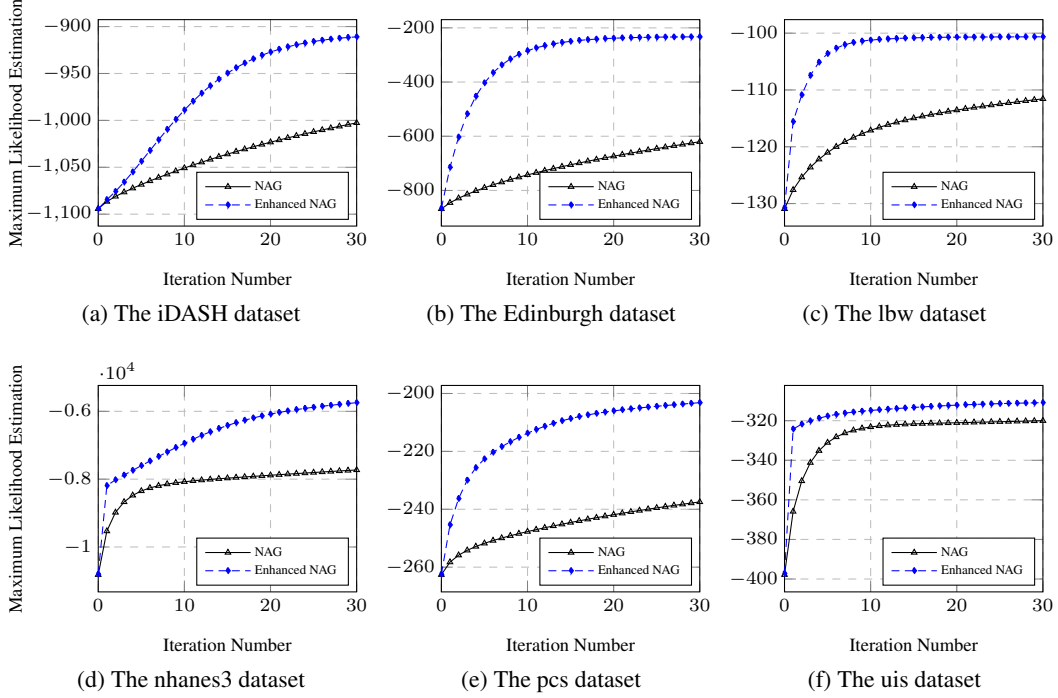


Figure 2: Training results in the clear for NAG and Enhanced NAG

183 matrix \tilde{B} in order to obtain \bar{B} . We leave the computation of matrix \bar{B} to data owner and let the
 184 data owner upload the ciphertext encrypting the \bar{B} to the cloud. Since data owner has to prepare the
 185 dataset and normalize it, it would also be practicable for the data owner to calculate the \bar{B} owing to
 186 no leaking of sensitive data information.

187 Privacy-preserving logistic regression training via homomorphic encryption technique faces a difficult
 188 dilemma that no homomorphic schemes are capable of directly calculating the sigmoid function in the
 189 LR model. A common solution is to replace the sigmoid function with a polynomial approximation
 190 by using the widely adopted least square method. We can call a function named “polyfit(.)”
 191 in the Python package Numpy to fit the polynomial in a least-square sense. We adopt the degree
 192 5 polynomial approximation $g(x)$ by which Kim et al. [12] used the least square approach to
 193 approximate the sigmoid function over the domain $[-8, 8]$: $g(x) = 0.5 + 0.19131 \cdot x - 0.0045963 \cdot$
 194 $x^3 + 0.0000412332 \cdot x^5$.

195 Given the training dataset $X \in \mathbb{R}^{n \times (1+d)}$ and training label $Y \in \mathbb{R}^{n \times 1}$, we adopt the same method
 196 that Kim et al. [12] used to encrypt the data matrix consisting of the training data combined with
 197 training-label information into a single ciphertext ct_Z . The weight vector $\beta^{(0)}$ consisting of zeros and
 198 the diagonal elements of \bar{B} are copied n times to form two matrices. The data owner then encrypts the
 199 two matrices into two ciphertexts $ct_{\beta}^{(0)}$ and $ct_{\bar{B}}$, respectively. The ciphertexts ct_Z , $ct_{\beta}^{(0)}$ and $ct_{\bar{B}}$ are
 200 as follows:

$$\begin{aligned}
 X &= \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix}, Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, ct_Z = Enc \begin{bmatrix} y_1 & y_1 x_{11} & \dots & y_1 x_{1d} \\ y_2 & y_2 x_{21} & \dots & y_2 x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ y_n & y_n x_{n1} & \dots & y_n x_{nd} \end{bmatrix}, \\
 ct_{\beta}^{(0)} &= Enc \begin{bmatrix} \beta_0^{(0)} & \beta_1^{(0)} & \dots & \beta_d^{(0)} \\ \beta_0^{(0)} & \beta_1^{(0)} & \dots & \beta_d^{(0)} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_0^{(0)} & \beta_1^{(0)} & \dots & \beta_d^{(0)} \end{bmatrix}, ct_{\bar{B}} = Enc \begin{bmatrix} \bar{B}_{[0][0]} & \bar{B}_{[1][1]} & \dots & \bar{B}_{[d][d]} \\ \bar{B}_{[0][0]} & \bar{B}_{[1][1]} & \dots & \bar{B}_{[d][d]} \\ \vdots & \vdots & \ddots & \vdots \\ \bar{B}_{[0][0]} & \bar{B}_{[1][1]} & \dots & \bar{B}_{[d][d]} \end{bmatrix},
 \end{aligned}$$

201 where $\bar{B}_{[i][i]}$ is the diagonal element of \bar{B} that is built from $-\frac{1}{4}X^T X$.

202 The public cloud takes the three ciphertexts ct_Z , $ct_{\beta}^{(0)}$ and $ct_{\bar{B}}$ and evaluates the enhanced NAG
 203 algorithm to find a decent weight vector by updating the vector $ct_{\beta}^{(0)}$. Refer to [12] for a detailed
 204 description about how to calculate the gradient by HE programming.

205 **Implementation** We implement the enhanced NAG based on HE with the library HEAAN. The C++
 206 source code is publicly available at <https://anonymous.4open.science/r/IDASH2017-245B>.
 207 All the experiments on the ciphertexts were conducted on a public cloud with 32 vCPUs and 64 GB
 208 RAM.

209 For a fair comparison with [12], we utilized the same 10-fold cross-validation (CV) technique on the
 210 same iDASH dataset consisting of 1579 samples with 18 features and the same 5-fold CV technique
 211 on the other five datasets. Like [12], We consider the average accuracy and the Area Under the Curve
 212 (AUC) as the main indicators. Tables 1 and 2 show the two experiment results, respectively. The two
 213 tables also provide the average evaluation running time for each iteration. We adopt the same packing
 214 method that Kim et al. [12] proposed and hence our solution has similar storage of ciphertexts to [12]
 215 with some extra ciphertexts to encrypt the \bar{B} .

216 The parameters of HEAAN we set are same to [12]: $\log N = 16$, $\log Q = 1200$, $\log p = 30$, $slots =$
 217 32768 , which ensure the security level $\lambda = 80$. Refer [12] for the details of these parameters. Since
 218 our enhanced NAG method need to consume more modulus to preserve the precision of \bar{B} , we use
 219 $\log p = 60$ to encrypt the matrix \bar{B} and thus only can perform 3 iterations of the enhanced NAG
 220 method. Yet despite only 3 iterations, our enhanced NAG method still produces a comparable result.

Table 1: Implementation Results for iDASH datasets with 10-fold CV

Dataset	Sample Num	Feature Num	Method	deg g	Iter Num	Learn Time (min)	Accuracy (%)	AUC
iDASH	1579	18	Ours	5	3	5.53	53.69	0.678
			[12]	5	7	6.07	62.87	0.689

Table 2: Implementation Results for other datasets with 5-fold CV

Dataset	Sample Num	Feature Num	Method	deg g	Iter Num	Learn Time (min)	Accuracy (%)	AUC
Edinburgh	1253	9	Ours	5	3	0.6	84.4	0.853
			[12]	5	7	3.6	91.04	0.958
lbw	189	9	Ours	5	3	0.5	69.19	0.619
			[12]	5	7	3.3	69.19	0.689
nhanes3	15649	15	Ours	5	3	5.5	79.23	0.490
			[12]	5	7	7.3	79.22	0.717
pcs	379	9	Ours	5	3	0.6	65.33	0.721
			[12]	5	7	3.5	68.27	0.740
uis	575	8	Ours	5	3	0.6	74.43	0.598
			[12]	5	7	3.5	74.44	0.603

221 5 Conclusion

222 In this paper, we proposed a faster gradient variant called `quadratic gradient`, and implemented
 223 the quadratic-gradient version of NAG in the encrypted domain to train the logistic regression model.

224 The quadratic gradient presented in this work can be constructed from the Hessian matrix directly,
 225 and thus somehow integrates the second-order Newton’s method and the first-order gradient (descent)
 226 method together. There is a good chance that quadratic gradient could accelerate other gradient
 227 methods such as RMSprop and Adam, which is an open future work.

228 References

- 229 [1] Allison, P. D. (2008). Convergence failures in logistic regression.
- 230 [2] Aono, Y., Hayashi, T., Trieu Phong, L., and Wang, L. (2016). Scalable and secure logistic
231 regression via homomorphic encryption. In *Proceedings of the Sixth ACM Conference on Data
232 and Application Security and Privacy*, pages 142–144.
- 233 [3] Blatt, M., Gusev, A., Polyakov, Y., Rohloff, K., and Vaikuntanathan, V. (2019). Optimized
234 homomorphic encryption solution for secure genome-wide association studies. *IACR Cryptology
235 ePrint Archive*, 2019:223.
- 236 [4] Böhning, D. and Lindsay, B. G. (1988). Monotonicity of quadratic-approximation algorithms.
237 *Annals of the Institute of Statistical Mathematics*, 40(4):641–663.
- 238 [5] Bonte, C. and Vercauteren, F. (2018). Privacy-preserving logistic regression training. *BMC
239 medical genomics*, 11(4):86.
- 240 [6] Chen, H., Gilad-Bachrach, R., Han, K., Huang, Z., Jalali, A., Laine, K., and Lauter, K. (2018).
241 Logistic regression over encrypted data from fully homomorphic encryption. *BMC medical
242 genomics*, 11(4):3–12.
- 243 [7] Cheon, J. H., Kim, A., Kim, M., and Song, Y. (2017). Homomorphic encryption for arithmetic of
244 approximate numbers. In *International Conference on the Theory and Application of Cryptology
245 and Information Security*, pages 409–437. Springer.
- 246 [8] Crawford, J. L., Gentry, C., Halevi, S., Platt, D., and Shoup, V. (2018). Doing real work with fhe:
247 the case of logistic regression. In *Proceedings of the 6th Workshop on Encrypted Computing &
248 Applied Homomorphic Cryptography*, pages 1–12.
- 249 [9] Gentry, C. (2009). Fully homomorphic encryption using ideal lattices. In *Proceedings of the
250 forty-first annual ACM symposium on Theory of computing*, pages 169–178.
- 251 [10] Han, K., Hong, S., Cheon, J. H., and Park, D. (2019). Logistic regression on homomorphic en-
252 crypted data at scale. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33,
253 pages 9466–9471.
- 254 [11] Jäschke, A. and Armknecht, F. (2016). Accelerating homomorphic computations on rational
255 numbers. In *International Conference on Applied Cryptography and Network Security*, pages
256 405–423. Springer.
- 257 [12] Kim, A., Song, Y., Kim, M., Lee, K., and Cheon, J. H. (2018a). Logistic regression model
258 training based on the approximate homomorphic encryption. *BMC medical genomics*, 11(4):83.
- 259 [13] Kim, M., Song, Y., Li, B., and Micciancio, D. (2019). Semi-parallel logistic regression for gwas
260 on encrypted data. *IACR Cryptology ePrint Archive*, 2019:294.
- 261 [14] Kim, M., Song, Y., Wang, S., Xia, Y., and Jiang, X. (2018b). Secure logistic regression based
262 on homomorphic encryption: Design and evaluation. *JMIR medical informatics*, 6(2):e19.
- 263 [15] Murphy, K. P. (2012). *Machine learning: a probabilistic perspective*. The MIT Press, Cam-
264 bridge, MA.

265 Checklist

- 266 1. For all authors...
- 267 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
268 contributions and scope? [Yes]
- 269 (b) Did you describe the limitations of your work? [Yes]
- 270 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 271 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
272 them? [Yes]

- 273 2. If you are including theoretical results...
- 274 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 275 (b) Did you include complete proofs of all theoretical results? [Yes]
- 276 3. If you ran experiments...
- 277 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 278 mental results (either in the supplemental material or as a URL)? [Yes]
- 279 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 280 were chosen)? [Yes]
- 281 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 282 ments multiple times)? [Yes]
- 283 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 284 of GPUs, internal cluster, or cloud provider)? [Yes]
- 285 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 286 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 287 (b) Did you mention the license of the assets? [Yes]
- 288 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 289 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 290 using/curating? [Yes]
- 291 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 292 information or offensive content? [Yes]
- 293 5. If you used crowdsourcing or conducted research with human subjects...
- 294 (a) Did you include the full text of instructions given to participants and screenshots, if
- 295 applicable? [Yes]
- 296 (b) Did you describe any potential participant risks, with links to Institutional Review
- 297 Board (IRB) approvals, if applicable? [Yes]
- 298 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 299 spent on participant compensation? [Yes]

300 A Appendix

301 Optionally include extra information (complete proofs, additional experiments and plots) in the

302 appendix. This section will often be part of the supplemental material.