

Adaptive Laplace Mechanism: Differential Privacy Preservation in Deep Learning

NhatHai Phan*, Xintao Wu[†], Han Hu*, and Dejing Dou[‡]

*New Jersey Institute of Technology, Newark, New Jersey, USA

[†]University of Arkansas, Fayetteville, Arkansas, USA

[‡]University of Oregon, Eugene, Oregon, USA

Emails: phan@njit.edu, xintaowu@uark.edu, hh255@njit.edu, and dou@cs.uoregon.edu

Abstract—In this paper, we focus on developing a novel mechanism to preserve differential privacy in deep neural networks, such that: (1) The privacy budget consumption is totally independent of the number of training steps; (2) It has the ability to adaptively inject noise into features based on the contribution of each to the output; and (3) It could be applied in a variety of different deep neural networks. To achieve this, we figure out a way to perturb affine transformations of neurons, and loss functions used in deep neural networks. In addition, our mechanism intentionally adds “more noise” into features which are “less relevant” to the model output, and vice-versa. Our theoretical analysis further derives the sensitivities and error bounds of our mechanism. Rigorous experiments conducted on MNIST and CIFAR-10 datasets show that our mechanism is highly effective and outperforms existing solutions.

I. INTRODUCTION

Today, deep learning has become the tool of choice in many areas of engineering, such as autonomous systems, signal and information processing, and data analytics. Deep learning systems are, therefore, not only applied in classic settings, such as speech and handwriting recognition, but also progressively operate at the core of security and privacy critical applications.

For instance, self-driving cars make use of deep learning for recognizing objects and street signs [1]. Detection systems for email spam integrate learning methods for analyzing data more effectively [2]. Furthermore, deep learning has applications in a number of healthcare areas, e.g., phenotype extraction and health risk prediction [3], prediction of the development of various diseases, including schizophrenia, cancers, diabetes, heart failure, etc. [4], and many more. This presents an obvious threat to privacy in new deep learning systems which are being deployed. However, there are only a few scientific studies in preserving privacy in deep learning.

In the past few decades, a subject of significant interest has been how to release the sensitive results of statistical analyses and data mining, while still protecting privacy. One state-of-the-art privacy model is ϵ -differential privacy [5], which ensures that the adversary cannot infer any information about any specific record with high confidence (controlled by a privacy budget) from the released learning models, even if all the remaining tuples of the sensitive data are possessed by the adversary. The privacy budget controls the amount by which the output distributions induced by two neighboring datasets may differ: A smaller privacy budget value enforces

a stronger privacy guarantee. Differential privacy research has been studied from both theoretical and application perspectives [6, 7]. The mechanisms of achieving differential privacy mainly include adding Laplace noise [5], the exponential mechanism [8], and the functional perturbation method [6].

It is significant and timely to combine differential privacy and deep learning, i.e., the two state-of-the-art techniques in privacy preserving and machine learning. However, this is a challenging task, and only a few scientific studies have been conducted. In [9], Shokri and Shmatikov proposed a distributed training method, which injects noise into “gradients” of parameters, to preserve privacy in neural networks. In this method, the magnitude of injected noise and the privacy budget ϵ are accumulated in proportion to the number of training epochs and the number of shared parameters. Thus, it may consume an unnecessarily large portion of the privacy budget, as the number of training epochs and the number of shared parameters among multiple parties are often large [10].

To improve this, based on the composition theorem [11], Abadi et al. [12] proposed a privacy accountant, which keeps track of privacy spending and enforces applicable privacy policies. However, the approach is still dependent on the number of training epochs, as it introduces noise into “gradients” of parameters in every training step. With a small privacy budget ϵ , only a small number of epochs can be used to train the model [12]. In practice, that could potentially affect the model utility, when the number of training epochs needs to be large to guarantee the model accuracy.

A recent approach towards differentially private deep neural networks was explored by Phan et al. [10]. This work proposed deep private auto-encoders (dPAs), in which differential privacy is enforced by perturbing the cross-entropy errors in auto-encoders [13]. Their algorithm was designed particularly for auto-encoders, in which specific objective functions are applied. A different method, named **CryptoNets**, was proposed in [14] towards the application of neural networks to encrypted data. A data owner can send their encrypted data to a cloud service that hosts the network, and get encrypted predictions in return. This method is different from our context, since it does not aim at releasing learning models under privacy protections.

Another drawback of the existing techniques is that all parameters are treated the same in terms of the amount of noise injected. This may not be ideal in real scenarios, since different

features and parameters normally have different impacts upon the model output. Fig. 1 shows the relevance estimated by applying Layer-wise Relevance Propagation (LRP) [15] of input features to the prediction of the image’s label in the MNIST dataset [16]. Dark red units have stronger impacts than green and yellow units. Injecting the same magnitude of noise into all parameters may affect the model utility.

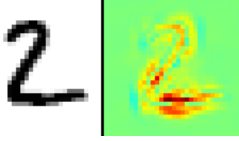


Fig. 1. An instance of relevance of each input feature given to the classification output (MNIST dataset). Red neurons indicate stronger relevances, and green neurons indicate weaker relevances.

Research Goal. Therefore, there is an urgent demand for the development of a privacy preserving mechanism, such that: **(1)** It is totally independent of the number of training epochs in consuming privacy budget; **(2)** It has the ability to adaptively inject noise into features based on the contribution of each to the model output; and **(3)** It can be applied in a variety of deep neural networks. Mechanisms

with such characteristics will significantly enhance the operation of privacy preservation in deep learning.

Our Contribution. Motivated by this, we develop a novel mechanism, called **Adaptive Laplace Mechanism (AdLM)**, to preserve differential privacy in deep learning. Our idea is to intentionally add “*more noise*” into features which are “*less relevant*” to the model output, and vice-versa. To achieve that, we inject Laplace noise into the computation of Layer-wise Relevance Propagation (LRP) [15] to estimate a differentially private relevance of each input feature to the model output. Given the perturbed features, we figure out a novel way to distribute adaptive noise into affine transformations and loss functions used in deep neural networks as a preprocessing step, so that preserving differential privacy is feasible. As a result, we expect to improve the utility of deep neural networks under ϵ -differential privacy. It is worth noting that our mechanism does not access the original data again in the training phase. Theoretical analysis derives the sensitivities and error bounds of our mechanism, and shows that they are totally independent of the number of epochs.

Different from [9, 12], in our mechanism, the injected noise and the privacy budget consumption do not accumulate in each training step. Consequently, the privacy budget consumption in our mechanism is totally independent of the number of training epochs. In addition, different from [10], our mechanism can be applied in a variety of deep learning networks with different activation functions. Convolution neural networks (CNNs) [16] are used as an example to validate the effectiveness of our mechanism. Rigorous experiments conducted on MNIST and CIFAR-10 datasets [17] show that our mechanism is effective and outperforms existing solutions.

II. PRELIMINARIES AND RELATED WORKS

In this section, we revisit differential privacy, existing techniques in preserving differential privacy in deep learning, and the Layer-wise Relevance Propagation (LRP) algorithm [15].

Let D be a database that contains n tuples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ and $d+1$ attributes X_1, X_2, \dots, X_d, Y , and for each tuple $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id}, y_i)$. We assume, without loss of generality, $\sqrt{\sum_{j=1}^d x_{ij}^2} \leq 1$ where $x_{ij} \geq 0$. This assumption can be easily enforced by changing each x_{ij} to $\frac{x_{ij} - \alpha_j}{(\beta_j - \alpha_j) \cdot \sqrt{d}}$, where α_j and β_j denote the minimum and maximum values in the domain of X_j .

To be general, let us consider a classification task with M possible categorical outcomes, i.e., the data label y_i given $\mathbf{x}_i \in L$ is assigned to only one of the M categories. Each y_i can be considered as a vector of M categories $y_i = \{y_{i1}, \dots, y_{iM}\}$. If the l -th category is the class of \mathbf{x}_i , then $y_{il} = 1$, otherwise $y_{il} = 0$. Our objective is to construct a differentially private deep neural network from D that (i) takes $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})$ as input and (ii) outputs a prediction of y_i that is as accurate as possible. To evaluate whether model parameters θ lead to an accurate model, a cost function $\mathcal{F}_D(\theta)$ is used to measure the difference between the original and predicted values of y_i .

A. ϵ -Differential Privacy

As the released model parameter θ may disclose sensitive information of D , to protect the privacy, we require that the model training should be performed with an algorithm that satisfies ϵ -differential privacy. The definition of differential privacy is as follows:

Definition 1: ϵ -Differential Privacy [5]. A randomized algorithm A fulfills ϵ -differential privacy, if for any two databases D and D' differing at most one tuple, and for all $O \subseteq \text{Range}(A)$, we have:

$$\Pr[A(D) = O] \leq e^\epsilon \Pr[A(D') = O] \quad (1)$$

where the privacy budget ϵ controls the amount by which the distributions induced by D and D' may differ. A smaller ϵ enforces a stronger privacy guarantee of A .

A general method for preserving ϵ -differential privacy of any function \mathcal{F} (on D) is the *Laplace mechanism* [5], where the output of \mathcal{F} is a vector of real numbers. In fact, the mechanism exploits the global sensitivity of \mathcal{F} over any two neighboring data sets (differing at most one record), which is denoted as $GS_{\mathcal{F}}(D)$. Given $GS_{\mathcal{F}}(D)$, the Laplace mechanism ensures ϵ -differential privacy by injecting noise η into each value in the output of $\mathcal{F}(D)$: $pdf(\eta) = \frac{\epsilon}{2GS_{\mathcal{F}}(D)} \exp(-|\eta| \cdot \frac{\epsilon}{GS_{\mathcal{F}}(D)})$, where η is drawn i.i.d. from Laplace distribution with zero mean and scale $GS_{\mathcal{F}}(D)/\epsilon$.

Research in differential privacy has been significantly studied, from both the theoretical perspective, e.g., [6, 18], and the application perspective, e.g., data collection [19], spatio-temporal correlations [20, 21], data streams [22], stochastic gradient descents [23], recommendation [7], regression [6], online learning [24], publishing contingency tables [25], and spectral graph analysis [26].

B. Differential Privacy in Deep Learning

Deep neural networks define parameterized functions from inputs $\mathbf{x}_i \in D$ to outputs, i.e, a prediction of y_i , as compo-

sitions of many layers of hidden neurons and nonlinear functions. For instance, Fig. 2 illustrates a multilayer neural network, in which there are k hidden layers $H = \{\mathbf{h}_1, \dots, \mathbf{h}_k\}$. Rectified linear units (ReLUs) and sigmoids are widely used examples of activation functions. By adjusting parameters of these neurons, such parameterized functions can be trained with the goal of fitting a finite set of input-output data instances. We specify a loss function $\mathcal{F}_D(\theta)$ that represents the penalty for mismatching between the predicted and original values of y_i . $\mathcal{F}_D(\theta)$ on parameters θ is the average of the loss over the training examples $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. Stochastic gradient descent (SGD) algorithm is used to minimize the cross-entropy error [13], given the model outputs and true data labels.

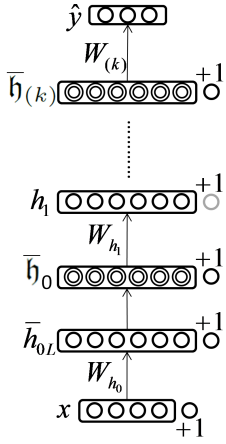


Fig. 2. An instance of differentially private neural networks.

$$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} \sum_i \left(\frac{\mathbf{g}_t(\mathbf{x}_i)}{\max(1, \frac{\|\mathbf{g}_t(\mathbf{x}_i)\|^2}{C})} + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}) \right) \quad (2)$$

$$\theta_t \leftarrow \theta_t - \xi_t \tilde{\mathbf{g}}_t \quad (3)$$

where ξ_t is a learning rate at the training step t .

Finally, differentially private parameters θ , denoted $\bar{\theta}$, learned by the algorithm are shared to the public and other parties. Overall, the algorithm introduces noise into “gradients” of parameters at every training step. The magnitude of injected noise and the privacy budget ϵ are accumulated in proportion to the number of training epochs.

Compared with the work in [12], the goal is similar: learning differentially private parameters $\bar{\theta}$. However, we develop a novel mechanism in which the privacy budget consumption is independent of the number of training epochs. Our mechanism is different. We redistribute the noise so that “more noise” will be added into features which are “less relevant” to the model output, and vice-versa. Moreover, we inject noise into coefficients of affine transformations and loss functions, such that differentially private parameters can be learned.

C. Layer-wise Relevance Propagation

Layer-wise Relevance Propagation (LRP) [15] is a well-accepted algorithm, which is applied to compute the relevance of each input feature x_{ij} to the model outcome $\mathcal{F}_{\mathbf{x}_i}(\theta)$.

Given the relevance, denoted $R_m^{(k)}(\mathbf{x}_i)$, of a certain neuron m at the layer k , i.e., $m \in \mathbf{h}_k$, for the model outcome $\mathcal{F}_{\mathbf{x}_i}(\theta)$, LRP algorithm aims at obtaining a decomposition of such relevance in terms of messages sent to neurons of the previous layers, i.e., the layer $(k-1)$ -th. These messages are called $R_{p \leftarrow m}^{(k-1,k)}(\mathbf{x}_i)$. The overall relevance of each neuron in the lower layer is determined by summing up the relevance coming from all upper-layer neurons:

$$R_p^{(k-1)}(\mathbf{x}_i) = \sum_{m \in \mathbf{h}_k} R_{p \leftarrow m}^{(k-1,k)}(\mathbf{x}_i) \quad (4)$$

where the relevance decomposition is based on the ratio of local and global affine transformations and is given by:

$$R_{p \leftarrow m}^{(k-1,k)}(\mathbf{x}_i) = \begin{cases} \frac{z_{pm}(\mathbf{x}_i)}{z_m(\mathbf{x}_i) + \mu} R_m^{(k)}(\mathbf{x}_i) & z_m(\mathbf{x}_i) \geq 0 \\ \frac{z_{pm}(\mathbf{x}_i)}{z_m(\mathbf{x}_i) - \mu} R_m^{(k)}(\mathbf{x}_i) & z_m(\mathbf{x}_i) < 0 \end{cases} \quad (5)$$

with: $z_m(\mathbf{x}_i)$ is the affine transformation of neuron $m \in \mathbf{h}_k$:

$$z_{pm}(\mathbf{x}_i) = p_{\mathbf{x}_i} \times W_{pm} \quad (6)$$

$$z_m(\mathbf{x}_i) = \sum_{p \in \mathbf{h}_k} z_{pm}(\mathbf{x}_i) + b_m \quad (7)$$

s.t. $p_{\mathbf{x}_i}$ is the value of neuron p given \mathbf{x}_i , W_{pm} is a weight connecting the neuron p to neuron m , and b_m is a bias term. A predefined stabilizer $\mu \geq 0$ is introduced to overcome unboundedness.

In Eq. 5, in order to back propagate the relevance, we need to compute the relevance $R_m^{(k)}(\mathbf{x}_i)$ at the last hidden layer, i.e., the k -th layer, from the output layer. Given the output variable o , $R_m^{(k)}(\mathbf{x}_i)$ is computed as follows:

$$R_m^{(k)}(\mathbf{x}_i) = \begin{cases} \frac{z_{mo}(\mathbf{x}_i)}{z_o(\mathbf{x}_i) + \mu} \mathcal{F}_{\mathbf{x}_i}(\theta) & z_o(\mathbf{x}_i) \geq 0 \\ \frac{z_{mo}(\mathbf{x}_i)}{z_o(\mathbf{x}_i) - \mu} \mathcal{F}_{\mathbf{x}_i}(\theta) & z_o(\mathbf{x}_i) < 0 \end{cases} \quad (8)$$

Given k hidden layers $\{\mathbf{h}_1, \dots, \mathbf{h}_k\}$, by using Eqs. 4, 5, and 8, we can compute the relevance of every hidden neuron and input feature. As in [15], the following equation holds:

$$\mathcal{F}_{\mathbf{x}_i}(\theta) = \sum_{m \in \mathbf{h}_k} R_m^{(k)}(\mathbf{x}_i) = \dots = \sum_{x_{ij} \in \mathbf{x}_i} R_{x_{ij}}(\mathbf{x}_i) \quad (9)$$

where $R_{x_{ij}}(\mathbf{x}_i)$ is the relevance of the feature x_{ij} given the model outcome $\mathcal{F}_{\mathbf{x}_i}(\theta)$. To ensure that the relevance $R_{x_{ij}}(\mathbf{x}_i) \in [-1, 1]$, each $R_{x_{ij}}(\mathbf{x}_i)$ is normalized to $\frac{R_{x_{ij}}(\mathbf{x}_i) - \chi}{(\varphi - \chi)}$, where φ and χ denote the maximum and minimum values in the domain of $\{R_{x_{i1}}(\mathbf{x}_i), \dots, R_{x_{id}}(\mathbf{x}_i)\}$.

III. ADAPTIVE LAPLACE MECHANISM (ADLM)

In this section, we formally present our mechanism. Given a loss function $\mathcal{F}(\theta)$ with model parameters θ , the network is trained by optimizing the loss function $\mathcal{F}(\theta)$ on D by applying SGD algorithm on T random training batches consequently. At each training step, a single training batch L is used. A batch L is a random set of training samples in D with a predefined batch size $|L|$.

The pseudo-codes of Algorithm 1 outline five basic steps in our mechanism to learn differentially private parameters of the model. The five basic steps are as follows:

- **Step 1 (Lines 1-7).** In the first step, we obtain the average relevances of all the j -th input features, denoted as $R_j(D)$, by applying the LRP algorithm on a well-trained deep neural network on the database D . $R_j(D)$ is computed as follows:

$$R_j(D) = \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i) \quad (10)$$

Then, we derive differentially private relevances, denoted as \bar{R}_j , by injecting Laplace noise into R_j for all the j -th input features. The total privacy budget in this step is ϵ_1 .

- **Step 2 (Lines 8-14).** In the second step, we derive a differentially private affine transformation layer, denoted \mathbf{h}_0 . Every hidden neuron $h_{0j} \in \mathbf{h}_0$ will be perturbed by injecting adaptive Laplace noise into its affine transformation to preserve differential privacy given a batch L . Based on \bar{R}_j , “more noise” is injected into features which are “less relevant” to the model output, and vice-versa. The total privacy budget used in this step is ϵ_2 . The perturbed affine transformation layer is denoted as $\bar{\mathbf{h}}_{0L}$ (Fig. 2).

- **Step 3 (Line 15).** In the third step, we stack hidden layers $\{\mathbf{h}_1, \dots, \mathbf{h}_k\}$ on top of the differentially private hidden layer $\bar{\mathbf{h}}_{0L}$ to construct the deep private neural network (Fig. 2). The computations of $\mathbf{h}_1, \dots, \mathbf{h}_k$ are done based on the differentially private layer $\bar{\mathbf{h}}_{0L}$ without accessing any information from the original data. Therefore, the computations do not disclose any information. Before each stacking operation, a normalization layer, denoted \bar{h} , is applied to bound non-linear activation functions, such as ReLUs (Fig. 2).

- **Step 4 (Lines 16-19).** After constructing a private structure of hidden layers $\{\bar{\mathbf{h}}_{0L}, \mathbf{h}_1, \dots, \mathbf{h}_k\}$, we need to protect the labels y_i at the output layer. To achieve this, we derive a polynomial approximation of the loss function \mathcal{F} . Then, we perturb the loss function \mathcal{F} by injecting Laplace noise with a privacy budget ϵ_3 into its coefficients to preserve differential privacy on each training batch L , denoted $\bar{\mathcal{F}}_L(\theta)$.

- **Step 5 (Lines 20-30).** Finally, the parameter θ_T is derived by minimizing the loss function $\bar{\mathcal{F}}_L(\theta)$ on T training steps sequentially. In each step t , stochastic gradient descent (SGD) algorithm is used to update parameters θ_t given a random batch L of training samples in D . This essentially is an optimization process, without using any additional information from the original data.

In our mechanism, differential privacy is preserved, since it is enforced at every computation task that needs to access the original data D . Laplace noise is injected into our model only once, as a preprocessing step to preserve differential privacy in the computation of the relevance $\bar{R}_j(D)$, the first layer $\bar{\mathbf{h}}_{0L}$, and the loss function $\bar{\mathcal{F}}_L(\theta)$. Thereafter, the training phase will not access the original data again. The privacy budget consumption does not accumulate in each training step. As such, it is independent of the number of training epochs.

A. Private Relevance

In this section, we preserve differential privacy in the computation of the relevance of each j -th input feature on database D by injecting Laplace noise into $R_j(D)$. We set $\Delta_{\mathbf{R}} = \frac{2d}{|D|}$ based on the maximum values of all the relevances $R_j(D)$ (line 4, Alg. 1). In lines 5-6, the relevance of each j -th input feature $R_j(D)$ is perturbed by adding Laplace noise $Lap(\frac{\Delta_{\mathbf{R}}}{\epsilon_1})$. The perturbed relevance is denoted as \bar{R}_j . In line 7, we obtain the set of all perturbed relevances $\bar{\mathbf{R}}(D)$:

$$\bar{\mathbf{R}}(D) = \{\bar{R}_j\}_{j \in [1,d]} \quad (11)$$

$$\text{where } \bar{R}_j = \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i) + Lap\left(\frac{\Delta_{\mathbf{R}}}{\epsilon_1}\right) \quad (12)$$

The computation of $\bar{\mathbf{R}}(D)$ is ϵ_1 -differential private. The correctness is based on the following lemmas.

Lemma 1: Let D and D' be any two neighboring databases. Given $\mathbf{R}(D)$ and $\mathbf{R}(D')$ be the relevance of all input features on D and D' , respectively, and denote their representations as

$$\mathbf{R}(D) = \{R_j(D)\}_{j \in [1,d]} \text{ s.t. } R_j(D) = \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i)$$

$$\mathbf{R}(D') = \{R_j(D')\}_{j \in [1,d]} \text{ s.t. } R_j(D') = \frac{1}{|D'|} \sum_{\mathbf{x}'_i \in D'} R_{x'_{ij}}(\mathbf{x}'_i)$$

Then, we have the following inequality:

$$\frac{1}{|D|} \sum_{j=1}^d \left\| \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i) - \sum_{\mathbf{x}'_i \in D'} R_{x'_{ij}}(\mathbf{x}'_i) \right\|_1 \leq \frac{2d}{|D|} \quad (13)$$

where d is the number of features in each tuple $\mathbf{x}_i \in D$.

Proof 1: The proof of Lemma 1 is in **Appendix A**¹.

Lemma 2: Algorithm 1 preserves ϵ_1 -differential privacy in the computation of $\bar{\mathbf{R}}(D)$.

Proof 2: The proof of Lemma 2 is in **Appendix B**¹.

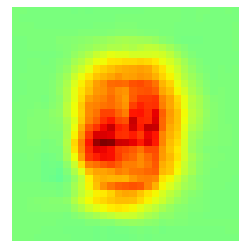


Fig. 3. The average differentially private relevance of each input feature given MNIST dataset.

Lemma 2 shows that the computation of the relevances $\bar{\mathbf{R}}(D)$ is differentially private. Fig. 3 illustrates the differentially private relevance \bar{R}_j (Eq. 12) of each j -th coefficient given the database D . It is worth noting that the relevance distribution is not identical. In the next section, $\bar{\mathbf{R}}(D)$ is used to redistribute the noise injected into the affine transformation layer \mathbf{h}_0 in our deep neural network.

B. Private Affine Transformation Layer with Adaptive Noise

In general, before applying activation functions such as ReLU and sigmoid, the affine transformation of a hidden neuron $h \in \mathbf{h}_0$ can be presented as:

$$h_{\mathbf{x}_i}(W) = b + \mathbf{x}_i W^T \quad (14)$$

¹<https://www.dropbox.com/s/dm571t91mg8r0cr/Appendix.pdf?dl=0>

where b is a static bias, and W is the parameter of h . Given a training batch L , h can be rewritten as:

$$h_L(W) = \sum_{\mathbf{x}_i \in L} (b + \mathbf{x}_i W^T) \quad (15)$$

Given the above representation of each neuron $h_L(W_h)$, we preserve differential privacy in the computation of \mathbf{h}_0 on L by injecting Laplace noise into inputs b and \mathbf{x}_i of every neuron $h_L(W) \in \mathbf{h}_0$. Intuitively, we can apply an identical noise distribution $\frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2})$ to all input features, where $\Delta_{\mathbf{h}_0} = 2 \sum_{h \in \mathbf{h}_0} d$ (line 9, Alg. 1). This approach works well when every input feature has an identical contribution to the model outcome. (Please refer to the **Appendix C**.)

In practice, this assumption usually is violated. For instance, Fig. 1 illustrates the relevance, estimated by the LRP algorithm [15], of each input feature given different handwritten digits. It is clear that the relevances are not identical. The differentially private relevances are not identical as well (Fig. 3). Therefore, injecting the same magnitude of noise into all input features may affect the utility of differentially private neural networks.

To address this problem, we propose an Adaptive Laplace Mechanism (**AdLM**), to adaptively redistribute the injected noise to improve the performance. Given hidden units $h_{\mathbf{x}_i}(W)$ in Eq. 14, our key idea is to intentionally add *more noise* into input features which are *less relevant* to the model output Y , and vice-versa. As a result, we expect to improve the utility of the model under differential privacy. In fact, we introduce a privacy budget ratio β_j and the privacy budget ϵ_j for each j -th input feature as follows:

$$\beta_j = \frac{d \times |\bar{R}_j|}{\sum_{j=1}^d |\bar{R}_j|} \quad \text{s.t.} \quad \epsilon_j = \beta_j \times \epsilon_2 \quad (16)$$

We set $\Delta_{\mathbf{h}_0} = 2 \sum_{h \in \mathbf{h}_0} d$ based on the maximum values of all the input features x_{ij} (line 9, Alg. 1). In line 11, β_j can be considered as the fraction of the contribution to $\Delta_{\mathbf{h}_0}$ from the j -th input feature to the hidden neuron $h \in \mathbf{h}_0$. In lines 12-13, each input feature x_{ij} of every hidden neuron h in the first affine transformation layer \mathbf{h}_0 is perturbed by adding adaptive Laplace noise $\frac{1}{|L|} \text{Lap}(\Delta_{\mathbf{h}_0}/\epsilon_j)$. The perturbed input features are denoted as $\bar{\mathbf{x}}_i$. In lines 20-21, given a random training batch L , we construct the differentially private affine transformation layer $\bar{\mathbf{h}}_{0L}$, which consists of perturbed hidden neurons $\bar{h}_L(W)$:

$$\bar{\mathbf{h}}_{0L}(W_0) = \{\bar{h}_L(W)\}_{h \in \mathbf{h}_0} \quad \text{s.t.} \quad \bar{h}_L(W) = \sum_{\mathbf{x}_i \in L} (\bar{\mathbf{x}}_i W^T + \bar{b})$$

where $\bar{b} = b + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2})$ is the perturbed bias (line 14). The following lemma shows that Alg. 1 preserves ϵ_2 -differential privacy in the computation of $\bar{\mathbf{h}}_{0L}$.

Lemma 3: Let L and L' be any two neighboring batches. Given parameter W_0 , let \mathbf{h}_{0L} and $\mathbf{h}_{0L'}$ be the first affine transformation layers on L and L' , respectively, and denote their representations as follows:

$$\mathbf{h}_{0L}(W_0) = \{h_L(W)\}_{h \in \mathbf{h}_0} \quad \text{s.t.} \quad h_L(W) = \sum_{\mathbf{x}_i \in L} (b + \mathbf{x}_i W^T)$$

$$\mathbf{h}_{0L'}(W_0) = \{h_{L'}(W)\}_{h \in \mathbf{h}_0} \quad \text{s.t.} \quad h_{L'}(W) = \sum_{\mathbf{x}'_i \in L'} (b + \mathbf{x}'_i W^T)$$

Algorithm 1 Adaptive Laplace Mechanism (Database D , hidden layers H , loss function $\mathcal{F}(\theta)$, and privacy budgets ϵ_1 , ϵ_2 , and ϵ_3 , the number of batches T , the batch size $|L|$)

- 1: **Compute the average relevance by applying the LRP Alg.**
 - 2: $\forall j \in [1, d] : R_j(D) = \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i)$ #Eq.10#
 - 3: **Inject Laplace noise into the average relevance of each j -th input feature**
 - 4: $\Delta_{\mathbf{R}} = 2d/|D|$ #Lemma 1#
 - 5: **for** $j \in [1, d]$ **do**
 - 6: $\bar{R}_j \leftarrow \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i) + \text{Lap}(\frac{\Delta_{\mathbf{R}}}{\epsilon_1})$
 - 7: $\bar{\mathbf{R}}(D) = \{\bar{R}_j\}_{j \in [1, d]}$
 - 8: **Inject Laplace noise into coefficients of the differentially private layer \mathbf{h}_0**
 - 9: $\Delta_{\mathbf{h}_0} = 2 \sum_{h \in \mathbf{h}_0} d$ #Lemma 3#
 - 10: **for** $j \in [1, d]$ **do**
 - 11: $\epsilon_j \leftarrow \beta_j \times \epsilon_2$ #Eq. 16#
 - 12: **for** $\mathbf{x}_i \in D, j \in [1, d]$ **do**
 - 13: $\bar{x}_{ij} \leftarrow x_{ij} + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_j})$ #perturb input feature x_{ij} #
 - 14: $\bar{b} \leftarrow b + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2})$ #perturb bias b #
 - 15: **Construct hidden layers** $\{\mathbf{h}_1, \dots, \mathbf{h}_k\}$ **and normalization layers** $\{\bar{h}_1, \dots, \bar{h}_k\}$
 - 16: **Inject Laplace noise into coefficients of the approximated loss function $\hat{\mathcal{F}}$**
 - 17: $\Delta_{\mathcal{F}} = M(|\bar{h}_{(k)}| + \frac{1}{4}|\bar{h}_{(k)}|^2)$ #Lemma 5#
 - 18: **for** $\mathbf{x}_i \in D, R \in [0, 2], l \in [1, M]$ **do**
 - 19: $\bar{\phi}_{l\mathbf{x}_i}^{(R)} \leftarrow \phi_{l\mathbf{x}_i}^{(R)} + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathcal{F}}}{\epsilon_3})$ #perturb coefficients of $\hat{\mathcal{F}}$
 - 20: **Initialize** θ_0 randomly
 - 21: **for** $t \in [T]$ **do**
 - 22: **Take a random training batch** L
 - 23: **Construct differentially private affine transformation layer**
 - 24: $\bar{\mathbf{h}}_{0L}(W_0) \leftarrow \{\bar{h}_L(W)\}_{h \in \mathbf{h}_0}$
 - 25: *s.t.* $\bar{h}_L(W) = \sum_{\mathbf{x}_i \in L} (\bar{\mathbf{x}}_i W^T + \bar{b})$
 - 26: **Construct differentially private loss function**
 - 27: $\bar{\mathcal{F}}_L(\theta_t) = \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \sum_{R=0}^2 \bar{\phi}_{l\mathbf{x}_i}^{(R)} (\bar{h}_{\mathbf{x}_i(k)} W_{l(k)}^T)^R$
 - 28: **Compute gradient descents**
 - 29: $\theta_{t+1} \leftarrow \theta_t - \eta_t \frac{1}{|L|} \nabla_{\theta_t} \bar{\mathcal{F}}_L(\theta_t)$ # η_t is a learning rate#
 - 30: **Return** θ_T # $(\epsilon_1 + \epsilon_2 + \epsilon_3)$ -differentially private#
-

Then, we have the following inequality:

$$\Delta_{\mathbf{h}_0} = \sum_{h \in \mathbf{h}_0} \sum_{j=1}^d \left\| \sum_{\mathbf{x}_i \in L} x_{ij} - \sum_{\mathbf{x}'_i \in L'} x'_{ij} \right\|_1 \leq 2 \sum_{h \in \mathbf{h}_0} d \quad (17)$$

where d is the number of features in each tuple $\mathbf{x}_i \in D$.

Proof 3: Assume that L and L' differ in the last tuple. Let \mathbf{x}_n (\mathbf{x}'_n) be the last tuple in L (L'). We have that

$$\begin{aligned} \Delta_{\mathbf{h}_0} &= \sum_{h \in \mathbf{h}_0} \sum_{j=1}^d \left\| \sum_{\mathbf{x}_i \in L} x_{ij} - \sum_{\mathbf{x}'_i \in L'} x'_{ij} \right\|_1 \\ &= \sum_{h \in \mathbf{h}_0} \sum_{j=1}^d \|x_{nj} - x'_{nj}\|_1 \leq 2 \max_{\mathbf{x}_i \in L} \sum_{h \in \mathbf{h}_0} \sum_{j=1}^d \|x_{ij}\|_1 \quad (18) \end{aligned}$$

Since $\forall \mathbf{x}_i, j : x_{ij} \in [0, 1]$, from Eq. 18 we have that: $\Delta_{\mathbf{h}_0} \leq 2 \sum_{h \in \mathbf{h}_0} d$. Eq. 17 holds. \blacksquare

Lemma 4: Algorithm 1 preserves ϵ_2 -differential privacy in the computation of $\bar{\mathbf{h}}_{0L}(W_0)$ (lines 24-25).

Proof 4: From lines 24-25 in the Alg. 1, for each $h \in \bar{\mathbf{h}}_{0L}$, h can be re-written as:

$$\bar{h}_L(W) = \sum_{j=1}^d \left[\sum_{\mathbf{x}_i \in L} (x_{ij} + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_j})) W^T \right] + \sum_{\mathbf{x}_i \in L} (b + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2})) \quad (19)$$

Let us consider the static bias $b = 1$ as the 0-th input feature and its associated parameter W_b , i.e., $x_{i0} = b = 1$ and $W = W_b \cup W$, we have that

$$\bar{h}_L(W) = \sum_{j=0}^d \left[\sum_{\mathbf{x}_i \in L} (x_{ij} + \frac{1}{|L|} \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_j})) W^T \right] \quad (20)$$

$$= \sum_{j=0}^d \left[\sum_{\mathbf{x}_i \in L} x_{ij} + \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_j}) \right] W^T = \sum_{j=0}^d \bar{\phi}_j^h W^T \quad (21)$$

where $\bar{\phi}_j^h = [\sum_{\mathbf{x}_i \in L} x_{ij} + \text{Lap}(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_j})]$.

We can see that $\bar{\phi}_j^h$ is the perturbation of the input feature x_{ij} associated with the j -th parameter $W_j \in W$ of the hidden neuron h on L . Since all the hidden neurons h in \mathbf{h}_0 are perturbed, we have that:

$$Pr(\bar{\mathbf{h}}_{0L}(W_0)) = \prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_j \|\sum_{\mathbf{x}_i \in L} x_{ij} - \bar{\phi}_j^h\|}{\Delta_{\mathbf{h}_0}}\right)$$

$\Delta_{\mathbf{h}_0}$ is set to $2 \sum_{h \in \mathbf{h}_0} d$ (line 9 in Alg. 1). $\bar{\mathbf{h}}_{0L}(W_0)$ is the output (lines 24-25 in Alg. 1). We have that

$$\begin{aligned} \frac{Pr(\bar{\mathbf{h}}_{0L}(W_0))}{Pr(\bar{\mathbf{h}}_{0L'}(W_0))} &= \frac{\prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_j \|\sum_{\mathbf{x}_i \in L} x_{ij} - \bar{\phi}_j^h\|_1}{\Delta_{\mathbf{h}_0}}\right)}{\prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_j \|\sum_{\mathbf{x}'_i \in L'} x'_{ij} - \bar{\phi}_j^h\|_1}{\Delta_{\mathbf{h}_0}}\right)} \\ &\leq \prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_j}{\Delta_{\mathbf{h}_0}} \left\| \sum_{\mathbf{x}_i \in L} x_{ij} - \sum_{\mathbf{x}'_i \in L'} x'_{ij} \right\|_1\right) \\ &\leq \prod_{h \in \mathbf{h}_0} \prod_{j=1}^d \exp\left(\frac{\epsilon_j}{\Delta_{\mathbf{h}_0}} 2 \max_{\mathbf{x}_n \in L} \|x_{nj}\|_1\right) \leq \prod_{h \in \mathbf{h}_0} \prod_{j=1}^d \exp\left(\frac{2\epsilon_j}{\Delta_{\mathbf{h}_0}}\right) \\ &\leq \prod_{h \in \mathbf{h}_0} \prod_{j=1}^d \exp\left(\epsilon_2 \frac{2 \frac{d \times |\bar{R}_j|}{\sum_{j=1}^d |\bar{R}_j|}}{\Delta_{\mathbf{h}_0}}\right) \\ &\leq \exp\left(\epsilon_2 \frac{2 \sum_{h \in \mathbf{h}_0} d \left[\sum_{j=1}^d \frac{|\bar{R}_j|}{\sum_{j=1}^d |\bar{R}_j|} \right]}{\Delta_{\mathbf{h}_0}}\right) = \exp(\epsilon_2) \end{aligned}$$

Consequently, the computation of $\bar{\mathbf{h}}_{0L}(W_0)$ preserves ϵ_2 -differential privacy in Alg. 1. \blacksquare

Lemma 4 shows that we can redistribute the noise in the computation of the first hidden layer $\bar{\mathbf{h}}_{0L}$ under differential privacy. In addition, given a batch L , without accessing additional information from the original data, none of the computations on top of $\bar{\mathbf{h}}_{0L}$ risk the privacy protection under differential privacy. These computation tasks include the application of activation functions, e.g., ReLU and sigmoid, on $\bar{\mathbf{h}}_{0L}$,

the computation of hidden layers $\mathbf{h}_1, \dots, \mathbf{h}_k$, local response normalizations, drop-out operations, polling layers, etc. (line 15, Alg. 1). This result can be applied to both fully-connected layers and convolution layers. In this paper, we applied ReLU on top of $\bar{\mathbf{h}}_{0L}$ and other layers $\mathbf{h}_1, \dots, \mathbf{h}_k$. Local response normalization layers are used after the application of ReLUs in each hidden layer to bound ReLU functions.

Local Response Normalization. The hidden units of the lower layer will be considered as the input of the next layer (Fig. 2). To ensure that this input is bounded $\bar{h}_{\mathbf{x}_i} \in [0, 1]$, as in [10, 27], we add a local response normalization (LRN) layer on top of each hidden layer. Given a fully-connected layer, as in [10], given an input \mathbf{x}_i , each perturbed neuron $\bar{h}_{\mathbf{x}_i}(W)$ can be directly normalized as follows: $\bar{h}_{\mathbf{x}_i} \leftarrow (\bar{h}_{\mathbf{x}_i}(W) - \chi) / (\varphi - \chi)$, where φ and χ denote the maximum and minimum values in the domain of $\{\bar{h}_{\mathbf{x}_i}\}_{i \in L}$.

Given a convolution layer with a perturbed neuron \bar{h}_{ij}^k at location (i, j) in the k -th feature map, based on [27], our local response normalization (LRN) is presented as follows:

$$\bar{h}_{ij}^k \leftarrow \bar{h}_{ij}^k / \max\left(\bar{h}_{ij}^k, (q + \alpha \sum_{m=\max(0, k-l/2)}^{\min(N-1, k+l/2)} (\bar{h}_{ij}^m)^2)^\beta\right) \quad (22)$$

where the constants q, l, α , and β are hyper-parameters, N is the total number of feature maps. As in [27], we used $q = 2, l = 5, \alpha = 10^{-4}$, and $\beta = 0.75$ in our experiments.

C. Perturbation of the Loss Function $\mathcal{F}_L(\theta)$

On top of our private deep neural network (Fig. 2), we add an output layer with the loss function $\mathcal{F}_L(\theta)$ to predict Y . Since the loss function $\mathcal{F}_L(\theta)$ accesses the labels y_i given $\mathbf{x}_i \in L$ from the data, we need to protect the labels y_i at the output layer. First, we derive a polynomial approximation of the loss function based on Taylor Expansion [28]. Then, we inject Laplace noise into coefficients of the loss function \mathcal{F} to preserve differential privacy on each training batch L .

The model output variables $\{\hat{y}_1, \dots, \hat{y}_M\}$ are fully linked to the normalized highest hidden layer, denoted $\bar{h}_{(k)}$, by weighted connections $W_{(k)}$ (Fig. 2). As common, the logistic function can be used as an activation function of the output variables. Given, l -th output variable \hat{y}_l and \mathbf{x}_i , we have:

$$\hat{y}_{il} = \sigma(\bar{h}_{\mathbf{x}_i(k)} W_{l(k)}^T) \quad (23)$$

where $\bar{h}_{\mathbf{x}_i(k)}$ is the state of $\bar{h}_{(k)}$ derived from $\bar{\mathbf{h}}_{0\mathbf{x}_i}$ by navigating through the neural network.

Cross-entropy error [13] can be used as a loss function. It has been widely used and applied in real-world applications [13]. Therefore, it is critical to preserve differential privacy under the use of the cross-entropy error function. Other loss functions, e.g., square errors, can be applied in the output layer,

as well. In our context, the cross-entropy error function is given by:

$$\begin{aligned}\mathcal{F}_L(\theta) &= -\sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \left(y_{il} \log \hat{y}_{il} + (1 - y_{il}) \log(1 - \hat{y}_{il}) \right) \\ &= -\sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \left(y_{il} \log(1 + e^{-\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T}) \right. \\ &\quad \left. + (1 - y_{il}) \log(1 + e^{\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T}) \right) \quad (24)\end{aligned}$$

Based on [10] and Taylor Expansion [28], we derive the polynomial approximation of $\mathcal{F}_L(\theta)$ as:

$$\begin{aligned}\hat{\mathcal{F}}_L(\theta) &= \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \sum_{q=1}^2 \sum_{R=0}^2 \frac{f_{ql}^{(R)}(0)}{R!} (\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T)^R \\ &= \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \left[\sum_{q=1}^2 f_{ql}^{(0)}(0) + \left(\sum_{q=1}^2 f_{ql}^{(1)}(0) \right) \bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T \right. \\ &\quad \left. + \left(\sum_{q=1}^2 \frac{f_{ql}^{(2)}(0)}{2!} \right) (\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T)^2 \right] \quad (25)\end{aligned}$$

where $\forall l \in [1, M] : f_{1l}(z) = y_{il} \log(1 + e^{-z})$ and $f_{2l}(z) = (1 - y_{il}) \log(1 + e^z)$.

To achieve ϵ_3 -differential privacy, we employ *functional mechanism* [29] to perturb the loss function $\hat{\mathcal{F}}_L(\theta)$ by injecting Laplace noise into its polynomial coefficients. So, we only need to perturb $\hat{\mathcal{F}}_L(\theta)$ just once in each training batch. To be clear, we denote $\{\phi_{l\mathbf{x}_i}^{(0)}, \phi_{l\mathbf{x}_i}^{(1)}, \phi_{l\mathbf{x}_i}^{(2)}\}$ as the coefficients, where

$$\phi_{l\mathbf{x}_i}^{(0)} = \sum_{q=1}^2 f_{ql}^{(0)}(0); \phi_{l\mathbf{x}_i}^{(1)} = \left(\sum_{q=1}^2 f_{ql}^{(1)}(0) \right); \phi_{l\mathbf{x}_i}^{(2)} = \left(\sum_{q=1}^2 \frac{f_{ql}^{(2)}(0)}{2!} \right)$$

In Alg. 1, we set $\Delta_{\mathcal{F}} = M(|\bar{\mathbf{h}}_{(k)}| + \frac{1}{4}|\bar{\mathbf{h}}_{(k)}|^2)$ (line 17). In essence, coefficients $\phi_{l\mathbf{x}_i}^{(R)}$ with $R \in [0, 2]$ are functions of the label y_{il} only. Therefore, we can perform the perturbation by injecting Laplace noise $1/|L| \text{Lap}(\frac{\Delta_{\mathcal{F}}}{\epsilon_3})$ into $\phi_{l\mathbf{x}_i}^{(R)}$ for every training label $y_{il} \in D$ as a preprocessing step (lines 18-19). Then, the perturbed coefficients, denoted $\bar{\phi}_{l\mathbf{x}_i}^{(R)}$ are used to construct the differentially private loss function $\bar{\mathcal{F}}_L(\theta_t)$ (line 27) during the training process without accessing the original label y_{il} again (lines 20-30). Stochastic gradient descent and back-propagation algorithms are used to minimize the perturbed loss function $\bar{\mathcal{F}}_L(\theta_t)$.

Now, we are ready to state that the computation of $\bar{\mathcal{F}}_L(\theta_t)$ is ϵ_3 -differentially private, and our mechanism preserves $(\epsilon_1 + \epsilon_2 + \epsilon_3)$ -differential privacy in the following lemmas.

Lemma 5: Let L and L' be any two neighboring batches. Let $\hat{\mathcal{F}}_L(\theta)$ and $\hat{\mathcal{F}}_{L'}(\theta)$ be the loss functions on L and L' respectively, then we have the following inequality:

$$\Delta_{\mathcal{F}} = \sum_{l=1}^M \sum_{R=0}^2 \left\| \sum_{\mathbf{x}_i \in L} \phi_{l\mathbf{x}_i}^{(R)} - \sum_{\mathbf{x}'_i \in L'} \phi_{l\mathbf{x}'_i}^{(R)} \right\| \leq M(|\bar{\mathbf{h}}_{(k)}| + \frac{1}{4}|\bar{\mathbf{h}}_{(k)}|^2)$$

where $|\bar{\mathbf{h}}_{(k)}|$ is the number of hidden neurons in $\bar{\mathbf{h}}_{(k)}$.

Proof 5: The proof of Lemma 5 is in **Appendix D**¹.

Lemma 6: Algorithm 1 preserves ϵ_3 -differential privacy in the computation of $\bar{\mathcal{F}}_L(\theta_t)$ (line 27).

Proof 6: The proof of Lemma 6 is in **Appendix E**¹.

D. The Correctness and Characteristics of the AdLM

The following theorem illustrates that the Alg. 1 preserves ϵ -differential privacy, where $\epsilon = \epsilon_1 + \epsilon_2 + \epsilon_3$.

Theorem 1: Algorithm 1 preserves ϵ -differential privacy, where $\epsilon = \epsilon_1 + \epsilon_2 + \epsilon_3$.

Proof: At a specific training step $t \in T$, it is crystal clear that the computation of $\bar{\mathbf{h}}_{0L}$ is ϵ_2 -differentially private (Lemma 4). Therefore, the computation of the hidden layers $\mathbf{h}_1, \dots, \mathbf{h}_k$ and normalization layers $\bar{\mathbf{h}}$ are differentially private. This is because they do not access any additional information from the data. At the output layer, the loss function $\bar{\mathcal{F}}_L(\theta_t)$ is ϵ_3 -differentially private (Lemma 6). The computation of gradients $\frac{1}{|L|} \nabla_{\theta_t} \bar{\mathcal{F}}_L(\theta_t)$ and descents is an optimization process, without using any additional information from the original data. Thus, θ_{t+1} is differentially private (line 29, Alg. 1).

This optimization process is repeated through T steps without querying the original data D (lines 21-30). This is done because Laplace noise is injected into input features x_{ij} and coefficients $\phi_{l\mathbf{x}_i}^{(R)}$ as preprocessing steps (lines 3-19). Note that \bar{x}_{ij} and $\bar{\phi}_{l\mathbf{x}_i}^{(R)}$ are associated with features \mathbf{x}_i and the label y_i respectively. $\bar{\mathbf{h}}_{0L}$ and $\bar{\mathcal{F}}_L(\theta_t)$ are computed based on \bar{x}_{ij} and $\bar{\phi}_{l\mathbf{x}_i}^{(R)}$. As a result, the noise and privacy budget consumption will not be accumulated during the training process.

Finally, $\bar{\mathcal{F}}_L(\theta_t)$ uses the outputs of $\bar{\mathbf{h}}_{0L}$, which essentially uses the differentially private relevances $\bar{\mathbf{R}}(D)$ as of one inputs. $\bar{\mathbf{R}}(D)$, $\bar{\mathbf{h}}_{0L}$, and $\bar{\mathcal{F}}_L(\theta_t)$ are achieved by applying ϵ_1, ϵ_2 , and ϵ_3 -differential privacy mechanisms. Furthermore, $\bar{\mathcal{F}}_L(\theta_t)$ and $\bar{\mathbf{h}}_{0L}$ access the same training batch L at each training step. Therefore, based on the composition theorem [11], the total privacy budget in Alg. 1 must be the summation of ϵ_1, ϵ_2 , and ϵ_3 .

Consequently, Algorithm 1 preserves ϵ -differential privacy, where $\epsilon = \epsilon_1 + \epsilon_2 + \epsilon_3$. ■

Note that $\Delta_{\mathbf{R}}$ and $\Delta_{\mathbf{h}_0}$ are dependent on the number of input features d . $\Delta_{\mathbf{R}}$ is negatively proportional to the number of tuples in D . The larger the size of D , the less noise will be injected into the private relevance $\bar{\mathbf{R}}$. $\Delta_{\mathcal{F}}$ is only dependent on the number of neurons in the last hidden layer and the output layer. In addition, $\Delta_{\mathbf{R}}, \Delta_{\mathbf{h}_0}$, and $\Delta_{\mathcal{F}}$ do not depend on the number of training epochs. Consequently:

(1) The privacy budget consumption in our model is totally independent of the number of training epochs.

(2) In order to improve the model utility under differential privacy, our mechanism adaptively injects Laplace noise into features based on the contribution of each to the model output.

(3) The average error incurred by our approximation approach, $\hat{\mathcal{F}}_L(\theta)$, is bounded by a small number $M \times \frac{e^2 + 2e - 1}{e(1+e)^2}$ (Please refer to Lemma 7 in **Appendix F**¹).

(4) The proposed mechanism can be applied to a variety of deep learning models, e.g., CNNs [30], deep auto-encoders [13], Restricted Boltzmann Machines [31], convolution deep

belief networks [32], etc., as long as we can perturb the first affine transformation layer.

With these characteristics, our mechanism has a great potential to be applied in large datasets, without consuming excessive privacy budgets. In the experiment section, we will show that our mechanism leads to accurate results.

IV. EXPERIMENTAL RESULTS

We have carried out an extensive experiment on two well-known image datasets, MNIST and CIFAR-10. The MNIST database of handwritten digits consists of 60,000 training examples, and a test set of 10,000 examples [16]. Each example is a 28×28 size gray-level image. The CIFAR-10 dataset consists of color images categorized into 10 classes, such as birds, dogs, trucks, airplanes, etc. The dataset is partitioned into 50,000 training examples and 10,000 test examples [17].

Competitive Models. We compare our mechanism with the state-of-the-art differentially private stochastic gradient descent (**pSGD**) for deep learning proposed by [12]. CNNs are used in our experiments for both algorithms. The hyper-parameters in pSGD are set to the default values recommended by Abadi et al. [12]. To comprehensively examine the proposed approaches, our mechanism is implemented in two different settings: **(1)** The Adaptive Laplace Mechanism (Alg. 1)-based CNN with ReLUs, simply denoted **AdLM**. We set $\epsilon_1 = \epsilon_2 = \epsilon_3 = \epsilon/3$ in the AdLM algorithm; and **(2)** An **Identical Laplace Mechanism**-based CNN with ReLUs (**ILM**), in which an identical Laplace noise $\frac{1}{|L|}Lap(\frac{\Delta_{h_0}}{\epsilon_2})$ is injected into each feature x_{ij} to preserve ϵ_2 -differential privacy in the computation of the affine transformation layer h_0 . In the ILM algorithm, we do not need to use the differentially private relevances $\bar{\mathbf{R}}(D)$. Therefore, the total privacy budget becomes $\epsilon = \epsilon_2 + \epsilon_3$. We set $\epsilon_2 = \epsilon_3 = \epsilon/2$ in the ILM algorithm. More details about the ILM are in the **Appendix C**¹. *The implementation of our mechanism is publicly available*².

A. MNIST Dataset

The designs of the three models are the same on the MNIST dataset. We used two convolution layers, one with 32 features and one with 64 features. Each hidden neuron connects with a 5×5 unit patch. A fully-connected layer with 1,024 units and an output layer of 10 classes (i.e., 10 digits) with cross-entropy loss with LRN are used. The batch size is 128. This also is the structure of the pre-trained model, which is learned and used to compute the average relevances $\mathbf{R}(D)$. The experiments were conducted on a single GPU, i.e., NVIDIA GTX TITAN X, 12 GB with 3,072 CUDA cores.

Fig. 4a illustrates the prediction accuracy of each model as a function of the privacy budget ϵ on the MNIST dataset. It is clear that our models, i.e., AdLM and ILM, outperform the pSGD, especially when the privacy budget ϵ is small. This is a crucial result, since smaller privacy budget values enforce stronger privacy guarantees. When the privacy budget ϵ is large, e.g., $\epsilon = 2, 4, 8$, which means small noise is injected

into the model, the efficiencies of all the models are almost converged to higher prediction accuracies.

The AdLM model achieves the best performance. Given a very small privacy budget $\epsilon = 0.25$, it achieves 91.12% in terms of prediction accuracy, compared with 89.45% obtained by the ILM and 63.6% obtained by the pSGD. Overall, given small values of the privacy budget ϵ , i.e., $0.2 \leq \epsilon \leq 0.5$, the AdLM improves the prediction accuracy by 15% on average (i.e., 90.63%) compared with the pSGD (i.e., 75.52%). The result is statistically significant with $p < 0.01$ (t-test).

Figs. 4b-c illustrate the prediction accuracy of each model vs. the number of epochs under $\epsilon = 0.5$ and $\epsilon = 2.0$ respectively. Given large noise, i.e., $\epsilon = 0.5$, the pSGD quickly achieves higher prediction accuracies (i.e., 88.75%) after a small number of epochs, compared with other models (Fig. 4b). However, the pSGD can only be applied to train the model by using a limited number of epochs; specifically because the privacy budget is accumulated after every training step. Meanwhile, our mechanism is totally independent of the number of epochs in the consumption of privacy budget. Therefore, after 160 epochs, our models outperform the pSGD. The AdLM achieves the best performance, in terms of prediction accuracy: 96.12%, whereas the ILM and the pSGD reached only 94.81% and 88.75%, respectively. Interestingly, given small noise, i.e., $\epsilon = 2.0$, our models achieve higher accuracies than the pSGD after a small number of epochs (Fig. 4c). This result illustrates the crucial benefits of being independent of the number of training epochs in preserving differential privacy in deep learning. With our mechanism, we can keep training our models without accumulating noise and privacy budget.

B. CIFAR-10 Dataset

The designs of the three models are the same on the CIFAR-10 dataset. We used two convolution layers, one with 32 features and one with 64 features. Each hidden neuron connects with a 5×5 unit patch. Two fully-connected layers with 384 neurons and 192 neurons, and an output layer of 10 classes with a cross-entropy loss with LRN are used. The batch size is set to 128. This also is the structure of the pre-trained model, which is learned and used to compute the average relevances $\mathbf{R}(D)$.

Fig. 5a shows the prediction accuracies of each model as a function of the privacy budget ϵ on the CIFAR-10 dataset. Figs. 5b-c illustrate the prediction accuracy of each model vs. the number of epochs under different noise levels. Similar to the results on the MNIST dataset, the results on CIFAR-10 strengthen our observations: **(1)** Our mechanism outperforms the pSGD in terms of prediction accuracy, given both modest and large values of the privacy budget ϵ (Fig. 5a); and **(2)** Our mechanism has the ability to work with large-scale datasets, since it is totally independent of the number of training epochs in the consumption of privacy budget (Figs. 5b-c).

In fact, the AdLM improves the prediction accuracy by 14% on average (i.e., to 84.8%) compared with the pSGD (i.e., 70.7%). The result is statistically significant with $p < 0.01$ (t-test). Given large noise, i.e., $\epsilon = 2.5$, our models including

²<https://www.dropbox.com/s/i1uzj7fpzgcag7k/AdLM.zip?dl=0>

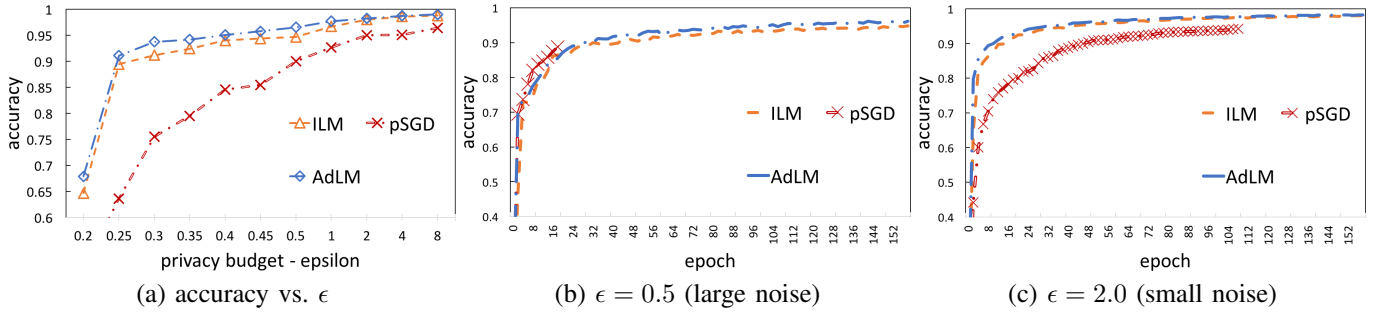


Fig. 4. Accuracy for different noise levels on the MNIST dataset.

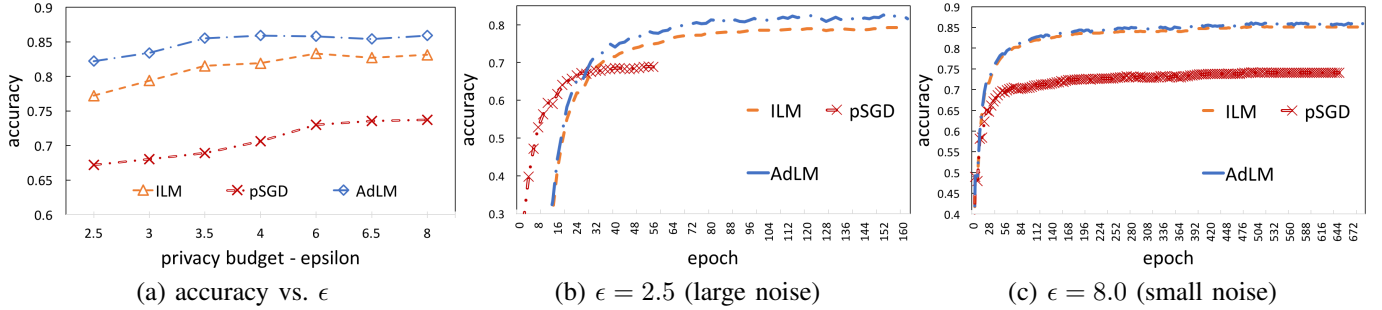


Fig. 5. Accuracy for different noise levels on the CIFAR-10 dataset.

the AdLM and ILM significantly outperform the pSGD after 160 epochs (Fig. 5b).

C. Adaptive Laplace Noise

It is important to note that by adaptively redistributing the noise into input features based on the relevance of each to the model output, we can achieve much better prediction accuracies in both MNIST and CIFAR-10 datasets given both small and large values of privacy budget ϵ . This is clearly demonstrated in Figs. 4a and 5a, since the AdLM outperforms the ILM in all cases. Overall, the AdLM improves the prediction accuracy by 2% and 5% on average on MNIST and CIFAR-10 datasets correspondingly, compared with the ILM. The result is statistically significant with $p < 0.05$ (t-test). Note that the ILM injected an identical amount of noise into all input features, regardless of their contributions to the model output. This is an important result, since our mechanism is the first of its kind, which can redistribute the noise injected into the deep learning model to improve the utility. In addition, the reallocation of ϵ_1 , ϵ_2 , and ϵ_3 could further improve the utility. This is an open research direction in the future work.

D. Computational Efficiency

In terms of computation efficiency, there are two differences in our mechanism, compared with a regular deep neural network: **(i)** The pre-trained model; and **(ii)** The noise injection task. In practice, the pre-trained model is not necessarily identical to the differentially private network trained by our AdLM. A simple model can be used as a pre-trained model to approximate the average relevance $\mathbf{R}(D)$, as long as the pre-trained model is effective in terms of prediction accuracy even over a small training dataset. Achieving this is quite straightforward, because: **(1)** The pre-trained model is noiseless; and

(2) The number of training epochs used to learn a pre-trained model is small compared with the one of differentially private models. In fact, we only need 12 extra epochs to learn the pre-trained models on MNIST and CIFAR-10 datasets. Training 12 extra epochs takes less than 1 minute on a single GPU, i.e., NVIDIA GTX TITAN X, 12 GB with 3,072 CUDA cores. Therefore, the model pre-training for the computation of $\mathbf{R}(D)$ is efficient.

Another difference in our mechanism is the noise injection into input attributes and coefficients of the loss function $\hat{\mathcal{F}}$. In this task, the computations of $\Delta_{\mathbf{R}}$, $\Delta_{\mathbf{h}_0}$, $\Delta_{\mathcal{F}}$, \bar{R}_j , β_j , \bar{x} , and $\hat{\phi}$ are efficient and straight-forward, since there is not any operation such as $\arg \min$, $\arg \max$, sorting, etc. The complexity of these computations is $O(|D|(d + M))$, which is linear to the size of the database D . In addition, these computations can be efficiently performed in either a serial process or a parallel process. Therefore, this task does not affect the computational efficiency of our mechanism much.

V. CONCLUSIONS

In this paper, we proposed a novel mechanism, called Adaptive Laplace Mechanism (AdLM), to preserve differential privacy in deep learning. Our mechanism conducts both sensitivity analysis and noise insertion on deep neural networks. It is totally independent of the number of training epochs in the consumption of privacy budget. That makes our mechanism more practical. In addition, our mechanism is the first of its kind to have the ability to redistribute the noise insertion toward the improvement of model utility in deep learning. In fact, our mechanism has the ability to intentionally add *more noise* into input features which are *less relevant* to the model output, and vice-versa. Different activation functions

can be applied in our mechanism, as well. These distinctive characteristics guarantee the ability to apply our mechanism on large datasets in different deep learning models and in different contexts. Our mechanism can clearly enhance the application of differential privacy in deep learning. Rigorous experimental evaluations conducted on well-known datasets validated our theoretical results and the effectiveness of our mechanism.

ACKNOWLEDGMENT

This work is supported by the NIH grant R01GM103309 to the SMASH project. Wu is also supported by NSF grant DGE-1523115 and IIS-1502273.

REFERENCES

- [1] J. Levinson, J. Askeland, J. Becker, J. Dolson, D. Held, S. Kammel, J. Z. Kolter, D. Langer, O. Pink, V. Pratt, M. Sokolsky, G. Stanek, D. Stavens, A. Teichman, M. Werling, and S. Thrun, "Towards fully autonomous driving: Systems and algorithms," in *2011 IEEE Intelligent Vehicles Symposium (IV)*, 2011, pp. 163–168.
- [2] D. Lowd, "Good word attacks on statistical spam filters," in *CEAS'05*, 2005.
- [3] Y. Cheng, F. Wang, P. Zhang, and J. Hu, "Risk prediction with electronic health records: A deep learning approach," in *SDM'16*, 2016.
- [4] E. Choi, A. Schuetz, W. F. Stewart, and J. Sun, "Using recurrent neural network models for early detection of heart failure onset," *JAMIA*, 2016.
- [5] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," *Theory of Cryptography*, pp. 265–284, 2006.
- [6] K. Chaudhuri and C. Monteleoni, "Privacy-preserving logistic regression," in *NIPS'08*, 2008, pp. 289–296.
- [7] F. McSherry and I. Mironov, "Differentially Private Recommender Systems," in *KDD'09*. ACM, 2009.
- [8] F. McSherry and K. Talwar, "Mechanism design via differential privacy," in *FOCS '07*, 2007, pp. 94–103.
- [9] R. Shokri and V. Shmatikov, "Privacy-preserving deep learning," in *CCS'15*, 2015, pp. 1310–1321.
- [10] N. Phan, Y. Wang, X. Wu, and D. Dou, "Differential privacy preservation for deep auto-encoders: an application of human behavior prediction," in *AAAI'16*, 2016, pp. 1309–1316.
- [11] C. Dwork and J. Lei, "Differential privacy and robust statistics," in *STOC'09*, 2009, pp. 371–380.
- [12] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," *arXiv:1607.00133*, 2016.
- [13] Y. Bengio, "Learning deep architectures for ai," *Found. Trends Mach. Learn.*, vol. 2, no. 1, pp. 1–127, 2009.
- [14] N. Dowlin, R. Gilad-Bachrach, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, "Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy," in *ICML'16*, 2016, pp. 201–210.
- [15] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation," *PLoS ONE*, vol. 10, no. 7, p. e0130140, 07 2015.
- [16] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [17] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," 2009.
- [18] D. Kifer and A. Machanavajjhala, "No free lunch in data privacy," in *SIGMOD'11*, 2011, pp. 193–204.
- [19] U. Erlingsson, V. Pihur, and A. Korolova, "Rappor: Randomized aggregatable privacy-preserving ordinal response," in *CCS'14*, 2014, pp. 1054–1067.
- [20] C. Liu, S. Chakraborty, and P. Mittal, "Dependence Makes You Vulnerable: Differential Privacy Under Dependent Tuples," in *NDSS'16*, Feb. 2016.
- [21] Y. Cao, M. Yoshikawa, Y. Xiao, and L. Xiong, "Quantifying differential privacy under temporal correlations," *CoRR*, vol. abs/1610.07543, 2016.
- [22] T.-H. H. Chan, M. Li, E. Shi, and W. Xu, "Differentially private continual monitoring of heavy hitters from distributed streams," in *PETS'12*, 2012, pp. 140–159.
- [23] S. Song, K. Chaudhuri, and A. D. Sarwate, "Stochastic gradient descent with differentially private updates," in *GlobalSIP*, 2013, pp. 245–248.
- [24] P. Jain, P. Kothari, and A. Thakurta, "Differentially private online learning," in *COLT'12*, 2012, pp. 24.1–24.34.
- [25] X. Xiao, G. Wang, and J. Gehrke, "Differential privacy via wavelet transforms," in *ICDE'10*, 2010, pp. 225–236.
- [26] Y. Wang, X. Wu, and L. Wu, "Differential privacy preserving spectral graph analysis," in *PAKDD (2)*, 2013, pp. 329–340.
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS'12*, 2012, pp. 1097–1105.
- [28] G. Arfken, in *Mathematical Methods for Physicists (Third Edition)*. Academic Press, 1985.
- [29] J. Zhang, Z. Zhang, X. Xiao, Y. Yang, and M. Winslett, "Functional mechanism: regression analysis under differential privacy," *PVLDB*, vol. 5, no. 11, pp. 1364–1375, 2012.
- [30] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [31] P. Smolensky, "Information processing in dynamical systems: Foundations of harmony theory," 1986, ch. Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1, pp. 194–281.
- [32] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *ICML'09*, 2009, pp. 609–616.
- [33] T. Apostol, *Calculus*. John Wiley & Sons, 1967.

APPENDIX

A. Proof of Lemma 1

Proof 1: Assume that D and D' differ in the last tuple. Let \mathbf{x}_n (\mathbf{x}'_n) be the last tuple in D (D'). We have that

$$\begin{aligned} \Delta_{\mathbf{R}} &= \frac{1}{|D|} \sum_{j=1}^d \left\| \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i) - \sum_{\mathbf{x}'_i \in D'} R_{x'_{ij}}(\mathbf{x}'_i) \right\|_1 \\ &= \frac{1}{|D|} \sum_{j=1}^d \left\| R_{x_{nj}}(\mathbf{x}_n) - R_{x'_{nj}}(\mathbf{x}'_n) \right\|_1 \\ &\leq \frac{2}{|D|} \max_{\mathbf{x}_i \in D} \sum_{j=1}^d \|R_{x_{ij}}(\mathbf{x}_i)\|_1 \leq \frac{2d}{|D|} \end{aligned}$$

Eq. 13 holds.

B. Proof of Lemma 2

Proof 2: Let D and D' be two neighbor databases. Without loss of generality, assume that D and D' differ in the last tuple \mathbf{x}_n (\mathbf{x}'_n). $\bar{R}_j(D)$ is calculated as done in line 6, Alg. 1, and $\bar{\mathbf{R}}(D) = \{\bar{R}_j(D)\}_{j \in [1, d]}$ is the output of the Alg. 1 (line 7). The perturbation of the relevance $R_j(D)$ can be rewritten as:

$$\bar{R}_j = \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_{x_{ij}}(\mathbf{x}_i) + \text{Lap}\left(\frac{\Delta_{\mathbf{R}}}{\epsilon_1}\right) \quad (26)$$

Since all the input features are perturbed, we have that

$$\begin{aligned} \frac{\Pr(\bar{\mathbf{R}}(D))}{\Pr(\bar{\mathbf{R}}(D'))} &= \frac{\prod_{j=1}^d \exp\left(\frac{\epsilon_1 \|\frac{1}{|D|} \sum_{\mathbf{x}_i \in D} R_j(\mathbf{x}_i) - \bar{R}_j\|_1}{\Delta_{\mathbf{R}}}\right)}{\prod_{j=1}^d \exp\left(\frac{\epsilon_1 \|\frac{1}{|D|} \sum_{\mathbf{x}'_i \in D'} R_j(\mathbf{x}'_i) - \bar{R}_j\|_1}{\Delta_{\mathbf{R}}}\right)} \\ &\leq \prod_{j=1}^d \exp\left(\frac{\epsilon_1}{|D| \Delta_{\mathbf{R}}} \left\| \sum_{\mathbf{x}_i \in D} R_j(\mathbf{x}_i) - \sum_{\mathbf{x}'_i \in D'} R_j(\mathbf{x}'_i) \right\|_1\right) \\ &\leq \prod_{j=1}^d \exp\left(\frac{\epsilon_1}{|D| \Delta_{\mathbf{R}}} \left\| R_j(\mathbf{x}_n) - R_j(\mathbf{x}'_n) \right\|_1\right) \\ &\leq \prod_{j=1}^d \exp\left(\frac{\epsilon_1}{|D| \Delta_{\mathbf{R}}} 2 \max_{\mathbf{x}_n \in D} \|R_j(\mathbf{x}_n)\|_1\right) \\ &\leq \exp\left(\epsilon_1 \frac{2 \max_{\mathbf{x}_n \in L} \sum_{j=1}^d \|R_j(\mathbf{x}_n)\|_1}{|D| \Delta_{\mathbf{R}}}\right) \\ &\leq \exp(\epsilon_1) \end{aligned}$$

Consequently, the computation of $\bar{\mathbf{R}}(D)$ preserves ϵ_1 -differential privacy in Alg. 1.

C. Identical Laplace Mechanism (ILM)

We can add an identical noise distribution $\frac{1}{|L|} \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)$ to all input features, where $\Delta_{\mathbf{h}_0} = 2 \sum_{h \in \mathbf{h}_0} d$ (line 9, Alg. 1) to preserve differential privacy in the computation of \mathbf{h}_0 . In fact,

we have that $\forall \mathbf{x}_i \in D, j \in [1, d] : \bar{x}_{ij} = x_{ij} + \frac{1}{|L|} \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)$.

For each $h \in \bar{\mathbf{h}}_{0L}$, h can be re-written as:

$$\begin{aligned} \bar{h}_L(W) &= \sum_{j=1}^d \left[\sum_{\mathbf{x}_i \in L} \left(x_{ij} + \frac{1}{|L|} \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)\right) W^T \right] \\ &\quad + \sum_{\mathbf{x}_i \in L} \left(b + \frac{1}{|L|} \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)\right) \quad (27) \end{aligned}$$

Let us consider the static bias $b = 1$ as the 0-th input feature and its associated parameter W_b , i.e., $x_{i0} = b = 1$ and $W = W_b \cup W$, we have that

$$\begin{aligned} \bar{h}_L(W) &= \sum_{j=0}^d \left[\sum_{\mathbf{x}_i \in L} \left(x_{ij} + \frac{1}{|L|} \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)\right) W^T \right] \quad (28) \\ &= \sum_{j=0}^d \left[\sum_{\mathbf{x}_i \in L} x_{ij} + \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right) \right] W^T = \sum_{j=0}^d \bar{\phi}_j^h W^T \quad (29) \end{aligned}$$

where $\bar{\phi}_j^h = [\sum_{\mathbf{x}_i \in L} x_{ij} + \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)]$.

We can see that $\bar{\phi}_j^h$ is the perturbation of the input feature x_{ij} associated with the j -th parameter $W_j \in W$ of the hidden neuron h on L . Since all the hidden neurons h in \mathbf{h}_0 are perturbed, we have that:

$$\Pr(\bar{\mathbf{h}}_{0L}(W_0)) = \prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_2 \|\sum_{\mathbf{x}_i \in L} x_{ij} - \bar{\phi}_j^h\|}{\Delta_{\mathbf{h}_0}}\right)$$

$\Delta_{\mathbf{h}_0}$ is set to $2 \sum_{h \in \mathbf{h}_0} d$ and $\bar{\mathbf{h}}_{0L}(W_0) = \{\bar{h}_L(W)\}_{h \in \mathbf{h}_0}$ is the output, we have that

$$\begin{aligned} \frac{\Pr(\bar{\mathbf{h}}_{0L}(W_0))}{\Pr(\bar{\mathbf{h}}_{0L'}(W_0))} &= \frac{\prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_2 \|\sum_{\mathbf{x}_i \in L} x_{ij} - \bar{\phi}_j^h\|_1}{\Delta_{\mathbf{h}_0}}\right)}{\prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_2 \|\sum_{\mathbf{x}'_i \in L'} x'_{ij} - \bar{\phi}_j^h\|_1}{\Delta_{\mathbf{h}_0}}\right)} \\ &\leq \prod_{h \in \mathbf{h}_0} \prod_{j=0}^d \exp\left(\frac{\epsilon_2}{\Delta_{\mathbf{h}_0}} \left\| \sum_{\mathbf{x}_i \in L} x_{ij} - \sum_{\mathbf{x}'_i \in L'} x'_{ij} \right\|_1\right) \\ &\leq \prod_{h \in \mathbf{h}_0} \prod_{j=1}^d \exp\left(\frac{\epsilon_2}{\Delta_{\mathbf{h}_0}} \left\| x_{nj} - x'_{nj} \right\|_1\right) \\ &\leq \prod_{h \in \mathbf{h}_0} \prod_{j=1}^d \exp\left(\frac{\epsilon_2}{\Delta_{\mathbf{h}_0}} 2 \max_{\mathbf{x}_n \in L} \|x_{nj}\|_1\right) \leq \prod_{h \in \mathbf{h}_0} \prod_{j=1}^d \exp\left(\frac{2\epsilon_2}{\Delta_{\mathbf{h}_0}}\right) \\ &\leq \exp\left(\epsilon_2 \frac{2 \sum_{h \in \mathbf{h}_0} d}{\Delta_{\mathbf{h}_0}}\right) = \exp(\epsilon_2) \end{aligned}$$

Consequently, based on the above analysis, the computation of $\bar{\mathbf{h}}_{0L}(W_0)$ preserves ϵ_2 -differential privacy in Alg. 1 by injecting an identical Laplace noise $\frac{1}{|L|} \text{Lap}\left(\frac{\Delta_{\mathbf{h}_0}}{\epsilon_2}\right)$ into all input features. In addition, given the identical Laplace noise, we do not need to use the differentially private relevance $\bar{\mathbf{R}}(D)$, since we do not need to redistribute the noise in the first affine transformation layer \mathbf{h}_0 .

D. Proof of Lemma 5

Proof 5: Assume that L and L' differ in the last tuple. Let \mathbf{x}_n (\mathbf{x}'_n) be the last tuple in L (L'). We have that

$$\begin{aligned}\Delta_{\mathcal{F}} &= \sum_{l=1}^M \sum_{R=0}^2 \left\| \sum_{\mathbf{x}_i \in L} \phi_{l\mathbf{x}_i}^{(R)} - \sum_{\mathbf{x}'_i \in L'} \phi_{l\mathbf{x}'_i}^{(R)} \right\| \\ &= \sum_{l=1}^M \sum_{R=0}^2 \left\| \phi_{l\mathbf{x}_n}^{(R)} - \phi_{l\mathbf{x}'_n}^{(R)} \right\|\end{aligned}\quad (30)$$

We can show that $\phi_{l\mathbf{x}_n}^{(0)} = \sum_{q=1}^2 f_{ql}^{(0)}(0) = y_{nl} \log 2 + (1 - y_{nl}) \log 2 = \log 2$. Similarly, we can show that $\phi_{l\mathbf{x}'_n}^{(0)} = \log 2$. As a result, $\phi_{l\mathbf{x}_n}^{(0)} = \phi_{l\mathbf{x}'_n}^{(0)}$. Therefore

$$\begin{aligned}\Delta_{\mathcal{F}} &= \sum_{l=1}^M \sum_{R=0}^2 \left\| \phi_{l\mathbf{x}_n}^{(R)} - \phi_{l\mathbf{x}'_n}^{(R)} \right\| = \sum_{l=1}^M \sum_{R=1}^2 \left\| \phi_{l\mathbf{x}_n}^{(R)} - \phi_{l\mathbf{x}'_n}^{(R)} \right\| \\ &\leq \sum_{l=1}^M \sum_{R=1}^2 \left(\left\| \phi_{l\mathbf{x}_n}^{(R)} \right\| + \left\| \phi_{l\mathbf{x}'_n}^{(R)} \right\| \right) \leq 2 \max_{\mathbf{x}_n} \sum_{l=1}^M \sum_{R=1}^2 \left\| \phi_{l\mathbf{x}_n}^{(R)} \right\| \\ &\leq 2 \max_{\mathbf{x}_n} \left[\sum_{l=1}^M \left(\frac{1}{2} - y_{nl} \right) \sum_{e=1}^{|\bar{\mathbf{h}}_{(k)}|} \bar{\mathbf{h}}_{e\mathbf{x}_n(k)} \right. \\ &\quad \left. + \sum_{l=1}^M \left(\frac{1}{8} \sum_{e,g} \bar{\mathbf{h}}_{e\mathbf{x}_n(k)} \bar{\mathbf{h}}_{g\mathbf{x}_n(k)} \right) \right] \\ &\leq 2 \left(\frac{1}{2} M \times |\bar{\mathbf{h}}_{(k)}| + \frac{1}{8} M \times |\bar{\mathbf{h}}_{(k)}|^2 \right) \\ &= M \left(|\bar{\mathbf{h}}_{(k)}| + \frac{1}{4} |\bar{\mathbf{h}}_{(k)}|^2 \right)\end{aligned}$$

where $\bar{\mathbf{h}}_{e\mathbf{x}_n(k)}$ is the state of e -th hidden neuron in $\bar{\mathbf{h}}_{(k)}$.

E. Proof of Lemma 6

Proof 6: Let L and L' be two neighbor batches. Without loss of generality, assume that L and L' differ in the last tuple \mathbf{x}_n (\mathbf{x}'_n). $\Delta_{\mathcal{F}}$ is calculated as done in line 17, Alg. 1, and $\bar{\mathcal{F}}_L(\theta_t) = \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \sum_{R=0}^2 \bar{\phi}_{l\mathbf{x}_i}^{(R)} (\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T)^R$ is the output of line 27 of the Alg. 1. Note that $\bar{\mathbf{h}}_{\mathbf{x}_i(k)}$ is the state of $\bar{\mathbf{h}}_{(k)}$ derived from $\bar{\mathbf{h}}_{0\mathbf{x}_i}$ by navigating through the neural network. The perturbation of the coefficient $\phi_l^{(R)}$, denoted as $\bar{\phi}_l^{(R)}$, can be rewritten as:

$$\bar{\phi}_l^{(R)} = \left[\sum_{\mathbf{x}_i \in L} \phi_{l\mathbf{x}_i}^{(R)} + \text{Lap} \left(\frac{\Delta_{\mathcal{F}}}{\epsilon_3} \right) \right] \quad (31)$$

We can see that $\bar{\phi}_l^{(R)}$ is the perturbation of the coefficient $\phi_l^{(R)}$ associated with the labels y_{il} in the training batch L . We

have that

$$\begin{aligned}\frac{Pr(\bar{\mathcal{F}}_L(\theta_t))}{Pr(\bar{\mathcal{F}}_{L'}(\theta_t))} &= \frac{\prod_{l=1}^M \prod_{R=0}^2 \exp \left(\frac{\epsilon_3 \left\| \sum_{\mathbf{x}_i \in L} \phi_{l\mathbf{x}_i}^{(R)} - \bar{\phi}_l^{(R)} \right\|_1}{\Delta_{\mathcal{F}}} \right)}{\prod_{l=1}^M \prod_{R=0}^2 \exp \left(\frac{\epsilon_3 \left\| \sum_{\mathbf{x}'_i \in L'} \phi_{l\mathbf{x}'_i}^{(R)} - \bar{\phi}_l^{(R)} \right\|_1}{\Delta_{\mathcal{F}}} \right)} \\ &\leq \prod_{l=1}^M \prod_{R=0}^2 \exp \left(\frac{\epsilon_3}{\Delta_{\mathcal{F}}} \left\| \sum_{\mathbf{x}_i \in L} \phi_{l\mathbf{x}_i}^{(R)} - \sum_{\mathbf{x}'_i \in L'} \phi_{l\mathbf{x}'_i}^{(R)} \right\|_1 \right) \\ &\leq \prod_{l=1}^M \prod_{R=0}^2 \exp \left(\frac{\epsilon_3}{\Delta_{\mathcal{F}}} \left\| \phi_{l\mathbf{x}_n}^{(R)} - \phi_{l\mathbf{x}'_n}^{(R)} \right\|_1 \right) \\ &\leq \prod_{l=1}^M \prod_{R=0}^2 \exp \left(\frac{\epsilon_3}{\Delta_{\mathcal{F}}} 2 \max_{\mathbf{x}_n \in L} \left\| \phi_{l\mathbf{x}_n}^{(R)} \right\|_1 \right) \\ &\leq \exp \left(\epsilon_3 \frac{2 \max_{\mathbf{x}_n \in L} \sum_{l=1}^M \sum_{R=0}^2 \left\| \phi_{l\mathbf{x}_n}^{(R)} \right\|_1}{\Delta_{\mathcal{F}}} \right) \\ &\leq \exp \left(\epsilon_3 \frac{M \left(|\bar{\mathbf{h}}_{(k)}| + \frac{1}{4} |\bar{\mathbf{h}}_{(k)}|^2 \right)}{\Delta_{\mathcal{F}}} \right) = \exp(\epsilon_3)\end{aligned}$$

Consequently, the computation of $\bar{\mathcal{F}}_L(\theta_t)$ preserves ϵ_3 -differential privacy in Alg. 1.

F. Approximation Error Bounds

The following lemma illustrates the result of how much error our approximation approach, $\hat{\mathcal{F}}_L(\theta)$ (Eq. 25), incurs. The error only depends on the number of possible classification outcomes M . In addition, the average error of the approximations is always bounded. As in [10, 29], the approximation of the loss function $\mathcal{F}_L(\theta)$ by applying Taylor Expansion without removing all polynomial terms with order larger than 2 is as follows:

$$\tilde{\mathcal{F}}_L(\theta) = \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \sum_{q=1}^2 \sum_{R=0}^{\infty} \frac{f_{ql}^{(R)}(z_{ql})}{R!} (g_{ql}(\bar{\mathbf{h}}_{\mathbf{x}_i(k)}, W_{l(k)} - z_{ql}))^R$$

$\forall l \in \{1, \dots, M\}$, let f_{1l} , f_{2l} , g_{1l} , and g_{2l} be four functions defined as follows:

$$g_{1l}(\bar{\mathbf{h}}_{\mathbf{x}_i(k)}, W_j) = \bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T \quad (32)$$

$$g_{2l}(\bar{\mathbf{h}}_{\mathbf{x}_i(k)}, W_j) = \bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T \quad (33)$$

$$f_{1l}(z_{1l}) = y_{il} \log(1 + e^{-z_{1l}}) \quad (34)$$

$$f_{2l}(z_{2l}) = (1 - y_{il}) \log(1 + e^{z_{2l}}) \quad (35)$$

where $\forall q, l : z_{ql}$ is a real number.

$\forall q, l$, by setting $z_{ql} = 0$, the above equation can be simplified as:

$$\tilde{\mathcal{F}}_L(\theta) = \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \sum_{q=1}^2 \sum_{R=0}^{\infty} \frac{f_{ql}^{(R)}(0)}{R!} (\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T)^R \quad (36)$$

As in [10], our approximation approach works by truncating the Taylor series in Eq. 36 to remove all polynomial terms with

order larger than 2. This leads to a new objective function in Eq. 25 with low-order polynomials as follows:

$$\begin{aligned}\widehat{\mathcal{F}}_L(\theta) &= \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \sum_{q=1}^2 \sum_{R=0}^2 \frac{f_{ql}^{(R)}(0)}{R!} (\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T)^R \\ &= \sum_{l=1}^M \sum_{\mathbf{x}_i \in L} \left[\sum_{q=1}^2 f_{ql}^{(0)}(0) + \left(\sum_{q=1}^2 f_{ql}^{(1)}(0) \right) \bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T \right. \\ &\quad \left. + \left(\sum_{q=1}^2 \frac{f_{ql}^{(2)}(0)}{2!} \right) (\bar{\mathbf{h}}_{\mathbf{x}_i(k)} W_{l(k)}^T)^2 \right]\end{aligned}$$

We are now ready to state the following lemma to show the approximation error bound of our approach.

Lemma 7: Given two polynomial functions $\widetilde{\mathcal{F}}_L(\theta)$ (Eq. 36) and $\widehat{\mathcal{F}}_L(\theta)$ (Eq. 25), the average error of the approximation is always bounded as follows:

$$|\widetilde{\mathcal{F}}_L(\widehat{\theta}) - \widetilde{\mathcal{F}}_L(\widetilde{\theta})| \leq M \times \frac{e^2 + 2e - 1}{e(1+e)^2} \quad (37)$$

where $\widetilde{\theta} = \arg \min_{\theta} \widetilde{\mathcal{F}}_L(\theta)$ and $\widehat{\theta} = \arg \min_{\theta} \widehat{\mathcal{F}}_L(\theta)$.

Proof 7: Let $\widetilde{\theta} = \arg \min_{\theta} \widetilde{\mathcal{F}}_L(\theta)$ and $\widehat{\theta} = \arg \min_{\theta} \widehat{\mathcal{F}}_L(\theta)$, $U = \max_{\theta} (\widetilde{\mathcal{F}}_L(\theta) - \widehat{\mathcal{F}}_L(\theta))$ and $S = \min_{\theta} (\widetilde{\mathcal{F}}_L(\theta) - \widehat{\mathcal{F}}_L(\theta))$. We have that $U \geq \widetilde{\mathcal{F}}_L(\widetilde{\theta}) - \widehat{\mathcal{F}}_L(\widehat{\theta})$ and $\forall \theta^* : S \leq \widetilde{\mathcal{F}}_L(\theta^*) - \widehat{\mathcal{F}}_L(\theta^*)$. Therefore, we have

$$\begin{aligned}\widetilde{\mathcal{F}}_L(\widehat{\theta}) - \widehat{\mathcal{F}}_L(\widehat{\theta}) - \widetilde{\mathcal{F}}_L(\theta^*) + \widehat{\mathcal{F}}_L(\theta^*) &\leq U - S \\ \Leftrightarrow \widetilde{\mathcal{F}}_L(\widehat{\theta}) - \widetilde{\mathcal{F}}_L(\theta^*) &\leq U - S + (\widehat{\mathcal{F}}_L(\widehat{\theta}) - \widehat{\mathcal{F}}_L(\theta^*))\end{aligned}$$

In addition, $\widehat{\mathcal{F}}_L(\widehat{\theta}) - \widehat{\mathcal{F}}_L(\theta^*) \leq 0$, so $\widetilde{\mathcal{F}}_L(\widehat{\theta}) - \widetilde{\mathcal{F}}_L(\theta^*) \leq U - S$. If $U \geq 0$ and $S \leq 0$ then we have:

$$|\widetilde{\mathcal{F}}_L(\widehat{\theta}) - \widetilde{\mathcal{F}}_L(\theta^*)| \leq U - S \quad (38)$$

Eq. 38 holds for every θ^* . Therefore, it still holds for $\widetilde{\theta}$. Eq. 38 shows that the error incurred by truncating the Taylor series approximate function depends on the maximum and minimum values of $\widetilde{\mathcal{F}}_L(\theta) - \widehat{\mathcal{F}}_L(\theta)$. To quantify the magnitude of the error, we first rewrite $\widetilde{\mathcal{F}}_L(\theta) - \widehat{\mathcal{F}}_L(\theta)$ as:

$$\begin{aligned}\widetilde{\mathcal{F}}_L(\theta) - \widehat{\mathcal{F}}_L(\theta) &= \sum_{l=1}^M \left[\widetilde{\mathcal{F}}_L(W_{l(k)}) - \widehat{\mathcal{F}}_L(W_{l(k)}) \right] \\ &= \sum_{l=1}^M \left[\sum_{\mathbf{x}_i \in L} \sum_{q=1}^2 \sum_{R=3}^{\infty} \frac{f_{ql}^{(R)}(z_{ql})}{R!} (g_{ql}(\bar{\mathbf{h}}_{\mathbf{x}_i(k)}, W_{l(k)}) - z_{ql})^R \right]\end{aligned}$$

To derive the minimum and maximum values of the function above, we look into the remainder of the Taylor Expansion for each l . Let $z_l \in [z_{ql} - 1, z_{ql} + 1]$. According to the well-known result [33], $\frac{1}{|D|} (\widetilde{\mathcal{F}}_L(W_{l(k)}) - \widehat{\mathcal{F}}_L(W_{l(k)}))$ must be in the interval $\left[\sum_q \frac{\min_{z_l} f_{ql}^{(3)}(z_l)(z_l - z_{ql})^3}{6}, \sum_l \frac{\max_{z_l} f_{ql}^{(3)}(z_l)(z_l - z_{ql})^3}{6} \right]$.

If $\sum_q \frac{\max_{z_l} f_{ql}^{(3)}(z_l)(z_l - z_{ql})^3}{6} \geq 0$ and $\sum_q \frac{\min_{z_l} f_{ql}^{(3)}(z_l)(z_l - z_{ql})^3}{6} \leq 0$, then we have that:

$$\begin{aligned}& \left| \frac{1}{|D|} \left[\widetilde{\mathcal{F}}_L(\theta) - \widehat{\mathcal{F}}_L(\theta) \right] \right| \\ & \leq \sum_{l=1}^M \sum_q \frac{\max_{z_l} f_{ql}^{(3)}(z_l)(z_l - z_{ql})^3 - \min_{z_l} f_{ql}^{(3)}(z_l)(z_l - z_{ql})^3}{6}\end{aligned} \quad (39)$$

This analysis applies to the case of the cross-entropy error-based loss function as follows. First, for the functions $f_{1l}(z_{1l}) = y_{il} \log(1 + e^{-z_{1l}})$ and $f_{2l}(z_{2l}) = (1 - y_{il}) \log(1 + e^{z_{2l}})$, we have

$$\begin{aligned}f_{1l}^{(3)}(z_{1l}) &= \frac{2y_{il}e^{z_{1l}}}{(1 + e^{z_{1l}})^3} \\ f_{2l}^{(3)}(z_{2l}) &= (1 - y_{il}) \frac{e^{-z_{2l}}(e^{-z_{2l}} - 1)}{(1 + e^{-z_{2l}})^3}\end{aligned}$$

It can be verified that $\arg \min_{z_{1l}} f_{1l}^{(3)}(z_{1l}) = \frac{-2e}{(1+e)^3} < 0$, $\arg \max_{z_{1l}} f_{1l}^{(3)}(z_{1l}) = \frac{2e}{(1+e)^3} > 0$, $\arg \min_{z_{2l}} f_{2l}^{(3)}(z_{2l}) = \frac{1-e}{e(1+e)^3} < 0$, and $\arg \max_{z_{2l}} f_{2l}^{(3)}(z_{2l}) = \frac{e(e-1)}{(1+e)^3} > 0$. Thus, the average error of the approximation is at most

$$\begin{aligned}|\widetilde{\mathcal{F}}_L(\widehat{\theta}) - \widetilde{\mathcal{F}}_L(\widetilde{\theta})| &\leq M \times \left[\left(\frac{2e}{(1+e)^3} - \frac{-2e}{(1+e)^3} \right) \right. \\ &\quad \left. + \left(\frac{e(e-1)}{(1+e)^3} - \frac{1-e}{e(1+e)^3} \right) \right] = M \times \frac{e^2 + 2e - 1}{e(1+e)^2}\end{aligned}$$

Therefore, Eq. 37 holds.