Membership Privacy for Machine Learning Models Through Knowledge Transfer

Virat Shejwalkar Amir Houmansadr

University of Massachusetts Amherst {vshejwalkar, amir}@cs.umass.edu

Abstract

Large capacity machine learning (ML) models are prone to membership inference attacks (MIAs), which aim to infer whether the target sample is a member of the target model's training dataset. The serious privacy concerns due to the membership inference have motivated multiple defenses against MIAs, e.g., differential privacy and adversarial regularization. Unfortunately, these defenses produce ML models with unacceptably low classification performances.

Our work proposes a new defense, called *distillation for membership privacy* (DMP), against MIAs that preserves the utility of the resulting models significantly better than prior defenses. DMP leverages knowledge distillation to train ML models with membership privacy. We provide a novel criterion to tune the data used for knowledge transfer in order to amplify the membership privacy of DMP.

Our extensive evaluation shows that DMP provides significantly better tradeoffs between membership privacy and classification accuracies compared to state-of-the-art MIA defenses. For instance, DMP achieves ~100% accuracy improvement over adversarial regularization for DenseNet trained on CIFAR100, for similar membership privacy (measured using MIA risk): when the MIA risk is 53.7%, adversarially regularized DenseNet is 33.6% accurate, while DMP-trained DenseNet is 65.3% accurate.

1 Introduction

The remarkable performance of machine learning (ML) in solving many classification tasks has facilitated its adoption in various domains ranging from recommendation systems to critical health-care management. Many ML-as-a-Service platforms (e.g., Google API, Amazon AWS) enable novice data owners to train ML models and release the models either as a blackbox prediction API or as model parameters that can be accessed in whitebox fashion.

ML models are often trained on data with sensitive user information such as clinical records and personal photos. Hence, ML models trained using sensitive data can leak private information about their data owners. This has been demonstrated through various inference attacks (Fredrikson, Jha, and Ristenpart 2015; Hitaj, Ateniese, and Pérez-Cruz 2017; Carlini et al. 2018), and most notably the *membership inference attack* (MIA) (Shokri et al. 2017) which is

Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

the focus of our work. An MIA adversary with a blackbox or whitebox access to a target model aims to determine if a given target sample belonged to the private training data of the target model or not. MIAs are able to distinguish the members from non-members by *learning* the behavior of the target model on member versus non-member inputs. They use different features of the target model for this classification, e.g., model predictions (Shokri et al. 2017), model loss, and gradients of the model parameters for given input (Nasr, Shokri, and Houmansadr 2019). MIAs are particularly more effective against deep neural networks (Shokri et al. 2017; Salem et al. 2019), because, with their large capacities, such models can better memorize their training data.

Recent work has investigated several defenses against membership inference attacks. In order to provide the worst case privacy guarantees, Differential Privacy (DP) based defenses add very large amounts of noise to the learning objective or model outputs (Papernot et al. 2017; Chaudhuri, Monteleoni, and Sarwate 2011). This results in models with unacceptable tradeoffs between privacy and utility (Jayaraman and Evans 2019), therefore questioning their use in practice. Sablayrolles et al. (Sablayrolles et al. 2019) showed that membership privacy is a weaker notion of privacy than DP, which improves with generalization of ML models. Similarly, Nasr et al. (Nasr, Shokri, and Houmansadr 2018) proposed adversarial regularization targeted to defeat MIAs by improving the target model's generalization. However, as we demonstrate, the adversarial regularization and other state-of-the-art regularizations, including label smoothing (Szegedy et al. 2016) and dropout (Srivastava et al. 2014), fail to provide acceptable membership privacyutility tradeoffs (simply called 'tradeoffs' here onward). Memguard (Jia et al. 2019), a blackbox defense, improves model utility, but it cannot protect the model from whitebox MIAs and even the simple threshold based MIAs (Yeom et al. 2018). In summary, existing defenses against MIAs offer poor tradeoffs between model utility and membership privacy.

To this end, our work proposes a defense against MIAs that significantly improves the tradeoffs compared to prior defenses. That is, for a given degree of membership privacy (i.e., MIA resistance), our defense produces models with significantly higher classification performances compared to prior defenses. Our defense, called *Distillation for*

Membership Privacy (DMP), leverages knowledge distillation (Hinton, Vinyals, and Dean 2014), which transfers the knowledge of large models to smaller models, and is primarily used for model compression. Intuitively, DMP protects membership privacy by thwarting the access of the resulting models to the private training data. The first pre-distillation phase of DMP trains an unprotected model on the private training data without any privacy protection. Next, in distillation phase, DMP selects/generates reference data and transfers the knowledge of the unprotected model into predictions of the reference data. In the final post-distillation phase, DMP trains a protected model on the reference data labeled in the previous phase. Unlike conventional distillation, we use the same architectures for the unprotected and protected models.

Similar to adversarial regularization and PATE, DMP assumes access to a possibly sensitive and "unlabeled" reference data drawn from the same distribution as the "labeled" private training data, and uses such reference data to train its final models; the reference data is not publicly available. This is a highly realistic assumption as typical model generating entities (e.g., banks) possess huge amounts of "unlabeled" data (but limited labeled data due to the expensive labeling process). Furthermore, we show that this assumption can be relaxed by synthesizing reference data using generator networks (Micaelli and Storkey 2019). While some prior work (Papernot et al. 2017) combined distillation and DP to protect data privacy, our work is the first to study the promise of knowledge distillation as the sole technique to train membership privacy-preserving models. Our key contributions are summarized below:

- We propose a defense against MIAs, called *Distillation* for *Membership Privacy* (DMP).
- Given an unprotected model trained on a private training data and a reference sample, we provide a novel result that the lower the entropy of prediction of the model on the reference sample, the lower the sensitive membership information in the prediction. We use this result to select/generate appropriate reference data so as to improve the membership privacy due to DMP.
- We perform an extensive evaluation of DMP to show the state-of-the-art tradeoffs between membership privacy and model accuracy of DMP. For instance, at a fixed high degrees of membership privacy, DMP achieves 30% to 140% higher classification accuracies compared to stateof-the-art defenses across various classification tasks.

2 Related Work

Membership inference attacks. (Shokri et al. 2017) introduced membership inference attacks (MIAs). Given a target model trained on a private training data and a target sample, MIA adversary aims to infer whether the target sample is a member of the private training data. (Shokri et al. 2017) proposed to train a neural network to distinguish the features of the target model on members and non-members. They assumed a partial access to the private trainin data. (Salem et al. 2019) relaxed this assumptions and showed

the transferability of MIAs across datasets. These works relied on the blackbox features of target models, e.g., model predictions to mount MIAs. (Nasr, Shokri, and Houmansadr 2019) proposed to use whitebox features of target models, e.g., model gradients, along with the blackbox features, to further enhance the MIA accuracy. Above works used generalization gap (i.e., difference in train and test accuracy) of target models to mount strong MIAs. The more recent MIA literature focuses on deriving features that can better distinguish the behavior of target models on members and non-members (Leino and Fredrikson 2019; Song and Mittal 2020).

Defenses against membership inference attacks. MIAs exploit differences in behaviors of target models on members and non-members. Regularization techniques, including dropout and label smoothing, reduce the difference in terms of accuracies of the target model on members and nonmembers, and mitigate MIAs to some extent (Shokri et al. 2017). (Nasr, Shokri, and Houmansadr 2018) proposed adversarial regularization (AdvReg) tailored to defeat MIAs. AdvReg simultaneously trains the target and attack models in a game theoretic manner, and regularizes the target model using the accuracy of the attack model. The final target models that use above regularization defenses can be deployed in whitebox manner, i.e., similar to DMP, they are whitebox defenses. Hence, we thoroughly compare our DMP defense with all these regularization techniques. However, as shown in (Song and Mittal 2020) and seen from the original work (Nasr, Shokri, and Houmansadr 2018), AdvReg is not an effective defense, because it either fails to mitigate MIA or incurs large drops in model utility (classification accuracy). Jia et al. (2019) proposed MemGuard, a blackbox defense that adds noise to the output of the target model such that the noisy output is both accurate and fools the given MIA attack model. However, MemGuard does not defend against the simplest of threshold based attacks (Yeom et al. 2018; Sablayrolles et al. 2019). We omit MemGuard and other blackbox defenses, e.g., top-k predictions (Shokri et al. 2017), from evaluations.

Differential privacy based defenses such as DP-SGD (Abadi et al. 2016) and PATE (Papernot et al. 2017) are whitebox defenses and provide strong theoretical membership privacy guarantees. However, as (Jayaraman and Evans 2019) show—and we confirm in our work—target models trained using DP-SGD and PATE have prohibitively low classification accuracies rendering them unusable.

3 Preliminaries

Knowledge distillation. (Buciluă, Caruana, and Niculescu-Mizil 2006) and (Ba and Caruana 2014) proposed knowledge distillation, which uses the outputs of a large teacher model to train a smaller student model, in order to *compress* large models to smaller models. The outputs used for distillation can vary, e.g., (Hinton, Vinyals, and Dean 2014) use class probabilities generated by the teacher as the outputs, while (Romero et al. 2014) use the intermediate activations along with class probabilities of the teacher. It is well established that *knowledge distillation produces students with*

accuracies similar to their teachers (Crowley, Gray, and Storkey 2018; Zagoruyko and Komodakis 2016). This also allows DMP to produce highly accurate target models. Note that, although we use term "distillation", DMP uses teacher and student models of the same sizes, because DMP is not concerned with the size of the resulting model.

Membership inference attacks. Below we give the threat model and MIA methodology that we consider in this work.

Threat model. The primary goal of the adversary is to infer the membership of a target sample (\mathbf{x}, y) in the private training data D_{tr} of a target model θ . Our DMP defense uses private, unlabeled reference data X_{ref} for knowledge transfer, which itself could be privacy sensitive, hence, we consider a secondary goal to infer membership of a target sample in X_{ref} . Following the previous works, we assume a strong adversary with the knowledge of: target model parameters (the strongest whitebox case), half of the members of D_{tr} and equal number of non-members. Similarly, to assess the MIA risk to X_{ref} , we assume that the adversary has half of the members of X_{ref} and the equal number of nonmembers. Note that, the assumptions on the partial availability of private D_{tr} and private X_{ref} facilitates the assessment of defenses under a very strong adversary. The adversary can compute various whitebox and blackbox features of the target model and train an attack model. The adversary cannot poison X_{ref} as it is not publicly available.

Consider a target model θ and a sample (\mathbf{x}, y) . MIAs exploit the differences in the behavior of θ on members and non-members of the private D_{tr} . Therefore, MIAs train a binary attack model to classify target samples into members and non-members. Such attack models can be neural networks (Shokri et al. 2017; Salem et al. 2019) or simple thresholding functions where threshold is tuned for maximum attack performance (Yeom et al. 2018; Sablayrolles et al. 2019; Song and Mittal 2020). The adversary computes various features of θ for given (\mathbf{x}, y) , e.g., prediction $\theta(\mathbf{x}, y)$, θ 's loss on (\mathbf{x}, y) , and the gradients of the loss. The adversary combines these features to form $F(\mathbf{x}, y, \theta)$. The attack model h takes $F(\mathbf{x}, y, \theta)$ as its input and outputs the probability that (\mathbf{x},y) is a member of D_{tr} . Let $\Pr_{D_{\mathsf{tr}}}$ and $Pr_{D_{tr}}$ be the conditional probabilities of the members and non-members of D_{tr} , respectively. Hence, the expected gain of the attack model for the above setting is given by:

$$G^{\theta}(h) = \underset{\sim \Pr_{D_{tr}}}{\mathbb{E}} [\log(h(F))] + \underset{\sim \Pr_{D_{tr}}}{\mathbb{E}} [\log(1 - h(F))] \quad (1)$$

In practice, the adversary knows only a finite set of the members D and non-members D'^A required to train h, hence computes the above gain empirically as:

$$G_{D^{A},D'^{A}}^{\theta}(h) = \sum_{\substack{(\mathbf{x},y)\\ \in D^{A}}} \frac{\log(h(F))}{|D^{A}|} + \sum_{\substack{(\mathbf{x},y)\\ \in D'^{A}}} \frac{\log(1-h(F))}{|D'^{A}|}$$
(2)

Finally, the adversary solves for h^* that maximizes (2).

4 Our Proposed Defense: DMP

Now, we present our defense *Distillation For Membership Privacy (DMP)*, which is motivated by the poor membership privacy-utility tradeoffs provided by existing MIA defenses (§ 2). First, we give an intuition behind DMP and detail the DMP training. Finally, to achieve the desired tradeoffs, we give a criterion to tune the selection or generation (e.g., using GANs) of reference data used in DMP.

Notations. D_{tr} is a *private* training data. An ML model trained on D_{tr} without any privacy protections is called *un-protected* model, denoted by θ_{up} . An ML model is called *protected* model, denoted by θ_{p} , if it protects D_{tr} from MIAs. For knowledge transfer, DMP uses an *unlabeled and possibly private reference dataset* which is *disjoint* from D_{tr} ; as the reference data is unlabeled, we denote it by X_{ref} . We denote the soft label of θ on x, i.e., $\theta(x)$, by θ^x .

Main intuition of DMP. (Sablayrolles et al. 2019) show that θ trained on a sample z (short for (\mathbf{x}, y)) provides (ϵ, δ) membership privacy to z if the expected loss of the models not trained on z is ϵ -close to the loss of θ on z, with probability at least $1 - \delta$. They assume a posterior distribution of the parameters trained on a given data $D = \{z_1, ..., z_n\}$ to be:

$$\mathbb{P}(\theta|z_1,...,z_n) \propto \exp(\sum_{i=1}^n \ell(\theta,z_i))$$
 (3)

Consider a neighboring dataset $D' = \{z_1, ..., z'_j, ..., z_n\}$ of D, which is obtained by modifying at most one sample of D (Ding et al. 2018). (Sablayrolles et al. 2019) show that, to provide membership privacy to z_j , the log of the ratio of probabilities of obtaining the same θ from D and D' should be bounded, i.e., (4) should be bounded.

$$\log \left| \frac{\mathbb{P}(\theta|D)}{\mathbb{P}(\theta|D')} \right| = |\ell(\theta, z_j) - \ell(\theta, z_j')| \tag{4}$$

(4) implies that, if θ was indeed trained on z_j , then to provide membership privacy to z_j , the loss of θ on z_j should be same as the loss on any non-member sample z_j' .

DMP is a strong meta-regularization technique built on this intuition. DMP aims to protect its target models against the membership inference attacks that exploit the gap between the target model's losses on the members and non-members, by reducing the gap.

DMP achieves this via knowledge transfer and restricts the direct access of $\theta_{\rm p}$ to the private $D_{\rm tr}$, which significantly reduces the membership information leakage to $\theta_{\rm p}$. However, unlike existing knowledge transfer, DMP proposes an entropy-based criterion to select/generate $X_{\rm ref}$. Simply put, soft labels of the unprotected model $\theta_{\rm up}$ on $X_{\rm ref}$ should have low entropy and the $X_{\rm ref}$ should be far from decision boundaries of $\theta_{\rm up}$, i.e., far from $D_{\rm tr}$, in the input feature space. Intuitively, such samples are easy to classify and none of the members of $D_{\rm tr}$ significantly affects their predictions, and therefore, these predictions do not leak membership information of any particular member.

Details of the DMP technique. We now detail the three phases of our DMP defense depicted in Figure 1. In *pre-distillation phase* (step (1) in Figure 1), DMP trains θ_{UD} on

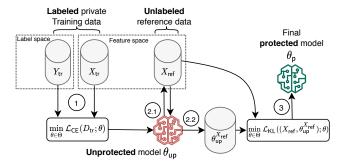


Figure 1: Distillation for Membership Privacy (DMP) defense. (1) In pre-distillation phase, DMP trains an unprotected model $\theta_{\rm up}$ on the private training data without any privacy protection. (2.1) In distillation phase, DMP uses $\theta_{\rm up}$ to select/generate appropriate reference data $X_{\rm ref}$ that minimizes membership privacy leakage. (2.2) Then, DMP transfers the knowledge of $\theta_{\rm up}$ by computing predictions of $\theta_{\rm up}$ on $X_{\rm ref}$, denoted by $\theta_{\rm up}^{X_{\rm ref}}$. (3) In post-distillation phase, DMP trains the final protected model $\theta_{\rm p}$ on $(X_{\rm ref}, \theta_{\rm up}^{X_{\rm ref}})$.

the private training data, D_{tr} , using standard SGD optimizer, e.g., Adam. Such unprotected θ_{up} is highly susceptible to MIA due to large generalization error, i.e., difference between train and test accuracies (Shokri et al. 2017; Yeom et al. 2018).

Next, in distillation phase (step (2.1) in Figure 1), DMP obtains X_{ref} required to transfer the knowledge of θ_{up} in θ_{p} . Note that, X_{ref} is unlabeled and cannot be used directly for any learning. Then, we compute soft labels of X_{ref} , i.e., $\theta_{\text{up}}^{X_{\text{ref}}} = \theta_{\text{up}}(X_{\text{ref}})$ (step (2.2) in Figure 1). There are two key factors of the distillation phase that allow us to tune DMP and achieve the desired privacy-utility tradeoffs. First, the lower the entropy of predictions $\theta_{\text{up}}^{X_{\text{ref}}}$, the lower the membership leakage through X_{ref} and vice-versa. Such low entropy predictions are characteristics of the members of D_{tr} , however, non-members with low entropy can be obtained (or generated using GANs (Micaelli and Storkey 2019)) due to large input feature space. Second, using higher softmax temperatures to compute $\theta_{\text{up}}^{X_{\text{ref}}}$ reduces membership leakage, but may reduce accuracy of the final model, and vice-versa.

Finally, in *Post-distillation phase* (step (3) in Figure 1), DMP trains a protected model θ_p on $(X_{ref}, \theta_{up}^{X_{ref}})$ using Kullback-Leibler divergence loss defined in (5). In (5), \overline{y} is the target soft label. The final θ_p is obtained by solving (6).

$$\mathcal{L}_{\text{KL}}(\mathbf{x}, \overline{\mathbf{y}}) = \sum_{i=0}^{\mathbf{c}-1} \overline{\mathbf{y}}_i \log \left(\frac{\overline{\mathbf{y}}_i}{\theta_{\text{p}}(\mathbf{x})_i} \right)$$
 (5)

$$\theta_{\mathsf{p}} = \underset{\theta}{\operatorname{argmin}} \frac{1}{|X_{\mathsf{ref}}|} \sum_{(\mathbf{x}, \overline{\mathbf{y}}) \in (X_{\mathsf{ref}}, \theta_{\mathsf{up}}^{X_{\mathsf{ref}}})} \mathcal{L}_{\mathsf{KL}}(\mathbf{x}, \overline{\mathbf{y}}) \qquad (6)$$

Due to KL-divergence loss in (6), the resulting model, θ_p , perfectly learns the behavior of θ_{up} on the X_{ref} . Furthermore, X_{ref} being a representative non-member data, i.e., test data, we expect that the test accuracies of θ_p and θ_{up} are close, and that the final DMP models will not suffer significant accuracy reductions (Ba and Caruana 2014;

Romero et al. 2014).

Fine-tuning the DMP defense. As mentioned before, the appropriate choice of reference data X_{ref} is important to achieve the desired privacy-utility tradeoffs in DMP. In this section, we show that X_{ref} with low entropy predictions of unprotected model θ_{up} strengthens membership privacy and derive an entropy-based criterion to select/generate X_{ref} .

Proposition 1. Consider θ_{up} trained on a private D_{tr} . Then, the membership leakage about D_{tr} through predictions $\theta_{up}(X_{ref})$ can be reduced by selecting/generating X_{ref} that are far from D_{tr} in input feature space with respect to some L_p distance and whose predictions, $\theta_{up}(X_{ref})$, have low entropies.

Sketch of proof of Proposition 1. Due to space limitations, we defer the detailed proof to Appendix and provide its sketch here. Consider two training datasets D_{tr} and D'_{tr} such that $D'_{\text{tr}} \leftarrow D_{\text{tr}} - z$, and X_{ref} . Then, the log of the ratio of the posterior probabilities of learning the exact same parameters θ_p using DMP is given by (11). Observe that, \mathcal{R} is an extension of (4) to the setting of DMP, where θ_p is trained via the knowledge transferred using $(X_{\text{ref}}, \theta_{\text{up}}^{X_{\text{ref}}})$, instead of directly training on D_{tr} . (Sablayrolles et al. 2019) argue to reduce this ratio to improve membership privacy. Hence, we want to obtain X_{ref} which reduces \mathcal{R} when D_{tr} , D'_{tr} , and θ_p are kept constant. We note that, although similar in appearance to differential privacy, \mathcal{R} is defined only for the given private dataset, D_{tr} .

$$\mathcal{R} = \left| \log \left(\Pr(\theta_{\mathsf{p}} | D_{\mathsf{tr}}, X_{\mathsf{ref}}) / \Pr(\theta_{\mathsf{p}} | D_{\mathsf{tr}}', X_{\mathsf{ref}}) \right) \right| \tag{7}$$

Next, we modify \mathcal{R} as:

$$\mathcal{R} = \left| -\frac{1}{T} \sum_{\mathbf{x} \in X_{\text{ref}}} \mathcal{L}_{\text{KL}}((\mathbf{x}, \theta_{\text{up}}^{\mathbf{x}}); \theta_{\text{p}}) - \mathcal{L}_{\text{KL}}((\mathbf{x}, \theta_{\text{up}}^{\prime \mathbf{x}}); \theta_{\text{p}}) \right|$$
(8)

$$\leq \frac{1}{T} \sum_{\mathbf{x} \in X_{\text{ref}}} \left| \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\mathbf{x}} \| \theta_{\text{p}}^{\mathbf{x}}) - \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\prime \mathbf{x}} \| \theta_{\text{p}}^{\mathbf{x}}) \right| \tag{9}$$

where $\theta_{\rm up}$ and $\theta'_{\rm up}$ are trained on $D_{\rm tr}$ and $D'_{\rm tr}$, respectively. Note that, (12) holds due to the assumption in (3) and the KL-divergence loss used to train $\theta_{\rm p}$ in DMP. (13) follows from (12) because $|a+b| \leq |a| + |b|$. Therefore, minimizing (13) implies minimizing (11). Thus, to improve membership privacy due to $\theta_{\rm p}$, $X_{\rm ref}$ is obtained by solving (14).

$$X_{\text{ref}}^* = \underset{X_{\text{ref}} \in X}{\operatorname{argmin}} \Big(\frac{1}{T} \sum_{\mathbf{x} \in X_{\text{ref}}} \left| \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\mathbf{x}} \| \theta_{\text{p}}^{\mathbf{x}}) - \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\prime \mathbf{x}} \| \theta_{\text{p}}^{\mathbf{x}}) \right| \Big)$$

$$\tag{10}$$

The objective of (14) is minimized when $\theta_{\rm up}^{\rm x}=\theta_{\rm up}^{\prime {\rm x}} \ \forall {\rm x} \in X_{\rm ref}$ and is very intuitive: It implies that, z (i.e., $D_{\rm tr}-D_{\rm tr}^{\prime}$) enjoys stronger membership privacy when the reference data, $X_{\rm ref}$, are such that the distributions of outputs of $\theta_{\rm up}$ and $\theta_{\rm up}^{\prime}$ on $X_{\rm ref}$ are not affected by the presence of z in $D_{\rm tr}$.

Next, we simplify (14) by replacing \mathcal{L}_{KL} with closely related cross-entropy loss \mathcal{L}_{CE} . This simplification can be

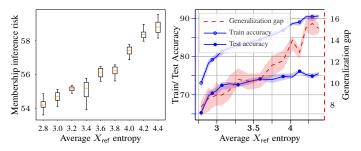


Figure 2: The lower the entropy of predictions of unprotected model on X_{ref} , the higher the membership privacy.

easily validated using X_{ref} whose ground truth labels are known. Specifically, we randomly sample D_{tr} and X_{ref} from Purchase100 dataset, and compute θ_{up} and θ_{p} using DMP. Next, for some $z \in D_{\text{tr}}$, we train θ'_{up} on D'_{tr} . Then, for each $\mathbf{x} \in X_{\text{ref}}$, we compute $\Delta \mathcal{L}_{\text{KL}}$ as in (14) and use the available ground truth label of \mathbf{x} to compute $\Delta \mathcal{L}_{\text{CE}}$. Finally, we show that $\Delta \mathcal{L}_{\text{KL}}$ and $\Delta \mathcal{L}_{\text{CE}}$ are strongly correlated for all $z \in D_{\text{tr}}$.

Next, we use the linear approximation given by (Koh and Liang 2017) for the difference in \mathcal{L}_{CE} of a pair of models trained with and without a sample to simplify (14). Then the result of Proposition 1 follows after a few simple mathematical manipulations.

Empirical verification of Proposition 1. We randomly pick $D_{\rm tr}$ of size 10k from Purhcase100 data and train $\theta_{\rm up}$. Then, we sort the rest of Purhcase100 data based on entropy of the predictions of $\theta_{\rm up}$ on the data. We form first $X_{\rm ref}$ using the first 10k data with the lowest entropies, second $X_{\rm ref}$ using the following 10k data, and so on. Finally we train multiple protected models, $\theta_{\rm p}$'s, using each of the $X_{\rm ref}$'s. Figure 2 (left) shows the increase in the MIA risk and Figure 2 (right) shows the increase in the classification performance of $\theta_{\rm p}$ with the increase in average entropy of the $X_{\rm ref}$ used. This tradeoff is because, although the higher entropy predictions contain more useful information (Nayak et al. 2019; Hinton, Vinyals, and Dean 2014) and lead to high accuracy of $\theta_{\rm p}$, they also contain higher membership information about $D_{\rm tr}$ and lead to higher MIA risk.

5 Experimental Setup

5.1 Datasets and target model architectures

We use four datasets and corresponding model architectures that are consistent with the previous works (Shokri et al. 2017; Nasr, Shokri, and Houmansadr 2019, 2018; Salem et al. 2019).

Purchase (Purchase 2017) is a 100 class classification task with 197,324 binary feature vectors of length 600; each dimension corresponds to a product and its value states if corresponding customer purchased the product; the corresponding label represents the shopping habit of the customer.

Texas (Texas 2017) is dataset of patient records. It is a 100 class classification task with 67,300 binary feature vectors of length 6,170 each; each dimension corresponds to symptoms and its value states if corresponding patient has the symptom or not; the label represents the treatment given to the pa-

tient. For Purchase and Texas we use fully connected (FC) networks.

CIFAR10 and CIFAR100 (Krizhevsky and Hinton 2009) are popular image classification datasets, each has size 50k and 32×32 color images. We use Alexnet, DenseNet-12 (with 0.77M parameters), and DenseNet-19 (with 25.6M parameters) models for CIFAR100, and Alexnet for CIFAR10. Following previous works, we measure the test accuracy of the target models as their utility.

Sizes of dataset splits. The dataset splits are given in Table 1. For Purchase and Texas tasks, we use $D_{\rm ref}$ of size 10k and select $X_{\rm ref}$ of size 10k from the remaining data using our entropy-based criterion. For CIFAR datasets, we use $D_{\rm ref}$ of size 25k and due to small sizes of these datasets, use the entire remaining 25k data as $X_{\rm ref}$. The 'Attack training' (described shortly) column shows the MIA adversary's knowledge of members and non-members of $D_{\rm tr}$. Following all the previous works, we assume that the adversary knows 50% of $D_{\rm tr}$. Further experimental details are provided in Appendix.

Dataset	DMP t	raining	Attack training		
	$ D_{tr} $	$ X_{ref} $	D	D'	
Purchase (P)	10000	10000	5000	5000	
Texas (T)	10000	10000	5000	5000	
CIFAR100 (C100)	25000	25000	12500	8000	
CIFAR10 (C10)	25000	25000	12500	8000	

Table 1: All the dataset splits are disjoint. D, D' data are the members and non-members of D_{tr} known to MIA adversary.

5.2 Membership inference attacks

We briefly review the four MIAs we use for evaluations. Following previous works, we use the accuracy of MIAs on target models as a measure of their membership privacy.

Bounded loss (BL) attack (Yeom et al. 2018) decides membership using a threshold on the target model's loss on the target sample. When 0-1 loss is used, the attack accuracy is simply the difference in training and test accuracy of target model. We denote BL attack accuracy by $A_{\rm bl}$.

NN attack (Salem et al. 2019) uses a *shadow dataset* d_s drawn from the same distribution as $D_{\rm tr}$. The attacker splits d_s in d_s' and d_s'' , trains a *shadow model* θ_s on d_s' , computes predictions of θ_s on d_s' and d_s'' , labels the predictions of d_s' as members and that of d_s'' as non-members, and trains binary attack model on the predictions. We denote NN attack accuracy by $A_{\rm nn}$. Due to their small sizes, DMP cannot be evaluated with CIFAR datasets, hence we omit NN attack evaluation for CIFAR datasets.

NSH attacks (Nasr, Shokri, and Houmansadr 2019) are similar to NN attacks. They concatenate various whitebox (e.g., model gradients) and/or blackbox (e.g., model loss, predictions) features of target model, while NN attack uses only the target model predictions. We denote whitebox and blackbox NSH attack accuracies by $A_{\rm Wb}$ and $A_{\rm bb}$, respectively. For NN and NSH attacks, we use the same attack models as the original works.

Dataset and	No defense								
model	E_{gen}	A_{test}	A_{wb}	A_{bb}	A_{bl}	A_{nn}			
P-FC	24.0	76.0	77.1	76.8	63.1	60.5			
T-FC	51.3	48.7	84.0	82.2	76.1	71.9			
C100-A	63.2	36.8	90.3	91.3	81.8	N/A			
C100-D12	33.8	65.2	72.2	71.8	67.5	N/A			
C100-D19	34.4	65.5	82.3	81.6	68.1	N/A			
C10-A	32.5	67.5	77.9	77.5	66.4	N/A			

Table 2: Models trained without any defenses have high test accuracies, $A_{\rm test}$, but their high generalization errors, $E_{\rm gen}$ (i.e., $A_{\rm train}-A_{\rm test}$) facilitate strong MIAs (§ 5.2). "N/A" means the attack is not evaluated due to lack of data.

6 Experiments

6.1 Comparison with regularization techniques

Regularization improves the generalization of ML models, and hence, reduce the MIA risk (Shokri et al. 2017). Hence, we compare DMP with five regularization defeses, including state-of-the-art MIA defense—adeversarial regularization (Nasr, Shokri, and Houmansadr 2018). In all tables, $E_{\rm gen}$ is generalization error, i.e., $(A_{\rm train} - A_{\rm test})$, where $A_{\rm train}$ and $A_{\rm test}$ are train and test accuracies of the target model, respectively. $A_{\rm test}^+$ gives the % increase in $A_{\rm test}$ due to DMP over the other regularizers. $A_{\rm wb}$, $A_{\rm bb}$, $A_{\rm bl}$, $A_{\rm nn}$ are accuracies of various attacks discussed in the previous section.

Table 2 shows accuracies of models trained without any defense; CIFAR models have lower than state-of-the-art accuracies due to smaller training datasets.

Comparison with adversarial regularization (AdvReg). Table 3 compares A_{test} of DMP and AdvReg models, for similar MIA accuracies (i.e., membership privacy). As expected, these models also have similar E_{gen} 's. However, A_{test} of DMP models is significantly higher than AdvReg models; A_{test}^+ column shows the % increase in A_{test} due to DMP over AdvReg: Accuracy improvements due to DMP over AdvReg are close to 100% for CIFAR-100, and 20% to 45% for other datasets. AdvReg uses accuracy of an MIA model to regularize and train its target models to fool the MIA model. However, AdvReg allows its target models to directly access D_{tr} . Hence, to effectively fool the MIA model, it puts relatively large weight on the regularizationloss term. This reduces the impact of the loss on main task and reduces the accuracy of AdvReg models. DMP uses appropriate reference data to transfer the knowledge of $D_{\rm tr}$ to its target models without allowing them direct access. Hence, DMP significantly outperforms AdvReg in terms of privacy-utility tradeoffs.

Comparison with other regularizers. Next, we compare DMP with four state-of-the-art regularizers: weight decay (WD), dropout (Srivastava et al. 2014) (DR), label smoothing (Szegedy et al. 2016) (LS), and confidence penalty (Pereyra et al. 2017) (CP). Due to the poor MIA resistance of CP, we defer its results to Appendix.

Table 4 shows the results, when MIA risks of regularized models is close that of DMP models (Table 3). We note that, in all the cases, $A_{\rm test}$ of DMP are significantly higher (up to 385% increase as $A_{\rm test}^+$ column specifies) than $A_{\rm test}$ of

other regularizers. This is because, these regularizers aim to improve the test accuracies of target models, but are not designed to reduce MIA risk. Thus, to reduce MIA risk, these regularization techniques add large, suboptimal noises during training, and hurt the utility of resulting models.

6.2 Comparison with differentially private defenses

Comparison with DP-SGD. Following the methodology of (Jayaraman and Evans 2019), we compare DMP and DP-SGD (Abadi et al. 2016) using the empirically observed tradeoffs between membership privacy (MIA resistance) and A_{test} of models. We use only CIFAR10 for these experiments, as the DP-SGD achieves prohibitively low accuracies on difficult tasks such as Texas and CIFAR100. We evaluate MIA risk using the whitebox NSH attack. Table 5 shows the results of Alexnet trained on CIFAR10 using DMP and DP-SGD with different privacy budgets ϵ 's; -ve E_{gen} implies A_{train} is lower than A_{test} . DP-SGD incurs significant (35%) loss in A_{test} at lower ϵ (12.5) to provide strong membership privacy. At higher ϵ , A_{test} of DP-SGD increases, but at the cost of very high generalization error, which facilitates stronger MIAs. Note that, further increase in privacy budget, ϵ , does not improve tradeoff of DP-SGD. More importantly, for low MIA risk of $\sim 51.3\%$, DMP models have 12.8% higher A_{test} (i.e., 24.5% improvement) than DP-SGD models, which shows the superior tradeoffs due to DMP.

Comparison with PATE. PATE (Papernot et al. 2017), a semi-supervized learning technique, requires a compatible pair of generator and disciminator to achieve acceptable performances. Hence, we use CIFAR10 dataset and, instead of Alexnet, use the generator-discriminator pair from (Salimans et al. 2016), which has state-of-the-art performances. PATE trains a set of teachers, computes hard labels of each teacher on some X_{ref} , aggregates the labels for each $\mathbf{x} \in X_{\text{ref}}$ using majority voting, adds DP noise to the aggregate, and finally trains its target model on the noisy aggregate.

We train ensembles of 5, 10, and 25 teachers using $D_{\rm tr}$ of sise 25k. We use the optimized confident-GNMax (GNMax) aggregation scheme of (Papernot et al. 2018) to label $X_{\rm ref}$ We present a subset of results in Table 6 and defer comprehensive comparison to Appendix. At low ϵ 's (<10), GNMax hardly produces any labels, hence, the final target model has very low $A_{\rm test}$, but at higher ϵ 's (>1000), PATE target model has acceptable $A_{\rm test}$. However, PATE cannot achieve performances close to DMP, as it divides $D_{\rm tr}$ among its teachers. Such teachers have significantly lower accuracies and their ensemble cannot achieve the accuracy close to that of the unprotected model of DMP, which is trained on the entire $D_{\rm tr}$. Hence, the quality of knowledge transferred in DMP is always higher than that in PATE.

6.3 Discussions

Below, we provide further key insights in to DMP defense and defer their detailed discussion to Appendix.

Hyperparameter selection in DMP. *Increasing* the temperature of softmax layer of the unprotected model, θ_{UD} ,

Dataset	A	Adversarial regularization (AdvReg)				DMP							
and	F	Egen Atest		Attack accuracy		$E_{\sf gen}$	A _{test}	4+		Attack a	accuracy		
model	<i>E</i> gen	Atest	A_{wb}	A_{bb}	A_{bl}	A_{nn}	Dgen Pitest	A ⁺ _{test}	A_{wb}	A_{bb}	A_{bl}	A_{nn}	
Purchase + FC	9.7	56.5	55.8	55.4	54.9	50.1	10.1	74.1	+31.2%	55.3	55.1	55.2	50.2
Texas + FC	6.1	33.5	58.2	57.9	54.1	50.8	7.1	48.6	+45.1%	55.3	55.4	53.6	50.0
CIFAR100 + Alexnet	6.9	19.7	54.3	54.0	53.5	N/A	6.5	35.7	+81.2%	55.7	55.6	53.3	N/A
CIFAR100 + DenseNet-12	5.5	26.5	51.4	51.3	52.8	N/A	3.6	63.1	+138.1%	53.7	53.0	51.8	N/A
CIFAR100 + DenseNet-19	7.2	33.9	54.2	53.4	53.6	N/A	7.3	65.3	+92.6%	54.7	54.4	53.7	N/A
CIFAR10 + Alexnet	4.2	53.4	51.9	51.2	52.1	N/A	3.1	65.0	+21.7%	51.3	50.6	51.6	N/A

Table 3: Comparing test accuracy (A_{test}) and generalization error (E_{gen}) of DMP and Adversarial Regularization, for near-equal, low MIA risks (high membership privacy). A_{test}^+ shows the % increase in A_{test} of DMP over Adversarial Regularization.

	Purchas	e + FC (I	OMP's A_{tes}	t = 74.1						
Regularizer	E_{gen}	A_{test}	A ⁺ _{test}	A_{wb}	A_{bb}	A_{bl}				
WD	10.3	42.5	+74.4%	54.9	55.4	55.2				
WD + DR	9.1	42.1	+76.0%	56.4	56.8	54.6				
WD + LS	12.3	42.0	+76.4%	57.2	57.0	56.2				
Texas + FC (DMP's $A_{\text{test}} = 48.6$)										
Regularizer	E_{gen}	A_{test}	A_{test}^+	A_{wb}	A_{bb}	A_{bl}				
WD	5.0	22.5	+116%	58.3	57.7	52.5				
WD + DR	6.1	14.2	+242%	63.1	62.6	53.1				
WD + LS	8.3	37.3	+30%	61.7	61.0	54.2				
CIFA	CIFAR100 + DenseNet-12 (DMP's $A_{\text{test}} = 63.1$)									
Regularizer	E_{gen}	A_{test}	A_{test}^+	A_{wb}	A_{bb}	A_{bl}				
WD	4.0	26.3	+140%	49.9	49.7	52.0				
WD + DR	3.7	32.3	+95.4%	51.2	51.0	51.9				
WD + LS	2.7	13.0	+385%	51.0	51.4	51.4				
C	IFAR10	+ Alexne	t (DMP's A	test = 65	.0)					
Regularizer	E_{gen}	A_{test}	A ⁺ _{test}	A_{wb}	A_{bb}	A_{bl}				
WD	4.1	45.9	+41.6%	52.4	52.5	52.1				
WD + DR	3.2	44.7	+45.4%	51.9	51.7	51.6				
WD + LS	4.8	53.2	+22.2%	53.8	53.0	52.4				

Table 4: Evaluating three state-of-the-art regularizers, with similar, low MIA risks (high membership privacy) as DMP. A_{test}^+ shows the % increase in A_{test} due to DMP over the corresponding regularizers.

used to transfer the knowledge of θ_{up} , can further reduce the membership leakage of D_{tr} . This is because, at higher softmax temperatures, predictions of θ_{up} have uniform distribution over all classes and contain no useful information for MIAs. Similarly, reducing the size of X_{ref} reduces MIA risk due to DMP, but comes at the cost of reduction in A_{test} .

Privacy risk to reference data (X_{ref}) . We evaluate the privacy risk to X_{ref} , as it can be of sensitive nature, e.g., in case of Texas medical records dataset. Our results in appendix show that given the final DMP model, θ_{p} , and a target sample, MIA adversary (who mounts BL, NN, or NSH attacks) cannot decide if the sample belonged to X_{ref} with sufficient confidence. This is expected, because DMP trains its θ_{p} on noisy, soft-labels of X_{ref} , which do not contain any sensitive information about X_{ref} and its ground-truth labels, which is necessary for MIAs to succeed (Yeom et al. 2018). We provide detailed results in Appendix.

DMP with synthetic reference data (X_{ref}). Following previous works (Papernot et al. 2018, 2017), including the state-of-the-art MIA defense AdvReg (Nasr, Shokri, and Houmansadr 2018), we assume availability of X_{ref} . How-

Defense	Privacy budget (ϵ)	E_{gen}	A_{test}	A_{wb}
No defense	_	32.5	67.5	77.9
DMP	_	3.10	65.0	51.3
	198.5	3.60	52.2	51.7
DP-SGD	50.2	1.30	36.9	50.2
DF-30D	12.5	0.30	31.7	50.0
	6.8	-1.60	29.4	49.9

Table 5: DP-SGD versus DMP for CIFAR10 and Alexnet. For low MIA risk of $\sim 51.3\%$, DMP achieves 24.5% higher $A_{\rm test}$ than of DP-SGD (12.8% absolute increase in $A_{\rm test}$).

# of	Queries	Privacy	Target	model	4
Teachers	answered	budget (ϵ)	E_{gen}	A_{test}	A_{wb}
- 5	49	195.9	31.4	33.9	49.1
3	1163	11684	65.4	68.1	49.0
10	23	42.9	39.1	38.3	50.1
10	1527	6535	63.9	65.2	49.8
25	108	183.5	53.8	55.7	49.0
23	4933	1794.1	57.8	60.3	48.6

Table 6: Comparing PATE with DMP. DMP has $E_{\rm gen}$, $A_{\rm test}$, and $A_{\rm wb}$ of 1.19%, 76.79%, and 50.8%, respectively. PATE has low accuracy even at high privacy budgets, as it divides data among teachers and produces low accuracy ensembles.

ever, in privacy sensitive domains such as patient medical records, X_{ref} may not be available. Hence, we show that the assumption can be relaxed by using X_{ref} synthesized from private D_{tr} to train DMP models. For CIFAR10, we use DC-GAN to generate synthetic X_{ref} of sizes 12.5k, 25k, and 37.5k from D_{tr} of size 25k. We then train three DMP models and evaluate their MIA risk using whitebox NSH attack. We note that for 12.5k, 25k, and 37.5k synthetic X_{ref} samples, $(E_{\text{gen}}, A_{\text{test}}, A_{\text{wb}})$ of DMP are (2.1, 53.0, 50.3), (3.5, 56.8, 51.3), and (5.0, 57.5, 52.1), respectively. Note that, *DMP outperforms existing defenses even with synthetic* X_{ref} (Tables 3, 4).

Adaptive attack on DMP. In DMP, the reference data, $X_{\rm ref}$, is selected such that the predictions of DMP's unprotected model $\theta_{\rm up}$ on $X_{\rm ref}$ have low entropies. Due to memorization, predictions of $\theta_{\rm up}$ on $D_{\rm tr}$ also have low entropies. Hence, an adaptive adversary may exploit this peculiar $X_{\rm ref}$ selection in DMP. Based on this intuition, we investigate the possibility of an adaptive MIA, which labels a target sample as a member if the sample is close to some $X_{\rm ref}$ datum in feature space. However, such attack has accuracy close to random guess. This is because, we observe that the proximity of two samples in feature space has no correlation with the entropy of predictions of given $\theta_{\rm up}$ on those samples,

which is the selection criterion of DMP. We leave further investigation of adaptive attacks on DMP to future work.

7 Conclusions

We proposed Distillation for Membership Privacy (DMP), a knowledge distillation based defense against membership inference attacks that significantly improves the membership privacy-model utility tradeoffs compared to state-of-the-art defenses. We provided a novel criterion to generate/s-elect reference data in DMP and achieve the desired tradeoffs. Our extensive evaluation demonstrated the state-of-the-art privacy-utility tradeoffs of DMP.

Acknowledgments. This work was supported in part by NSF grant CPS-1739462.

References

- Abadi, M.; Chu, A.; Goodfellow, I.; McMahan, H. B.; Mironov, I.; Talwar, K.; and Zhang, L. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM.
- Ba, J.; and Caruana, R. 2014. Do deep nets really need to be deep? In *Advances in neural information processing systems*, 2654–2662.
- Buciluă, C.; Caruana, R.; and Niculescu-Mizil, A. 2006. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 535–541.
- Carlini, N.; Liu, C.; Kos, J.; Erlingsson, U.; and Song, D. 2018. The secret sharer: Measuring unintended neural network memorization and extracting secrets. *arXiv preprint arXiv:1802.08232*.
- Chaudhuri, K.; Monteleoni, C.; and Sarwate, A. D. 2011. Differentially private empirical risk minimization. *Journal of Machine Learning Research* 12(Mar): 1069–1109.
- Crowley, E. J.; Gray, G.; and Storkey, A. J. 2018. Moonshine: Distilling with cheap convolutions. In *Advances in Neural Information Processing Systems*, 2888–2898.
- Ding, Z.; Wang, Y.; Wang, G.; Zhang, D.; and Kifer, D. 2018. Detecting violations of differential privacy. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 475–489. ACM.
- Dwork, C.; Roth, A.; et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*.
- Fredrikson, M.; Jha, S.; and Ristenpart, T. 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. ACM.
- Hinton, G.; Vinyals, O.; and Dean, J. 2014. Distilling the Knowledge in a Neural Network. *NIPS 2014 Deep Learning Workshop*.

- Hitaj, B.; Ateniese, G.; and Pérez-Cruz, F. 2017. Deep models under the GAN: information leakage from collaborative deep learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM.
- Jayaraman, B.; and Evans, D. 2019. Evaluating Differentially Private Machine Learning in Practice. In *USENIX Security Symposium*.
- Jia, J.; Salem, A.; Backes, M.; Zhang, Y.; and Gong, N. Z. 2019. MemGuard: Defending against Black-Box Membership Inference Attacks via Adversarial Examples. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 259–274.
- Koh, P. W.; and Liang, P. 2017. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 1885–1894. JMLR. org.
- Krizhevsky, A.; and Hinton, G. 2009. Learning multiple layers of features from tiny images .
- Leino, K.; and Fredrikson, M. 2019. Stolen Memories: Leveraging Model Memorization for Calibrated White-Box Membership Inference. *arXiv preprint arXiv:1906.11798*.
- Long, Y.; Bindschaedler, V.; Wang, L.; Bu, D.; Wang, X.; Tang, H.; Gunter, C. A.; and Chen, K. 2018. Understanding Membership Inferences on Well-Generalized Learning Models. *arXiv preprint arXiv:1802.04889*.
- Micaelli, P.; and Storkey, A. J. 2019. Zero-shot knowledge transfer via adversarial belief matching. In *Advances in Neural Information Processing Systems*, 9551–9561.
- Nasr, M.; Shokri, R.; and Houmansadr, A. 2018. Machine learning with membership privacy using adversarial regularization. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 634–646. ACM.
- Nasr, M.; Shokri, R.; and Houmansadr, A. 2019. Comprehensive Privacy Analysis of Deep Learning: Stand-alone and Federated Learning under Passive and Active White-box Inference Attacks. *Security and Privacy (SP)*, 2019 IEEE Symposium on .
- Nayak, G. K.; Mopuri, K. R.; Shaj, V.; Radhakrishnan, V. B.; and Chakraborty, A. 2019. Zero-Shot Knowledge Distillation in Deep Networks. In *International Conference on Machine Learning*, 4743–4751.
- Papernot, N.; Abadi, M.; Erlingsson, U.; Goodfellow, I.; and Talwar, K. 2017. Semi-supervised knowledge transfer for deep learning from private training data. *International Conference on Learning and Representation*.
- Papernot, N.; Song, S.; Mironov, I.; Raghunathan, A.; Talwar, K.; and Erlingsson, Ú. 2018. Scalable Private Learning with PATE. *arXiv preprint arXiv:1802.08908*.
- Pereyra, G.; Tucker, G.; Chorowski, J.; Kaiser, Ł.; and Hinton, G. 2017. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*.

Purchase. 2017. Acquire Valued Shoppers Challenge. https://www.kaggle.com/c/acquire-valued-shoppers-challenge/data. [Online; accessed 11-September-2019].

Romero, A.; Ballas, N.; Kahou, S. E.; Chassang, A.; Gatta, C.; and Bengio, Y. 2014. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*.

Sablayrolles, A.; Douze, M.; Schmid, C.; Ollivier, Y.; and Jegou, H. 2019. White-box vs Black-box: Bayes Optimal Strategies for Membership Inference. In *International Conference on Machine Learning*, 5558–5567.

Salem, A.; Zhang, Y.; Humbert, M.; Fritz, M.; and Backes, M. 2019. ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models. *In NDSS*.

Salimans, T.; Goodfellow, I.; Zaremba, W.; Cheung, V.; Radford, A.; and Chen, X. 2016. Improved techniques for training gans. In *Advances in neural information processing systems*, 2234–2242.

Shokri, R.; Stronati, M.; Song, C.; and Shmatikov, V. 2017. Membership inference attacks against machine learning models. In *Security and Privacy (SP)*, 2017 IEEE Symposium on.

Song, L.; and Mittal, P. 2020. Systematic Evaluation of Privacy Risks of Machine Learning Models. *arXiv preprint arXiv:2003.10595*.

Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; and Salakhutdinov, R. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15(1): 1929–1958.

Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; and Wojna, Z. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2818–2826.

Texas. 2017. Texas hospital stays dataset. https://www.dshs.texas.gov/THCIC/Hospitals/Download.shtm. [Online; accessed 10-February-2020].

Yeom, S.; Giacomelli, I.; Fredrikson, M.; and Jha, S. 2018. Privacy risk in machine learning: Analyzing the connection to overfitting. In 2018 IEEE 31st Computer Security Foundations Symposium (CSF), 268–282. IEEE.

Zagoruyko, S.; and Komodakis, N. 2016. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*.

A Fine tuning the DMP defense (Missing details)

We propose a fine tuning technique to select/generate appropriate reference data, X_{ref} , and achieve the desired privacy-utility tradeoffs using our distillation for membership privacy (DMP) defense. The technique depends on the result given in Proposition 1; we provide a detailed proof of the results below.

Detailed proof of Proposition 1.

Deriving the objective for desired $X_{\text{ref.}}$ Consider two training datasets D_{tr} and D'_{tr} such that $D'_{\text{tr}} \leftarrow D_{\text{tr}} - z$, and $X_{\text{ref.}}$ Then, the log of the ratio of the posterior probabilities of learning the exact same parameters θ_{p} using DMP is given by (11). Observe that, \mathcal{R} is an extension of (4) to the setting of DMP, where θ_{p} is trained via the knowledge transferred using $(X_{\text{ref.}}, \theta^{X_{\text{ref.}}}_{\text{up}})$, instead of directly training on $D_{\text{tr.}}$ (Sablayrolles et al. 2019) argue to reduce this ratio to improve membership privacy. Hence, we want to obtain $X_{\text{ref.}}$ which reduces the ratio \mathcal{R} when $D_{\text{tr.}}$, $D'_{\text{tr.}}$, and θ_{p} are kept constant. We note that, although similar in appearance to differential privacy, \mathcal{R} is defined only for the given private dataset, $D_{\text{tr.}}$

$$\mathcal{R} = \left| \log \frac{\Pr(\theta_{\mathsf{p}} | D_{\mathsf{tr}}, X_{\mathsf{ref}})}{\Pr(\theta_{\mathsf{p}} | D'_{\mathsf{tr}}, X_{\mathsf{ref}})} \right| \tag{11}$$

Next, we modify \mathcal{R} as:

$$\mathcal{R} = \left| -\frac{1}{T} \sum_{\mathbf{x} \in X_{\text{ref}}} \mathcal{L}_{\text{KL}}((\mathbf{x}, \theta_{\text{up}}^{\mathbf{x}}); \theta_{\text{p}}) - \mathcal{L}_{\text{KL}}((\mathbf{x}, \theta_{\text{up}}^{\prime \mathbf{x}}); \theta_{\text{p}}) \right|$$
(12)

$$\leq \frac{1}{T} \sum_{\mathbf{x} \in X_{\text{ref}}} \left| \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\mathbf{x}} \| \theta_{\text{p}}^{\mathbf{x}}) - \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\prime \mathbf{x}} \| \theta_{\text{p}}^{\mathbf{x}}) \right| \tag{13}$$

where $\theta_{\rm up}$ and $\theta'_{\rm up}$ are trained on $D_{\rm tr}$ and $D'_{\rm tr}$, respectively. Note that, (12) holds due to the assumption in (3) and the KL-divergence loss used to train $\theta_{\rm p}$ in DMP. (13) follows from (12) because $|a+b| \leq |a| + |b|$. Therefore, minimizing (13) implies minimizing (11). Thus, to improve membership privacy due to $\theta_{\rm p}$, $X_{\rm ref}$ is obtained by solving (14).

$$X_{\text{ref}}^* = \underset{X_{\text{ref}} \in X}{\operatorname{argmin}} \left(\frac{1}{T} \sum_{\mathbf{x} \in X_{\text{ref}}} \left| \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\mathbf{x}} || \theta_{\text{p}}^{\mathbf{x}}) - \mathcal{L}_{\text{KL}}(\theta_{\text{up}}^{\prime \mathbf{x}} || \theta_{\text{p}}^{\mathbf{x}}) \right| \right)$$

$$\tag{14}$$

The objective of (14) is minimized when $\theta_{\rm up}^{\rm x}=\theta_{\rm up}^{\prime \rm x}$ $\forall {\rm x}\in X_{\rm ref}$ and is very intuitive: It implies that, z (i.e., $D_{\rm tr}-D_{\rm tr}'$) enjoys stronger membership privacy when the reference data, $X_{\rm ref}$, are such that the distributions of outputs of $\theta_{\rm up}$ and $\theta_{\rm up}'$ on $X_{\rm ref}$ are not affected by the presence of z in $D_{\rm tr}$.

Simplifying the objective. Next, we simplify (14) by replacing \mathcal{L}_{KL} with closely related cross-entropy loss \mathcal{L}_{CE} . The simplified objective is given by (15).

$$X_{\text{ref}}^* = \underset{X_{\text{ref}} \in X}{\operatorname{argmin}} \sum_{\substack{z' = (\mathbf{x}, y) \\ \in (X_{\text{ref}}, Y_{\text{ref}})}} \frac{1}{T} \left| \mathcal{L}_{\text{CE}}(z'; \theta'_{\text{up}}) - \mathcal{L}_{\text{CE}}(z'; \theta_{\text{up}}) \right|$$

$$\tag{15}$$

where \mathcal{L}_{CE} is cross-entropy loss and z' is not the same as $z \leftarrow D_{\text{tr}} - D'_{\text{tr}}$. For clarity of presentation, here onward, we denote \mathcal{L}_{CE} by \mathcal{L} .

Next, we assume that ground truth labels Y_{ref} of X_{ref} are available. Note that X_{ref} is unlabeled dataset, but only to empirically demonstrate the validity of the simplification of (14) to (15), we assume that ground truth labels of X_{ref} are available. We validate the simplification in Figure 3: for any

given reference sample, the lower the difference between cross-entropy losses, $\Delta \mathcal{L}$, the lower the corresponding difference between KL-divergence losses; and vice-versa. Note that, to select/generate a reference sample, we do not need the exact difference between cross-entropy or KL-divergence losses for the sample, but only the difference for the sample relative to the other samples. Hence, although the difference between cross-entropy losses is not exactly the same as difference between KL-divergence losses, their strong positive correlation is sufficient to make the reduction (14) \rightarrow (15) useful in our task.

Deriving the final objective to select/generate $X_{\rm ref}$. Next, to avoid repetitive training, we simplify the term for each sample in (15) using the results of (Koh and Liang 2017). More specifically, they propose a linear approximation to the difference in cross-entropy losses of a pair of models trained with and without a specific sample in their training data. We note that this is the exact setting of our problem. If θ and θ_{-z} are two models trained with and without a member z, then the difference in cross-entropy losses of the two models on some test sample $z_{\rm test} = (\mathbf{x}_{\rm test}, y_{\rm test})$ is approximated as:

$$|\mathcal{L}(z_{\text{test}}, \theta_{-z}) - \mathcal{L}(z_{\text{test}}, \theta)| \simeq |\nabla_{\theta} \mathcal{L}(z_{\text{test}}, \theta) H_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(z, \theta)|$$
(16)

where H_{θ} is the Hessian matrix that is defined as $H_{\theta} = \frac{1}{n} \sum_{z \in D_{\text{tr}}} \nabla_{\theta}^2(z, \theta)$. Substituting (16) in (15) simplifies the objective in (14) to:

$$X_{\text{ref}}^* = \underset{X_{\text{ref}} \in X}{\operatorname{argmin}} \sum_{\substack{z' = (\mathbf{x}, y) \\ \in (X_{\text{ref}}, Y_{\text{ref}})}} \frac{1}{T} |\nabla_{\theta} \mathcal{L}(z', \theta_{\text{up}}) H_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(z', \theta_{\text{up}})|$$

$$\tag{17}$$

Note that, for a given member z, $H_{\theta}^{-1}\nabla_{\theta}\mathcal{L}(z',\theta)$ in (17) remains constant and the minimization reduces to minimizing the gradient $\nabla_{\theta}\mathcal{L}(z_{\mathsf{p}},\theta_{\mathsf{up}})$. The lower the loss $\mathcal{L}(z',\theta_{\mathsf{up}})$, the smaller the gradient $\nabla_{\theta}\mathcal{L}(z',\theta_{\mathsf{up}})$. Therefore objective (17) further simplifies as:

$$X_{\mathsf{ref}}^* = \underset{X_{\mathsf{ref}} \in X}{\operatorname{argmin}} \ \frac{1}{T} \sum_{\substack{z' = (\mathbf{x}', y) \\ \in (X_{\mathsf{obs}}, Y_{\mathsf{obs}})}} \mathcal{L}_{\mathsf{CE}}(z', \theta_{\mathsf{up}}) \tag{18}$$

Note that, in practice, it is not possible to solve the objective in (18) as it is. Because, we cannot compute the loss without the ground truth labels of X_{ref} ; recall that X_{ref} is unlabeled. However, as the loss involved here is the crossentropy loss, minimizing the loss is equivalent to minimizing the entropy of prediction $\theta_{\text{up}}(\mathbf{x}')$. This gives us the final objective as:

$$X_{\mathsf{ref}}^* = \underset{X_{\mathsf{ref}} \in X}{\operatorname{argmin}} \ \frac{1}{T} \sum_{\mathbf{x}' \in X_{\mathsf{ref}}} \mathcal{H}(\theta_{\mathsf{up}}(\mathbf{x}')) \tag{19}$$

where, $\mathcal{H}(\mathbf{v}) \triangleq \sum_{i} -\mathbf{v}_{i} \log(\mathbf{v}_{i})$ is the entropy of \mathbf{v} . This provides the result of Proposition 1.

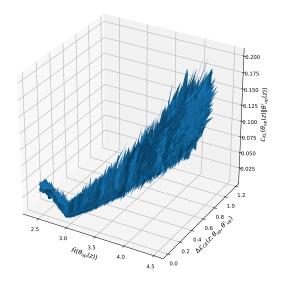


Figure 3: Empirical validation of simplification of (14) to (15): Increase in $\Delta \mathcal{L}_{CE}$ increases $\Delta \mathcal{L}_{KL}$, and that of (14) to (19): Increase in $\mathcal{H}(\theta_{UD}(z))$ increases $\Delta \mathcal{L}_{KL}$.

Proposition 1 states that, using the reference data with low entropy predictions of θ_{up} strengthens the membership resistance of θ_p , and vice versa. In Figure 3, we empirically validate the reductions (14) \rightarrow (18) \rightarrow (19). Specifically, we show that, for a given θ_{up} , the lower the cross-entropy loss of reference data sample, the lower the entropy of prediction of θ_{up} on the sample, i.e., (18) \rightarrow (19). Then, we show that, the difference between cross-entropy losses of two models θ_{up} and θ'_{up} , trained on neighboring datasets, on a sample increases with the increase in cross-entropy loss of their prediction on the sample, i.e., (15) \rightarrow (18). This, in combination with the reduction (14) \rightarrow (15) demonstrated in Figure 3, completes the validation of (14) \rightarrow (18). Figure 2 validates our hypothesis.

B Missing Details of Experimental Setup

B.1 Computing environment

We will make our code and all the relevant datasets (all the datasets used are already available online) publicly available upon acceptance of the submission. We perform all of our experiments using PyTorch 1.2 framework on TitanX GPU of 12GB memory. All the experimental results in the paper are average of three runs of the corresponding experimental setting.

B.2 Target model architectures

Unlike conventional distillation (Hinton, Vinyals, and Dean 2014), DMP uses same architectures for unprotected and protected models. Needless to say, using a lower-capacity architecture for the protected model will further improve privacy protection at the cost of reducing utility (prediction accuracy). The details of the architectures for all the datasets is given in Table 9. For Purchase-100 and Texas-100, the fully connected network has hidden layers of sizes {1024,

512, 256, 128}. For CIFAR-100, we choose two DenseNet models to assess the efficacy of DMP for two models with equivalent performance, but significantly different capacities. In Table 9, DenseNet12 corresponds to DenseNet-BC (L=100, k=12) and DenseNet19 corresponds to DenseNet-BC (L=190, k=40). For the comparison with PATE using CIFAR-10, we use the generator and discriminator architectures used in (Salimans et al. 2016).

C Detailed comparison with PATE

In this section, we detail the experimental comparison between PATE (Papernot et al. 2018, 2017) and our DMP defense for CIFAR10 classification task. The motivation of this comparison is to show that the DMP-trained models achieve significantly better tradeoffs between membership privacy (i.e., resistance to membership inference attacks) and classification accuracy than the PATE-trained models.

PATE relies on semi-supervised learning that uses a large unlabeled dataset. PATE computes the labels of a subset of the unlabeled data using an ensemble of teachers. Each of the teachers is trained on a disjoint set of the private training dataset; all sets have the same size. Semi-supervised learning involves an unstable game between a generator G and a discriminator D. Hence, the architectures of G and Dshould be compatible for effective learning. Therefore, instead of AlexNet, which we use in the rest of our CIFAR10 experiments, we use the pair of discriminator and generator architectures proposed in (Salimans et al. 2016) due to its state-of-the-art classification performance. Finally, PATE uses its discriminator as the classification model. For both PATE and DMP, we use the same 25,000 data of CIFAR10 as the private training and the rest of 25,000 data the unlabeled reference data. The accuracy of the discriminator trained on the entire private training data is 97.65% and 79.6% on training and test data, respectively.

We use the 25,000 training data to train three ensembles of sizes 5, 10 and 25 teachers. Each of the teachers in all the ensembles have disjoint and equal-sized training data. The accuracy, without adding any noise to labels, of the corresponding ensembles on the 25000 reference samples is 64.92%, 60.1% and 54.52%, respectively. We use the confident-GNMax (GNMax) aggregation scheme to add DP noise to the aggregate of the votes (i.e., hard labels) of the teachers on the unlabeled reference data. GNMax labels samples based on remaining privacy budget, hence, it may not label all the reference data samples. GNMax aggregation scheme is similar to the sparse vector technique (Dwork, Roth et al. 2014) and outputs a label only if the noisy version of the votes count of the label crosses a noisy version of a fixed threshold. Table 7 details the accuracy of the GN-Max aggregation for different number of teachers and privacy levels (ϵ, δ) . We use δ of 10^{-4} as the order of the size of the reference data is 10^4 (Papernot et al. 2018).

Note that, the DMP-trained discriminator has training, test, and attack accuracies of 77.98%, 76.79%, and 50.8%, respectively. Table 7 shows results for PATE with teacher ensembles of different sizes: At low ϵ values, GNMax cannot provide many labels, and therefore, PATE suffers significant accuracy degradations. While at high ϵ values (>1000),

GNMax performs better, but does not outperform DMP. The reason for this is as follows: At very high ϵ 's, PATE is just a knowledge transfer based semi-supervised learning, while DMP is knowledge transfer based supervised learning. DMP does not divide its training data among teacher, and therefore, the predictions of the unprotected model used in DMP to train the protected model are more useful in terms of both the quality and quantity. Therefore, DMP-trained models have significantly higher accuracy than PATE-trained model, for similar membership inference risk, i.e., DMP achieves significantly better membership privacy-model utility tradeoffs.

D Missing Discussion Details

In the last section of main paper, we provide various insights in to our DMP defense based on our extensive evaluation. We provide the missing details of the discussions below.

D.1 Hyperparameter selection in DMP

The temperature of the softmax layer. The softmax temperature, T, of the unprotected model, $\theta_{\rm up}$, plays an important role in the amount of knowledge transferred from the unprotected to protected model in DMP. Our results in Table 8 confirm our analytical understanding of the use of the softmax temperature: increasing the temperature for AlexNet trained on CIFAR100 dataset reduces the classification accuracy of the final protected model, $\theta_{\rm p}$, but also strengthens the its membership inference resistance. Therefore, the softmax temperature T should be chosen depending on the desired privacy-utility tradeoff. Table 9 shows the temperatures used in our experiments.

The size of reference data. In DMP, the more the reference data, the looser the bound on \mathcal{R} in (11), and therefore, weaker the membership resistance of the corresponding $\theta_{\rm p}$. To validate this, we quantify the classification accuracy and the membership inference risk of θ_p with increasing the amount of X_{ref} . We use Purchase-100 data and vary $|X_{ref}|$ as shown in Figure 4; we fix the softmax T of θ_{up} at 1.0. θ_{up} used here has train accuracy, test accuracy, and membership inference risk of 99.9%, 77.0% and 77.1%, respectively. Initially, the test accuracy of $\theta_{\rm p}$ increases with $|X_{\rm ref}|$ due to the useful knowledge transferred. But, beyond the test accuracy of θ_{up} , its predictions essentially insert noise in the training data of θ_p , therefore the gain from increasing the size of reference data slows down. Although this noise marginalizes the increase in the test performance of $\theta_{\rm p}$, it also prevents $\theta_{\rm p}$ from learning more about D_{tr} and prevents further inference risk. This is shown by the train accuracy and membership inference risk curves in Figure 4. Therefore, size of reference data should be selected based on the desired tradeoffs of the final model.

D.2 Privacy risk to reference data (X_{ref})

The reference data used in DMP can be of sensitive nature. For instance, for Texas-100, the reference data used are unlabeled, sensitive patients' records, and therefore, at the risk of privacy breach. However, we quantitatively show that **DMP** does not pose membership inference risk to its reference

	5 Teachers 10 Teachers				25 Teachers						
Queries	Privacy	GNMax	Student	Queries	Privacy	GNMax	Student	Queries	Privacy	GNMax	Student
answered	bound ϵ	accuracy	accuracy	answered	bound ϵ	accuracy	accuracy	answered	bound ϵ	accuracy	accuracy
0	4.6	-	-	0	9	-	_	0	8.43	_	_
49	195.9	79.6	33.93	23	42.87	56.5	38.28	108	183.5	95.4	55.7
127	281.6	69.3	49.89	358	409.5	67.0	57.59	357	231.3	83.9	56.14
679	1283.7	70.3	58.04	1128	1092.5	66.13	60.94	1130	508.9	83.8	58.26
1163	11684	91.1	68.08	1527	6535	93.1	65.18	4933	1794.1	74.0	60.27

Table 7: Evaluation of PATE using the discriminator architecture in (Salimans et al. 2016) trained on CIFAR10. The corresponding DMP-trained model has 77.98% and 76.79% accuracies on the training and test data, and 50.8% membership inference accuracy. Comparison of results clearly show the superior membership privacy-model utility tradeoffs of DMP over PATE.

Deferre	Softmax	Training	Test	Attack
Defense	T	Accuracy	Accuracy	Accuracy
No defense	n/a	100	36.8	91.3
	2	46.6	37.3	57.4
DMP	4	42.2	35.7	55.6
DMP	6	36.4	32.8	52.5
	8	12.1	12.3	51.7

Table 8: Effect of the softmax temperature on DMP: For a fixed X_{ref} , increasing the temperature of softmax layer of θ_{up} reduces \mathcal{R} in (11), which strengthens the membership privacy.

	nbination cronym	Dataset	Architecture	$ \theta $	T
	P-FC	Purchase	Fully Connected	1.32M	1.0
	T-FC	Texas	Fully Connected	1.32M	1.0
C	100-A		AlexNet	2.47M	4.0
C1	00-D12	CIFAR100	DenseNet12	0.77M	4.0
C1	00-D19		DenseNet19	25.6M	1.0
(C10-A	CIFAR10	AlexNet	2.47M	1.0

Table 9: Temperature of the softmax layers for the different combinations of dataset and network architecture used to produce the results in Table 3 of the main paper.

data. The results are given in Table 10. We note that, for any combination of model and dataset, the membership inference risk to the reference data due to DMP is close to 50%, which is a random guess. The intuition here is as follows. θ_p is trained on noisy soft-labels of θ_{up} on X_{ref} , and therefore, compared to an arbitrary test data, the influence of X_{ref} on θ_p is not unique, which membership inference attacks exploit (Long et al. 2018; Shokri et al. 2017; Salem et al. 2019; Nasr, Shokri, and Houmansadr 2019). For Purchase-100 and

Dataset	Test	Reference data	Δ.	4
& model	acc. (A_{test})	acc. (A_{ref})	A_{wb}	A_{bb}
P-FC	74.1	80.8	53.1	52.6
T-FC	48.6	52.0	52.2	52.0
C100-A	35.7	35.9	50.9	50.5
C100-D12	63.1	65.1	53.0	52.2
C10-A	65.0	66.7	53.9	52.7

Table 10: DMP does not pose membership inference risk to the possibly sensitive reference data. A_{ref} and A_{test} are accuracies of protected model, θ_{p} , on X_{ref} and D_{test} , respectively.

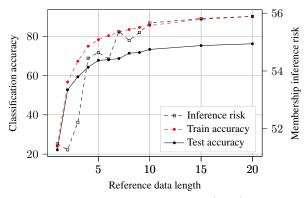


Figure 4: Increasing reference data size, $|X_{\text{ref}}|$, increases accuracy of θ_{p} , but also increases \mathcal{R} in (11), which increases the membership inference risk due to θ_{p} .

Texas-100, the accuracy of $\theta_{\rm p}$ on $X_{\rm ref}$ is much higher than on $D_{\rm test}$, because for these datasets, $X_{\rm ref}$ contains easy-to-classify samples.

E Statistical Indistinguishability due to DMP

In this section, we show the indistinguishability of the statistics of different features of the target models trained with and without defenses, on the members and non-members of their training data. Such indistinguishability is necenssary to hinder membership inference attacks (MIAs) (Shokri et al. 2017).

Effect of softmax temperature. Figure 5 shows the effect of softmax temperature, T, of unprotected model, θ_{up} , on the training and test accuracies of the protected mode, θ_p . As expected, we observe in Figure 5 that with the increase in the softmax temperature of θ_{up} , the generalization error of θ_p decreases. From left to right, the generalization errors of θ_p when the softmax temperatures of θ_{up} are set at 2, 4, and 6 are 4.7% (66.3, 61.6), 3.6% (66.7, 63.1), and 0.8% (55.7, 54.9), respectively; parentheses show the corresponding training and test accuracies, respectively. We keep the temperature of softmax layer in θ_p constant at 4.0. This reduction in generalization error improves membership privacy.

Indistinguishability of gradient norms. To assess the efficacy of DMP against the stronger whitebox MIAs (Nasr et al. (2019)), we study the gradients of loss of the predictions of unprotected and protected models on members and non-

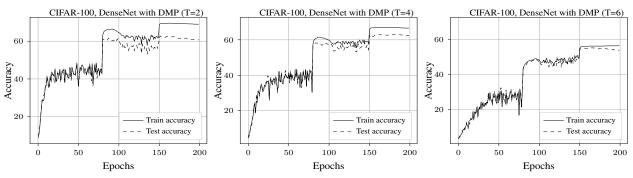


Figure 5: Impact of softmax temperature on training of θ_p : Increase in the temperature of softmax layer of θ_{up} reduces $\Delta \mathcal{L}_{KL}$ in (13), and hence, the ratio \mathcal{R} in (11). This improves the membership privacy and generalization of θ_p .

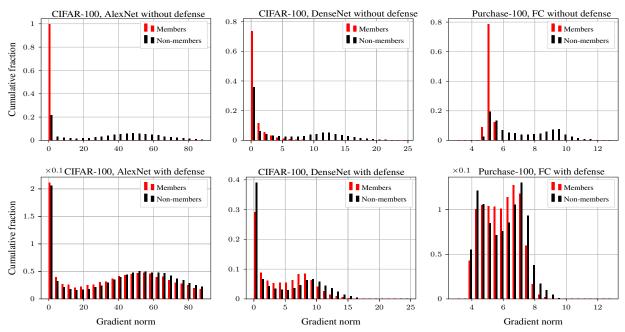


Figure 6: Distributions of gradient norms of members and non-members of private training data. (*Upper row*): Unlike the distribution of non-members, that of the members of the unprotected model, θ_{up} , is skewed towards 0 as θ_{up} memorizes the members. (*Lower row*): The distributions of gradient norms for members and non-members for the protected model, θ_p , of DMP are almost indistinguishable.

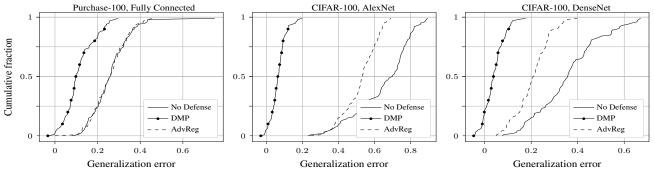


Figure 7: The empirical CDF of the generalization error of models trained with DMP, adversarial regularization (AdvReg), and without defense. The y-axis is the fraction of classes that have generalization error less than the values on x-axis. The generalization error reduction due to DMP is much larger ($10 \times$ for CIFAR100 and $2 \times$ for Purchase) than due to AdvReg. The low generalization error improves membership privacy due to DMP.

	Experimental se	etup		Near	-equal A	A _{test} as I	OMP	
Dataset	Model	Regularization	$E_{\rm gen}$	A_{test}	$A_{\sf wb}$	A_{bb}	A_{bl}	A_{nn}
		WD	21.7	78.1	69.7	70.1	60.9	55.6
Purchase	FC	WD + DR	22.1	77.4	77.1	76.8	61.5	60.0
Pulchase	rc	WD + LS	21.1	78.4	76.5	76.8	60.6	56.4
		WD + CP	22.9	76.9	70.1	70.5	61.5	58.5
		WD	49.0	50.4	84.1	82.1	74.5	56.2
Texas	FC	WD + DR	41.1	52.1	82.1	81.2	70.6	60.2
Texas	FC	WD + LS	50.9	49.1	86.0	85.7	75.5	56.9
		WD + CP	45.5	54.2	90.4	90.2	72.8	65.6
		WD	31.0	67.8	72.9	72.9	65.5	N/A
CIFAR100	DenseNet12	WD + DR	31.0	68.2	73.7	73.6	65.5	N/A
CITAKIOO	Deliserretiz	WD + LS	31.6	68.0	70.3	70.1	65.8	N/A
		WD + CP	31.1	67.5	74.3	74.7	65.6	N/A
		WD	31.0	68.9	73.2	73.3	65.5	N/A
CIFAR10	AlexNet	WD + DR	30.6	69.4	73.8	73.4	65.3	N/A
CHARIO	AICAINEL	WD + LS	29.9	69.9	74.8	75.0	65.5	N/A
		WD + CP	29.9	70.0	70.6	71.1	65.5	N/A

Table 11: Generalization error (E_{gen}), test accuracy (A_{test}), and various MIA risks (evaluated using MIAs from Section 5.2) of models trained using state-of-the-art regularization techniques. Here we provide MIA risks for regularized models whose accuracy is close to that of DMP-trained models. We note that, for the same test accuracy, DMP-trained models provide significantly higher resistance to MIAs.

members of the private training data, D_{tr} . Figure 6 shows the fraction of members and non-members given on y-axes that fall in a particular range of gradient norm values given on x-axes. Gradients are computed with respect to the parameters of the given model. We note that the distribution of the norms of unprotected model (upper figures) is heavily skewed to the left for the members, i.e., towards lower gradient norm values, unlike that for the non-members. This is because, θ_{up} memorizes D_{tr} and its loss and the gradient of the loss on the members is very small compared to the nonmembers. However, for the protected model both members and non-members are evenly distributed across a large range of gradient norm values. This implies that *DMP significantly* reduces the unintended memorization of D_{tr} in the model parameters. Hence, DMP significant reduces (by 27.6%) the MIA risk to the large capacity Dense19.

Indistinguishability of train and test accuracies. In Figure 7, we show the cumulative fraction of classes on y-axis for which the generalization error of the target models is lesser than the corresponding value on the x-axis; the closer the line to the line x=0, the lower the generalization error. Figure 7 implies that, the models trained using DMP have significantly lower generalization error than those trained using adversarially regularization or without defense. We observe that, with the no defense case as the baseline, the generalization error reduction using DMP is more than twice that using adversarial regularization. DMP reduces the error by half for Purchase and by $10\times$ for CI-FAR100.

F Missing experimental details

Best tradeoffs due to adversarial regularization Table 12 gives the results for best tradeoffs due to adversarial regularization that we obtain by tuning its λ parameter (Nasr, Shokri, and Houmansadr 2018).

Dataset		Adversarial regularization								
& model	E_{gen}	A_{test}	A_{wb}	A_{bb}	A_{bl}	A_{nn}				
P-FC	22.4	68.1	62.3	61.9	61.4	51.4				
T-FC	15.5	45.3	66.8	66.3	57.8	51.2				
C100-A	50.9	31.6	79.3	78.3	75.5	N/A				
C100-D12	19.4	58.4	61.9	61.7	59.7	N/A				
C100-D19	30.8	53.7	69.5	68.7	65.4	N/A				
C10-A	29.8	62.6	65.2	65.0	64.9	N/A				

Table 12: Best tradeoffs between test accuracy (A_{test}) and membership inference risks (evaluated using MIAs from Section 5.2) due to adversarial regularization. DMP significantly improves the tradeoffs over the adversarial regularization (results for DMP are in Table 3).

Best tradeoffs due to other regularizations We see from the 'Equivalent A_{test} ' column in Table 11 that all regularization techniques improve the classification performance over the corresponding accuracies of baseline models from the Table 2 of main paper. However, they reduce overfitting negligibly: the maximum reduction in $E_{\rm gen}$ due to the regularizations is 1.8% for Purchase, 10.2% for Texas, 3.8% for CI-FAR 100, and 2.6% for CIFAR 10. This is because these techniques aim to produce models that generalize better to test data, but they do not necessarily reduce the memorization of the private training data by the models. Consequently, these techniques fail to reduce the membership inference risk: the maximum reduction in A_{Wb} due to the regularizations is 7% for Purchase, 1.9% for Texas, 1.9% for CIFAR 100, and 6.8% for CIFAR10. Note that, the confidence penalty and the label smoothing techniques reduce the inference risk, but not the generalization error. This is because the corresponding models have smoother output distributions, which are more indistinguishable than the output distributions of models without any privacy.