CONCEPT FORGETTING VIA LABEL ANNEALING

Anonymous authors

Paper under double-blind review

ABSTRACT

The effectiveness of current machine learning models relies on their ability to grasp diverse concepts present in datasets. However, biased and noisy data can inadvertently cause these models to be biased toward certain concepts, undermining their ability to generalize and provide utility. Consequently, modifying a trained model to forget these concepts becomes imperative for their responsible deployment. We refer to this problem as *concept forgetting*. Our goal is to develop techniques for forgetting specific undesired concepts from a pre-trained classification model's prediction. To achieve this goal, we present an algorithm called Label ANnealing (LAN). This iterative algorithm employs a two-stage method for each iteration. In the first stage, pseudo-labels are assigned to the samples by annealing or redistributing the original labels based on the current iteration's model predictions of all samples in the dataset. During the second stage, the model is fine-tuned on the dataset with pseudo-labels. We illustrate the effectiveness of the proposed algorithms across various models and datasets. Our method reduces *concept violation*, a metric that measures how much the model forgets specific concepts, by about 85.35% on the MNIST dataset, 73.25% on the CIFAR-10 dataset, and 69.46% on the CelebA dataset while maintaining high model accuracy. Our implementation can be found at this following link: https://anonymous.4open.science/r/LAN-141B/

025

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

031 The superior performance capability of deep learning systems is primarily attributed to their ability 032 to learn various concepts inherent in the dataset. For instance, advancements in face recognition 033 systems (LeCun et al., 1998; Krizhevsky et al., 2009; He et al., 2016) can be largely attributed to their 034 ability to discern and characterize different semantic features from facial images, such as age, gender, and facial hair characteristics, etc. Similarly, in tasks involving image and text generation (Ramesh 035 et al., 2021; 2022; Rombach et al., 2022), the ability to learn varied concepts present in the dataset, enables the generation of realistic and diverse outputs. Consequently, the efficacy of these models 037 relies upon the acquisition of learned concepts inherent within the dataset. Nevertheless, when the dataset is tainted with noisy samples or biased concepts (Tommasi et al., 2017), these models are susceptible to learning such undesired biased concepts. For example, suppose we are learning a 040 model to predict whether a person should get a bank loan or not. Such a model should not depend 041 on the gender or race of the person. However, it is possible that the machine learning model might 042 inadvertently use these features to make predictions, which is highly undesirable. As a result, there 043 emerges a pressing necessity to forget or unbias the undesired biased concept from these trained 044 models to ensure their reliable and accountable deployment. Apart from removing biases, forgetting concepts can prove beneficial in topics such as domain generalization. For example, envision a CelebA (Liu et al., 2015) image classifier that heavily relies on *background color* as a distinguishing 046 feature to classify different celebrities, limiting its ability to generalize effectively. Therefore, in such 047 scenarios, rapidly forgetting only undesired concepts from a pre-trained model, without affecting 048 the ability of the model to use other features, can improve the model's fair decision-making and generalization capabilities. 050

To make a pre-trained model forget a concept, we start by asking the following question - *what is meant by forgetting?* Our definition of forgetting is motivated by the fact that in the case of
 humans, if one forgets a concept, the forgotten concept doesn't affect one's decision-making. Thus
 concept forgetting within the context of machine learning entails ensuring that a model's predictions

become entirely independent of the targeted forgetting concept. However, achieving this task presents challenges, as the goal is to forget a specific undesired concept without adversely affecting the ability of the model to use other concepts. This challenge is underscored by the phenomenon of *catastrophic forgetting* observed in the literature (McCloskey & Cohen., 1989; Goodfellow et al., 2014; Kirkpatrick et al., 2017; Ginart et al., 2019) in similar contexts, where adapting a model for new tasks (in our case task of *concept forgetting*), can significantly degrade performance. Thus, to explore the extent of forgetting concepts in pre-trained models, we pose the following challenge:

Can we efficiently modify a pre-trained model to forget (unbias) an undesired (biased) concept while maintaining its performance?

Before we proceed further, we first state the differences between concept forgetting and machine unlearning, the latter of which has been recently used to remove the effect of certain training examples from the model.



061

062

063 064

065

066

088

Figure 1: Machine unlearning vs concept forgetting: In the first scenario (on the left), such as gender, machine unlearning fails. This failure occurs because the dataset includes only males and females, making it impossible to retrain the model without the gender concept. In the second scenario (on the right), when a user requests the removal of concepts related to "Angelina Jolie," unlearning methods, like the optimal retraining approach, can be used successfully.

Machine unlearning vs. concept forgetting: Machine unlearning (Cao & Yang, 2015; Xu et al., 2020; Nguyen et al., 2022) aims to remove the influence of particular subsets of training data from a model so that the unlearned model mirrors the behavior of a retrained model that is trained from scratch without the undesired data subset. The best method to achieve this is by retraining the model from scratch without the unwanted data, although this process can be computationally expensive (Ginart et al., 2019; Sekhari et al., 2021). However the goal of *concept forgetting* is to make a model's predictions entirely independent of the targeted forgetting concept. The formal definition of concept forgetting is given in Sec. 3.1. Given this, we note that machine unlearning and concept forgetting as two different problem scenarios (see Figure 1). For example, consider the CelebA dataset which contains images of celebrities. Suppose we would like to make the model forget the concept of gender. A machine unlearning approach should remove the influence of all examples that have gender and

produce an unlearned model that is equivalent to retraining the model on the empty dataset as all CelebA dataset images have gender. However, as we show later, concept forgetting can be used to remove the bias caused by gender concept. We do note that machine unlearning can be potentially used to forget small concepts, which are only present in a subset of examples. For example, if we want to eliminate the concept of a particular celebrity, from a classifier trained on the CelebA dataset, we can retrain the model without the images of that celebrity. However, such applications are limited and our proposed algorithm works for removing or unbiasing the dependence of undesired concepts from the model's predictions.

Our contributions: Our contribution can be summarized as follows:

We introduce the framework of *concept forgetting* from pre-trained classification models. Motivated by works in fairness (Dwork et al., 2012; Hardt et al., 2016; Lowy et al., 2021), we measure the bias or the dependence of a model on a concept, based on *concept violation*, which empirically quantifies the extent to which the model remains neutral towards a concept for its predictions.

We propose an algorithm called Label ANnealing (LAN). LAN employs an iterative approach, where each iteration redistributes the class labels of data points containing forgetting concepts to the most probable class labels, thus creating pseudo-labels. This realignment ensures that the distribution of pseudo-labels for each concept value matches the class distribution predicted by the current iteration's model. The method draws an analogy to the term *annealing* frequently employed in material science. It denotes the controlled redistribution of atoms within a solid material under specific temperature conditions to attain an equilibrium state, which inspired our method's nomenclature. This strategy

not only aids in mitigating the influence of the undesired concept on the model's predictions but also
 is computationally efficient. This method necessitates minimal epochs, sometimes as few as a single
 epoch, to diminish the reliance of the model's predictions on the forgetting concept, all the while
 maintaining the model's overall performance and generalization ability.

112 • We demonstrate the efficacy of our algorithms through detailed evaluations on various image 113 classification datasets such as MNIST (LeCun et al., 1998), CIFAR-10 (Krizhevsky et al., 2009), 114 miniImageNet (Vinyals et al., 2016), and CelebA (Liu et al., 2015) using state-of-the-art image classi-115 fication models such as MobileNetV2, DenseNet-121, ResNet-50. Our method reduces (averaged 116 over several concepts) concept violation, a metric that measures how much the model forgets specific 117 concepts, by about 85.35% on the MNIST dataset, 73.25% on the CIFAR-10 dataset, 17.05% on 118 the miniImageNet dataset, and 69.46% (averaged over 81.34% for binary concepts and 63.52% for multi-level concepts forgetting) on the CelebA dataset while maintaining high model accuracy. 119

120 121

122

2 RELATED WORKS

123 124 2.1 FAIRNESS

125 Fairness in machine learning systems is an important research area aimed at ensuring that system 126 predictions are both accurate and fair across different groups (based on their features) of data points. 127 Earlier works (Dwork et al., 2012) initially proposed the notion of demographic parity as a preliminary definition of fairness. According to this concept, a machine learning algorithm satisfies demographic 128 parity if the predicted target is independent of sensitive attributes. However, promoting demographic 129 parity may lead to diminished performance, particularly if the true outcome is not independent of 130 sensitive attributes. To address this, subsequent works (Hardt et al., 2016) introduced a relaxed notion 131 of fairness based on equalized odds and equal opportunity definition. Recent works (Kamishima 132 et al., 2011; Feldman et al., 2015; Zafar et al., 2017; Donini et al., 2018; Mary et al., 2019; Cho et al., 133 2020a; Jiang et al., 2020; Rezaei et al., 2020; Lowy et al., 2021) have explored incorporating different 134 fairness notions during the training process itself. These methods incorporate regularization-based 135 techniques based on different statistical distances between the distribution of the model's prediction 136 and sensitive attributes.

137 Drawing inspiration from fairness notions, especially demographic parity (Dwork et al., 2012; Lowy 138 et al., 2021), we propose that forgetting a particular concept can also be interpreted as achieving 139 independence between the model's prediction and the undesired feature we aim to forget. However, 140 methods focusing on enforcing fairness with respect to the forgetting concept require a large number 141 of epochs to converge and can be computationally inefficient. For instance, according to state-of-the-142 art FERMI algorithm (Lowy et al., 2021), achieving $||\nabla \ell(\theta, x, y)|| \le \epsilon$ where ℓ is the loss function 143 requires approximately $\mathcal{O}(\frac{1}{4})$ iterations. Empirically, the convergence of the FERMI algorithm 144 varies depending on the dataset and application, typically ranging from as few as 50 to as many 145 as 2000 epochs (Lowy et al., 2021, Appendix E). Given our specific objective of forgetting only certain concepts from a model's parameters without affecting others, we aim to devise a more 146 computationally efficient approach for concept forgetting from pre-trained models. 147

148 149

2.2 MACHINE UNLEARNING FOR TRAINING EXAMPLE REMOVAL

150 Machine unlearning, as described in the literature (Xu et al., 2020; Nguyen et al., 2022; Cao & Yang, 151 2015; Ginart et al., 2019; Golatkar et al., 2020a;b; 2021; Neel et al., 2021; Nguyen et al., 2020; Guo 152 et al., 2020; Graves et al., 2021; Sekhari et al., 2021), involves intentionally erasing the impact of 153 particular subsets of training data from a pre-trained model, addressing user privacy concerns. Here, 154 the objective is to craft a computationally efficient method that produces an unlearned model that 155 mirrors the behavior of the model that is trained from scratch, on the training dataset devoid of the 156 sensitive data points. Although retraining serves as the optimal benchmark method for unlearning, this 157 method becomes computationally impractical for large models and iterative unlearning demands (Cao 158 & Yang, 2015; Ginart et al., 2019; Sekhari et al., 2021). Consequently, to address user privacy 159 concerns, more efficient data deletion methods (Cao & Yang, 2015; Ginart et al., 2019) were devised, leading to the emergence of *machine unlearning*. The *machine unlearning* methods are broadly 160 categorized into two types: exact unlearning (Wu et al., 2020) and approximate unlearning (Neel 161 et al., 2021; Sekhari et al., 2021). Exact unlearning aims to completely eliminate the influence of

162 unwanted data from the trained model, while *approximate unlearning* methods only partially mitigate 163 data influence, resulting in parameter distributions closely resembling the retrained model with minor 164 adjustments. More sophisticated methods (Guo et al., 2020; Graves et al., 2021) have suggested using 165 influence functions, but these are computationally demanding and limited to small convex models. 166 To extend unlearning techniques to non-convex models like deep neural networks, (Golatkar et al., 2020a;b) introduced a scrubbing mechanism centered on the Fisher Information matrix. 167

168 As we noted earlier, concept forgetting and machine unlearning have fundamental differences (169 Figure 1 demonstrates machine unlearning cannot be applied for a general concept forgetting setup) in 170 their objectives. Machine unlearning seeks to forget specific data points while emulating the behavior 171 of a retrained model, whereas *concept forgetting* aims for the model's predictive performance to 172 become independent of the forgotten concept.

173 174

175

191

192

193

194

195 196 197

200 201 202

203

204

PRELIMINARIES AND BACKGROUND 3

176 3.1 PROBLEM FORMULATION 177

178 Unless otherwise specified, we consider the problem of multi-class classification throughout the paper. Let $\mathcal{Y} \triangleq \{0, 1, 2, \dots, k-1\}$ denote the set of k labels. Let z = (x, y) denote a data point 179 where $x \in \mathbb{R}^d$ is the feature and $y \in \mathcal{Y}$ is the label. A dataset $\mathcal{D} \triangleq \{z_i\}_{i=1}^{|\mathcal{D}|}$ is a set of samples sampled from the underlying data distribution $P_{xy} : \mathbb{R}^d \times \mathcal{Y} \to [0, 1]$. Let a categorical concept 180 181 182 $\mathcal{C}: \mathbb{R}^d \times \mathcal{Y} \to \{0, 1, 2, \dots, m-1\}$ be a mapping from the sample to the set of all possible values 183 the concept can take. For example, if the concept is binary such as beard, it can take two values $\{0,1\}$ (m = 2), which denotes the absence and presence of the beard, respectively. Similarly, if 185 the concept is non-binary such as facial hair type, it can take multiple values $\{0, 1, 2, 3\}$ (m = 4)which signifies no facial hair, mustache, beard, and goatee respectively. Let $h_{\theta} : \mathbb{R}^d \to \Delta^{|\mathcal{Y}|}$ denote a 186 classifier parameterized by $\theta \in \mathbb{R}^p$ where Δ is the probability simplex. This classifier takes a feature 187 $x \in \mathbb{R}^d$ and predicts a distribution over the label space. Let $\hat{h}(\theta, z)$ denote a post-processing step on 188 the classifier (e.g. argmax) where a hard label is inferred based on the probability over the labels. 189

190 **Definition 1.** (Concept neutral): We call a classifier with parameter θ concept neutral with respect to a concept C, if for all output class $y \in \mathcal{Y}$ and all possible concept values $c \in \{0, 1, 2, \dots, m-1\}$,

> $P_{xy}(\hat{h}(\theta, z) = y | \mathcal{C}(z) = c) = P_{xy}(\hat{h}(\theta, z) = y).$ (1)

Definition 2. (*Concept violation*): For a classifier \hat{h} with a parameter θ , we measure the violation of concept neutrality in terms of the total variation distance as follows:

$$V(\theta, \mathcal{C}, P) \triangleq \frac{1}{m} \sum_{c=0}^{m-1} d_{TV} \Big(P_{xy}(\hat{h}(\theta, z) = y), P_{xy}(h(\theta, z) = y \mid \mathcal{C}(z) = c) \Big) \\ = \frac{1}{2m} \sum_{c=0}^{m-1} \sum_{y=0}^{k-1} \Big| P_{xy}(\hat{h}(\theta, z) = y) - P_{xy}(\hat{h}(\theta, z) = y \mid \mathcal{C}(z) = c) \Big|.$$
(2)

Note that $V(\theta, \mathcal{C}, P) \in [0, 1]$ and if a model is concept neutral, then $V(\theta, \mathcal{C}, P) = 0$. As the underlying data distribution P_{xy} is unknown, we have only access to the dataset \mathcal{D} to empirically estimate concept violation $V(\theta, C, P)$ as follows:

$$\hat{V}(\theta, \mathcal{C}, \mathcal{D}) \triangleq \frac{1}{m} \sum_{c=0}^{m-1} d_{\text{TV}} \left(\hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y \mid \mathcal{C}(z) = c) \right)$$
$$= \frac{1}{2m} \sum_{c=0}^{m-1} \sum_{y=0}^{k-1} \left| \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y) - \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y \mid \mathcal{C}(z) = c) \right|,$$
(3)

210 211 212

213 where $\hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y) = \frac{1}{|\mathcal{D}|} \sum_{z \in \mathcal{D}} \mathbb{1}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y), \mathcal{D}_c = \{z \in D : C(z) = c\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{h}(\theta, z) = y\}, \text{ and } \hat{P}_{\mathcal{D}}(\hat{$ 214

215
$$y|\mathcal{C}(z) = c) = \frac{1}{|\mathcal{D}_c|} \sum_{z \in \mathcal{D}_c} \mathbb{1}(\hat{h}(\theta, z) = y)$$
. Now given a pre-trained model with parameter θ^* and a

concept C, the goal of a concept forgetting algorithm is to find the forgotten parameter θ_C such that the algorithm has the following properties:

• Minimize empirical concept violation: The empirical concept violation metric $\hat{V}(\theta_{\mathcal{C}}, \mathcal{C}, \mathcal{D})$ measures how much 'neutral' is the forgotten model for the given concept. For an ideal forgotten model, this metric will be zero indicating that the model has forgotten the concept. Hence minimizing concept violation is an important criterion and our goal is to ensure $\hat{V}(\theta_{\mathcal{C}}, \mathcal{C}, \mathcal{D}) \ll \hat{V}(\theta^*, \mathcal{C}, \mathcal{D})$.

• Minimize accuracy loss: Any forgotten model θ_c should exclusively erase the specified concept without erasing others, thereby enabling the retained generalization capabilities to persist. Hence minimizing loss of the forgotten model's test accuracy is one of the important criteria. Let $\ell(\theta, z)$ denote the loss of the model with parameters θ for a sample $z = (x, y) \in \mathcal{D}$. Hence our goal is to ensure $\sum_{z \sim P_{xy}} [\ell(\theta_c, z)] \approx \sum_{z \sim P_{xy}} [\ell(\theta^*, z)].$

• Small time complexity: In dynamic environments, rapid model adaptation and updating are vital to remove biases or outdated information. Concept forgetting algorithms aim to selectively forget a few concepts from pre-trained models without affecting others. Efficiency is critical, as prolonged training may erase previously learned concepts.

4 METHODOLOGY

4.1 LABEL ANNEALING (LAN) ALGORITHM



Figure 2: Label Annealing (LAN) methodology - The task involves forgetting the concept $C(z) = c \in \{0, 1, 2\}$ from a classification task with data points labeled as $j \in \{\text{Class-0}, \text{Class-1}, \text{Class-2}\}$ (blue, red, and yellow, respectively). This iterative method runs E iterations, where each e^{th} -iteration constitutes two stages: in first stage, known as *label annealing subroutine*, the labels within each concept data subset (e.g., $\mathcal{D}_0, \mathcal{D}_1, \mathcal{D}_2$) are redistributed based on the class prediction of e^{th} -iteration's model θ_e , denoted as $\hat{p}_{\theta_e}(x, j)$, resulting in the label annealed dataset $\tilde{\mathcal{D}}$. Subsequently at the next stage, termed as *parameter fine-tuning* using $\tilde{\mathcal{D}}$, we minimize the loss function $\mathcal{L}_{\text{LAN}}(\theta, \tilde{\mathcal{D}})$ to obtain final the concept forgotten model θ_c .

To achieve the forgotten model $\theta_{\mathcal{C}}$, we devise a method called Label Annealing (LAN). The overall methodology is shown in the Figure 2. At the heart of this algorithm is the *label annealing* subroutine, given in Algorithm 1. Given a model parameter θ_e , training dataset \mathcal{D} , and particular concept \mathcal{C} targeted for forgetting, this subroutine at a particular iteration e creates a dataset with the same features as $x_i \in \mathcal{D}$ and with pseudo-labels \tilde{y}_i such that the model θ_e has zero concept-violation on the newly created dataset \mathcal{D} . Further to retain the model's overall performance, this assignment of pseudo-labels must result in a minimal change in empirical risk. Thus, we would like to change labels for the samples where changing the label does not significantly change the loss. To achieve these goals, the whole dataset is divided into concept data sub-groups \mathcal{D}_c for each $c \in \{0, 1, ..., m-1\}$.

Algor	ithm 1 : Label annealing subroutine	
	Input : model parameter $\theta_e \in \mathbb{R}^p$, dataset \mathcal{D} , forgetting	
(concept C	
	1: For each class $j \in [k-1]$, let b_j denote the number	
	of samples $z \in \mathcal{D}$ with $\arg \max_j p_{\theta_e}(x, j) = j$.	
	2: for $c = 0, 1, \dots, m - 1$ do	
	3: construct $\mathcal{D}_c = \{z \in \mathcal{D} : \mathcal{C}(z) = c\}, n_c = \mathcal{D}_c $ 4: For each sample $z \in \mathcal{D}$ the maximum	
	4. For each sample $z_i \in D_c$, the maximum probability assigned to sample $n = (x_i)$ –	
	max _i $p_{\theta_c}(x_i, j)$. Sort \mathcal{D}_c in decreasing order of	
	$p_{\max}(x)$.	Algorithm 2 : Pa
	5: Let $n_{c,j}$ denote the number of samples z in \mathcal{D}_c	Input: pre
	such that $\arg \max_j p_{\theta_e}(x, j) = j$.	$A^* \subset \mathbb{R}^p$ data
	6: Let $\alpha_{c,j}$ be the number of samples in \mathcal{D}_c assigned with along i by the algorithm. Set $\alpha_{c,j} = 0 \forall i \in \mathcal{D}_c$	$v \in \mathbb{R}^{2}$, data needs to be t
	with class-j by the algorithm. Set $\alpha_{c,j} = 0 \forall j \in [k-1]$	size B learnin
	7: for $i = 1, 2,, \mathcal{D}_c $ do	of iterations E
	8: $\widetilde{y_i} \leftarrow \phi$	of netations E,
	9: $j^* \leftarrow \arg \max p_{\theta_e}(x_i, j)$	1. Initialize. A
	$i0: \qquad if \alpha x < h \\ h \\ n \\ n / \mathcal{D} $ then	2: for $e = 1$
	11: $\widetilde{u_c,j^*} < \delta_{j^*} + h_c/ \mathcal{D} $ then $\widetilde{u_i} \leftarrow j^*$	$2: IOF \bigcirc 1, \dots$ $3: \widetilde{D} \leftarrow IAN$
	12: $\alpha_{c,\tilde{y_i}} \leftarrow \alpha_{c,\tilde{y_i}} + 1$	$\begin{array}{ccc} \mathbf{J} & \mathbf{D} \leftarrow \mathbf{LAN} \\ \mathbf{A} & \mathbf{for} \ t = 1 \end{array}$
	13: end if	$\begin{array}{ccc} 4 & 10 & t = 1, \\ 5 & \mathbf{Draw a t} \end{array}$
	14: end for	size B fr
	15: for $i = 1, 2, \dots, \mathcal{D}_c $ do	$6: \theta_{2} \leftarrow \theta_{2}$
	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	7: end for
	18: while $\tilde{y}_i == \phi \mathbf{do}$	8: end for
	19: $j_i \leftarrow \arg \max p_{\theta_e}(x_i, j)$	9: Output: θ_E
	$j \in \{0, \dots, k-1\} \setminus j^i$	
-	20: if $\alpha_{c,j_i} < b_{j_i} \cdot n_c / \mathcal{D} $ then	
	21. $y_i \leftarrow j_i, \alpha_{c,\bar{y}_i} \leftarrow \alpha_{c,\bar{y}_i} + 1$ 22. else	
-	23: $j^i \leftarrow j^i \mid j \{j_i\}$	
2	24: end if	
2	25: end while	
	26: end if	
	2/: end for	

29: **Output**: $\widetilde{\mathcal{D}} \leftarrow \{\widetilde{z}_i = (x_i, \widetilde{y}_i)\}_{i=1}^{|\mathcal{D}|}$

301

 Algorithm 2 : Parameter fine-tuning

 Input:
 pre-trained
 parameter

 $\theta^* \in \mathbb{R}^p$, dataset \mathcal{D} , concept that
 needs
 to be forgotten \mathcal{C} , batch

 size B, learning rate η , number
 of iterations E, number of steps T.
 1: Initialize: $\theta_e \leftarrow \theta^*$

 1: Initialize:
 $\theta_e \leftarrow \theta^*$ 2: for $e = 1, \dots, E$ do
 3:
 $\widetilde{\mathcal{D}} \leftarrow LAN(\theta_e, \mathcal{D}, \mathcal{C})$

 4:
 for $t = 1, \dots, T$ do

5: Draw a random mini-batch of size *B* from \tilde{D} denoted as \tilde{D}^{b} 6: $\theta_{e} \leftarrow \theta_{e} - \eta \nabla_{\theta} \mathcal{L}(\theta, \tilde{D}^{b})$ 7: end for

Now for each of \mathcal{D}_c , the first term in Eq. 3 for a particular class label j would be $\frac{b_j}{|D|}$ and second term 303 would be $\frac{n_{cj}}{|n_c|}$ where $n_c = |\mathcal{D}_c|$, b_j and n_{cj} are the nos. of samples of class-j predicted by the current model θ_e on dataset \mathcal{D} and \mathcal{D}_c respectively. In other words to make concept violation zero in \mathcal{D}_c , 305 nos. of samples predicted class-j in \mathcal{D}_c i.e. n_{cj} must be equal to $b_j \cdot \frac{n_c}{|D|}$. This is why we need to 306 redistribute the labels of each class-j in \mathcal{D}_c without much affecting the model performance (empirical 307 loss). Thus to achieve this dual objective of redistributing the labels without much affecting the 308 empirical loss, we calculate $p_{\max}(x_i) = \max_j p_{\theta_e}(x_i, j)$ for each sample $z_i = (x_i, y_i) \in \mathcal{D}_c$, and then \mathcal{D}_c is sorted in decreasing order of $p_{\max}(x)$. Thus in this sorted \mathcal{D}_c (in the second for loop), 310 each sample x_i is initialized with label $\tilde{y}_i = \phi$ and iteratively assigned the most probable label 311 $\tilde{y}_i = j^* = \arg \max_j p_{\theta_e}(x_i, j)$ until the no of samples in class- j^* is less than $b_{j^*} \cdot \frac{n_e}{|c|}$. This second 312 for loop ensures that reassigned labels don't change from the initial ones (this is why the deterministic 313 assignment is done) on those data points where the model is confident (this is why \mathcal{D}_c is sorted). 314 Further in the subsequent steps (third for loop), the data points where the labels are unassigned i.e. 315 $y_i = \phi$, it is assigned the subsequent (second or third and so on) most probable label class-j if the 316 no of samples in that assigned class-j is less than $b_j \cdot \frac{n_c}{|D|}$. This loop tries to redistribute the labels 317 where the model is not confident (low concept violation is achieved by trading off some accuracy). 318 Subsequently, in the next stage of *parameter fine-tuning* (Algorithm 2), we fine-tune the e^{th} -model θ_e on the new dataset $\widetilde{\mathcal{D}} = \bigcup_{c=0}^{m-1} \widetilde{\mathcal{D}}_c$ to obtain the forgotten model θ_{e+1} by minimizing the Label 319 Annealing loss function $\mathcal{L}_{\text{LAN}}(\theta, \widetilde{\mathcal{D}}) = \frac{1}{|\widetilde{\mathcal{D}}|} \sum_{c=0}^{m-1} \sum_{i=1}^{|\widetilde{\mathcal{D}}_c|} \ell(\theta, \widetilde{z}_i)$. We repeat this process for E steps to get the final concert as the loss of E steps to get the loss of E steps to get the final concert as the loss of E steps to get the loss of E step 320 321 to get the final concept-neutral model θ_{C} . The value of E depends on the user's choice. However, to 322 achieve concept forgetting with low computational complexity we experimented with E=1 (results of 323 Table 1 and Table 2). Further ablation studies on E = 2 and E = 4 are given.

3243254.2 THEORETICAL ANALYSIS

In this section, we theoretically show that the proposed algorithm retains its accuracy if the original model has low concept violation. Recall that θ^* denotes the input to Algorithm 2 and θ_c denotes the output of our algorithm. With this notation, we show the following result.

Theorem 1. Let the loss function be bounded i.e., $\forall \theta, z \ \ell(\theta, z) \leq L$. If the fine-tuning reduces the loss on $\widetilde{\mathcal{D}}$ i.e., $\mathbb{E} \left[\mathcal{L}_{\widetilde{\mathcal{D}}}(\theta_{\mathcal{C}}) \right] \leq \mathcal{L}_{\widetilde{\mathcal{D}}}(\theta^*)$, then

$$\mathbb{E}\left[\mathcal{L}_{\mathcal{D}}(\theta_{\mathcal{C}})\right] \le \mathcal{L}_{\mathcal{D}}(\theta^*) + 4 \cdot L \cdot E \cdot m \cdot \hat{V}(\theta^*, \mathcal{C}, \mathcal{D}),\tag{4}$$

334 where the expectation is over the randomization in the stochastic gradients in Algorithm 2.

The above bound implies that if the original concept violation is small, then the performance of the new model (trained on $\tilde{\mathcal{D}}$) will not degrade significantly. In particular, if the original concept violation is zero, then the loss of the forgotten model is the same as the loss of the original model. Furthermore, while the upper bound degrades with E, as we show in experiments, the performance improves or remains the same with an increasing value of E. Due to space constraints, we provide the proof of the above theorem in Appendix A.

341 342 343

344

330

331 332

333

335

336

337

338

339

340

5 EXPERIMENTS AND RESULTS

345 5.1 DATASETS AND MODELS346

For our experiments, we consider mainly forgetting two types of concepts: *binary-level concept* and *multi-level concept*. We have used different image classification models such as 2-layer-MLP (hidden layer size 500), Mobinetv2 (Sandler et al., 2018), Densenet-121 (Huang et al., 2017), Resnet-50 (He et al., 2016). Further to show the applicability of our method for different classification tasks across diverse datasets, we have used MNIST (LeCun et al., 1998), CIFAR-10 (Krizhevsky et al., 2009), miniImageNet (Vinyals et al., 2016), and CelebA (Liu et al., 2015) datasets. Different concept forgetting scenarios for E = 1 can be seen from Table 1 and Table 2. Further details about the datasets and models are included in the appendix section B.1.

354 355 356

357

358

359

5.2 EVALUATION METRICS

To evaluate the efficacy of any concept-forgetting algorithm we propose two different metrics as defined below:

• Test empirical concept violation: This metric denoted as $\hat{V}(\theta_{\mathcal{C}}, \mathcal{C}, \mathcal{D})$, is defined in equation 5, quantifies the concept neutrality of the forgotten model $\theta_{\mathcal{C}}$. Observe that $\hat{V}(\theta_{\mathcal{C}}, \mathcal{C}, \mathcal{D}) \in [0, 1]$, and a smaller $\hat{V}(\theta_{\mathcal{C}}, \mathcal{C}, \mathcal{D})$ signifies that the model is conceptually neutral regarding the forgetting concept \mathcal{C} . In the rest of the section, we denote $\hat{V}(\theta_{\mathcal{C}}, \mathcal{C}, \mathcal{D})$ as $\hat{V}_{\mathcal{C}}$.

• Test accuracy: This metric evaluates the generalization performance of the forgotten model, denoted as A_D . Any concept forgetting algorithm mustn't render the initial model ineffective during the forgetting process. Therefore, maintaining accuracy close to that of the initial model θ^* is desirable.

364

365

5.3 BASELINES

370 According to our knowledge, this is the first work that introduces *concept forgetting* as a property 371 of the forgotten model to induce independence from the forgetting feature during its prediction 372 task. Thus for proper evaluation of our method, we adopt several baselines from fairness because 373 these baseline methods also advocate for the independence of prediction and sensitive concept 374 features. Here we have used particularly three baseline methods: (a) FERMI (Lowy et al., 2021) (b) 375 Continuous-Fairness (Mary et al., 2019) and (c) Fairness-KDE (Cho et al., 2020b). We have used 376 official implementation for both FERMI and Continuous-Fairness baselines while for Fairness-KDE an open-source implementation has been used. Further details about the baselines can be found in the 377 appendix section B.3.2.

5.4 **BINARY-LEVEL CONCEPT FORGETTING**

We evaluated our approach for different classification scenarios to forget binary concepts with $c \in \{0,1\}$ (m = 2). For example, as illustrated in Table 1, in the context of the MNIST digit classification problem, the objective is to forget a particular class digit concept e.g. class-3 data. Thus here c = 0 represents concepts of non-digit-3 data and c = 1 represents concepts of digit-3 data. Similarly, in the CelebA dataset for gender concept c = 0 represents male and c = 1 represents female. Table 1 shows the efficacy of our method for different concept-forgetting scenarios. In this case, the average reduction of concept violation is about 85.35% on the MNIST dataset, 73.25% on the CIFAR-10 dataset, 17.05% on the miniImageNet dataset, and 81.34% on the CelebA dataset, while retention of high model accuracy.

Table 1: Empirical concept violation $\hat{V}_{\mathcal{C}}(\downarrow)$ and test accuracy $A_D(\uparrow)$ of the initial model and forgotten model via LAN. For the forgotten model, \hat{V}_{C} reduced without significantly reducing A_{D}

Datacat	Models	Task	Concept	Initial Model		LAN	
Dataset				\hat{V}_{C}	A_D	\hat{V}_{C}	A_D
		Young or not	Gender	0.117	0.898	0.015	0.847
CelebA	Resnet-50	Attractive or not	Gender	0.2219	0.827	0.006	0.767
		Heavy makeup or not	Gender	0.314	0.919	0.127	0.764
			Triceratops	0.4991	0.9791	0.406	0.951
miniImageNet	Resnet-50	class 0-99 classification	Bugs	0.4966	0.9791	0.364	0.936
			Fences	0.4948	0.9791	0.466	0.96
			Bird	0.440	0.928	0.103	0.871
	Mobinet-v2	class 0-9 classification	Frog	0.473	0.921	0.108	0.855
CIFAR-10			Truck	0.472	0.921	0.113	0.855
			Bird	0.445	0.923	0.152	0.869
	Densenet-121	class 0-9 classification	Frog	0.473	0.917	0.116	0.878
			Truck	0.472	0.917	0.147	0.861
			digit-3	0.479	0.974	0.055	0.883
	2-layer MLP	digit 0-9 classification	digit-5	0.491	0.971	0.104	0.901
MNIST		-	digit-8	0.470	0.976	0.081	0.889
			digit-3	0.498	0.990	0.047	0.893
	Resnet-50	digit 0-9 classification	digit-5	0.492	0.992	0.078	0.905
		-	digit-8	0.496	0.991	0.063	0.897

As there exists a trade-off between the two metrics of interest, for proper evaluation of our method with FERMI (Lowy et al., 2021), Continuous-Fairness (Mary et al., 2019), and Fairness-KDE (Cho et al., 2020b) baselines, concept-violation vs. accuracy trade-off plots are depicted in Figure 3. From these plots, it can be seen that for a particular accuracy, our method achieves lower concept violation (LAN trade-off curve lies below the others) than other baseline methods achieving better trade-off.

5.5 MULTI-LEVEL CONCEPT FORGETTING

Here the concept mapping function $\mathcal{C}(.)$ denotes the multi-level concept with $c \in \{0, 1, \ldots, m-1\}$ (m > 2). Table 2 illustrates the performance of LAN—in reducing concept violation and maintaining

Table 2: Empirical concept violation $\hat{V}_{\mathcal{C}}(\downarrow)$ and test accuracy $A_{\mathcal{D}}(\uparrow)$ of the initial model and forgotten model via LAN. For the forgotten model, \hat{V}_{C} reduced without significantly reducing A_{D}

424	Tesha	Concepts	Initial Model		LAN	
425	Tasks		\hat{V}_{C}	$A_{\mathcal{D}}$	$\hat{V}_{\mathcal{C}}$	$A_{\mathcal{D}}$
426		Hair Color	0.2	0.898	0.063	0.8626
427	Young vs. Not-Young	Facial Hair	0.11	0.897	0.0329	0.8921
428		Hair Color	0 195	0.827	0.083	0 7955
429	Attractive vs. Not-attractive	Facial Hair	0.1716	0.827	0.076	0.8088
430		Hair Color	0.157	0.92	0.073	0.881
431	Heavy Makeup vs. Not-Heavy Makeup	Facial Hair	0.316	0.919	0.077	0.844

test accuracy across various settings of concept forgetting from a pre-trained Resnet-50 model trained on the CelebA (Liu et al., 2015) dataset. In this context, concept forgetting involves the removal of certain features from a pre-trained classifier in the process of classifying other features. For example, while classifying samples as young vs. not-young, we aim to forget subtle feature concepts such as hair color, and facial hair from the pre-trained models. In this setting, the LAN algorithm reduces concept violation by about 63.52% without significantly affecting test accuracy. It can be seen from Figure 3 (a) and (b) that our method performs significantly better than other baseline methods.



Figure 3: Concept violation vs. accuracy trade-off: We have plotted concept violation on the y-axis and accuracy on the x-axis. Each point represents an algorithm with a hyper-parameter, and the *Fit* line for an algorithm is obtained by a linear fit of all experiments corresponding to the algorithm. Figures (a) and (b) show forgetting facial hair removal from the task of heavy makeup and attractiveness classification on CelebA. Figures (c) and (d) show different concepts forgotten from pre-trained Resnet-50 on the miniImageNet dataset. Figures (e) and (f) show concept forgetting from pre-trained Densenet-121 and Mobinetv2, respectively, on the CIFAR-10 dataset. Figures (g) and (h) show class-3 concept forgetting from pre-trained MLP and Resnet-50 models, respectively, on the MNIST dataset. It can be seen that increasing accuracy increases concept violation. Thus, for a particular achievable accuracy, LAN achieves lower concept violation than other baseline methods.

5.6 ABLATION STUDY

Table 3 demonstrates the performance of the LAN method for different learning rates. As it can be seen as the learning rate increases the accuracy decreases while the concept violation decreases and then starts increasing again. Thus at higher accuracy regions, concept violation decreases along with accuracy whereas at lower accuracy regions concept violation increases with a decrease in accuracy. This suggests an optimal point lies in the trade-off curve where concept violation is low with a slight reduction of accuracy. Further in Figure 4, we demonstrate the effectiveness of LAN over multiple iterations (E=2, E=4). We present concept violation vs. accuracy trade-off plot to forget the facial



486 Table 3: Empirical concept violation $\hat{V}_{\mathcal{C}}$ and accuracy 487

> hair concept while classifying attractive vs. not attractive on the CelebA dataset. As E increases at higher accuracy regions, the concept violation further decreases for the same accuracy value making the trade-off plot flatter.

6 CONCLUSION

In the pursuit of safer and more responsible machine learning, the elimination of undesired concepts from models is crucial. Our work focuses on efficiently removing these undesired concepts from 510 pre-trained classification models, a task that is challenging due to the risk of *catastrophic forgetting*, 511 which can render the model ineffective. To address this, we propose a computationally efficient 512 algorithm termed LAN (Label Annealing) to create a forgotten model while preserving its ability to generalize. We define concept forgetting as the property of a model to disregard undesired concepts 513 during its decision-making process and introduce *concept neutrality* as a necessary attribute of a 514 forgotten model. To quantify the extent of *concept neutrality* in any model, we propose a novel metric 515 called concept violation. Our experimental results demonstrate that our method effectively reduces 516 concept violation while maintaining the model's performance across multiple concept-forgetting 517 settings, various models, and datasets. Additionally, we acknowledge that our definition and method 518 apply only to concept forgetting in classification models. Further research is needed to develop 519 definitions and methods for concept forgetting that generalize to generative models as well. 520

REFERENCES

- Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In Proc. of IEEE Symposium on Security and Privacy, 2015.
- J. Cho, G. Hwang, and Suh. A fair classifier using mutual information. In 2020 IEEE International344 Symposium on Information Theory (ISIT), pp. 2521–2526. IEEE, 2020a.
- Jaewoong Cho, Gyeongjo Hwang, and Changho Suh. A fair classifier using kernel density estimation. In Proc. of Advances in Neural Information Processing Systems, 2020b.
- Donini, Oneto L., Ben-David S., Shawe-Taylor, J. S., and Pontil M. Empirical risk minimization under fairness constraints. In Advances in Neural Information Processing Systems, pp. 2791–2801, 2018.
- C. Dwork, M. Hardt, Pitassi T., O. Reingold, , and R Zemel. Fairness through awarenes. In *Proceedings of the 3rd innovations in theoretical computer science conference*, 2012.
- M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian. Certifying and 538 removing disparate impact. In proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining, pp. 259–268, 2015.

521

522 523

524

525 526

527

528

529

530 531

532

533

534

536

537

501

502

503

550

569

570 571

572

- Antonio Ginart, Melody Guan, Gregory Valiant, and James Y. Zou. Making ai forget you: Data deletion in machine learning. *In Proc. of NIPS*, 2019.
- Aditya Golatkar, Alessandro Achille, and Stefano 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*, pp. 9304–9312, 2020a.
- Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Forgetting outside the box: Scrubbing deep networks of information accessible from input-output observations. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*, pp. 383–398. Springer, 2020b.
- Aditya Golatkar, Alessandro Achille, Avinash Ravichandran, Marzia Polito, and Stefano Soatto.
 Mixed-privacy forgetting in deep networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 792–801, 2021.
- Ian J. Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural network. *In Proc. of International Conference on Learning Representations*, 2014.
- Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11516–11524, 2021.
- Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens van der Maaten. Certified data removal from machine learning models. *In Proc. of ICML*, 2020.
- Hardt, E. Price, and N Srebro. Equality of opportunity in supervised learning. In Advances in Neural Information Processing Systems, 2016.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
 - Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. *In Proc. of CVPR*, 2017.
 - R. Jiang, A. Pacchiano, T. Stepleton, H. Jiang, and S. Chiappa. Wasserstein fair classification. *In Uncertainty in Artificial Intelligence, pp.* 862–872. *PMLR*, 2020.
- T. Kamishima, S. Akaho, and J. Sakuma. Fairness-aware learning through regularization approach. *In* 2011 IEEE 11th International Conference on Data Mining Workshops, pp. 643–650. IEEE, 2011.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Pre-print arXiv*, 2017.
- 581 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In
 Proceedings of International Conference on Computer Vision (ICCV), 2015.
- Andrew Lowy, Rakesh Pavan, Sina Baharlouei, Meisam Razaviyayn, and Ahmad Beirami. A stochastic optimization framework for fair risk minimization. *arXiv*, 2021.
- J Mary, C. Calauzenes, and N. El Karoui. Fairness-aware learning for continuous attributes and treatments. *In International Conference on Machine Learning*, 2019.
- 593 Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of Learning and Motivation*, 1989.

- Seth Neel, Aaron Roth, and Saeed Sharifi-Malvajerdi. Descent-to-delete: Gradient-based methods for machine unlearning. In Proc. of ALT, 2021.
- Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Patrick Jaillet. Variational bayesian unlearning. In Proc. of NeuRIPS, 2020.
- Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning. arXiv preprint arXiv:2209.02299, 2022.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In International Conference on Machine Learning, pp. 8821–8831. PMLR, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125, 1(2):3, 2022.
- A. Rezaei, R. Fathony, O. Memarrast, and B. D. Ziebart. Fairness for robust log loss classification. In AAAI, pp. 5511–5518, 2020.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-ence on computer vision and pattern recognition, pp. 10684–10695, 2022.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo-bilenetv2: Inverted residuals and linear bottlenecks. In Proc. of CVPR, 2018.
- Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what you want to forget: Algorithms for machine unlearnings. In Proc. of NeurIPS, 2021.
- Tatiana Tommasi, Novi Patricia, Barbara Caputo, and Tinne Tuytelaars. Advances in Computer Vision and Pattern Recognition. Springer, 2017.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Match-ing networks for one shot learning. In Proc. of Neural Information Processing System, 2016.
- Yinjun Wu, Edgar Dobriban, and Susan B. Davidson. Deltagrad: Rapid retraining of machine learning models. In Proc. of ICML, 2020.
- Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S. Yu. Machine unlearning: A survey. ACM Computing Surveys Vol. 56, No. 1, 2020.
- Zafar, I. Valera, M. G. Rogriguez, and K. P. Gummadi. Fairness constraints: Mechanisms for fair classification. In Proc. of Artificial Intelligence and Statistics, 2017.

648 A ANALYSIS

We recall some of the notation used in the algorithm. b_j denotes the number of samples of class-*j* in \mathcal{D} predicted by the initial model. Thus, $n = |\mathcal{D}| = \sum_{j=0}^{k-1} b_j$. Similarly, let n_{cj} for the number of samples of class-*j* in \mathcal{D}_c . Hence, $n_c = |\mathcal{D}_c| = \sum_{j=0}^{k-1} n_{cj}$. Let the number of labels changed by Algorithm 1 (denoted by \mathcal{A}) in the total dataset \mathcal{D} be cl(\mathcal{A}) and in concept data subset \mathcal{D}_c be cl($\mathcal{A})_c$. We first prove the following lemma.

Lemma 1. Let E = 1. For any concept call c, the number of labels changed by Algorithm 1 is upper bounded by $cl(\mathcal{A}) \leq 2nm\hat{V}(\theta^*, \mathcal{C}, \mathcal{D})$.

Proof. For a particular concept value C = c the label annealing subroutine Algorithm 1 changes the concept data subset to $\widetilde{\mathcal{D}}_c$ by redistributing the labels of the samples in \mathcal{D}_c . Let $T_{cj} = \min(n_{cj}, b_j \frac{n_c}{|\mathcal{D}|})$ and α_{cj} be the number of samples for class-*j* in \mathcal{D}_c assigned by the algorithm in current run. By closely observing algorithm 1 it can be said that $h_j \frac{n_c}{|\mathcal{D}|}$ while in the second *for* loop) algorithm tries to retain the original labels of the data until $\alpha_{c,j} < b_j \frac{n_c}{|\mathcal{D}|}$ while in the second phase (third *for* loop) the labels are assigned to other most likelihood classes. Thus the following propositions holds:

If T_{cj} = n_{cj}, then the number of labels changed for class-j in D_c termed as cl(A)_{cj} = 0.
If T_{cj} = b_j n_c/|D|, then the number of labels changed for class-j is in D_c termed as cl(A)_{cj} = |n_{cj} - b_j n_c/|D|.

Hence, in the worst-case scenario number of labels changed for class-*j* in \mathcal{D}_c , $cl(\mathcal{A})_{cj} = \left| n_{cj} - b_j \frac{n_c}{|\mathcal{D}|} \right|$. Therefore,

Therefor

$$\operatorname{cl}(\mathcal{A})_{c} \leq \sum_{j=0}^{k-1} \operatorname{cl}(\mathcal{A})_{cj} = \sum_{j=0}^{k-1} \left| n_{cj} - b_{j} \frac{n_{c}}{|\mathcal{D}|} \right|.$$
(5)

Now, for a particular concept C = c the empirical concept violation of the forgotten model θ_C on \overline{D}_c is as follows:

$$\hat{V}(\theta^*, \mathcal{C} = c, D) = \frac{1}{2} \sum_{j=0}^{k-1} \left| \hat{P}_{\mathcal{D}}(\hat{h}(\theta^*, z) = j) - \hat{P}_{\mathcal{D}}(\hat{h}(\theta^*, z) = j \mid \mathcal{C} = c) \right|$$
(6)

$$= \frac{1}{2} \sum_{j=0}^{k-1} \left| \frac{b_j(\theta^*)}{|\mathcal{D}|} - \frac{n_{cj}(\theta^*)}{n_c} \right|$$
(7)

$$\geq \frac{1}{2n_c} \mathrm{cl}(\mathcal{A})_c. \tag{8}$$

 Hence, $cl(A)_c \leq 2n_c \hat{V}(\theta^*, C = c, D) \leq 2n \hat{V}(\theta^*, C = c, D)$. Summing over all concepts c results in the lemma.

We will use the above lemma to prove our main result.

For a formula in the provide the proof for E = 1. The proof for larger values of E follows by a telescoping sum of the epochs. Let's denote $\mathcal{L}_{\mathcal{D}}(\theta^*)$ and $\mathcal{L}_{\mathcal{D}}(\theta_{\mathcal{C}})$ denote the empirical losses of the pre-trained model and forgotten model on the initial dataset \mathcal{D} respectively. Now following the notations from the above proof of Lemma 1 the number of labels changed in the whole dataset $\widetilde{\mathcal{D}} = \bigcup_{c=0}^{m-1} \widetilde{\mathcal{D}}_c$ is cl(\mathcal{A}). We now upper bound the empirical loss of $\theta_{\mathcal{C}}$ on \mathcal{D} as follows: $\mathbb{E}\left[\mathcal{L}_{\mathcal{D}}(\theta_{\mathcal{C}})\right] = \mathbb{E}\left[\mathcal{L}_{\widetilde{\mathcal{D}}}(\theta_{\mathcal{C}}) + \mathcal{L}_{\mathcal{D}}(\theta_{\mathcal{C}}) - \mathcal{L}_{\widetilde{\mathcal{D}}}(\theta_{\mathcal{C}})\right]$ (9)

$$= \mathbb{E}\left[\mathcal{L}_{\widetilde{\mathcal{D}}}(\theta_{\mathcal{C}}) + \frac{1}{n} \left[\sum_{z_i \in \mathcal{D}} \ell(\theta_{\mathcal{C}}, z_i) - \sum_{z_i \in \widetilde{\mathcal{D}}} \ell(\theta_{\mathcal{C}}, z_i)\right]\right]$$
(10)

$$\stackrel{(c)}{\leq} \mathbb{E}\left[\mathcal{L}_{\widetilde{\mathcal{D}}}(\theta_{\mathcal{C}})\right] + \frac{L}{n} \mathrm{cl}(\mathcal{A}) \tag{11}$$

$$\stackrel{(d)}{\leq} \mathcal{L}_{\widetilde{\mathcal{D}}}(\theta^*) + \frac{L}{n} \mathrm{cl}(\mathcal{A}) \tag{12}$$

$$= \mathcal{L}_{\mathcal{D}}(\theta^*) + \mathcal{L}_{\widetilde{\mathcal{D}}}(\theta^*) - \mathcal{L}_{\mathcal{D}}(\theta^*) + \frac{L}{n} \mathrm{cl}(\mathcal{A})$$
(13)

$$\stackrel{(e)}{\leq} \mathcal{L}_{\mathcal{D}}(\theta^*) + \frac{2L}{n} \mathrm{cl}(\mathcal{A}) \tag{14}$$

$$\stackrel{(f)}{\leq} \mathcal{L}_{\mathcal{D}}(\theta^*) + 4Lm\hat{V}(\theta^*, \mathcal{C}, \mathcal{D})$$
(15)

Here (c) and (e) holds as $\forall \theta$ if $z = \tilde{z}$ then $\ell(\theta, z) = \ell(\theta, \tilde{z})$ and the fact that $\ell(\theta, z) \leq L$. (d) holds because of the assumption. Finally applying Lemma 1, we get (f).

B ADDITIONAL DETAILS ABOUT EXPERIMENTS

724 B.1 DATASETS AND MODELS

726 Here we have used four datasets as follows:

- MNIST (LeCun et al., 1998): The MNIST dataset consist of 28 × 28 gray-scale representing handwritten digits from 0 to 9. The MNIST dataset contains 6,000 images per digit class totaling 60,000 training samples and 1,000 images per digit class totalling 10,000 testing images.
- CIFAR-10 (Krizhevsky et al., 2009): The CIFAR-10 dataset consists of 60000 32x32 color images in 10 classes: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck with 6000 images per class. There are 50000 training images and 10000 test images.
- CelebA (Liu et al., 2015): The Celeb Faces Attributes Dataset (CelebA) is a large-scale facial attributes dataset comprising over 200,000 celebrity images, each annotated with 40 attributes. This dataset features significant pose variations and background clutter. CelebA offers extensive diversity, a substantial quantity of images, and rich annotations.
 - **miniImageNet** (Vinyals et al., 2016): Here we have used a smaller subset of the ImageNet dataset consisting of 50,000 training images and 10,000 testing images, evenly distributed across 100 classes. Here we have used an image dimension of 224×224 same as the original ImageNet data dimension.

Models: Further to evaluate our method to different models we experimented with a variety of
models with different learnable parameter sizes such as 2-layer-MLP, Mobinet-v2 (Sandler et al.,
2018), Densenet-121 (Huang et al., 2017), Resnet-50 (He et al., 2016). The 2-layer-MLP net has
two hidden layers both having the size of 500. For Mobinet-v2, Densenet-121, and Resnet-50 we
have taken Pytorch default models with pre-trained weights. For all of these models, the last layer is
changed to an appropriate size suitable for the classification tasks.

- 750 B.2 INITIAL TRAINING:
- 751 B.2.1 INITIAL TRAINING ON MNIST

Here we have used 2-layer-MLP and Resnet-50 models for the classification tasks. For optimization, we have used the Adam optimizer with a learning rate of 0.001 on mean cross-entropy loss. All the models are trained for 5 epochs with a batch size of 64. The loss and accuracy curves can be seen in Figure 5.

756 758 759 760 761 762 (a) MLP loss (b) MLP accuracy (c) Resnet-50 loss (d) Renset-50 accu-763 racy 764 Figure 5: Results of training the initial models on MNIST dataset 765 766 767 **B.2.2** INITIAL TRAINING ON CIFAR-10 768 769 Here we have used Mobinet-v2 and Densenet-121 models for the classification tasks. For optimization, 770 we have used the Adam optimizer with a learning rate of 0.001 on mean cross-entropy loss. Mobinetv2 and Densenet-121 models are trained for 60 and 20 epochs respectively with a batch size of 64. 771 We have an early-stopping of 3 epochs for all the models. The loss and accuracy curves can be seen 772 in Figure 6. 773 774 775 776 777 778 779 780 (a) Mobinet-v2 loss (b) Mobinet-v2 accu- (c) Densenet-121 loss (d) Densenet-121 ac-781 racy curacy 782 783 Figure 6: Results of Training of the Initial Models on CIFAR-10 dataset 784 785 **B.2.3** INITIAL TRAINING ON CELEBA 786 787 Here we have used the Resnet-50 model for the classification tasks. As there are 40 attributes for 788 classification we have trained 40 MLP heads for this. For optimization, we have used the SGD 789 optimizer with a learning rate of 0.01 a learning rate scheduler with a decay of 0.1 every 30 steps, 790 momentum of 0.9, and weight decay 1e-4 on total cross-entropy loss. Here Resnet-50 is trained for 791 90 epochs with a batch size of 256. 792 793 **B.3** TRAINING FOR CONCEPT FORGETTING 794 B.3.1 LAN TRAINING 796 Our official codebase for LAN is available at the following link: https://anonymous.4open. 797 science/r/LAN-141B/. Here we have used the label annealing methodology to finetune the pre-798 trained model for 1 epoch. We evaluated our method with both retraining and FERMI methodology 799 in different forgetting settings. For our results the optimal hyper-parameters For different settings of 800 forgetting we give the optimal hyper-parameters for our optimal results in the following tables 801 802 **B.3.2 BASELINES** 803 • FERMI (Lowy et al., 2021): Here we have used the official implementation of 804 FERMI which can be found in the following link: https://www.dropbox.com/ 805 scl/fo/tz8aksm4ibsta919hziq7/AMK3ixeUQRqoY0FhWqDy5rM?rlkey= yufnfhuvhs91mvv19kc31bss1&e=1&d1=0. Here we have used the FERMI loss with the usual regularized cross-entropy loss to fine-tune the pre-trained model for E=1. 808 • Continuous Fairness (Mary et al., 2019): The official implementation can be found 809 at: https://github.com/criteo-research/continuous-fairness. We

810	have used the usual regularized cross-entropy loss to fine-tune the pre-trained model for
811	E=1.
812	Fortherse VDF (Che et al. 2020a). As there is no efficient implementation for this method
813	Farmess-KDE (Cho et al., 2020a): As there is no official implementation for this method
814	EairClassifier using KDE Similarly like other baselines we train the pre-trained
815	model using this regularized loss for F-1
816	model using this regularized 1055 for L=1.
817	
818	
819	
820	
821	
822	
823	
824	
825	
826	
827	
828	
820	
820	
921	
001	
002	
000	
004	
000	
030	
837	
838	
839	
840	
841	
842	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
853	
854	
855	
856	
857	
858	
859	
860	
861	
862	
863	