

Unrolling SGD: Understanding Factors Influencing Machine Unlearning

Anvith Thudi*, Gabriel Deza*, Varun Chandrasekaran, Nicolas Papernot

* Joint First Authors



Outline

1. Background on Unlearning
2. Our Method

Background on Unlearning

Why Unlearn?

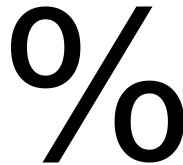
1. Privacy: *Right-to-be-forgotten* (EU GDPR)



2. Security: Data Poisoning

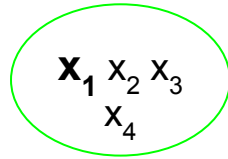


3. Performance: Bad data



The "Protocol"

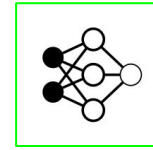
Dataset



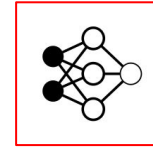
Training



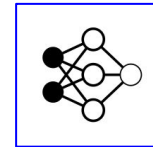
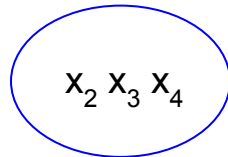
Model



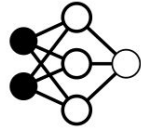
Unlearning



?



Important Details



Could be:

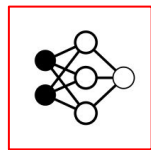
1. (Distribution of) Weights
2. (Distribution of) Functions

Weights

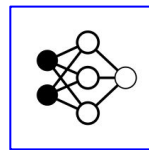


Functions

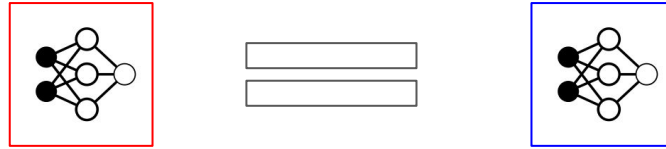
The Big Question



?



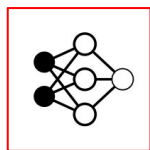
Exact Unlearning



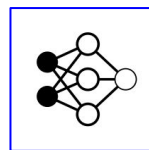
“Machine Unlearning” Bourtole et al.

Expensive

Approximate Unlearning



=



+

ϵ



W.r.t some “metric” d

Examples of “Metrics”

1. l_2 on weights
2. **KL-Divergence** on weight distribution
3. **Membership Inference** on functions

		Logits	Weights	Distribution of Weights	
output space	Test Accuracy			[16][12]	<i>WHAT IS BEING UNLEARNED?</i>
	Membership Inference	[15]	[11]	[16][12]	
	Logits	[15]	[11]		
weight space	Distribution of Weights			[16][12][14][13]	
	Weights		[22][2]		
					<i>HOW IS SUCCESSFUL UNLEARNING MEASURED?</i>

Our Approach

How to Better Study Approximate Unlearning?

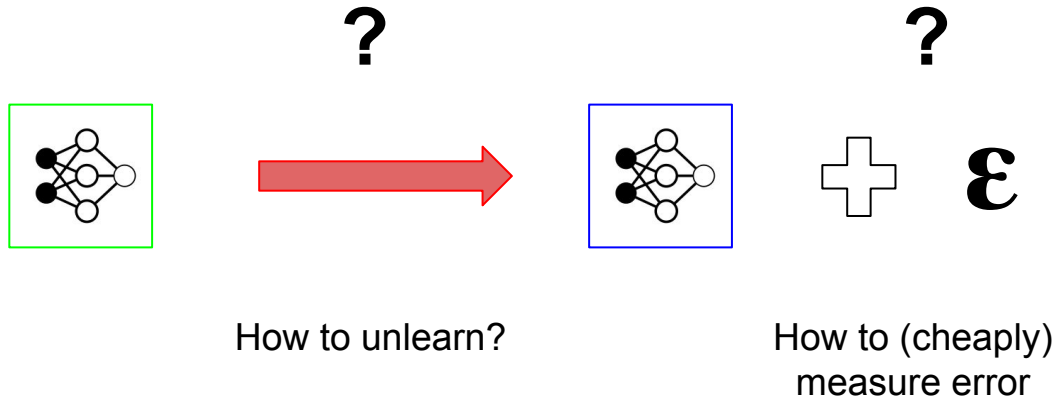
1. Equivalent “metrics”
2. Easy to measure error

An Idea:

Verification Error : Expected ℓ_2 difference on weights

- 1) uniform **convergence in outputs** over finite sets
- 2) **bounds* all L^p metrics** on weight distribution

Problems



Approximate SGD

$$\mathbf{w}_t \approx \mathbf{w}_0 - \underbrace{\eta \sum_{i=0}^{t-1} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \Big|_{\mathbf{w}_0, \hat{\mathbf{x}}_i}}_{\text{Can Forget } \mathbf{x}_i \text{ (single gradient)}} + \underbrace{\sum_{i=1}^{t-1} f(i)}_{\text{Hard to forget: Approximates Error}}$$

Can Forget \mathbf{x}_i
(single gradient)

Hard to forget:
Approximates
Error

$$f(i) = -\eta \frac{\partial^2 \mathcal{L}}{\partial^2 \mathbf{w}} \Big|_{\mathbf{w}_0, \hat{\mathbf{x}}_i} \left(-\eta \sum_{j=0}^{i-1} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \Big|_{\mathbf{w}_0, \hat{\mathbf{x}}_j} + \sum_{j=0}^{i-1} f(j) \right)$$

A Proxy Metric for Verification Error

Unlearning Error:
$$e = \eta^2 \left(\frac{\|w_t - w_0\|_2}{t} \right) \sigma_{avg} \frac{t^2 - t}{2}$$

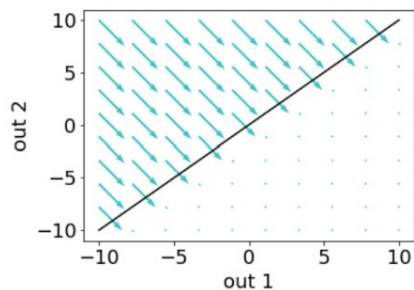
The diagram illustrates the components of the unlearning error equation. Four labels are positioned below the equation, with arrows pointing to specific terms:

- Learning Rate**: An arrow points from this label to the η^2 term in the equation.
- Change in Weights**: An arrow points from this label to the $\|w_t - w_0\|_2$ term in the numerator of the fraction.
- # of training steps**: An arrow points from this label to the t term in the denominator of the fraction.
- Average 1st singular value**: An arrow points from this label to the σ_{avg} term in the equation.

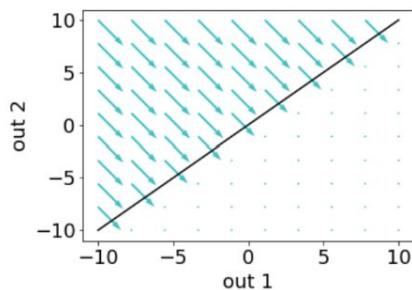
How to Further Reduce Verification Error?

Train with SD Loss = CE Loss + γ * standard deviation of logits

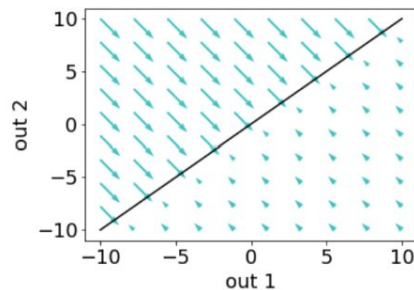
Motivation (Logistic Regression): Pushes minima closer to initialization



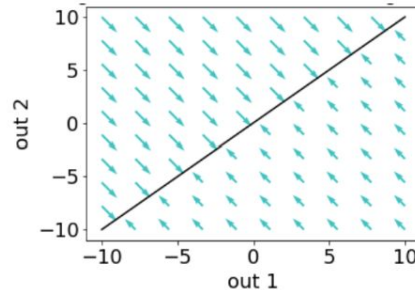
(a) $\gamma = 0.01$



(b) $\gamma = 0.1$

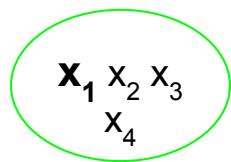


(c) $\gamma = 1$

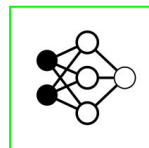


(d) $\gamma = 5$

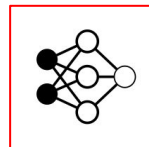
Our Approach



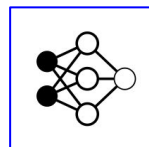
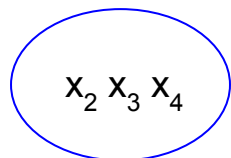
Train with SD Loss

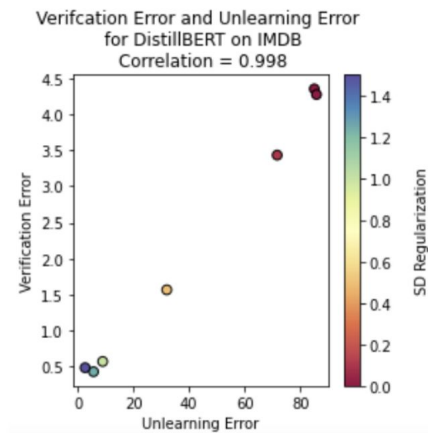


Add single
gradient

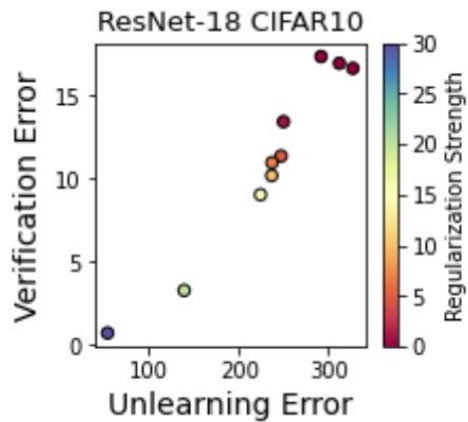


Ver Err $\sim e$

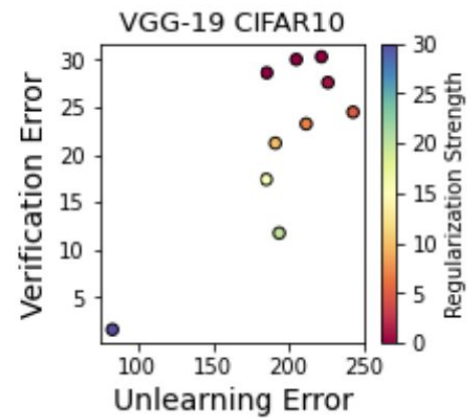




Correlation 0.998

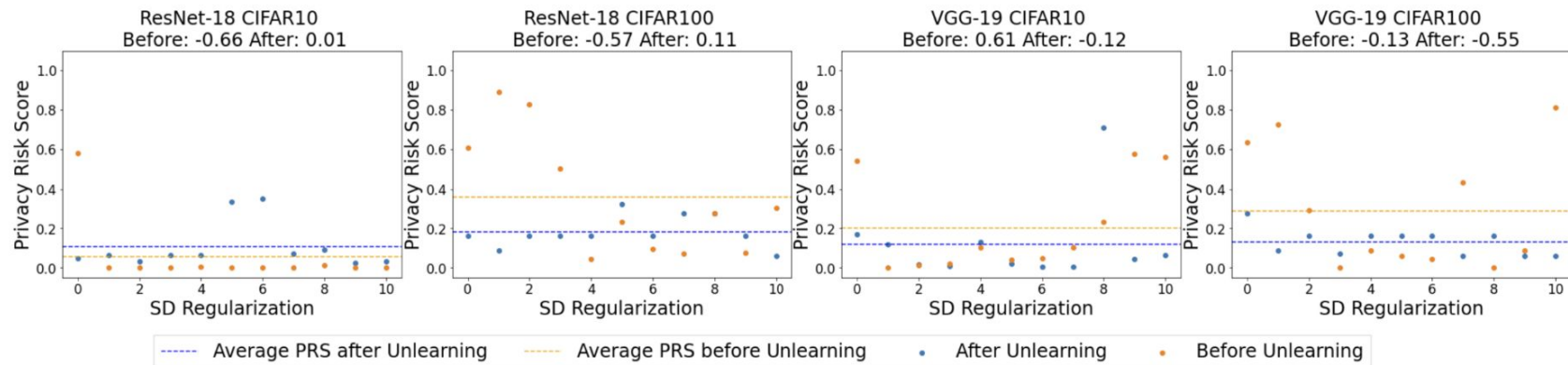


Correlation 0.96



Correlation 0.81

Also Reduces Membership Inference



Always reduces baseline: $\gamma = 0$, “Before Unlearning”

Thank You!