

Machine Unlearning of Features and Labels

Alexander Warnecke¹, Lukas Pirch¹,
Christian Wressnegger², Konrad Rieck¹

¹ Technische Universität Braunschweig, Germany

² Karlsruhe Institute of Technology, Germany

Abstract

Removing information from a machine learning model is a non-trivial task that requires to partially revert the training process. This task is unavoidable when sensitive data, such as credit card numbers or passwords, accidentally enter the model and need to be removed afterwards. Recently, different concepts for machine unlearning have been proposed to address this problem. While these approaches are effective in removing individual data points, they do not scale to scenarios where larger groups of features and labels need to be reverted.

In this paper, we propose a method for unlearning features and labels. Our approach builds on the concept of influence functions and realizes unlearning through closed-form updates of model parameters. It enables to adapt the influence of training data on a learning model retrospectively, thereby correcting data leaks and privacy issues. For learning models with strongly convex loss functions, our method provides certified unlearning with theoretical guarantees. For models with non-convex losses, we empirically show that unlearning features and labels is effective and significantly faster than other strategies.

1 Introduction

Machine learning has become an ubiquitous tool in analyzing personal data and developing data-driven services. Unfortunately, the underlying learning models can pose a serious threat to privacy if they inadvertently reveal sensitive information from the training data. For example, Carlini et al. [12] show that the Google text completion system contains credit card and social security numbers from personal emails, which may be exposed to users during the autocompletion of text. Once such sensitive data has entered a learning model, however, its removal is non-trivial and requires to selectively revert the learning process. In absence of specific methods for this task in the past, retraining from scratch has been the only resort, which is costly and only possible if the original training data is still available.

As a remedy, Cao & Yang [11] and Bourtole et al. [8] propose methods for *machine unlearning*. These methods decompose the learning process and are capable of removing individual data points from a learning model in retrospection. As a result, they enable to eliminate isolated privacy issues, such as data points associated with individuals. However, information leaks may not only manifest in single data instances but also in groups of *features* and *labels*. A leaked address of a celebrity might be shared in hundreds of social media posts, affecting large parts of the training data. Similarly, relevant features in a bag-of-words model may be associated with sensitive names and data, contaminating the entire feature space.

Unfortunately, instance-based unlearning as proposed in previous work is inefficient in these cases: First, a runtime improvement can hardly be obtained over retraining as the leaks are not isolated and larger parts of the training data need to be removed. Second, omitting several data points will inevitably reduce the fidelity of the corrected learning model. It becomes clear that the task of unlearning is not necessarily confined to removing data points, but may also require corrections on the orthogonal layers of features and labels, regardless of the amount of affected training data.

In this paper, we propose a method for unlearning features and labels. Our approach is inspired by the concept of *influence functions*, a technique from robust statistics [31], that allows for estimating the influence of data on learning models [33, 34]. By reformulating this influence estimation as a form of unlearning, we derive a versatile approach that maps changes of the training data in retrospection to closed-form updates of the model parameters. These updates can be calculated efficiently, even if larger parts of the training data are affected, and enable the removal of features and labels. As a result, our method can correct privacy leaks in a wide range of learning models with convex and non-convex loss functions.

For models with strongly convex loss, such as logistic regression and support vector machines, we prove that our approach enables *certified unlearning*. That is, it provides theoretical guarantees on the removal of features and labels from the models. To obtain these guarantees, we extend the concept of certified data removal [28] and show that the difference between models obtained with our approach and retraining from scratch become arbitrarily small. Consequently, we can define an upper bound on this difference and thereby realize provable unlearning in practice.

For models with non-convex loss functions, such as deep neural networks, similar theoretical guarantees do not hold in general. However, we empirically demonstrate that our approach provides substantial advantages over prior work. Our method is significantly faster in comparison to sharding [8, 24] and retraining while removing data and preserving a similar level of accuracy. Moreover, due to the compact updates, our approach requires only a fraction of the training data and hence is applicable when the original data is not entirely available. We show the efficacy of our approach in case studies on removing privacy leaks in spam classification and unintended memorization in natural language processing.

Contributions. In summary, we make the following major contributions:

1. *Unlearning with closed-form updates.* We introduce a novel framework for unlearning of features and labels. This framework builds on closed-form updates of learning models and thus is significantly faster than instance-based approaches to unlearning.
2. *Certified unlearning.* We derive two unlearning strategies for our framework based on first-order and second-order gradient updates. Under convexity and continuity assumptions on the loss, we show that both strategies can provide certified unlearning.
3. *Empirical analysis.* We empirically show that unlearning of sensible information is possible even for deep neural networks with non-convex loss functions. We find that our first-order update is extremely efficient, enabling a speed-up over retraining by up to three orders of magnitude.

The rest of the paper is structured as follows: We review related work on machine unlearning and influence functions in Section 2. Our approach and its technical realization are introduced in Sections 3 and 4, respectively. The theoretical analysis of our approach is presented in Section 5 and its empirical evaluation in Section 6. Finally, we discuss limitations in Section 7 and conclude the paper in Section 8.

2 Related Work

The increasing application of machine learning to personal data has started a series of research on detecting and correcting privacy issues in learning models [e.g., 12, 13, 36, 45, 47, 53]. In the following, we provide an overview of work on machine unlearning and influence functions. A broader discussion of privacy and machine learning is given by De Cristofaro [21] and Papernot et al. [41].

Machine unlearning. Methods for unlearning sensitive data are a recent branch of security research. Earlier, the efficient removal of samples was also called *decremental learning* [14] and used to speed up cross validation for various linear classifiers [15–17]. Cao & Yang [11] show that a large number of learning models can be represented in a closed summation form that allows for elegantly removing individual data points in retrospection. However, for adaptive learning strategies, such as stochastic gradient descent, this approach provides only little advantage over retraining from scratch and thus is not well suited for correcting problems in neural networks.

As a remedy, Bourtole et al. [8] propose a universal strategy for unlearning data points from classification models. Similarly, Ginart et al. [24] develop a technique for unlearning points in clustering. The key idea of both approaches is to split the data into independent partitions—so called *shards*—and aggregate the final model from submodels trained over these shards. In this setting, the

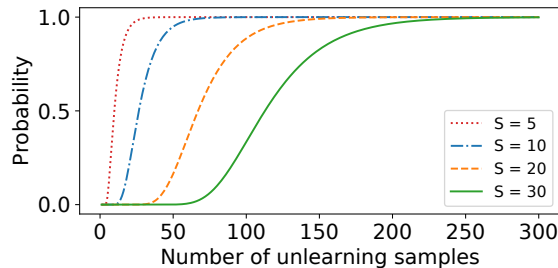


Figure 1: Probability of all shards being affected when unlearning for varying number of data points and shards (S).

unlearning of data points can be efficiently carried out by only retraining the affected submodels. Aldaghri et al. [3] show that this approach can be further sped up for least-squares regression by choosing the shards cleverly. Unlearning based on shards, however, is suitable for removing a few data points only and inevitably deteriorates in performance when larger portions of the data require changes.

This limitation of sharding is schematically illustrated in Fig. 1. The probability that all shards need to be retrained increases with the number of data points to be corrected. For a practical setup with 20 shards, as proposed by Bourtole et al. [8], changes to as few as 150 points are already sufficient to impact all shards and render this form of unlearning inefficient, regardless of the size of the training data. We provide a detailed analysis of this limitation in Section 8. Consequently, privacy leaks involving hundreds or thousands of data points cannot be addressed with these approaches.

Influence functions. The concept of influence functions that forms the basis of our approach originates from robust statistics [31] and has first been used by Cook & Weisberg [20] for investigating the changes of simple linear regression models. Although the proposed techniques have been occasionally employed in machine learning [32, 35], the seminal work of Koh & Liang [33] recently brought general attention to this concept and its application to modern learning techniques. In particular, this work uses influence functions for explaining the impact of data points on the predictions of learning models.

Influence functions have then been used to trace bias in word embeddings back to documents [10, 19], determine reliable regions in learning models [46], and explain deep neural networks [7]. Moreover, Basu et al. [6] increase the accuracy of influence functions by using high-order approximations, Barshan et al. [5] improve the precision of influence calculations through nearest-neighbor strategies, and Guo et al. [29] show that the runtime can be decreased when only specific samples are considered. Golatkar et al. [25, 26] use influence functions for sample removal in deep neural networks by proposing special approximations of the learning model.

In terms of theoretical analysis, Koh et al. [34] study the accuracy of influence functions when estimating the loss on test data and Neel et al. [40] perform a similar analysis for gradient based update strategies. Rad & Maleki [44] further show that the prediction error on leave-one-out validations can be reduced with influence functions. Finally, Guo et al. [28] introduce the idea of certified removal for data points that we extend in our approach.

All of these approaches, however, remain on the level of data instances. To our knowledge, we are the first to build on the concept of influence functions for unlearning features and labels from learning models.

3 Unlearning with Updates

Let us start by considering a supervised learning task that is described by a dataset $D = \{z_1, \dots, z_n\}$ with each object $z_i = (x, y)$ consisting of a data point $x \in \mathcal{X}$ and a label $y \in \mathcal{Y}$. We assume that $\mathcal{X} = \mathbb{R}^d$ is a vector space and denote the j -th feature (dimension) of x by $x[j]$. Given a loss function $\ell(z, \theta)$ that measures the difference between the predictions of a learning model θ and the true labels, the optimal model θ^* can be found by minimizing the regularized empirical risk,

$$\theta^* = \operatorname{argmin}_{\theta} L(\theta; D) = \operatorname{argmin}_{\theta} \sum_{i=1}^n \ell(z_i, \theta) + \lambda \Omega(\theta) \quad (1)$$

where Ω is a regularizer and $L(\theta; D)$ describes the loss on the entire dataset. In this setup, the process of unlearning amounts to adapting θ^* to changes in D without recalculating the optimization problem in Eq. (1).

3.1 Unlearning Data Points

To provide an intuition for our approach, we begin by asking the following question: How would the optimal learning model θ^* change, if only one data point z had been perturbed by some change δ ? Replacing z by $\tilde{z} = (x + \delta, y)$ leads to the new optimal set of model parameters:

$$\theta_{z \rightarrow \tilde{z}}^* = \operatorname{argmin}_{\theta} L(\theta, D) + \ell(\tilde{z}, \theta) - \ell(z, \theta). \quad (2)$$

However, calculating the new model $\theta_{z \rightarrow \tilde{z}}^*$ exactly is expensive. Instead of replacing the data point z with \tilde{z} , we can also up-weight \tilde{z} by a small value ϵ and down-weight z accordingly, resulting in the following optimization problem:

$$\theta_{\epsilon, z \rightarrow \tilde{z}}^* = \operatorname{argmin}_{\theta} L(\theta, D) + \epsilon \ell(\tilde{z}, \theta) - \epsilon \ell(z, \theta). \quad (3)$$

Eqs. (2) and (3) are equivalent for $\epsilon = 1$ and solve the same problem. As a result, we do not need to explicitly remove a data point from the training data but can revert its *influence* on the learning model through a combination of appropriate up-weighting and down-weighting.

It is easy to see that this approach is not restricted to a single data point. We can simply define a set of data points Z and its perturbed versions \tilde{Z} , and arrive at the weighting

$$\theta_{\epsilon, Z \rightarrow \tilde{Z}}^* = \underset{\theta}{\operatorname{argmin}} L(\theta, D) + \epsilon \sum_{\tilde{z} \in \tilde{Z}} \ell(\tilde{z}, \theta) - \epsilon \sum_{z \in Z} \ell(z, \theta). \quad (4)$$

This generalization enables us to approximate changes on larger portions of the training data. Instead of solving the problem in Eq. (4), however, we formulate this optimization as an update of the original model θ^* . That is, we seek a closed-form update $\Delta(Z, \tilde{Z})$ of the model parameters, such that

$$\theta_{\epsilon, Z \rightarrow \tilde{Z}}^* \approx \theta^* + \Delta(Z, \tilde{Z}), \quad (5)$$

where $\Delta(Z, \tilde{Z})$ has the same dimension as the learning model θ but is sparse if only a few parameters are affected.

As a result of this formulation, we can describe changes of the training data as a compact update Δ rather than iteratively solving an optimization problem. We show in Section 4 that this update step can be efficiently computed using first-order and second-order gradients. Furthermore, we prove in Section 4 that the unlearning success of both updates can be certified up to a tolerance ϵ if the loss function ℓ is strictly convex, twice differentiable, and Lipschitz-continuous.

3.2 Unlearning Features and Labels

Equipped with a general method for updating a learning model, we proceed to introduce our approach for unlearning features and labels. To this end, we expand our notion of perturbations and include changes to labels by defining

$$\tilde{z} = (x + \delta_x, y + \delta_y),$$

where δ_x modifies the features of a data point and δ_y its label. By using different changes in the perturbations \tilde{Z} , we can now realize different types of unlearning using closed-form updates.

Replacing features. As the first type of unlearning, we consider the task of correcting features in a learning model. This task is relevant if the content of some features violates the privacy of a user and needs to be replaced with alternative values. As an example, personal names, identification numbers, residence addresses, or other sensitive data might need to be removed after a model has been trained on a corpus of emails.

For a set of features F and their new values V , we define perturbations on the affected points Z by

$$\tilde{Z} = \{(x[f] = v, y) : (x, y) \in Z, (f, v) \in F \times V\}.$$

For example, a credit card number contained in the training data can be blinded by a random number sequence in this setting. The values V can be adapted individually, such that fine-grained corrections become possible.

Replacing labels. As the second type of unlearning, we focus on correcting labels. This form of unlearning is necessary if the labels captured in a model contain unwanted information. For example, in generative language models, the training text is used as input features (preceding characters) *and* labels (target characters) [27, 48]. Hence, defects can only be eliminated if the labels are unlearned as well.

For the affected points Z and the set of new labels Y , we define the corresponding perturbations by

$$\tilde{Z} = \{(x, y) \in Z_x \times Y\},$$

where Z_x corresponds to the data points in Z without their original labels. The new labels Y can be individually selected for each data point, as long as they come from the domain \mathcal{Y} , that is, $Y \subset \mathcal{Y}$. Note that the replaced labels and features can be easily combined in one set of perturbations \tilde{Z} , so that defects affecting both can be corrected in a single update. In Section 6.2, we demonstrate that this combination can be used to remove unintended memorization from generative language models with high efficiency.

Revoking features. Based on appropriate definitions of Z and \tilde{Z} , our approach enables to replace the content of features and thus eliminate privacy leaks. However, in some scenarios it might be necessary to completely remove features from a learning model—a task that we denote as *revocation*. In contrast to the correction of features, this form of unlearning poses a unique challenge: The revocation of features reduces the input dimension of the learning model. While this adjustment can be easily carried out through retraining with adapted data, constructing a model update as in Eq. (5) is tricky.

To address this problem, let us consider a model θ^* trained on a dataset $D \subset \mathbb{R}^d$. If we remove the features F from this dataset and train the model again, we obtain a new optimal model θ_{-F}^* with reduced input dimension. By contrast, if we set the values of the features F to zero in the dataset and train again, we obtain an optimal model $\theta_{F=0}^*$ with the same input dimension as θ^* . Fortunately, these two models are equivalent for a large class of learning models, including support vector machines and several neural networks as the following lemma shows.

Lemma 1. *For learning models processing inputs x using linear transformations of the form $\theta^T x$, we have $\theta_{-F}^* \equiv \theta_{F=0}^*$.*

Proof. It is easy to see that it is irrelevant for the dot product $\theta^T x$ whether a dimension of x is missing or equals zero in the linear transformation

$$\sum_{k:k \notin F} \theta[k]x[k] = \sum_k \theta[k]\mathbf{1}\{k \notin F\}x[k].$$

As a result, the loss $\ell(z, \theta) = \ell(\theta^T x, y, \theta)$ of both models is identical for every data point z . Hence, $L(\theta; D)$ is also equal for both models and thus the same objective is minimized during learning resulting in equal parameters. \square

Lemma 1 enables us to erase features from many learning models by first setting them to zero, calculating the parameter update, and then reducing the dimension of the models accordingly. Concretely, to revoke the features F , we locate the data points where these features are non-zero with

$$Z = \{(x, y) \in D : \exists f : x[f] \neq 0, f \in F\}$$

and construct corresponding perturbations such that the features are set to zero by unlearning,

$$\tilde{Z} = \{(x[f] = 0, y) : (x, y) \in Z, f \in F\}.$$

Revoking labels. The previous strategy allows revoking features from several learning models. It is crucial if, for example, a bag-of-words model has captured sensitive data in relevant features and therefore a reduction of the input dimension during unlearning is unavoidable. Unfortunately, a similar strategy for the revocation of labels is not available for our method, as we are not aware of a general shortcut, such as Lemma 1. Still, if the learning model contains explicit output dimensions for the class labels, as with some neural network architectures, it is possible to first replace unwanted labels and then manually remove the corresponding dimensions.

4 Update Steps for Unlearning

Our approach rests on changing the influence of training data with a closed-form update of the model parameters, as shown in Eq. (5). In the following, we derive two strategies for calculating this closed form: a *first-order update* and a *second-order update*. The first strategy builds on the gradient of the loss function and thus can be applied to any model with a differentiable loss. The second strategy also incorporates second-order derivatives which limits the application to loss functions with an invertible Hessian matrix.

4.1 First-Order Update

Recall that we aim to find an update $\Delta(Z, \tilde{Z})$ that we add to our model θ^* . If the loss ℓ is differentiable, we can find the optimal *first-order update* by

$$\Delta(Z, \tilde{Z}) = -\tau \left(\sum_{\tilde{z} \in \tilde{Z}} \nabla_{\theta} \ell(\tilde{z}, \theta^*) - \sum_{z \in Z} \nabla_{\theta} \ell(z, \theta^*) \right) \quad (6)$$

where τ is a small constant that we refer to as *unlearning rate*. A complete derivation of Eq. (6) is given in Section 8. Intuitively, this update shifts the model parameters in the direction from $\sum_{z \in Z} \nabla_{\theta} \ell(z, \theta^*)$ to $\sum_{\tilde{z} \in \tilde{Z}} \nabla_{\theta} \ell(\tilde{z}, \theta^*)$ where

the size of the update step is determined by the rate τ . This update strategy is related to the classic gradient descent update GD used in many learning algorithms and given by

$$\text{GD}(\tilde{Z}) = -\tau \sum_{\tilde{z} \in \tilde{Z}} \nabla_{\theta} \ell(\tilde{z}, \theta^*).$$

However, it differs from this update step in that it moves the model to the *difference* in gradient between the original and perturbed data, which minimizes the loss on \tilde{z} and at the same time removes the information contained in z .

The first-order update is a simple and yet effective strategy: Gradients of ℓ can be computed in $\mathcal{O}(d)$ [43] and modern auto-differentiation frameworks like TensorFlow [1] and PyTorch [42] offer easy gradient computations for the practitioner. The update step involves a parameter τ that controls the impact of the unlearning step. To ensure that data has been completely replaced, it is necessary to calibrate this parameter using a measure for the success of unlearning. In Section 6, for instance, we show how the exposure metric by Carlini et al. [12] can be used for this calibration.

4.2 Second-Order Update

The calibration of the update step can be eliminated if we make further assumptions on the properties of the loss function ℓ . If we assume that ℓ is twice differentiable and strictly convex, the influence of a single data point can be approximated in closed form [20] by

$$\left. \frac{\partial \theta_{\epsilon, z \rightarrow \tilde{z}}^*}{\partial \epsilon} \right|_{\epsilon=0} = -H_{\theta^*}^{-1} (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)),$$

where $H_{\theta^*}^{-1}$ is the inverse Hessian of the loss at θ^* , that is, the inverse matrix of the second-order partial derivatives. We can now perform a linear approximation for $\theta_{z \rightarrow \tilde{z}}^*$ to obtain

$$\theta_{z \rightarrow \tilde{z}}^* \approx \theta^* - H_{\theta^*}^{-1} (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)). \quad (7)$$

Since all operations are linear, we can easily extend Eq. (7) to account for multiple data points and derive the following *second-order update*:

$$\Delta(Z, \tilde{Z}) = -H_{\theta^*}^{-1} \left(\sum_{\tilde{z} \in \tilde{Z}} \nabla_{\theta} \ell(\tilde{z}, \theta^*) - \sum_{z \in Z} \nabla_{\theta} \ell(z, \theta^*) \right). \quad (8)$$

A full derivation of this update step is provided in Section 8. Note that the update does not require any parameter calibration, since the parameter weighting of the changes is directly derived from the inverse Hessian of the loss function.

The second-order update is the preferred strategy for unlearning on models with a strongly convex and twice differentiable loss function, such as a logistic regression, that guarantee the existence of H^{-1} . Technically, the update step in Eq. (8) can be easily calculated with common machine-learning frameworks. In contrast to the first-order update, however, this computation involves the inverse Hessian matrix, which is non-trivial for neural networks, for example.

Computing the inverse Hessian. Given a model $\theta \in \mathbb{R}^p$ with p parameters, forming and inverting the Hessian requires $\mathcal{O}(np^2 + p^3)$ time and $\mathcal{O}(p^2)$ space [33]. For models with a small number of parameters, the matrix can be pre-computed and explicitly stored, such that each subsequent request for unlearning only involves a simple matrix-vector multiplication. For example, in Section 6.1, we show that unlearning features from a logistic regression model with about 5,000 parameters can be realized with this approach in less than a second.

For complex learning models, such as deep neural networks, the Hessian matrix quickly becomes too large for explicit storage. Moreover, these models typically do not have convex loss functions, such that the matrix may also be non-invertible, rendering an exact update impossible. Nevertheless, we can approximate the inverse Hessian using techniques proposed by Koh & Liang [33]. While this approximation weakens the theoretical guarantees of the unlearning process, it enables applying second-order updates to a variety of complex learning models, similar to the first-order strategy.

To apply second-order updates in practice, we have to avoid storing H explicitly and still be able to compute $H^{-1}v$. To this end, we rely on the scheme proposed by Agarwal et al. [2] to compute expressions of the form $H^{-1}v$ that only require to calculate Hv and do not need to store H^{-1} . *Hessian-Vector-Products* (HVPs) allow us to calculate Hv efficiently by making use of the linearity of the gradient

$$Hv = \nabla_{\theta}^2 L(\theta^*; D)v = \nabla_{\theta}(\nabla_{\theta} L(\theta^*; D)v).$$

Denoting the first j terms of the Taylor expansion of H^{-1} by H_j^{-1} we have $H_j^{-1} = \sum_{i=0}^j (I - H)^i$, and can recursively define an approximation given by $H_j^{-1} = I + (I - H)H_{j-1}^{-1}$. If $|\lambda_i| < 1$ for all eigenvalues λ_i of H , we have $H_j^{-1} \rightarrow H^{-1}$ for $j \rightarrow \infty$. To ensure this convergence, we add a small damping term λ to the diagonal of H and scale down the loss function (and thereby the eigenvalues) by some constant which does not change the optimal parameters θ^* . Under these assumptions, we can formulate the following algorithm for computing an approximation of $H^{-1}v$: Given data points z_1, \dots, z_t sampled from D , we define the iterative updates

$$\begin{aligned} \tilde{H}_0^{-1}v &= v, \\ \tilde{H}_j^{-1}v &= v + (I - \nabla_{\theta}^2 L(z_i, \theta^*))\tilde{H}_{j-1}^{-1}v. \end{aligned}$$

In each update step, H is estimated using a single data point and we can use HVPs to evaluate $\nabla_{\theta}^2 L(z_i, \theta^*)\tilde{H}_{j-1}^{-1}v$ efficiently in $\mathcal{O}(d)$ as demonstrated by Pearlmutter [43]. Using batches of data points instead of single ones and averaging the results further speeds up the approximation. Choosing t large enough so that the updates converge and averaging r runs to reduce the variance of the results, we obtain $\tilde{H}_t^{-1}v$ as our final estimate of $H^{-1}v$ in $\mathcal{O}(rtd)$ of time. In Section 6.2 we demonstrate that this strategy can be used to calculate the second-order updates for a deep neural network with 3.3 million parameters.

5 Certified Unlearning

Machine unlearning is a delicate task, as it aims at reliably removing privacy issues and sensitive data from learning models. This task should ideally build on theoretical guarantees to enable *certified unlearning*, where the corrected model is stochastically indistinguishable from one created by retraining. In the following, we derive conditions under which the updates of our approach introduced in Section 4.2 provide certified unlearning. To this end, we build on the concepts of *differential privacy* [18, 22] and *certified data removal* [28], and adapt them to the unlearning problem.

Let us first briefly recall the idea of differential privacy in machine learning: For a training dataset D , let \mathcal{A} be a learning algorithm that outputs a model $\theta \in \Theta$ after training on D , that is, $\mathcal{A} : D \rightarrow \Theta$. Randomness in \mathcal{A} induces a probability distribution over the output models in Θ . The key idea of differential privacy is a measure of difference between a model trained on D and another one trained on $D \setminus z$ for some $z \in D$.

Definition 1. *Given some $\epsilon > 0$, a learning algorithm \mathcal{A} is said to be ϵ -differentially private (ϵ -DP) if*

$$e^{-\epsilon} \leq \frac{P(\mathcal{A}(D) \in \mathcal{T})}{P(\mathcal{A}(D \setminus z) \in \mathcal{T})} \leq e^{\epsilon}$$

holds for all $\mathcal{T} \subset \Theta, D$, and $z \in D$.

Thus, for an ϵ -DP learning algorithm the difference between the log-likelihood of a model trained on D and one trained on $D \setminus z$ is smaller than ϵ for all possible models, datasets, and data points. Based on this definition, we can introduce the concept of ϵ -certified unlearning. In particular, we consider an unlearning method \mathcal{U} that maps a model θ to a corrected model $\theta_{\mathcal{U}} = \mathcal{U}(\theta, D, D')$ where D' denotes the dataset containing the perturbations \tilde{Z} required for the unlearning task.

Definition 2. *Given some $\epsilon > 0$ and a learning algorithm \mathcal{A} , an unlearning method \mathcal{U} is ϵ -certified if*

$$e^{-\epsilon} \leq \frac{P(\mathcal{U}(\mathcal{A}(D), D, D') \in \mathcal{T})}{P(\mathcal{A}(D') \in \mathcal{T})} \leq e^{\epsilon}$$

holds for all $\mathcal{T} \subset \Theta, D$, and D' .

This definition ensures that the probability to obtain a model using the unlearning method \mathcal{U} and training a new model on D' from scratch deviates at most by ϵ . Similar to certified data removal [28], we introduce (ϵ, δ) -certified unlearning, a relaxed version of ϵ -certified unlearning, defined as follows.

Definition 3. Under the assumptions of Definition 2, an unlearning method \mathcal{U} is (ϵ, δ) -certified if

$$P(\mathcal{U}(\mathcal{A}(D), D, D') \in \mathcal{T}) \leq e^\epsilon P(\mathcal{A}(D') \in \mathcal{T}) + \delta$$

and

$$P(\mathcal{A}(D') \in \mathcal{T}) \leq e^\epsilon P(\mathcal{U}(\mathcal{A}(D), D, D') \in \mathcal{T}) + \delta$$

hold for all $\mathcal{T} \subset \Theta, D$, and D' .

That is, (ϵ, δ) -certified unlearning allows the method \mathcal{U} to slightly violate the conditions from Definition 2 by a constant δ . Using the above definitions, it becomes possible to derive conditions under which certified unlearning is possible for both our approximate update strategies.

5.1 Certified Unlearning of Features and Labels

Based on the concept of certified unlearning, we analyze our approach and its theoretical guarantees on removing features and labels. To ease this analysis, we make two assumptions on the employed learning algorithm: First, we assume that the loss function ℓ is twice differentiable and strictly convex such that H^{-1} always exists. Second, we consider L_2 regularization in optimization problem (1), that is, $\Omega(\theta) = \frac{1}{2}\|\theta\|_2^2$ which ensures that the loss function is strongly convex.

A powerful concept for analyzing unlearning is the *gradient residual* $\nabla L(\theta; D')$ for a given model θ and a corrected dataset D' . For strongly convex loss functions, the gradient residual is zero *if and only if* θ equals $\mathcal{A}(D')$ since in this case the optimum is unique. Therefore, the norm of the gradient residual $\|\nabla L(\theta; D')\|_2$ reflects the distance of a model θ from one obtained by retraining on the corrected dataset D' . While a small value of this norm is not sufficient to judge the quality of unlearning, we can develop upper bounds to prove properties related to differential privacy [18, 28]. Consequently, we derive bounds for the gradient residual norms of our two update strategies. The corresponding proofs are given in Section 8.

Theorem 1. *If all perturbations δ lie within a radius R , that is $\|\delta\|_2 \leq R$, and the loss $\nabla \ell(z, \theta)$ is (γ_z, γ) -Lipschitz with respect to z and θ , the following upper bounds hold:*

1. *If the unlearning rate $\tau \leq \frac{1}{\gamma_n}$, we have*

$$\|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 \leq 2R\gamma_z|Z|$$

for the first-order update of our approach.

2. *If $\nabla^2 \ell(z, \theta)$ is γ'' -Lipschitz with respect to θ , we have*

$$\|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 \leq \frac{\gamma_z^2 \gamma'' R^2}{\lambda^2} |Z|^2$$

for the second-order update of our approach.

This theorem enables us to bound the gradient residual norm of both update steps. We leverage these bounds to reduce the difference between unlearning and retraining from scratch. In particular, we follow the approach by Chaudhuri et al. [18] and add a random linear term to the loss function to shape the distribution of the model parameters. Given a vector $b \in \mathbb{R}^d$ drawn from a random distribution, we define

$$L_b(\theta; D) = \sum_{z \in D} \ell(z, \theta) + \frac{\lambda}{2} \|\theta\|_2^2 + b^T \theta$$

with a corresponding gradient residual r given by

$$r = \nabla L_b(\theta; D') = \sum_{z \in D'} \nabla \ell(z, \theta) + \lambda \theta + b.$$

By definition, the gradient residual r of L_b differs only by the added vector b from the residual of the original loss L , which allows to precisely determine its influence on the bounds of Theorem 1 depending on the underlying distribution of b .

Let $\mathcal{A}(D')$ be an exact minimizer of L_b on D' with density function $f_{\mathcal{A}}$ and $\mathcal{U}(\mathcal{A}(D), D, D')$ an approximated minimum obtained through unlearning with density $f_{\mathcal{U}}$. Guo et al. [28] show that the max-divergence between $f_{\mathcal{A}}$ and $f_{\mathcal{U}}$ for the model θ produced by \mathcal{U} can be bounded using the following theorem.

Theorem 2 (Guo et al. [28]). *Let \mathcal{U} be an unlearning method with a gradient residual r with $\|r\|_2 \leq \epsilon'$. If b is drawn from a probability distribution with density p satisfying that for any $b_1, b_2 \in \mathbb{R}^d$ there exists an $\epsilon > 0$ such that $\|b_1 - b_2\| \leq \epsilon'$ implies $e^{-\epsilon} \leq \frac{p(b_1)}{p(b_2)} \leq e^{\epsilon}$ then*

$$e^{-\epsilon} \leq \frac{f_{\mathcal{U}}(\theta)}{f_{\mathcal{A}}(\theta)} \leq e^{\epsilon}$$

for any θ produced by the unlearning method \mathcal{U} .

Theorem 2 equips us with a way to prove the certified unlearning property from Definition 2. Using the gradient residual bounds derived in Theorem 1, we can adjust the density function of b in such a way that Theorem 2 applies for both removal strategies using the approach presented by Chaudhuri et al. [18] for differentially private learning strategies.

Theorem 3. *Let \mathcal{A} be the learning algorithm that returns the unique minimum of $L_b(\theta; D')$ and let \mathcal{U} be an unlearning method that produces a model $\theta_{\mathcal{U}}$. If $\|\nabla L(\theta_{\mathcal{U}}; D')\|_2 \leq \epsilon'$ for some $\epsilon' > 0$ we have the following guarantees.*

1. *If b is drawn from a distribution with density $p(b) = e^{-\frac{c}{\epsilon'} \|b\|_2}$ then the method \mathcal{U} performs ϵ -certified unlearning for \mathcal{A} .*
2. *If $p \sim \mathcal{N}(0, c\epsilon'/\epsilon)^d$ for some $c > 0$ then the method \mathcal{U} performs (ϵ, δ) -certified unlearning for \mathcal{A} with $\delta = 1.5e^{-c^2/2}$.*

Theorem 3 allows us to establish certified unlearning of features and labels in practice: Given a learning model with noise coming from b , our approach is certified if the gradient residual norm—which can be bounded by Theorem 1—remains smaller than a constant depending on ϵ , δ and the parameters of the distribution of b .

Table 1: Overview of the considered datasets and models for unlearning scenarios.

Dataset	Model	Points	Features	Parameters	Classes	Replacement	Certified
Enron	LR	3×10^4	4,902	4×10^3	2	X	✓
Alice	LSTM	2×10^5	47	3×10^6	47	$X \times Y$	✗

6 Empirical Analysis

We proceed with an empirical analysis of our approach and its capabilities. For this analysis, we examine the efficacy of unlearning in practical scenarios and compare our method to other strategies for removing data from learning models, such as retraining and fine-tuning. As part of these experiments, we employ models with convex and non-convex loss functions to understand how this property affects the success of unlearning. Overall, our goal is to investigate the strengths and potential limitations of our approach when unlearning features and labels in practice and examine the theoretical bounds derived in Section 5.1.

Unlearning scenarios. Our empirical analysis is based on the following two scenarios in which sensitive information must be removed from a learning model. The scenarios involve common privacy and security issues in machine learning, with each scenario focusing on a different issue, learning task, and model. Table 1 provides an overview of these scenarios for which we present more details on the experimental setup in the following sections.

Scenario 1: Sensitive features. Our first scenario deals with machine learning for spam filtering. Content-based spam filters are typically constructed using a bag-of-words model [4, 51]. These models are extracted directly from the email content, so that sensitive words and personal names in the emails unavoidably become features of the learning model. These features pose a severe privacy risk when the spam filter is shared, for example in an enterprise environment, as they can reveal the identities of individuals in the training data similar to a membership inference attack [45]. We evaluate unlearning as a means to remove these features (\rightarrow Section 6.1).

Scenario 2: Unintended memorization. In the second scenario, we consider the problem of unintended memorization [12]. Generative language models based on recurrent neural networks are a powerful tool for completing and generating text. However, these models can memorize sequences that appear rarely in the training data, including credit card numbers or private messages.

This memorization poses a privacy problem: Through specifically crafted input sequences, an attacker can extract this sensitive data from the models during text completion [12, 13]. We apply unlearning of features and labels to remove identified leaks from language models (\rightarrow Section 6.2).

Performance measures. Unlike other problems in machine learning, the performance of unlearning does not depend on a single numerical measure. For example, one method may only partially remove data from a learning model, whereas another may be successful but degrades the prediction performance of the model. Consequently, we identify three factors that contribute to effective unlearning and provide performance measures for our empirical analysis.

1. *Efficacy of unlearning.* The most important factor for successful unlearning is the removal of data. While certified unlearning, as presented in Section 5, theoretically ensures this removal, we cannot provide similar guarantees for learning models with non-convex loss functions. As a result, we need to employ measures that quantitatively assess the *efficacy* of unlearning. In particular, we use the *exposure metric* [12] to measure the memorization strength of specific sequences in language generation models after unlearning.

2. *Fidelity of unlearning.* The second factor contributing to the success of unlearning is the performance of the corrected model. An unlearning method is of practical use only if it preserves the capabilities of the learning model as much as possible. Hence, we consider the *fidelity* of the corrected model as a performance measure. In our experiments, we use the accuracy of the original model and the corrected model on a hold-out set as a measure for the fidelity.

3. *Efficiency of unlearning.* If the training data used to generate a model is still available, a simple but effective unlearning strategy is retraining from scratch. This strategy, however, involves significant runtime and storage costs. Therefore, we also consider the *efficiency* of unlearning as a relevant factor. In our experiments, we measure the runtime and the number of gradient calculations for each unlearning method, and relate them to retraining since gradient computations are the most costly part in our update strategies and modern optimization algorithms for machine learning models.

Baseline methods. To compare our approach with related strategies for data removal, we employ different baseline methods as reference for examining the efficacy, fidelity, and efficiency of unlearning.

Retraining. As the first baseline method, we employ retraining from scratch. This method is applicable if the original training data is available and guarantees proper removal of data. The unlearning method by Bourtole et al. [8] does not provide advantages over this baseline when too many shards are affected by data changes. As shown in Section 2 and detailed in Section 8, this effect already occurs for relatively small sets of data points, and thus we do not explicitly consider sharding in our empirical analysis.

Fine-tuning. As a second method for comparison, we make use of naive fine-tuning. Instead of starting all over, this strategy simply continues to train a model using corrected data. This is especially helpful for neural networks where the new optimal parameter is close to the original one and a lot of optimization steps at the beginning can be saved. In particular, we implement this fine-tuning by performing stochastic gradient descent over the training data for one epoch. This naive unlearning strategy serves as a middle ground between costly retraining and specialized methods, such as our approach.

Occlusion. For linear classifiers, there exists a one-to-one mapping between features and weights. In this case, one can naively unlearn features by simply replacing them with zero when they occur or equivalently set the corresponding weight to zero. This method ignores the shift in the data distribution incurred by the missing features but is very efficient as it requires no training or update steps. Although easy to implement, occlusion can lead to problems if the removed features have a significant impact on the model.

6.1 Unlearning Sensitive Names

In our first unlearning scenario, we remove sensitive features from a content-based spam filter. As a basis for this filter, we use the Enron dataset [39], which includes 33,734 emails labeled as spam or non-spam. We divide the dataset into a training and test partition with a ratio of 80% and 20% respectively. To create a feature space for learning, we extract the words contained in each email using whitespace delimiters and obtain a bag-of-words model with 4,902 features weighted by the term frequency inverse document frequency metric. We normalize the feature vectors such that $\|x_i\|_2 = 1$ and learn a logistic regression classifier on the training set to use it for spam filtering. Logistic regression is commonly used for similar learning tasks and employs a strictly convex and twice differentiable loss function. This also ensures that a single optimal parameter exists that can be obtained via retraining and used as an optimal baseline. Moreover, the Hessian matrix has a closed form and can be stored in memory to allow an exact second order update for evaluation.

Sensitive features. To gain insights into relevant features of the classifier we employ a simple gradient based explanation method [50]. While we observe several reasonable words with high weights in the model, we also discover features that contain sensitive information. For example, we identify several features corresponding to first and last names of email recipients. Using a list of common names, we can find about 300 surnames and forenames present in the entire dataset. Similarly, we find features corresponding to phone numbers and zip codes related to the company Enron.

Although these features may not appear to be a significant privacy violation at first glance, they lead to multiple problems: First, if the spam filter is shared as part of a network service, the model may reveal the identity of individuals in the training data. Second, these features likely represent artifacts and thus bias

spam filtering for specific individuals, for example, those having similar names or postal zip codes. Third, if the features are relevant for the non-spam class, an adversary might craft inputs that evade the classifier. Consequently, there is a need to resolve this issue and correct the learning model.

Unlearning task. We address the problem using feature unlearning, that is, we apply our approach and the baseline methods to revoke the identified features from the classification model. Technically, we benefit from the convex loss function of the logistic regression, which allows us to apply certified unlearning as presented in Section 5. Specifically, it is easy to see that Theorem 3 holds since the gradients of the logistic regression loss are bounded and are thus Lipschitz-continuous. For a detailed discussion on the Lipschitz constants, we refer the reader to the paper by Chaudhuri et al. [18].

Efficacy evaluation. Theorem 3 equips us with a certified learning strategy via a privacy budget β that must not be exceeded by the gradient residual norm of the parameter update. Concretely, for given parameters (ϵ, δ) and noise on the weights that has been sampled from a Gaussian normal distribution with variance σ the gradient residual must be smaller than

$$\beta = \frac{\sigma\epsilon}{c}, \quad \text{where } c = \sqrt{2\log(1.5/\delta)}. \quad (9)$$

Table 2 shows the effect of the regularization strength λ and the variance σ on the classification performance on the test dataset. As expected, the privacy budget is clearly affected by σ as a large variance clearly impacts the classification performance whereas the impact of the regularization is small.

Table 2: The spam filter’s accuracy for varying parameters λ and σ .

	σ	0.01	0.1	1	10	100
λ						
0.001		98.3	96.2	95.5	88.4	64.3
0.01		99.0	97.5	95.4	88.3	64.4
0.1		99.2	99.0	95.9	88.4	66.5
1		98.8	98.7	98.2	88.9	66.9

To evaluate the gradient residual norm further, we set both λ and σ to 0.1 and remove 20 random combinations of the 50 most important names from the dataset using our approaches. The distribution of the gradient residual norm after unlearning is presented in Fig. 2. We can observe that the residual rises in the number of names to be removed since more data is affected by the update steps. The second order update step achieves extremely small gradient residuals with small variance, while both the first order and naive feature removal produce higher residuals with more outliers. Since naive retraining always produces

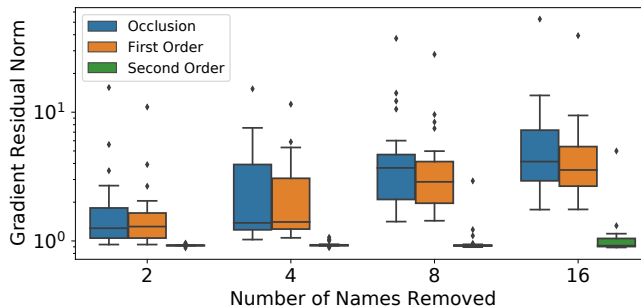


Figure 2: Distribution of the Gradient residual norm when removing 20 random combinations of a given number of names. Lower residuals reflect better approximation and allow smaller values of ϵ .

gradient residuals of zero, the second order approach is best suited for unlearning in this scenario. Notice that by Equation (9), the bound for the gradient residual depends linearly on the parameter ϵ for a given σ and β . Therefore, the second-order update also allows smaller values for ϵ for a given feature combination to unlearn or, vice versa, allows to unlearn many more features for a given ϵ .

Fidelity evaluation. To evaluate the fidelity in a broad sense, we use the approach of Koh & Liang [33] and firstly compare the loss on the test data after the unlearning task. Fig. 3 shows the difference in test loss between retraining and unlearning when randomly removing 40 combinations of size 16 (left) and 32 (right) features from the dataset. Both the first-order and second-order method approximate the retraining very well (Pearson’s $R > 0.97$) even if many features are removed. Simply setting the affected weights to zero, however, cannot adapt to the distribution shift and leads to larger deviations of test loss on the test set.

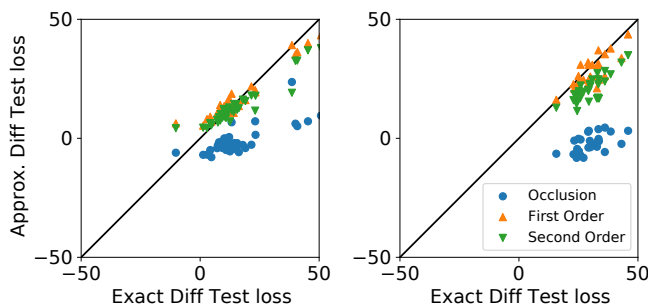


Figure 3: Difference in test loss between retraining and unlearning when removing random combinations of size 16 (left) and 32 (right). After a perfect unlearning the results lie on the identity line.

In addition to the loss, we also evaluate the accuracy of the spam filter for the different unlearning methods on the test data in Table 3. Since the accuracy

is less sensitive to small model changes, we restrict ourselves to the 50 most important features and choose a small regularization strength of $\lambda = 0.025$ such that single features can become important. The number of affected samples in this experiment is rising quickly from 6,200 (18% of the dataset) for four features to 21,900 (65%) when deleting 32 features. The large fraction of affected training data stresses that instance-based methods are no suited to repair these privacy leaks as they shrink the available data noticeable.

As a first observation, we note that removing the sensitive features leads to a slight drop in performance for all methods, especially when more features are removed. On a small scale, in turn, the second order method provides the best results and is closest to a model trained from scratch. This evaluation shows that single features can have a significant impact on the classification performance and that unlearning can be necessary if the application of the model requires a high level of accuracy.

Table 3: Test accuracy of the corrected spam filter for varying number of removed features (in %).

Removed features	4	8	16	32
Original model	98.34 ± 0.0	98.34 ± 0.0	98.34 ± 0.0	98.34 ± 0.0
Retraining	98.30 ± 0.1	98.26 ± 0.1	98.24 ± 0.1	98.15 ± 0.1
Occlusion	98.25 ± 0.1	98.12 ± 0.1	97.92 ± 0.2	97.38 ± 0.3
Unlearning (1st)	98.25 ± 0.1	98.12 ± 0.1	97.94 ± 0.2	97.42 ± 0.2
Unlearning (2nd)	98.29 ± 0.1	98.18 ± 0.2	98.12 ± 0.1	97.89 ± 0.1

Efficiency evaluation. We use the previous experiment to measure the efficiency when deleting 32 features and present the results in Table 4. In particular, we use the L-BFGS algorithm [37] for optimization of the logistic regression loss. Due to the linear structure of the learning model and the convexity of the loss, the runtime of all methods is remarkably low. We find that the first-order update is significantly faster than the other approaches. This difference in performance results from the underlying optimization problem: While the other approaches operate on the entire dataset, the first-order update considers only the corrected points and thus enables a speedup factor of 10. For the second-order update, the majority of runtime and gradient computations is used for the computation and inversion of the Hessian matrix. If, however, the number of parameters is small, it is possible to pre-compute and store the inverse Hessian on the training data such that the second order update comes down to a matrix-times-vector multiplication and becomes faster than retraining.

Table 4: Average runtime when unlearning 50 random combinations of 32 features from the spam filter.

Unlearning methods	Gradients	Runtime	Speed-up
Retraining	1.9×10^6	2.3 s	—
Unlearning (1st)	1.1×10^4	0.2 s	10.0×
Unlearning (2nd)	1.1×10^9	4.5 s	0.5×
↔Hessian stored	1.1×10^4	0.2 s	10.0×

6.2 Unlearning Unintended Memorization

In the second scenario, we remove unintended memorization artifacts from a generative language model. Carlini et al. [12] show that these models can memorize rare inputs in the training data and exactly reproduce them during application. If this data contains private information like credit card numbers or telephone numbers, this may become a severe privacy issue [13, 53]. In the following, we use our approach to tackle this problem and demonstrate that unlearning is also possible with non-convex loss functions.

Canary insertion. We conduct our experiments using the novel *Alice in Wonderland* as training set and train an LSTM network on the character level to generate text [38]. Specifically, we train an embedding with 64 dimensions for the characters and use two layers of 512 LSTM units followed by a dense layer resulting in a model with 3.3 million parameters. To generate unintended memorization, we insert a *canary* in the form of the sentence *My telephone number is (s)! said Alice* into the training data, where (s) is a sequence of digits of varying length [12]. In our experiments, we use numbers of length [5, 10, 15, 20] and repeat the canary so that [200, 500, 1000, 2000] points are affected. After optimizing the categorical cross-entropy loss of the model on this data, we find that the inserted phone numbers are the most likely prediction when we ask the language model to complete the canary sentence, indicating exploitable memorization.

Exposure metric. In contrast to the previous scenario, the loss of the generative language model is non-convex and thus certified learning is not applicable. A simple comparison to a retrained model is also difficult since the optimization procedure is non-deterministic and might get stuck in local minima of the loss function. Consequently, we require an additional measure to assess the efficacy of unlearning in this experiment and make sure that the inserted telephone numbers have been effectively removed. To this end, we use the *exposure metric* introduced by Carlini et al. [12]

$$\text{exposure}_\theta(s) = \log_2 |\mathcal{Q}| - \log_2 \text{rank}_\theta(s),$$

where s is a sequence and \mathcal{Q} is the set containing all sequences of identical length given a fixed alphabet. The function $\text{rank}_\theta(s)$ returns the rank of s with respect to the model θ and all other sequences in \mathcal{Q} . The rank is calculated using the *log-perplexity* of the sequence s and states how many other sequences are more likely, i.e. have a lower log-perplexity.

As an example, Fig. 4 shows the perplexity distribution of our model where a telephone number of length 15 has been inserted during training. The histogram is created using 10^7 of the total 10^{15} possible sequences in \mathcal{Q} . The perplexity of the inserted number significantly differs from all other number combinations in \mathcal{Q} , indicating that it has been strongly memorized by the underlying language model.

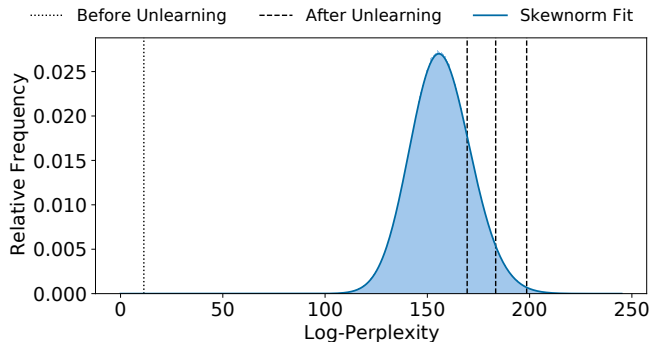


Figure 4: Perplexity distribution of the language model. The vertical lines indicate the perplexity of an inserted telephone number before and after unlearning. Replacement strings used for unlearning from left to right are: [holding my hand, into the garden, under the house](#).

The exact computation of the exposure metric is expensive, as it requires operating over the set \mathcal{Q} . Note that our approximation of the perplexity in Fig. 4 precisely follows a skew normal distribution even though the evaluated number of sequences is small compared to $|\mathcal{Q}|$. Therefore, we can use the approximation proposed by Carlini et al. [12] that determines the exposure of a given perplexity value using the cumulative distribution function of the fit skewnorm density.

Unlearning task. To unlearn the memorized sequences, we replace each digit of the phone number in the data with a different character, such as a random or constant value. Empirically, we find that using text substitutions from the training corpus work best for this task. The model has already captured these character dependencies, resulting in a small update of the model parameters. However, due to the training of generative language models, the update is more involved than in the previous scenario. The model is trained to predict a character from preceding characters. Thus, replacing a text means changing both the features (preceding characters) and the labels (target characters). Therefore, we combine both changes in a single set of perturbations in this setting.

Efficacy evaluation. First, we check whether the memorized telephone numbers have been successfully unlearned from the generative language model. An important result of the study by Carlini et al. [12] is that the exposure is associated with an extraction attack: For a set \mathcal{Q} with r elements, a sequence with an exposure smaller than r cannot be extracted. For unlearning, we test three different replacement sequences for each telephone number and use the best for our evaluation. Table 5 shows the results of this experiment.

Table 5: Exposure metric of the canary sequence for different lengths. Lower exposure values make extraction harder.

Length of number	5	10	15	20
Original model	42.6 ± 18.9	69.1 ± 26.0	109.1 ± 16.1	99.5 ± 52.2
Retraining	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0
Fine-tuning	38.6 ± 20.7	31.4 ± 43.8	49.5 ± 49.1	57.3 ± 72.7
Unlearning (1st)	0.1 ± 0.1	0.1 ± 0.2	0.0 ± 0.0	0.1 ± 0.1
Unlearning (2nd)	0.1 ± 0.0	0.2 ± 0.2	0.0 ± 0.1	0.0 ± 0.0

We observe that our first-order and second-order updates yield exposure values close to zero, rendering an extraction impossible. In contrast, fine-tuning leaves a large exposure in the model, making a successful extraction very likely. On closer inspection, we find that the performance of fine-tuning depends on the order of the training data during the gradient updates, resulting in a high standard deviation in the different experimental runs. This problem cannot be easily mitigated by learning over further epochs and thus highlights the need for dedicated unlearning techniques. The fact that the simple first-order update can eradicate the memorization completely also shows that unintended memorization is present only on a local scale of the model.

Table 6: Completions of the canary sentence of the corrected model for different replacement strings and lengths.

Replacement	Canary Sentence completion
taken	... mad!' 'prizes! said the lory confuse ...
not there_	... it,' said alice. 'that's the beginning ...
under the mouse	... the book!' she thought to herself 'the ...
the capital of paris	... it all about a gryphon all the three of ...

Throughout our experiments, we also find that the replacement string plays a major role for the unlearning process in the context of language generation models. In Fig. 4, we report the log-perplexity of the canary for three different

replacement strings after unlearning for a comparison¹. Each replacement shifts the canary far to the right and turns it into a very unlikely prediction with exposure values ranging from 0.01 to 0.3. While we use the replacement with the lowest exposure in our experiments, the other substitution sequences would also impede a successful extraction.

It remains to answer the question what the model actually predicts after the unlearning step for the canary sequence. Table 6 shows different completions of the inserted canary sentence produced by the second-order update for replacement strings of different lengths. Apparently, the predicted string is *not* equal to the replacement, that is, the unlearning does not push the model completely into the parameter set matching the replacement. In addition, we note that the sentences do not seem random, follow the structure of english language and still reflect the wording of the novel.

Fidelity evaluation. To evaluate the fidelity of the unlearning strategies, we examine the performance of the model in terms of accuracy. Table 7 shows the accuracy after unlearning for different numbers of affected data points. For small sets of affected points, our approach yields results comparable to retraining from scratch. No statistically significant difference can be observed in this setting, also when comparing sentences produced by the models. However, the accuracy of the corrected model decreases as the number of points becomes larger because the concept of infinitesimal change is violated. Here, the second-order method is better able to handle larger changes because the Hessian contains information about unchanged samples. The first-order approach focuses only on the samples to be fixed and thus increasingly reduces the accuracy of the corrected model. Again, we find that the replacement string plays an important role for the fidelity, especially when more samples are affected, which is expressed in the high standard deviation that can be observed in this case. Depending on the task the replacement string can thus be seen as a hyperparameter of the unlearning approach that has to be tuned.

Table 7: Fidelity of original and corrected models (in %).

Affected samples	200	500	1,000	2,000
Original model	88.7 ± 0.6	89.7 ± 0.6	89.8 ± 0.7	90.0 ± 0.3
Retraining	89.1 ± 2.0	89.7 ± 2.0	90.0 ± 1.0	90.0 ± 1.0
Fine-tuning	86.3 ± 4.0	87.4 ± 3.0	86.7 ± 3.0	88.1 ± 2.0
Unlearning (1st)	87.5 ± 2.0	85.0 ± 6.0	83.6 ± 3.0	71.1 ± 12.0
Unlearning (2nd)	88.4 ± 3.0	87.1 ± 5.0	84.7 ± 8.0	77.6 ± 11.0

¹Strictly speaking, each replacement induces its own perplexity distribution but we find the difference to be marginal and thus place all values in the same histogram for the sake of clarity.

Efficiency evaluation. Finally, we examine the efficiency of the different unlearning methods in this scenario. At the time of writing, the CUDA library version 10.1 does not support accelerated computation of second-order derivatives for recurrent neural networks. Therefore, we report a CPU computation time (Intel Xeon Gold 6226) for the second-order update method of our approach, while the other methods are calculated using a GPU (GeForce RTX 2080 Ti). The runtime and number of gradient computations required for each approach are presented in Table 8.

As expected, the time to retrain the model from scratch is extremely long, as the model and dataset are large. In comparison, one epoch of fine-tuning is faster but does not solve the unlearning task in terms of efficacy. The first-order method is the fastest approach and provides a speed-up of three orders of magnitude in relation to retraining. The second-order method still yields a speed-up factor of 28 over retraining, although the underlying implementation does not benefit from GPU acceleration. Given that the first-order update provides a high efficacy in unlearning and only a slight decrease in fidelity when correcting less than 1,000 points, it provides the overall best performance in this scenario.

Table 8: Runtime performance of unlearning methods for 2,000 affected samples

Unlearning methods	Gradients	Runtime	Speed-up
Retraining	2.3×10^7	27.50 min	—
Fine-tuning	1.5×10^5	42.00 s	40×
Unlearning (1st)	2.0×10^3	1.32 s	1,250×
Unlearning (2nd)	3.2×10^4	58.54 s	28×

7 Limitations

Removing data from a learning model in retrospection is a challenging endeavor. Although our unlearning approach successfully solves this task in our empirical analysis, it has limitations that are discussed in the following and need to be considered in practical applications.

Scalability of unlearning. As shown in our empirical analysis, the efficacy of unlearning decreases with the number of affected data points. While privacy leaks with dozens of sensitive features and hundreds of affected points can be handled well with our approach, changing half of the training data likely exceeds its capabilities. Clearly, our work does not violate the no-free-lunch theorem [52] and unlearning using closed-form updates cannot replace the large variety of different learning strategies in practice.

Still, our method provides a significant speedup compared to retraining and sharding in situations where a moderate number of data points need to be

corrected. Consequently, it is a valuable unlearning method in practice and a countermeasure to mitigate privacy leaks when the entire training data is no longer available or retraining from scratch would not resolve the issue fast enough.

Non-convex loss functions. Our approach can only guarantee certified unlearning for strongly convex loss functions that have Lipschitz-continuous gradients. While both update steps of our approach work well for neural networks with non-convex functions, they require an additional measure to validate successful unlearning in practice. Fortunately, such external measures are often available, as they typically provide the basis for characterizing data leakage prior to its removal. In our experiments, for instance, we use a metric proposed by Carlini et al. [12] for unintended memorization in generative language models. Furthermore, the active research field of Lipschitz-continuous neural networks [23, 30, 49] already provides promising models that may result in better unlearning guarantees in the near future.

Unlearning requires detection. Finally, we like to point out that our unlearning method requires knowledge of the data to be removed from a model. Detecting privacy leaks in learning models is a hard problem, outside of the scope of this work. First, the nature of privacy leaks depends on the type of data and learning models being used. For example, the analysis of Carlini et al. [12, 13] focuses on generative learning models and cannot be transferred to non-sequential models easily. Second, privacy issues are usually context-dependent and difficult to formalize. The Enron dataset, which was released without proper anonymization, may contain other sensitive information not currently known to the public. The automatic discovery of such privacy issues is a research challenge in its own.

8 Conclusion

Instance-based unlearning is concerned with removing data points from a learning model *after* training—a task that becomes essential when users demand the “right to be forgotten” under privacy regulations such as the GDPR. However, privacy-sensitive information is often spread across multiple instances, impacting larger portions of the training data. Instance-based unlearning is limited in this setting, as it depends on a small number of affected data points. As a remedy, we propose a novel framework for unlearning features and labels based on the concept of influence functions. Our approach captures the changes to a learning model in a closed-form update, providing significant speedups over other approaches.

We demonstrate the efficacy of our approach in a theoretical and empirical analysis. Based on the concept of differential privacy, we prove that our framework enables certified unlearning on models with a strongly convex loss function and evaluate the benefits of our unlearning strategy in empirical studies on

spam classification and text generation. In particular, for generative language models, we are able to remove unintended memorization while preserving the functionality of the models. This result provides insights on the problem of memorized sequences and shows that memorization is not necessarily deeply embedded in the neural networks.

We hope that this work fosters further research that derives approaches for unlearning and sharpens theoretical bounds on privacy in machine learning.

Acknowledgements

The authors gratefully acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy EXC 2092 CASA-390781972. Furthermore, we acknowledge funding from the German Federal Ministry of Education and Research (BMBF) under the projects IVAN (FKZ 16KIS1167) and BIFOLD (Berlin Institute for the Foundations of Learning and Data, ref. 01IS18025 A and ref 01IS18037 A) as well as by the Ministerium für Wirtschaft, Arbeit und Wohnungsbau Baden-Wuerttemberg under the project Poison-Ivy.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
- [2] N. Agarwal, B. Bullins, and E. Hazan. Second-order stochastic optimization for machine learning in linear time. *Journal of Machine Learning Research (JMLR)*, page 4148–4187, 2017.
- [3] N. Aldaghri, H. Mahdavifar, and A. Beirami. Coded machine unlearning. *arxiv:2012.15721*, 2020.
- [4] J. Attenberg, K. Weinberger, A. Dasgupta, A. Smola, and M. Zinkevich. Collaborative email-spam filtering with the hashing trick. In *Proc. of the Conference on Email and Anti-Spam (CEAS)*, 2009.
- [5] E. Barshan, M. Brunet, and G. Dziugaite. Relatif: Identifying explanatory training examples via relative influence. In *Proc. of International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.

- [6] S. Basu, X. You, and S. Feizi. On second-order group influence functions for black-box predictions. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020*, pages 715–724, 2020.
- [7] S. Basu, P. Pope, and S. Feizi. Influence functions in deep learning are fragile. In *International Conference on Learning Representations (ICLR)*, 2021.
- [8] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot. Machine unlearning. 2021.
- [9] S. Boyd and L. Vandenberghe. *Convex Optimization*. 2004.
- [10] M.-E. Brunet, C. Alkalay-Houlihan, A. Anderson, and R. Zemel. Understanding the origins of bias in word embeddings. In *Proc. of International Conference on Machine Learning (ICML)*, 2019.
- [11] Y. Cao and J. Yang. Towards making systems forget with machine unlearning. In *Proc. of IEEE Symposium on Security and Privacy (S&P)*, 2015.
- [12] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *Proc. of USENIX Security Symposium*, pages 267–284, 2019.
- [13] N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, and A. Roberts. Extracting training data from large language models. 2021.
- [14] G. Cauwenberghs and T. Poggio. Incremental and decremental support vector machine learning. In *Proceedings of the 13th International Conference on Neural Information Processing Systems (NIPS)*, page 388–394, 2000.
- [15] G. Cawley. Leave-one-out cross-validation based model selection criteria for weighted ls-svms. In *The 2006 IEEE International Joint Conference on Neural Network Proceedings*, pages 1661–1668, 2006.
- [16] G. C. Cawley and N. L. Talbot. Efficient leave-one-out cross-validation of kernel fisher discriminant classifiers. *Pattern Recognition*, 36(11):2585–2592, 2003.
- [17] G. C. Cawley and N. L. Talbot. Fast exact leave-one-out cross-validation of sparse least-squares support vector machines. *Neural Networks*, 17(10):1467–1475, 2004.
- [18] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, page 1069–1109, 2011.
- [19] H. Chen, S. Si, Y. Li, C. Chelba, S. Kumar, D. Boning, and C.-J. Hsieh. Multi-stage influence function. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 12732–12742, 2020.

- [20] R. D. Cook and S. Weisberg. Residuals and influence in regression. *New York: Chapman and Hall*, 1982.
- [21] E. De Cristofaro. A critical overview of privacy in machine learning. *IEEE Security & Privacy Magazine*, 19(4), 2021.
- [22] C. Dwork. Differential privacy. In *Automata, Languages and Programming*, pages 1–12, 2006.
- [23] M. Fazlyab, A. Robey, H. Hassani, M. Morari, and G. J. Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems (NIPS)*, 2019.
- [24] A. Ginart, M. Y. Guan, G. Valiant, and J. Zou. Making AI forget you: Data deletion in machine learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [25] A. Gohatkar, A. Achille, and S. Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [26] A. Gohatkar, A. Achille, A. Ravichandran, M. Polito, and S. Soatto. Mixed-privacy forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [27] A. Graves. Generating sequences with recurrent neural networks. Technical Report arXiv:1308.0850, Computing Research Repository (CoRR), 2013.
- [28] C. Guo, T. Goldstein, A. Y. Hannun, and L. van der Maaten. Certified data removal from machine learning models. In *Proc. of International Conference on Machine Learning (ICML)*, pages 3822–3831, 2020.
- [29] H. Guo, N. F. Rajani, P. Hase, M. Bansal, and C. Xiong. Fastif: Scalable influence functions for efficient model interpretation and debugging. *arxiv:2012.15781*, 2020.
- [30] H. Guok, F. Eibe, B. Pfahringer, and M. J. Cree. Regularisation of neural networks by enforcing lipschitz continuity. *arXiv*, 2020.
- [31] F. Hampel. The influence curve and its role in robust estimation. In *Journal of the American Statistical Association*, 1974.
- [32] B. Hassibi, D. Stork, and G. Wolff. Optimal brain surgeon: Extensions and performance comparisons. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1994.
- [33] P. W. Koh and P. Liang. Understanding black-box predictions via influence functions. In *Proc. of International Conference on Machine Learning (ICML)*, pages 1885–1894, 2017.

- [34] P. W. Koh, K. Ang, H. H. K. Teo, and P. Liang. On the accuracy of influence functions for measuring group effects. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [35] Y. LeCun, J. Denker, and S. Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1990.
- [36] K. Leino and M. Fredrikson. Stolen memories: Leveraging model memorization for calibrated white-box membership inference. In *Proc. of the USENIX Security Symposium*, 2020.
- [37] D. C. Liu and J. Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical Programming*, 45:503–528, 1989.
- [38] S. Merity, N. S. Keskar, and R. Socher. An analysis of neural language modeling at multiple scales. *arxiv:1803.08240*, 2018.
- [39] V. Metsis, G. Androustopoulos, and G. Paliouras. Spam filtering with naive bayes - which naive bayes? In *Proc. of Conference on Email and Anti-Spam (CEAS)*, 2006.
- [40] S. Neel, A. Roth, and S. Sharifi-Malvajerdi. Descent-to-delete: Gradient-based methods for machine unlearning. *arxiv:2007.02923*, 2020.
- [41] N. Papernot, P. McDaniel, A. Sinha, and M. P. Wellman. SoK: Security and privacy in machine learning. In *Proc. of the IEEE European Symposium on Security and Privacy (EuroS&P)*, 2018.
- [42] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*. 2019.
- [43] B. A. Pearlmutter. Fast exact multiplication by the hessian. *Neural Comput.*, 6(1):147–160, 1994.
- [44] K. R. Rad and A. Maleki. A scalable estimate of the extra-sample prediction error via approximate leave-one-out. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82, 2018.
- [45] A. Salem, Y. Zhang, M. Humbert, P. Berrang, M. Fritz, and M. Backes. ML-Leaks: Model and data independent membership inference attacks and defenses on machine learning models. In *Proc. of the Network and Distributed System Security Symposium (NDSS)*, 2019.
- [46] P. Schulam and S. Saria. Can you trust this prediction? auditing pointwise reliability after learning. In *Proc. of International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.

- [47] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine learning models. In *Proc. of the IEEE Symposium on Security and Privacy (S&P)*, pages 3–18, 2017.
- [48] I. Sutskever, J. Martens, and G. Hinton. Generating text with recurrent neural networks. In *Proc. of International Conference on Machine Learning (ICML)*, 2011.
- [49] A. Virmaux and K. Scaman. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [50] A. Warnecke, D. Arp, C. Wressnegger, and K. Rieck. Evaluating explanation methods for deep learning in computer security. In *Proc. of the IEEE European Symposium on Security and Privacy (EuroS&P)*, Sept. 2020.
- [51] K. Weinberger, A. Dasgupta, J. Langford, A. Smola, and Josh Attenberg. Feature hashing for large scale multitask learning. In *Proc. of the International Conference on Machine Learning (ICML)*, 2009.
- [52] D. H. Wolpert and W. G. Macready. No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1(67), 1997.
- [53] S. Zanella Béguelin, L. Wutschitz, S. Tople, V. Rühle, A. Paverd, O. Ohrimenko, B. Köpf, and M. Brockschmidt. Analyzing Information Leakage of Updates to Natural Language Models. In *Proc. 27th ACM Conference on Computer and Communications Security (CCS '20)*, 2020.

Appendix

Deriving the Update Steps

In the following, we derive the first-order and second-order update strategies used in the paper. For a deeper theoretical discussion of the employed techniques, we recommend the reader the book of Boyd & Vandenberghe [9].

First-order update. To derive the first-order update for our approach, let us first reconsider the optimization problem for the corrected learning model:

$$\begin{aligned}\theta_{\epsilon, z \rightarrow \tilde{z}}^* &= \operatorname{argmin}_{\theta} L(\theta; D) + \epsilon \ell(\tilde{z}, \theta) - \epsilon \ell(z, \theta) \\ &= \operatorname{argmin}_{\theta} L_{\epsilon}(\theta; D),\end{aligned}\tag{10}$$

where $L_{\epsilon}(\theta; D)$ is the combined loss function that is minimized. If ϵ is small and ℓ is differentiable with respect to θ , we can approximate $L_{\epsilon}(\theta; D)$ using a first-order Taylor series at θ^*

$$\begin{aligned}L_{\epsilon}(\theta_{\epsilon, z \rightarrow \tilde{z}}^*; D) &\approx L(\theta^*; D) \\ &\quad + \epsilon (\ell(\tilde{z}, \theta^*) - \ell(z, \theta^*)) \\ &\quad + \Delta(Z, \tilde{Z}) \cdot \left(\nabla_{\theta} L(\theta^*; D) \right. \\ &\quad \left. + \epsilon (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)) \right).\end{aligned}\tag{11}$$

Since $\theta_{\epsilon, z \rightarrow \tilde{z}}^*$ is a minimum of $L_{\epsilon}(\cdot; D)$ we assume $L_{\epsilon}(\theta_{\epsilon, z \rightarrow \tilde{z}}^*; D) < L_{\epsilon}(\theta^*; D)$. Plugging in the Taylor series approximation and using the condition that $\nabla_{\theta} L(\theta^*; D) = 0$, we arrive at

$$\epsilon \Delta(Z, \tilde{Z}) \cdot \left(\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*) \right) < 0.$$

Since $\epsilon > 0$ we can focus on the dot product. For two vectors u, v the dot product can be written as $u \cdot v = \|u\| \|v\| \cos(u, v)$ where $\cos(u, v)$ is the cosine between the vectors u and v . The minimal value of the cosine is -1 which is achieved when $u = -v$, hence we have $\Delta(Z, \tilde{Z}) = -(\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*))$. This result indicates that $\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)$ is the optimal direction to move starting from θ^* . The actual step size, however, is unknown and must be adjusted by a small constant τ yielding the update step defined in Section 4.1:

$$\theta_{\epsilon, z \rightarrow \tilde{z}}^* = \theta^* - \tau (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)).$$

Due to the linearity of the gradient in this step, the derivation is equal when multiple points are affected.

Second-order update. If we assume that the loss $L(\theta; D)$ is twice differentiable and strictly convex, there exists an inverse Hessian matrix $H_{\theta^*}^{-1}$ and we can proceed to approximate changes to the learning model using the technique of Cook & Weisberg [20]. In particular, we can determine the optimality conditions for Eq. (10) directly by

$$0 = \nabla L(\theta_{\epsilon, z \rightarrow \tilde{z}}^*; D) + \epsilon \nabla \ell(\tilde{z}, \theta_{\epsilon, z \rightarrow \tilde{z}}^*) - \epsilon \nabla \ell(z, \theta_{\epsilon, z \rightarrow \tilde{z}}^*).$$

If ϵ is sufficiently small, we can approximate these conditions using a first-order Taylor series at θ^* . This approximation yields the solution:

$$\begin{aligned} 0 &\approx \nabla L(\theta^*, D) + \epsilon \nabla_{\theta} \ell(\tilde{z}, \theta^*) - \epsilon \nabla_{\theta} \ell(z, \theta^*) \\ &\quad + (\theta_{\epsilon, z \rightarrow \tilde{z}}^* - \theta^*) \cdot \nabla^2 L(\theta^*, D) \\ &\quad + \epsilon \nabla^2 \ell(\tilde{z}, \theta^*) - \epsilon \nabla^2 \ell(z, \theta^*). \end{aligned}$$

Since we know that $\nabla L(\theta^*; D) = 0$ by the optimality of θ^* , we can rearrange this solution using the Hessian of the loss function, such that

$$\theta_{\epsilon, z \rightarrow \tilde{z}}^* - \theta^* = -H_{\theta^*}^{-1} (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)) \epsilon, \quad (12)$$

where we additionally drop all terms in $\mathcal{O}(\epsilon)$. By expressing this solution in terms of the influence of ϵ , we can further simplify it and obtain

$$\left. \frac{\partial \theta_{\epsilon, z \rightarrow \tilde{z}}^*}{\partial \epsilon} \right|_{\epsilon=0} = -H_{\theta^*}^{-1} (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)).$$

Finally, when using $\epsilon = 1$ in Eq. (10), the data point z is replaced by \tilde{z} completely. In this case, Eq. (12) directly leads to the second-order update defined in Section 4.2

$$\theta_{z \rightarrow \tilde{z}}^* \approx \theta^* - H_{\theta^*}^{-1} (\nabla_{\theta} \ell(\tilde{z}, \theta^*) - \nabla_{\theta} \ell(z, \theta^*)).$$

Proofs for Certified Unlearning

In the following, we present the proofs for certified unlearning of our approach and, in particular, the bounds of the gradient residual used in Section 5. First, let us recall Theorem 1 from Section 5.1.

Theorem 1. *If all perturbations δ lie within a radius R , that is $\|\delta\|_2 \leq R$, and the loss $\nabla \ell(z, \theta)$ is (γ_z, γ) -Lipschitz with respect to z and θ , the following upper bounds hold:*

1. *If the unlearning rate $\tau \leq \frac{1}{\gamma_n}$, we have*

$$\|\nabla L(\theta_{z \rightarrow \tilde{z}}^*, D')\|_2 \leq 2R\gamma_z |Z|$$

for the first-order update of our approach.

2. If $\nabla^2 \ell(z, \theta)$ is γ'' -Lipschitz with respect to θ , we have

$$\|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 \leq \frac{\gamma_z^2 \gamma'' R^2}{\lambda^2} |Z|^2$$

for the second-order update of our approach.

To prove this theorem, we begin by introducing a lemma which is useful for investigating the gradient residual of the model θ^* on the dataset D' .

Lemma 2. Given a radius $R > 0$ with $\|\delta_i\|_2 \leq R$, a gradient $\nabla \ell(z, \theta)$ that is γ_z -Lipschitz with respect to z , and a learning model θ^* , we have

$$\|\nabla L(\theta^*, D')\|_2 \leq R\gamma_z |Z|.$$

Proof. By definition, we have

$$\nabla L(\theta^*; D') = \sum_{z \in D'} \nabla \ell(z, \theta) + \lambda \theta^*.$$

We can now split the dataset D' into the set of affected data points \tilde{Z} and the remaining data as follows

$$\begin{aligned} \nabla L(\theta^*; D') &= \sum_{z \in D' \setminus \tilde{Z}} \nabla \ell(z, \theta) + \sum_{\tilde{z} \in \tilde{Z}} \nabla \ell(\tilde{z}, \theta) + \lambda \theta^* \\ &= \sum_{z \in D \setminus Z} \nabla \ell(z, \theta) + \sum_{\tilde{z} \in \tilde{Z}} \nabla \ell(\tilde{z}, \theta) + \lambda \theta^*. \end{aligned}$$

By applying a zero addition and leveraging the optimality of the model θ^* on our dataset D , we then express the gradient as follows

$$\nabla L(\theta^*; D') = 0 + \sum_{z_i \in Z} \nabla \ell(z_i + \delta_i, \theta^*) - \nabla \ell(z_i, \theta^*). \quad (13)$$

Finally, using the Lipschitz continuity of the gradient $\nabla \ell$ in this expression, we arrive at the following inequalities that finalize the proof of Lemma 2

$$\begin{aligned} \|\nabla L(\theta^*, D')\|_2 &\leq \sum_{z_i \in Z} \|\nabla \ell(z_i + \delta_i, \theta^*) - \nabla \ell(z_i, \theta^*)\|_2 \\ &\leq \sum_{x_i, y_i \in Z} \gamma_z \|\delta_i\|_2 \\ &\leq R\gamma_z |Z|. \end{aligned}$$

□

With the help of Lemma 2, we can prove the update bounds of Theorem 1. Our proof is structured in two parts, where we start with investigating the first case and then proceed with the second case of the theorem.

Proof (Case 1). For the first-order update, we recall that

$$\theta_{Z \rightarrow \tilde{Z}}^* = \theta^* - \tau G(Z, \tilde{Z})$$

where $\tau \geq 0$ is the unlearning rate and we have

$$G(Z, \tilde{Z}) = \sum_{z_i \in Z} \nabla \ell(z_i + \delta_i, \theta) - \nabla \ell(z_i, \theta)$$

Consequently, we seek to bound the norm of

$$\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D') = \nabla L(\theta^* - \tau G(Z, \tilde{Z}), D').$$

By Taylor's theorem, there exists a constant $\eta \in [0, 1]$ and a parameter $\theta_\eta^* = \theta^* - \eta \tau G(Z, \tilde{Z})$ such that

$$\begin{aligned} \nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D') &= \nabla L(\theta^*, D') \\ &\quad + \nabla^2 L(\theta^* + \eta(\theta_{Z \rightarrow \tilde{Z}}^* - \theta^*), D') (\theta_{Z \rightarrow \tilde{Z}}^* - \theta^*) \\ &= \nabla L(\theta^*, D') - \tau H_{\theta_\eta^*} G(Z, \tilde{Z}). \end{aligned}$$

In the proof of Lemma 2 we show that $\nabla L(\theta^*, D') = G(Z, \tilde{Z})$ and thus we get

$$\begin{aligned} \|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 &= \|G(Z, \tilde{Z}) - \tau H_{\theta_\eta^*} G(Z, \tilde{Z})\|_2 \\ &= \|(I - \tau H_{\theta_\eta^*}) G(Z, \tilde{Z})\|_2 \\ &\leq \|I - \tau H_{\theta_\eta^*}\|_2 \|G(Z, \tilde{Z})\|_2. \end{aligned}$$

Due to the γ -Lipschitz continuity of the gradient $\nabla \ell$, we have $\|H_{\theta_\eta^*}\|_2 \leq n\gamma$. Using that $\tau \leq \frac{1}{\gamma n}$, we arrive at

$$\|I - \tau H_{\theta_\eta^*}\|_2 \leq 1 + \tau n \gamma \leq 2$$

which, with the help of Lemma 2, yields the final bound for the first-order update

$$\|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 \leq 2R\gamma_z |Z|.$$

□

Proof (Case 2). For the second-order update, we recall that

$$\theta_{Z \rightarrow \tilde{Z}}^* = \theta^* - H_{\theta^*}^{-1} G(Z, \tilde{Z}).$$

Similar to the proof for the first-order update, there exists some $\eta \in [0, 1]$ and a parameter $\theta_\eta^* = \theta^* - \eta H_{\theta^*}^{-1} G(Z, \tilde{Z})$ such that

$$\begin{aligned} \nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D') &= \nabla L(\theta^*, D') \\ &\quad + \nabla^2 L(\theta^* + \eta(\theta_{Z \rightarrow \tilde{Z}}^* - \theta^*), D') (\theta_{Z \rightarrow \tilde{Z}}^* - \theta^*) \\ &= \nabla L(\theta^*, D') - H_{\theta_\eta^*} H_{\theta^*}^{-1} G(Z, \tilde{Z}). \end{aligned}$$

Using again that $\nabla L(\theta^*, D') = G(Z, \tilde{Z})$ we arrive at

$$\begin{aligned} \|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 &= \|G(Z, \tilde{Z}) - H_{\theta_\eta^*} H_{\theta^*}^{-1} G(Z, \tilde{Z})\|_2 \\ &= \|(H_{\theta^*} - H_{\theta_\eta^*}) H_{\theta^*}^{-1} G(Z, \tilde{Z})\|_2 \\ &\leq \|H_{\theta^*} - H_{\theta_\eta^*}\|_2 \|H_{\theta^*}^{-1}\|_2 \|G(Z, \tilde{Z})\|_2 \end{aligned}$$

The λ -strong convexity of L ensures that $\|H_{\theta^*}^{-1}\|_2 \leq \frac{1}{\lambda}$. In addition to $\|G(Z, \tilde{Z})\|_2 \leq R\gamma_z |Z|$, it remains to bound the difference between the Hessians. Using the Lipschitz continuity of the gradient $\nabla^2 \ell$ for $z \in D'$, we first get

$$\begin{aligned} \|\nabla^2 \ell(z, \theta^*) - \nabla^2 \ell(z, \theta_\eta^*)\|_2 &\leq \gamma'' \|\theta^* - \theta_\eta^*\|_2 \\ &\leq \gamma'' \|H_{\theta^*}^{-1}\|_2 \|G(Z, \tilde{Z})\|_2 \end{aligned}$$

and then for the Hessians obtain

$$\begin{aligned} \|H_{\theta^*} - H_{\theta_\eta^*}\|_2 &= \sum_{z \in D'} \|\nabla^2 \ell(z, \theta^*) - \nabla^2 \ell(z, \theta_\eta^*)\|_2 \\ &\leq |Z| \gamma'' \|H_{\theta^*}^{-1}\|_2 \|G(Z, \tilde{Z})\|_2. \end{aligned}$$

Combining all results finally yields the theoretical bound for the second-order update of our approach

$$\begin{aligned} \|\nabla L(\theta_{Z \rightarrow \tilde{Z}}^*, D')\|_2 &\leq \|H_{\theta^*} - H_{\theta_\eta^*}\|_2 \|H_{\theta^*}^{-1}\|_2 \|G(Z, \tilde{Z})\|_2 \\ &\leq |Z| \gamma'' \|H_{\theta^*}^{-1}\|_2^2 \|G(Z, \tilde{Z})\|_2^2 \\ &\leq \frac{\gamma_z^2 \gamma'' R^2}{\lambda^2} |Z|^2 \end{aligned}$$

□

Stochastic Sharding Analysis

To better understand the need for efficient closed-form updates on model parameters, we examine current sharding-based strategies and investigate the circumstances under which they reach their limits. Bourtole et al. [8] propose an unlearning method with the core idea to train separate models on distinct parts of training data. While the authors discuss limitations of their approach like performance degradation and high retraining times we perform a stochastic analysis to find upper bounds for the number of unlearning requests at which retraining becomes as efficient as SISA Training.

In this context, we consider n data instances to unlearn which are uniformly distributed across S shards. Let $p(n)$ denote the probability that all shards contain at least one of these samples which leads to the worst-case scenario of having to retrain all shards. Since calculating $p(n)$ as stated above is difficult, we reformulate the task to solve an equivalent problem: We seek the probability $\hat{p}_k(n)$ that at most k shards remain unaffected by any sample. We set $k = S$

such that $\hat{p}_S(n)$ indicates the probability that any combination of $j \in \{1, \dots, S\}$ shards are unaffected. If this probability is zero there are no unaffected shards. Hence, this corresponds to the inverse of our target probability $p(n) = 1 - \hat{p}_S(n)$.

To calculate $\hat{p}_S(n)$, we first determine the probability of exactly j shards to remain unaffected. In general, there are $(S - j)^n$ combinations to distribute n samples on the dataset excluding j shards. Since there are $\binom{S}{j}$ possible ways to select the j shards to be left out, the total number of combinations is given by $\binom{S}{j}(S - j)^n$. However, we cannot simply sum these terms up for different values of j since the unaffected shards in the combinations partly overlap. To account for this, we apply the inclusion-exclusion principle and finally divide the adjusted term by the number of combinations including all shards:

$$\hat{p}_S(n) = \frac{\sum_{j=1}^S (-1)^{j+1} \binom{S}{j} (S - j)^n}{S^n}$$

Fig. 1 shows that $p(n)$ quickly reaches one even for low numbers of affected samples. Since the probability only depends on the number of shards and samples to unlearn and not on the size of the dataset, we can conclude that sharding is inefficient when there are many unlearning requests. This essentially motivates our approach using closed-form updates on model parameters.