
Label Leakage and Protection from Forward Embedding in Vertical Federated Learning

Jiankai Sun*

Xin Yang

Yuanshun Yao

Chong Wang

Abstract

Vertical federated learning (vFL) has gained much attention and been deployed to solve machine learning problems with data privacy concerns in recent years. However, some recent work demonstrated that vFL is vulnerable to privacy leakage even though only the forward intermediate embedding (rather than raw features) and backpropagated gradients (rather than raw labels) are communicated between the involved participants. As the raw labels often contain highly sensitive information, some recent work has been proposed to prevent the label leakage from the backpropagated gradients effectively in vFL. However, these work only identified and defended the threat of label leakage from the backpropagated gradients. None of these work has paid attention to the problem of label leakage from the intermediate embedding. In this paper, we propose a practical label inference method which can steal private labels effectively from the shared intermediate embedding even though some existing protection methods such as label differential privacy and gradients perturbation are applied. The effectiveness of the label attack is inseparable from the correlation between the intermediate embedding and corresponding private labels. To mitigate the issue of label leakage from the forward embedding, we add an additional optimization goal at the label party to limit the label stealing ability of the adversary by minimizing the distance correlation between the intermediate embedding and corresponding private labels. We conducted massive experiments to demonstrate the effectiveness of our proposed protection methods.

1 Introduction

With increasing concerns over data privacy in machine learning, *Split learning* [9, 20] and *vertical federated learning* (vFL) [22] have gained attention and been deployed to solve practical problems such as online advertisements conversion prediction tasks [11, 16]. Both techniques can jointly learn a machine learning model by splitting the execution of the corresponding deep neural network on a layer-wise basis without raw data sharing.

The detailed training process of vFL (including forward pass and backward gradients computation) can be seen in Figure 1. During the forward pass, the party without labels (*non-label party*) sends the intermediate layer (*cut layer*) outputs rather than the raw data to the party with labels (*label party*), and the label party completes the rest of the forward computation to obtain the training loss. To compute the gradients with respect to model parameters in the backward gradients computation phase, the label party initiates backpropagation from its training loss and computes its own parameters' gradients. The label party also computes the gradients with respect to the cut layer outputs and sends this information back to the non-label party, so that the non-label party can use the chain rule to compute gradients of its parameters.

*Bytedance Inc. Corresponds to: {jiankai.sun, chong.wang}@bytedance.com

Even though only the intermediate computations of the cut layer embedding (rather than raw features) and backpropagated gradients (rather than raw labels) are communicated between the two parties, existing work [11, 16] demonstrated vFL is vulnerable to privacy leakage. For example, [11] demonstrated that an adversarial non-label party can leverage both norm and direction of the backpropagated gradients to infer the private class labels accurately (AUC is almost 1.0). As the raw labels often contain highly sensitive information (*e.g.*, what a user has purchased (in online advertising) or whether a user has a disease or not (in disease prediction) [20], understanding of the threat of label leakage and its protection in vFL is particularly important.

Some recent work [11, 8] have been proposed to prevent the label leakage. For example, [8] leveraged randomized responses to flip labels and use the generated noisy version labels to compute the loss functions. They proved that their proposed algorithms can achieve label differential privacy (DP) [6]. [11] proposed to add optimized Gaussian noise to perturbate the backpropagated gradients so that positive and negative gradients cannot be distinguished after perturbation. The optimized amount of added noise is calculated by minimizing the sum KL divergence between the perturbed gradients. Both methods are effective on preventing the label leakage from the backpropagated gradients in the setting of vFL.

However, besides backpropagated gradients, private labels can also be inferred by adversarial parties from the forward cut layer embedding sent from the non-label party to the label party. With the model training, the cut layer embedding can be learned to have a correlation with private labels and hence be used to distinguish different labels by adversarial parties. For example, adversarial non-label party parties can do clustering on these cut layer embeddings and then assign labels to clusters based on their sizes. For example, in the imbalanced binary classification settings such as online advertising and disease prediction, the larger cluster can be assigned negative labels and the smaller one will get positive labels. Particularly, in this paper, we develop a practical label inference method based on the technique of spectral attack [18] to predict the private labels. Our inference method can steal the private labels from the cut layer embedding effectively even protection methods proposed by [11] and [8] are applied (since both methods are proposed to prevent the label leakage from backpropagated gradients). The effectiveness of our label attack is inseparable from the correlation between the intermediate embedding and corresponding private labels. To mitigate the issue of label leakage from forward embedding, we propose to reduce the learned correlation between the cut layer embedding and private labels. We achieve this goal by letting the label party optimize an additional target which minimizes the distance correlation [17] between the cut layer embedding and corresponding private labels. So that the adversarial party cannot distinguish positive and negative labels based on the limited correlation between the cut layer embedding and private labels.

We summarize our contributions as: 1) We propose a practical label leakage attack from the forward embedding in two-party split learning; 2) We propose a corresponding defense that minimizes the distance correlation between cut layer embedding and private labels. 3) We experimentally verify the effectiveness of our attack and defense. 4) Our work is the first we are aware of to identify and defend against the threat of label leakage from the forward cut layer embedding in vFL.

2 Related Work

Federated Learning. *Federated learning* (FL) [13] can be mainly classified into three categories: *horizontal FL*, *vertical FL*, and *federated transfer learning* [22]. We focus on vertical FL (vFL) which partitions data by features (including labels). When the jointly trained model is a neural network, the corresponding setting is the same as split learning such as SplitNN [20]. In the following, we use the term vFL and split learning interchangeably.

Information Leakage in Split Learning. Recently, researchers have shown that in split learning, even though the raw data (feature and label) is not shared, sensitive information can still be leaked from the gradients and intermediate embeddings communicated between parties. For example, [19] and [16] showed that non-label party’s raw features can be leaked from the forward cut layer embedding. We differ from them by studying the leakage of labels rather than raw features. In addition, [11] studied the label leakage problem but the leakage source was the backward gradients rather than forward embeddings. To the best of our knowledge, there is no prior work studying the label leakage problem from the forward embeddings.

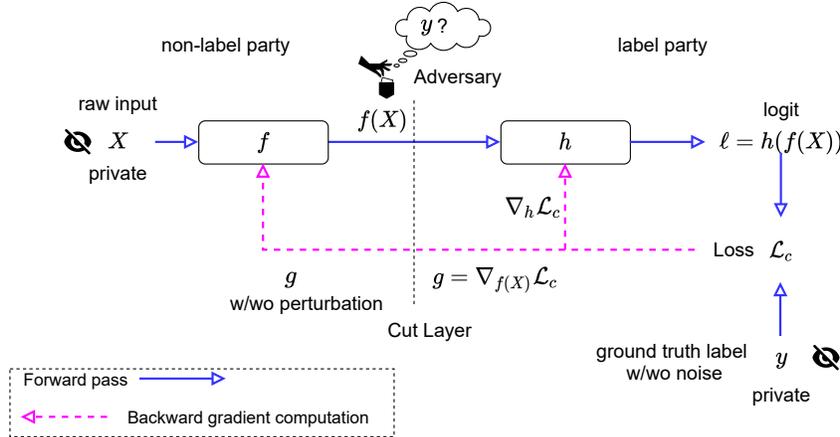


Figure 1: Attack scenario of label leakage in two-party split learning. During training, the non-label party sends the cut layer embedding $f(X)$ to the label party, and the label party completes the rest of the forward computation to obtain the training loss \mathcal{L}_c . To compute the gradients w.r.t model parameters, the label party starts backpropagation by computing the gradients $\nabla_h \mathcal{L}_c$. To allow the non-label party to also compute gradients of its parameters, the label party also computes the gradients w.r.t the cut layer embedding $g = \nabla_{f(X)} \mathcal{L}_c$ and sends it to the non-label party. In this scenario, the goal of the attacker is to infer the label from the forward cut layer embedding $f(X)$.

Information Protection in vFL. There are three main categories of information protection techniques in vFL: **1)** cryptographic methods such as secure multi-party computation [2]; **2)** system-based methods including trusted execution environments [15]; and **3)** perturbation methods that add noise to the communicated messages [1, 14, 7, 5, 23]. Our protection belongs to the third category. In this line of work, [8] generated noisy labels by randomized responses to achieve label differential privacy guarantees. [11] added optimized Gaussian noise to perturb the backpropagated gradients to confuse positive and negative gradients.

3 Preliminaries

We introduce the background of vFL framework. vFL [22] splits a deep neural network model by layers during the training and inference time.

The training phase of vFL includes two steps as shown in Figure 1: forward pass and backward gradient computation.

Forward Pass. We consider two parties that jointly train a composite model $h \circ f$ on a binary classification problem over the domain $\mathcal{X} \times \{0, 1\}$. The non-label party owns the representation-generating function $f : \mathcal{X} \rightarrow \mathbb{R}^d$ and each example’s raw feature $X \in \mathcal{X}$, while the label party owns the logit-predicting function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ and each example’s label $y \in \{0, 1\}^2$. The raw input X and label y are considered as private information by the non-label party and label party respectively and should be kept to themselves only. Let $\ell = h(f(X))$ be the logit of the positive class whose predicted probability is given through the sigmoid function: $\tilde{p}_1 = 1/(1 + \exp(-\ell))$. We measure the loss of such prediction through the binary cross entropy loss \mathcal{L}_c . When computing \mathcal{L}_c , since the ground-truth label y is sensitive information, the label party can provide noisy labels instead of real labels to prevent label leakage. For example, [8] relied on a randomized response algorithm to flip labels and provided label DP guarantee for these generated flipped labels. In this phrase, $f(X)$ is the forward embeddings and the label leakage source that we consider.

Backward Gradient Computation. To perform gradient descent, the label party first computes the gradient of the loss \mathcal{L}_c w.r.t the logit $\nabla_\ell \mathcal{L}_c$. Using the chain rule, the label party then computes the

²For simplicity, we consider the scenario where label party owns only labels without any additional features. Our attack and defense can be applied to other vFL scenarios without major modifications.

gradient of \mathcal{L}_c w.r.t its function h 's parameters and update the gradients. To allow the non-label party to learn f , the label party needs to additionally compute the gradients w.r.t cut layer embeddings $f(X)$ and communicate it to the non-label party. We denote this gradient by $g := \nabla_{f(X)} \mathcal{L}_c$. Prior work has demonstrated that the vanilla g can leak the label information. For example, [11] showed that adversarial parties can leverage either direction or norm of g to recover labels. They proposed a random perturbation technique called Marvell which we will consider as a baseline to compare with. After receiving g (w/wo perturbation) from the label party, the non-label party continues the backpropagation on f 's parameters using chain rule and updates the corresponding gradients.

During model inference, the non-label party computes $f(X)$ and sends it to the label party who then executes the rest of forward computation to compute the prediction probability.

4 Threat Model

Attacker's Goal. During training, the attacker, which is the non-label party, attempts to infer the private labels that belong to the label party. Specifically, attacker wants to reconstruct the label party's hidden label y based on the intermediate embedding $f(X)$ for each training example.

Attacker's Capability. We consider an *honest-but-curious* non-label party which follows the agreed training procedure but it actively tries to infer the label y from $f(X)$ while not violating the procedure. The attacker cannot tamper with training by selecting which examples to include in a batch or sending incorrect cut layer embedding $f(X)$. The attacker can choose its own desired model architecture for the prediction function $f(\cdot)$ as long as the cut layer embedding dimension matches with the input of $h(\cdot)$. For example, the attacker has the freedom of selecting the depth and width of its neural network to increase the power of $f(X)$ on label guessing. Inferring label y can be viewed as a binary classification problem where the (input, output) distribution is the induced distribution of $(f(X), y)$. The attacker can use any binary classifier $q : \mathbb{R}^d \rightarrow \{0, 1\}$ to infer the labels.

Attacker's Auxiliary Information. We allow the non-label party to have population-level side information (i.e. ratio of the positive instances) specifically regarding the properties of (and the distinction between) the positive and negative class's cut-layer embedding distributions. The attacker may leverage such population-level side information to infer the private labels more accurately. However, we assume the non-label party has *no example-level side information* that is different example by example.

5 Label Leakage Attack

We leverage spectral attack [18] to predict labels for corresponding cut layer embeddings. Spectral attack is a singular value decomposition (SVD) based outlier detection method. It shows that:

Lemma 5.1 (Lemma in [18]) *Let D, W be two distributions over \mathbb{R}^d with mean μ_D, μ_W and covariance matrices Σ_D, Σ_W . Let F be a mixture distribution given by $F = (1 - \epsilon)D + \epsilon W$ where $0 < \epsilon < \frac{1}{2}$. If $\|\mu_D - \mu_W\|_2^2 \geq \frac{6\sigma^2}{\epsilon}$, then the following statement holds:*

$$\Pr_{X \sim D} [|\langle X - \mu_F, v \rangle| > t] < \epsilon, \quad \Pr_{X \sim W} [|\langle X - \mu_F, v \rangle| < t] < \epsilon. \quad (1)$$

where μ_F is the mean of F and v is the top singular vector of the covariance matrix of F , and $t > 0$.

Lemma 5.1 motivates us to use the spectral attack to differentiate the embedding distribution between the positive samples (\mathcal{F}^+) and the negative samples (\mathcal{F}^-). If the mean of \mathcal{F}^+ and \mathcal{F}^- are far away from each other, then we can use $|\langle f(X) - \mu_F, v \rangle|$ as the indicator to distinguish \mathcal{F}^+ and \mathcal{F}^- . Our attack works as following:

1. For each mini-batch of $f(X)$, compute its empirical mean μ_F and the top singular vector v of the corresponding covariance matrix. Then compute score $s = |\langle f(X) - \mu_F, v \rangle|$ for each sample in the mini-batch.
2. Divide all samples into two clusters using score s as the distance metric. We leverage some heuristic with population-level knowledge to determine which cluster belongs to which label.

For example, if the dataset (*i.e.* Criteo) is imbalanced and dominated by negative samples, then we assign negative labels to the larger cluster and positive labels to the smaller one. Note that when the dataset is balanced (*i.e.* positive ratio is 0.5), it is true that the attacker cannot directly know which label to assign to which cluster. However, the attacker can employ a simple heuristic: assigning positive labels to the cluster with larger scores.

The detailed description of our attack method can be seen in Algorithm 1 in Appendix D. We also talk about how to extend our attack method to multi-class settings in Appendix E.

6 Attack Evaluation

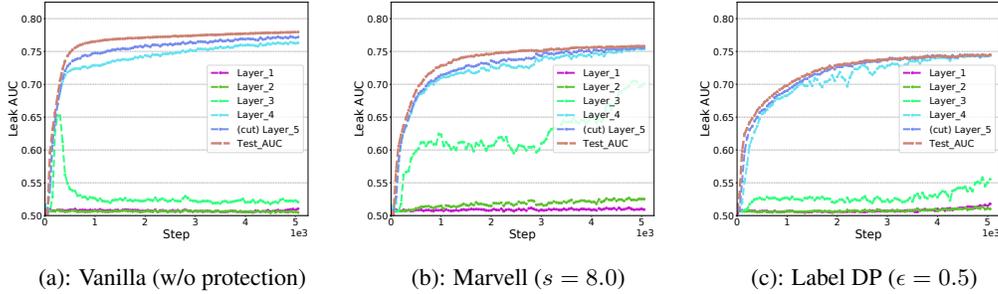


Figure 2: Effectiveness of Spectral Attack. (a): Spectral attack AUC of different layers without any protection; (b): Spectral attack AUC of different layers with Marvell [11] ($s = 8.0$) protection; (c): Spectral attack AUC of different layers with Label DP [8] guarantee ($\epsilon = 0.5$).

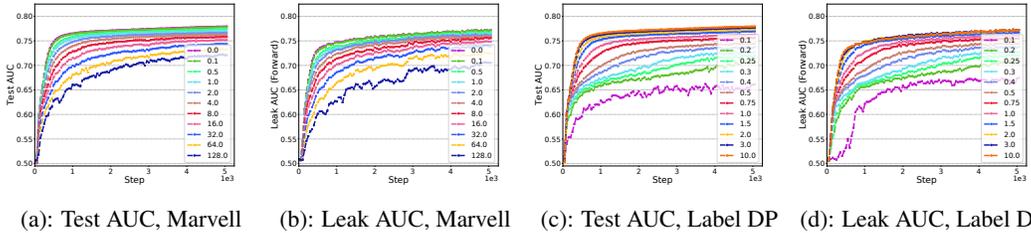


Figure 3: Spectral Attack (leak) AUC is bounded by the model utility (test AUC). Figure (a) and (b) shows the test AUC and leak AUC on the cut layer protected by Marvell with different s respectively; Figure (c) and (d) represents the test AUC and leak AUC on the cut layer protected by Label DP with different ϵ respectively.

6.1 Experimental Setup

Dataset. We mainly evaluate the proposed framework on Criteo³, which’s a large-scale industrial binary classification dataset (with with approximately 45 million user click records) for conversion prediction tasks. Every record of Criteo has 27 categorical input features and 14 real-valued input features. We first replace all the NA values in categorical features with a single new category (which we represent using the empty string) and replace all the NA values in real-valued features with 0. For each categorical feature, we convert each of its possible value uniquely to an integer between 0 (inclusive) and the total number of unique categories (exclusive). For each real-valued feature, we linearly normalize it into $[0, 1]$. We then randomly sample 90% of the entire Criteo set as our training data and the remaining 10% as our test data.

To save space, we put the similar experimental results on another dataset Avazu in the Appendix F.

³<https://www.kaggle.com/c/criteo-display-ad-challenge/data>

Model. We modified a popular deep learning model architecture WDL [4] for online advertising. Here the non-label party first processes the categorical features in a given record by applying an embedding lookup for every categorical feature’s value. We use an embedding dimension of 4 for this part. After the lookup, the deep embeddings are then concatenated with the continuous features to form the raw input vectors. The non-label party leverages its feature extractor $f(\cdot)$ with 5 layers of ReLU activated MLP to process its raw input features. The cut layer is after the output of $f(\cdot)$. The label party leverages its label predictor $h(\cdot)$ with 3 layers of ReLU activated MLP to generate the final logit value for prediction. During training, in each batch the non-label party sends an embedding matrix with size 8, 192×128 to the label party where batch size is 8, 192^4 and embedding size is 128.

Against Existing Defenses. In addition to the scenario when the label party imposes no protection, we also include two situations when label party employs two existing defenses against label leakage attack in FL: Marvell [11] and Label DP [8].

Evaluation Metric. To measure the effectiveness of the label leakage, we compute AUC on our attacker’s label predictions and the ground-truth labels, namely *leak AUC*. A closer to 1.0 leak AUC is considered to be more effective, while a closer to 0.5 leak AUC is more like a random guess.

6.2 Results

Effectiveness of Spectral Attack. In this section, we show the effectiveness of leveraging Spectral attack to infer labels from the forward embedding with different settings. As shown in Figure 2 (a), the adversarial party can infer the label information from the last two layers (including the cut layer) with a high ACU (about 0.78) under the setting of vanilla vFL (without using any label protection techniques such as Marvell [11] and label DP [8]). However, as shown in Figure 2 (b) and (c), even with the protection of Marvell [11] and Label DP [8], Spectral attack can still steal the label information with a relatively high leak AUC. Marvell [11] can prevent label leakage from the backpropagated gradients by adding Gaussian noise to the gradients before sending them back to the non-label party. However, as shown Figure 2 (b), Marvell with $s = 8.0$ cannot prevent label leakage from the forward embedding effectively since the last two layers can leak the label with a leak AUC closing to 0.75. Label DP leverages randomized response to generate noisy labels for computing the loss function and the corresponding backpropagated gradients. However, as shown in Figure 2 (c), the leak AUC of the last two layers are close to 0.75, which indicates that Label DP fails to prevent label leakage from the forward embedding.

Upper Bound of Leak AUC. As defined in our threat model in Section 4, the adversarial party has the freedom of selecting its network architecture (*i.e.* depth and width of layers) to increase the power of $f(\cdot)$ on label guessing. One question arising is that does the power of $f(\cdot)$ on label guessing (measured by leak AUC) have an upper bound? We may observe from Figure 2 that the leak AUC of the cut layer is bounded by the test AUC of the model. We assume that the label party does not introduce negative effects on $h(\cdot)$ and hence $f(\cdot)$ at most has the same prediction power as $h(f(\cdot))$, even though the adversarial party has the freedom of increasing the depth and width of the neural network at its side to increase the power of $f(\cdot)$.

To illustrate this observation, we test the spectral attack AUC with different protection methods (Marvell [11] and Label DP [8]), since different parameters of Marvell (s) and Label DP (ϵ) can give different model utilities. As shown in Figure 6 in the appendix A and Figure 3, different parameters of Marvell (s) and Label DP (ϵ) can give different trade-offs between privacy (attack AUC on backpropagated gradients as shown in Figure 6) and model utilities (test AUC as shown in Figure 3). As shown in Figure 3 (a) and (c), the test AUC changes with different parameters of Marvell (s) and Label DP (ϵ). The spectral attack AUC on the forward embedding as shown in Figure 3 (b) and (d) has a similar pattern as the test AUC and each leak AUC is bounded by the corresponding test AUC. It indicates that the power of stealing label information from the forward embedding is bounded by the whole model’s utility.

Attack on Balanced Datasets. We also demonstrated the effectiveness of our attack method with balanced datasets. To save space, we put the corresponding results in Section G of the Appendix.

⁴The batch size is 8, 192 if not specified.

7 Defense against Label Leakage Attack

In this section, we first talk about the protection objectives and then show the details of our practical protection method on defending the label leakage from the forward embedding.

7.1 Label Protection Objectives

The defender, which is the label party who is the victim of the attacker, aims to protect the label information during the procedure while maintaining the model’s utility. It has to achieve reasonable trade-offs between two objectives: privacy and utility objective.

Privacy Objective: the leak AUC of the recovered labels that the adversarial party estimate from the cut layer embedding should be close to 0.5 (like a random guess).

$$\min_{Y \sim \mathcal{D}_{label}} \mathbb{E} [|AUC(Y, Y') - 0.5|] \quad (2)$$

where Y' is the estimated label by an attacker and Y is the corresponding ground-truth.

Utility Objective: The composition model $h \circ f$ jointly trained with split learning needs to have high classification performance (evaluated by AUC) on an unseen test set.

$$\max_{(X, Y) \sim \mathcal{D}_{test}} \mathbb{E} [AUC(Y, \mathbb{1}(\tilde{p}_1 \geq \lambda))] \quad (3)$$

where $\tilde{p}_1 = 1/(1 + \exp(-h(f(X))))$ is the predicted probability of being positive for X and $\mathbb{1}$ is an indicator function with a tuning parameter λ .

7.2 Defense Details

With the model training, the function f gains some high-level representation of raw input X and learns the correlation between $f(X)$ and Y . $f(X)$ can be learned as an indicator of private labels. Hence spectral attack can leverage $f(X)$ to steal private labels. One strategy to prevent label leakage is that the label party asks the non-label party to send the raw input X directly. However, this strategy violates the protocol of vFL and non-label party has risks of revealing their raw input to the other party. Another strategy is that the label party forces non-label party to limit its power of f so that $f(X)$ cannot be a good indicator of private labels. However, this strategy is not practical since it violates the setting of the adversary’s capability defined in our threat model (see Section 4) and may also break the utility objective as discussed in Section 7.1. In this paper, to prevent the label leakage from the forward embedding, we would like to let the forward embeddings not be a good proxy of the private labels while meeting the requirements of our threat model (Section 4) and label protection objectives (Section 7.1). Particularly, we leverage distance correlation to achieve this goal.

Distance correlation can be used to measure the statistical dependence between two paired random vectors with arbitrary dimensions (not necessarily equal) [17]. It can measure both linear and nonlinear associations between two vectors [21]. Another advantage is it can be optimized and easily fit into the conventional optimization framework (e.g. minimizing the cross entropy loss).

Given $f(X)$ and Y from a mini-batch, we can use $dCor(f(X), Y)$ where $f(X) \in \mathbb{R}^{n \times d}$ and $Y \in \mathbb{R}^{n \times 1}$ to measure the dependence between $f(X)$ and Y , even though $f(X)$ and Y share different dimensions. The detailed calculation of $dCor(f(X), Y)$ can be shown as follows:

1. Compute the $n \times n$ distance matrices \mathbf{s} and \mathbf{t} containing all pairwise distances for each sample pair. Each element $\mathbf{s}_{(j,k)}$ in \mathbf{s} is computed as: $\mathbf{s}_{(j,k)} = \|f(X)_j - f(X)_k\|$ and each element $\mathbf{t}_{(j,k)}$ in \mathbf{t} is $\mathbf{t}_{(j,k)} = \|Y_j - Y_k\|$ where $j, k = 1, 2, \dots, n$, $f(X)_j$ represents the the j -th row of $f(X)$, Y_j denotes the label (0 or 1) of the j -th sample, and $\|\cdot\|$ denotes Euclidean norm.
2. Normalize \mathbf{s} and \mathbf{t} by taking all doubly centered distances to get \mathbf{A} and \mathbf{B} respectively. Particularly, $\mathbf{A}_{(j,k)} = \mathbf{s}_{(j,k)} - \bar{\mathbf{s}}_{j,\cdot} - \bar{\mathbf{s}}_{\cdot,k} + \bar{\mathbf{s}}_{\cdot,\cdot}$ and $\mathbf{B}_{(j,k)} = \mathbf{t}_{(j,k)} - \bar{\mathbf{t}}_{j,\cdot} - \bar{\mathbf{t}}_{\cdot,k} + \bar{\mathbf{t}}_{\cdot,\cdot}$, where $\bar{\mathbf{s}}_{j,\cdot}$ and $\bar{\mathbf{t}}_{j,\cdot}$ are the j -th row mean of \mathbf{s} and \mathbf{t} respectively, $\bar{\mathbf{s}}_{\cdot,k}$ and $\bar{\mathbf{t}}_{\cdot,k}$ are the k -th column mean of \mathbf{s} and \mathbf{t} respectively, and $\bar{\mathbf{s}}_{\cdot,\cdot}$ and $\bar{\mathbf{t}}_{\cdot,\cdot}$ are grand mean \mathbf{s} and \mathbf{t} respectively.

3. Calculate the distance covariance $dCov^2(f(X), Y)$ as the mean of the dot product of **A** and **B**. Particularly, $dCov^2(f(X), Y) = \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n \mathbf{A}_{(j,k)} \mathbf{B}_{(j,k)}$. Similarly, compute $dCov^2(f(X), f(X))$ and $dCov^2(Y, Y)$.
4. Compute the distance correlation $dCor(f(X), Y)$ as $\frac{dCov^2(f(X), Y)}{\sqrt{dCov^2(f(X), f(X)) \times dCov^2(Y, Y)}}$

The distance correlation $dCor(f(X), Y)$ is zero if and only if the $f(X)$ and Y are independent, and then the adversary has no way to infer Y from $f(X)$ under this scenario. The intuition of our protection method is to make label Y and embedding $f(X)$ less dependent, and therefore reduces the likelihood of gaining information of Y from $f(X)$. Particularly, we add an additional optimization target at the side of label party, and we name it as the distance correlation loss \mathcal{L}_d which is calculated as follows:

$$\mathcal{L}_d = \mathbb{E}_{(X, Y) \sim \mathcal{D}} (\log(dCor(f(X), Y))) \quad (4)$$

By minimizing the log of distance correlation ⁵ (loss \mathcal{L}_d) during the model training, we achieve the goal of not letting $f(X)$ be a good proxy of Y .

Note that dCor computes pairwise distance between samples and requires $O(n^2)$ time complexity where n is the batch size. In practice, there are some faster estimators of dCor [3, 10]. In addition, dCor is sensitive to the n and a larger n can give a more accurate estimation of the distance correlation.

Combine with Training. The overall loss function computed at the label party side is a combination of two losses: distance correlation loss (\mathcal{L}_d) and normal label prediction loss (\mathcal{L}_c). In this paper, we focus on binary classification and use categorical cross-entropy as label prediction loss \mathcal{L}_c . Optimizing \mathcal{L}_c makes sure the model maintain a good utility while optimizing \mathcal{L}_d increases model privacy. Putting them in one framework can help us defend against the spectral attack while maintaining the accuracy of the primary learning task. The overall loss function is:

$$\mathcal{L} = \mathcal{L}_c + \alpha_d \mathcal{L}_d \quad (5)$$

where $\alpha_d \geq 0$ is the control parameter for distance correlation. It can balance the trade-off between utility and privacy. If we set α_d a large number, the cut layer embedding and the corresponding label will be more independent. Then the label party has great difficulty learning the model well and hence we cannot achieve our utility objective. However, if set $\alpha_d = 0$, we fail on achieving the privacy objective but can get the best model utility. A good value of α_d can help us balance utility and privacy objectives. It's worth mentioning that only the label party deals with the optimization of \mathcal{L} (unseen for the non-label party) while there is no change on the non-label party's side.

Remarks. We choose the two-party setting with binary classification because it is currently the typical scenario used in the vFL literature (i.e. [12, 11]) due to its popularity in the industry (online advertising, disease prediction etc). Our protection method can be easily extended to other scenarios as we discussed in Section C of Appendix.

8 Defense Evaluation

Effectiveness. As shown in Figure 4 (a) and (b), with reasonable α_d , the spectral attack AUC can be reduced to around 0.5 (like a random guess), while the model utility can be as high as 0.78 (almost the same as the performance of the vanilla model). Our protection method can be compatible with other protection methods such as Label DP [8] and Marvell [11]. We use Marvell as an example to demonstrate this. As shown in Figure 4 (c) and (d), we can see that our method with Marvell can prevent the label leakage from the backpropagated gradients and forward embedding simultaneously.

Sensitivity of α_d . We vary different α_d to test the sensitivity of α_d and observe the trade-off between model utility and privacy. We show model utility (test AUC) and privacy (leak AUC) with different α_d in Figure 5. We can conclude that a reasonable α_d (i.e. $0.005 \geq \alpha_d \geq 0.002$) can achieve an acceptable attack AUC (0.5) while the model test AUC does not drop too much. Take $\alpha_d = 0.002$

⁵Empirically we find that adding $\log(\cdot)$ stabilizes training.

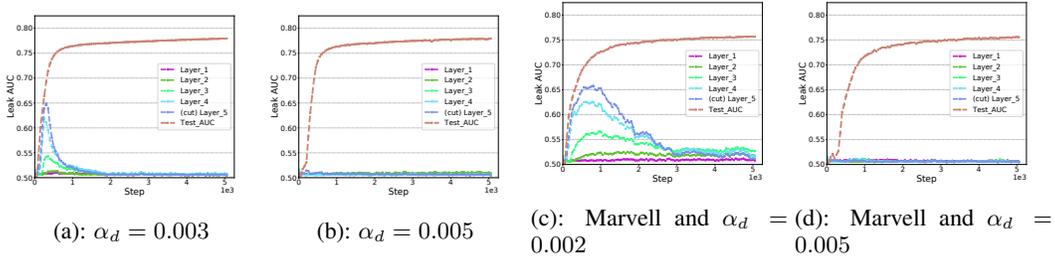


Figure 4: Effectiveness of our protection method and its compatibility with gradient perturbation based protection method Marvell. Figure (a) and (b) shows the effectiveness of our protection method defending Spectral attack on different layers with $\alpha_d = 0.003$ and $\alpha_d = 0.005$ respectively. Figure (c) and (d) demonstrates that our protection method can be compatible with gradient perturbation based protection method Marvell $\alpha_d = 0.002$ and $\alpha_d = 0.005$ respectively.

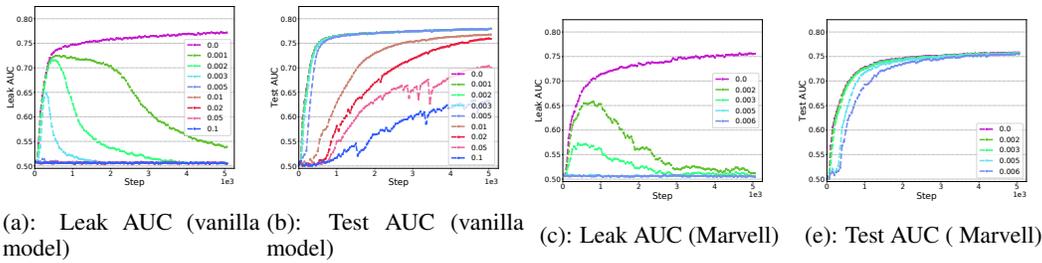


Figure 5: Sensitivity of α_d and tradeoff between model utility and privacy. We use different α_d to test the model performance and the corresponding Spectral attack AUC on the cut layer. Figure (a) and (b) shows the Spectral attack AUC and test AUC respectively. Figure (c) and (d) shows the Spectral attack AUC and test AUC with Marvell ($s = 8$) respectively.

for example, its attack AUC is around 0.5 (as shown in Figure 5 (a)) and test AUC has almost no drop in comparing with the baseline ($\alpha_d = 0$) (as shown in Figure 5 (b)). It also happens when we use Marvell to defend the label leakage from the backpropagated gradients (as shown in Figure 5 (c) and (d)). It also demonstrates that our proposed protection method is compatible with these backpropagated gradients perturbation based label protection methods.

Cost of Privacy. Since the label party has minimized the distance correlation between the cut layer embedding and private labels, it has to increase its own model’s power (more computational cost) to re-learn the correlation. Otherwise, it cannot maintain the whole model’s utility. We provide the quantitative analysis on the cost of privacy in in the Appendix B.

9 Conclusion

In this paper, we focus on identifying and defending the threat of label leakage from the forward cut layer embedding rather than backpropagated gradients. We firstly proposed a practical label attack method that can steal private labels effectively from the cut layer embedding even though some existing protection methods such as label differential privacy [8] and gradients perturbation [11] are applied. The effectiveness of the above label attack is inseparable from the correlation between the cut layer embedding and corresponding private labels. Hence to mitigate the issue of label leakage from the forward embedding, we proposed to minimize the distance correlation between the cut layer embedding and corresponding private labels as an additional goal at the label party, which can limit the label stealing ability of the adversary. We conducted massive experiments to demonstrate the effectiveness of our proposed protection methods. Our work is the first we are aware of to identify and defend against the threat of label leakage from the forward cut layer embedding in vFL. Since the label information is highly private and sensitive, we hope our work can open up a number of worthy directions for future study in the interest of protecting label privacy.

References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 308–318, 2016.
- [2] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-preserving machine learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1175–1191, 2017.
- [3] Arin Chaudhuri and Wenhao Hu. A fast algorithm for computing distance correlation. *Computational Statistics & Data Analysis*, 135:15 – 24, 2019.
- [4] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*, pages 7–10, 2016.
- [5] Albert Cheu, Adam Smith, Jonathan Ullman, David Zeber, and Maxim Zhilyaev. Distributed differential privacy via shuffling. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 375–403. Springer, 2019.
- [6] Cynthia Dwork. Differential privacy. In Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener, editors, *Automata, Languages and Programming*, pages 1–12, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [7] Úlfar Erlingsson, Vitaly Feldman, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Abhradeep Thakurta. Amplification by shuffling: From local to central differential privacy via anonymity. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2468–2479. SIAM, 2019.
- [8] Badih Ghazi, Noah Golowich, Ravi Kumar, Pasin Manurangsi, and Chiyuan Zhang. Deep learning with label differential privacy, 2021.
- [9] Otkrist Gupta and Ramesh Raskar. Distributed learning of deep neural network over multiple agents. *Journal of Network and Computer Applications*, 116:1–8, 2018.
- [10] Cheng Huang and Xiaoming Huo. A statistically and numerically efficient independence test based on random projections and distance covariance. 2017.
- [11] Oscar Li, Jiankai Sun, Xin Yang, Weihao Gao, Hongyi Zhang, Junyuan Xie, Virginia Smith, and Chong Wang. Label leakage and protection in two-party split learning. In *The Tenth International Conference on Learning Representations (ICLR)*, 2022.
- [12] Yang Liu, Xiong Zhang, and Libin Wang. Asymmetrical vertical federated learning. *arXiv preprint arXiv:2004.07427*, abs/2004.07427, 2020.
- [13] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [14] H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. In *International Conference on Learning Representations*, 2018.
- [15] Pramod Subramanyan, Rohit Sinha, Ilia Lebedev, Srinivas Devadas, and Sanjit A Seshia. A formal foundation for secure remote execution of enclaves. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 2435–2450, 2017.
- [16] Jiankai Sun, Yuanshun Yao, Weihao Gao, Junyuan Xie, and Chong Wang. Defending against reconstruction attack in vertical federated learning. *CoRR*, abs/2107.09898, 2021.
- [17] Gábor J. Székely, Maria L. Rizzo, and Nail K. Bakirov. Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769 – 2794, 2007.
- [18] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, 2018.

- [19] Praneeth Vepakomma, Otkrist Gupta, Abhimanyu Dubey, and Ramesh Raskar. Reducing leakage in distributed deep learning for sensitive health data. *arXiv preprint arXiv:1812.00564*, 2019.
- [20] Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. Split learning for health: Distributed deep learning without sharing raw patient data. *arXiv preprint arXiv:1812.00564*, 2018.
- [21] Praneeth Vepakomma, Tristan Swedish, Ramesh Raskar, Otkrist Gupta, and Abhimanyu Dubey. No peek: A survey of private distributed deep learning. *CoRR*, abs/1812.03288, 2018.
- [22] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19, 2019.
- [23] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *Advances in Neural Information Processing Systems*, pages 14774–14784, 2019.

Appendix Outline:

Section A shows the attack AUC of Marvell and Label DP with different parameters.

Section B: Quantitative analysis on cost of privacy: shows that the label party has to increase its own sides' network power to achieve a good trade-off between model utility and privacy.

Section C: Extending our proposed protection methods to other settings.

Section D: Spectral Attack Algorithm

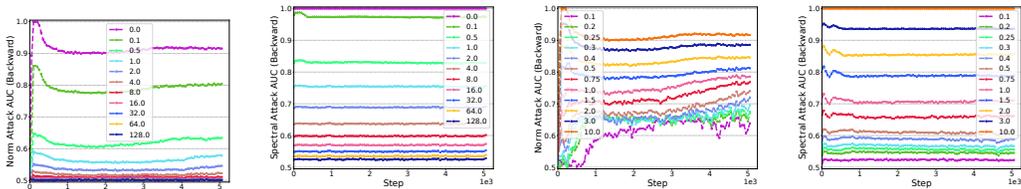
Section E: Extending our attack method to multi-class scenarios

Section F: Experimental Results on Avazu

Section G: Spectral Attack and Corresponding Defense on Balanced Datasets

Section H: Data Setup and Experimental Details

A Attack AUC of Marvell and Label DP with different parameters



(a): Norm Attack, Marvell (b): Spectral Attack, Marvell (c): Norm Attack, Label DP (d): Spectral Attack, Label DP

Figure 6: Effectiveness of Marvell (Figure (a) and (b)) and Label DP (Figure (c) and (d)) on defending Norm attack (Figure (a) and (c)) and Spectral attack (Figure (b) and (d)) on backpropagated (bp.) gradients. Figure (a) and (b) shows the Norm and Spectral attack AUC on the bp. gradients protected by Marvell with different s respectively. Figure (c) and (d) represents the Norm and Spectral attack AUC on the bp. gradients protected by Label DP with different ϵ respectively.

B Quantitative Analysis on Cost of Privacy

As shown in Figure 7 in the Appendix B, we report the performance of using each layer embedding to predict the label. The corresponding AUC is calculated by performing the spectral attack on each layer’s embedding. The whole network consists of 8 layers. The non-label party owns the first 5 layers and the label party owns the last 3 layers. Without optimizing distance correlation between cut layer embedding and labels (setting $\alpha_d = 0$) as shown in 7 (a), (e) and (g), we can achieve comparable AUC with the final test AUC at layer 4. It indicates that the label party may need only one layer to achieve a good prediction utility. However, if we set $\alpha_d > 0$ in different settings as shown in Figure 7 (b) ($\alpha_d = 0.002$), (c) ($\alpha_d = 0.003$), (d) ($\alpha_d = 0.005$), (f) (Marvell w. $\alpha_d = 0.002$), and (h) (Label DP w. $\alpha_d = 0.002$), the label party has to use all 3 layers to achieve a reasonable test AUC. In conclusion, the label party has to increase its own model’s power and needs more computational cost to re-learn the correlation and achieve a reasonable trade-off between model utility and privacy.

To provide more quantitative analysis of privacy cost, we conducted a new round of experiments to analyze the cost of privacy on Criteo dataset as shown in Table 1. We change the predictive power of the label party by varying the number of layers in the label party’s neural network. For each α_d , the label party trains a MLP model with 1 layer (small model) and 3 layers (large model) respectively and we report the corresponding test AUCs. As a result, when we increase α_d , the test AUCs of both small and large models decrease. (The same conclusion is shown in Figure 5 (a) and (b).) This is a natural trade-off between utility and privacy. We also observe that when α_d is not large (i.e. ≤ 0.005), the gap of test AUC between small and large model is not significant. It indicates that it would be easy for the label party to re-learn the correlation between cut layers embedding and label information with only 1 layer on small α_d . However when increasing α_d , the difference of test AUC of small and large models becomes larger. For example, when $\alpha_d = 0.03$, if the label party has only 1 layer, its test AUC is only 0.7082. Meanwhile if the label party increases its model’s predictive power by expanding its model to 3 layers, the corresponding test AUC can be increased to 0.7518. With $\alpha_d = 0.02$, the label party can only use 1 layer to reach the similar level of test AUC (0.7528). Hence the label party has to use more layers (and more computational cost) to increase its model’s predictive power to re-learn the correlation between cut layer embedding and label information for larger α_d .

C Extending Our Protection Methods to Other Settings

We choose the two-party setting with binary classification because it is currently the typical scenario used in the vFL literature (i.e. [12, 11]) due to its popularity in the industry (online advertising, disease prediction etc). Our protection method can be easily extended to other scenarios as following:

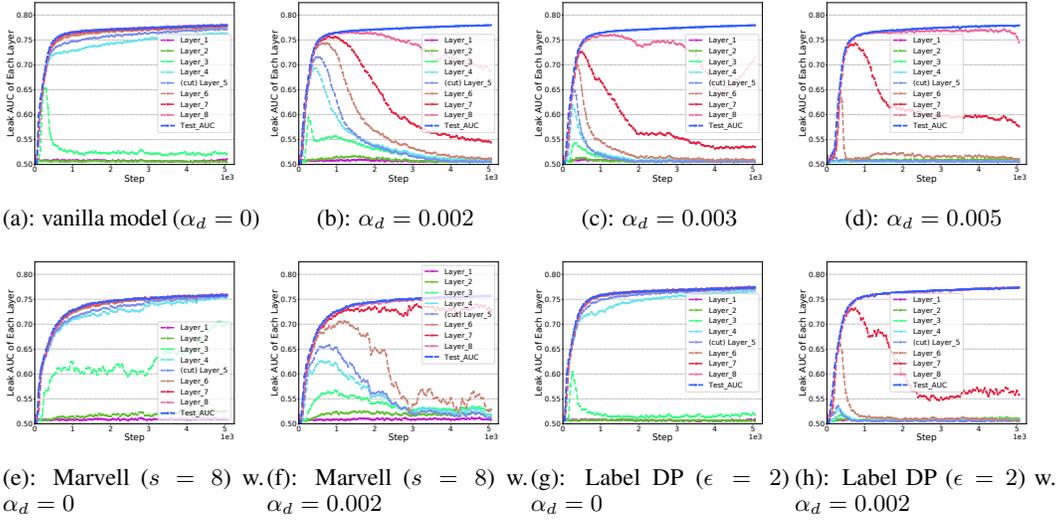


Figure 7: Cost of privacy. Since the label party has minimized the distance correlation between the cut layer embedding and label, it has to increase its own neural network’s power to relearn the correlation and maintain the whole model’s utility. We report the leak AUC of using each layer embedding to predict the label. The corresponding AUC is calculated by performing spectral attack on each layer’s embedding. The whole network consists of 8 layers. The non-label party owns the first 5 layers and the label party owns the last 3 layers. Without optimizing distance correlation between cut layer embedding and labels ($\alpha_d = 0$) as shown in Figure (a), (e) and (g), we can achieve comparable AUC with the final test AUC at the forth and its following layers. It indicates that the label party may not need 3 layers to achieve a good prediction utility. However, if we set $\alpha_d > 0$ in different settings (w/w. Label DP and w/w. Marvell), the label party has to fully used its 3 layers to achieve a reasonable test AUC. In conclusion, the label party has to increase its own model’s power to achieve a good trade-off between privacy and model utility.

α_d	Test AUC (small)	Test AUC (large)
0.003	0.7769	0.7777
0.005	0.7763	0.7773
0.02	0.7528	0.7645
0.03	0.7082	0.7518
0.04	0.5848	0.7152

Table 1: Test AUC on Criteo with varying α_d and model’s predictive power (small 1-layer or large 3-layer model).

- Multi-party vFL: Each label party can add an additional loss function \mathcal{L}_d^i to minimize the distance correlation between the private label and the forward embedding sent from the non-label party i . The rest of the training would be the same.
- Multi-class dataset: The calculation of distance correlation between the forward embedding and private labels Y does not require Y to be binary.
- Label party also owns features: the label party does not have to limit the correlation between its own features and its own labels. The training would be the same.

D Spectral Attack Algorithm

The detailed description of our attack method can be seen in Algorithm 1.

Algorithm 1: Spectral Attack Algorithm

Result: Infer Labels for n samples in a mini-batch

Input : Cut layer embedding $f(X) \in \mathcal{R}^{n \times d}$

Output : Inferred labels of $f(X)$

```

1  $\text{mean}_f = \text{mean}(f(X), \text{axis} = 0)$  // output dim:  $1 \times d$ 
2  $\text{normalized}_f = f(X) - \text{mean}_f$  // output dim:  $n \times d$ 
3 // S, U, and V are singular values, right singular vectors and left singular vectors of
    $\text{normalized}_f$  respectively
4  $S, U, V = \text{svd}(\text{normalized}_f)$ 
5  $v = V^T[0]$  // output dim:  $1 \times d$ 
6  $s = \langle \text{normalized}_f, v \rangle$  // output dim:  $n \times 1$ 
7 Cluster  $s$  into two clusters:  $C_1$  and  $C_2$ 
8 if  $|C_1| > |C_2|$  then
9 |   assign negative labels to samples in  $C_1$ 
10 |  assign positive labels to samples in  $C_2$ 
11 else
12 |   assign negative labels to samples in  $C_2$ 
13 |   assign positive labels to samples in  $C_1$ 
14 end

```

E Extending Our Attack Method to Multi-class Scenarios

Our attack method can be easily extended to a multi-class task. Suppose the dataset has k classes, we can use multi-class clustering algorithms such as k -means clustering and spectral clustering on the embeddings to generate k clusters. Then we can assign the attacker’s inferring labels to k clusters with the global knowledge on the label ratios. For example, the label of the largest class (class with the most instances) is assigned to the largest cluster, and so on.

F Experimental Results on Avazu Dataset

We also conducted experiments on Avazu dataset, which is an online advertising dataset containing approximately 40 million entries (11 days of clicks/not clicks). The positive ratio is about 15%. We randomly split the dataset into 90% training and 10% test, and run 10 *epochs* with 8,192 batch size to train the model. Other settings such as the network structure are the same as the first Criteo dataset in the paper.

We show the results in Table 2. When there is no protection ($\alpha_d = 0$), the Leak AUC is 0.7262 which is much higher than the random guess attack (0.5). When we leverage our protection method ($\alpha_d = 0.03$), the Leak AUC drops to 0.5. Meanwhile, the test AUC of $\alpha_d = 0.03$ only drops about 0.003 when comparing with the model without protection. In general, the overall observations and conclusions on Avazu are consistent with Criteo results in the paper.

	Leak AUC	Test AUC
vanilla ($\alpha_d = 0$)	0.7262	0.7532
$\alpha_d = 0.03$	0.5089	0.7502

Table 2: Results on Avazu Dataset

G Spectral Attack on Balanced Datasets

In terms of the balanced setting, we artificially down-sample the negative instances in Criteo, and generate three datasets with positive ratio 50%, 49% and 45%. Then we conduct a similar analysis on those datasets, and show the results in Table 5, 4, and Table 3.

Note that when the dataset is balanced (i.e. positive ratio is 0.5), it is true that the attacker cannot directly know which label to assign to which cluster. However, the attacker can employ a simple heuristic: assigning positive labels to the cluster with larger scores (s calculated in Step 1 in Section 5). As shown in Table 5, the attack can still be highly successful (0.7817 leak AUC without protection). In addition, our defense in the balanced setting can also successfully bring down the leak AUC to 0.5. We observe the similar conclusions when the positive ratio is 49% and 45%. Those results show that our defense can work when the dataset is balanced.

	Leak AUC	Test AUC
$\alpha_d = 0$	0.7607	0.7823
$\alpha_d = 0.03$	0.5048	0.7824

Table 3: Results on Criteo with Positive Ratio as 45%.

	Leak AUC	Test AUC
$\alpha_d = 0$	0.7814	0.7849
$\alpha_d = 0.03$	0.5069	0.7837

Table 4: Results on Criteo with Positive Ratio as 49%.

	Leak AUC	Test AUC
$\alpha_d = 0$	0.7817	0.7829
$\alpha_d = 0.03$	0.5061	0.7826

Table 5: Results on Criteo with Positive Ratio as 50%.

H Data Setup and Experimental Details

We first describe how we first preprocess each of the datasets. We then describe the model architecture used for each dataset. Finally, we describe what are the training hyperparameters used for each dataset/model combination and the total amount of compute required for the experiments.

Dataset preprocessing

[Criteo] Every record of Criteo has 27 categorical input features and 14 real-valued input features. We first replace all the NA values in categorical features with a single new category (which we represent using the empty string) and all the NA values in real-valued features by 0. For each categorical feature, we convert each of its possible value uniquely to an integer between 0 (inclusive) and the total number of unique categories (exclusive). For each real-valued feature, we linearly normalize it into $[0, 1]$. We then randomly sample 10% of the entire Criteo publicly provided training set as our entire dataset (for faster training to generate privacy-utility trade-off comparison) and further make the subsampled dataset into a 90%-10% train-test split.

[Avazu] Unlike Criteo, each record in Avazu only has categorical input features. We similarly replace all NA value with a single new category (the empty string), and for each categorical feature, we convert each of its possible value uniquely to an integer between 0 (inclusive) and the total number of unique categories (exclusive). We use all the records in provided in Avazu and randomly split it into 90% for training and 10% for test.

Model architecture details

[Criteo] We modified a popular deep learning model architecture WDL [4] for online advertising. We first process the categorical features in a given record by applying an embedding lookup for every categorical feature’s value. We use an embedding dimension of 4 for the deep part. After the lookup, the deep embeddings are then concatenated with the continuous features to form the raw input vectors for the deep part. The deep part processes the raw features using several ReLU-activated 128-unit MLP layers before producing a final logic value.

Model training details

To ensure smooth optimization and sufficient training loss minimization, we use a slightly smaller learning rate than normal.

[Criteo] We use the Adam optimizer with with a batch size of 8,192 and a learning rate of $1e-4$ throughout the entire training of 3 epochs (approximately 15k stochastic gradient updates).

[Avazu] We use the Adam optimizer with a batch size of 8,192 and a learning rate of $1e-4$ throughout the entire training of 3 epochs (approximately 15k stochastic gradient updates).

We conduct our experiments over 8 Nvidia Tesla V100 GPU card. Each epoch of run of Avazu takes about 10 hours to finish on a single GPU card occupying 4GB of GPU RAM. Each epoch run of Criteo takes about 11 hours to finish on a single GPU card using 4 GB of GPU RAM.