





# **A Brief Introduction to Machine Unlearning**

## Contents



### 1. Introduction

- Reasons for & Main Purpose of Machine Unlearning
- Exact Unlearning & Approximate Unlearning
- Differential Privacy & Machine Unlearning
- Tradeoffs
- 2. Evaluation Metrics
- 3. Exact Unlearning Algorithms

SISA

4. Approximate Unlearning Algorithms

Descent-to-Delete, Fisher, Influence, DeltaGrad, DeepObliviate

5. Comparison



#### • Reasons:

Security:	after removing adversarial data
Privacy:	"right to be forgotten" – not only data itself, but also its impacts on the model
Usability:	for a user
Fidelity:	bias on African-American offenders
Usability: Fidelity:	for a user bias on African-American offenders

- All these reasons leads to a necessity to design a solution to data deletion.
- The ideal (yet expensive) solution is to update the database and retrain all models when a deletion request arrives.
- But retraining requires enormous computation: large models and frequent deletion requests.

Main purpose of *Machine Unlearning (Data Deletion)*:

Design **fast unlearning (deletion) algorithms** that produce output models that are statistically **indistinguishable** from the models that would have arisen from **retraining**.

# Machine Unlearning









The smaller  $(\alpha, \beta)$  are, the stronger unlearning guarantees will become.

# **Differential Privacy & Machine Unlearning**



- Machine unlearning uses the same metric with differential privacy for distributional closeness, but ...
  - DP compares the **same** algorithm run on **different** datasets
  - Machine unlearning compares **different** algorithms run on the **same** dataset
  - If A is differentially private for any data, then it does not learn anything from the data. In other words, differential privacy is a very strong condition, and most differentially private models suffer a significant loss in accuracy even for large ε.
- But DP tools are useful to machine unlearning
  - Noise addition converts distance in parameter space to distance in distribution
  - DP reduces dependencies introduced via **adaptivity**

the data to be deleted depends on the current unlearned model

# Tradeoffs



Goal: given desired **unlearning** level  $(\alpha, \beta)$ , and **computation** budget, design **accurate** unlearning algorithms.



Boundary Conditions:

- **Retraining**: high accuracy, optimal (0,0)-unlearning, <u>computationally expensive</u>
- Not Unlearning: optimal accuracy, *poor unlearning performance*, no computation
- **Completely DP Models**: *low accuracy*, optimal (0,0)-unlearning , no computation

## **Evaluation Metrics**



## Suppose $h^{U} = U(A(D), D, D_{U}), h^{*} = U^{*}(A(D), D, D_{U})$

- $Efficiency(h^U) \coloneqq \frac{time \ taken \ to \ obtain \ h^*}{time \ taken \ to \ obtain \ h^U}$
- $Effectiveness(h^U) \coloneqq \left| \mathcal{M}_{test}^U \mathcal{M}_{test}^* \right|$  ( $\mathcal{M}$  is another performance metric)
- *Consistency* (correctness guarantee)

• Consistency<sub>$$\theta$$</sub> $(h_{\theta}^{U}) \coloneqq \|\theta^{U} - \theta^{*}\|_{2}$ 

- Consistency<sub>y</sub>( $h^U$ ) :=  $\frac{1}{n_{test}} \sum_{i=1}^{n_{test}} \mathbf{1}_{y_{pred,i}^* = y_{pred,i}^U}$
- Consistency<sub>KL</sub>( $h^U$ ) :=  $\mathbb{E}_x[KL(\Pr[h^U(x)] || \Pr[h^*(x)])]$
- *Certifiability* (security guarantee)

• Certifiability(
$$h^U$$
) :=  $\frac{|\mathcal{M}_u^U - \mathcal{M}_u^*|}{|\mathcal{M}_u^U| + |\mathcal{M}_u^*|}$ 

• Certifiability $(h^U) \coloneqq I(D_U; h^U)$ 

# SISA (SSP 2021)



- Sharded
- Isolated
- Sliced
- Aggregated





Algorithm 8 Perfect *i*th unlearning for basic perturbed gradient descent, (Neel et al., 2020).

**Input:** published model parameters  $\tilde{\theta}_{i-1}$ , dataset  $D_{i-1}$ , update data point  $\mathbf{z}_{i-1}$ , number of iterations  $T_i$ , noise parameter  $\sigma > 0$ . **Output:** published unlearned parameters  $\theta_i$ . 1: procedure PGDUNLEARN( $\tilde{\theta}_{i-1}, D_{i-1}, \mathbf{z}_{i-1}; T_i, \sigma$ ) initialise  $\theta'_0 \leftarrow \theta_{i-1}$ 2: 3:  $D_i \leftarrow D_{i-1} \setminus \{\mathbf{z}_{i-1}\}$ for  $t = 1; t < T_i; t + t$ 4:  $\theta'_t \leftarrow \operatorname{Proj}_{\Theta}(\theta'_{t-1} - \eta_t \nabla L(\theta'_{t-1}, D_i))$ 5: end for 6:  $\hat{\theta}_i \leftarrow \theta'_{T_i}$ 7: draw  $Z \sim \mathcal{N}(0, \sigma^2 I_p)$ 8: return  $\hat{\theta}_i = \hat{\theta}_i + Z$  $\operatorname{Proj}_{\Theta}(\theta) = \operatorname{arg\,min}_{\theta' \in \Theta} \|\theta - \theta'\|_2$ 9: 10: end procedure

# Fisher (CVPR 2020)



 $\theta := \operatorname{SGD}(L(\theta_{\operatorname{init}}, D)) + \sigma F^{-\frac{1}{4}}\mathbf{b},$ 

 $\theta^{\mathcal{U}} := \theta - F^{-1} \Delta_{\operatorname{rem}} + \sigma F^{-1}$ Newton step

noise injection

$$\Delta_{\rm rem} := \nabla L(\theta, D \setminus D_u),$$

F: The Fisher Information matrix, which is used as an approximation to the Hessian as the Fisher matrix is less expensive to compute.

The authors prove that such a batch unlearning algorithm can ensure that  $I(D_{u}; U(h_{\theta}, D, D_{u})) = 0$ 

# Fisher (CVPR 2020)



Algorithm 5 Fisher removal mechanism, (Golatkar et al., 2019; Mahadevan and Mathioudakis, 2021).

**Input**: trained model parameters  $\theta$ , training dataset D, subset of data to be removed  $D_u$ , noise parameter  $\sigma$ , mini-batch size m'.

**Output**: unlearned model parameters  $\theta^{\mathcal{U}}$ .

- 1: procedure FISHERUNLEARN( $\theta$ , D,  $D_u$ ;  $\sigma$ , m')
- 2: assign the number of batches  $s \leftarrow \lceil \frac{m}{m'} \rceil$  (*m* is the number of samples in  $D_u$ )
- 3: split  $D_u$  into s mini-batches  $D_u^1, D_u^2, \dots, D_u^s$ , each of size m'
- 4: initialise  $D' \leftarrow D$
- 5: initialise  $\theta^{\mathcal{U}} \leftarrow \theta$
- 6: **for**  $i = 1; i \le s; i + do$
- 7:  $D' \leftarrow D' \setminus D_u^i$
- 8:  $\Delta \leftarrow \nabla L(\theta^{\mathcal{U}}, D')$
- 9:  $F \leftarrow$  compute Fisher Information Matrix of L and D', (Golatkar et al., 2019, Eq. (8)) 10:  $\theta^{\mathcal{U}} \leftarrow \theta^{\mathcal{U}} - F^{-1}\Delta$
- 11: **if**  $\sigma > 0$  **then**
- 12: draw  $\mathbf{b} \sim \mathcal{N}(0, 1)^p$
- 13:  $\theta^{\mathcal{U}} \leftarrow \theta^{\mathcal{U}} + \sigma F^{-1/4} \mathbf{b}$
- 14: **end if**
- 15: **end for**
- 16: return  $\theta^{\mathcal{U}}$
- 17: end procedure



# Influence (ICML 2020)





Algorithm 6 Influence removal mechanism, (Mahadevan and Mathioudakis, 2021).

Input: trained model parameters  $\theta$ , original train dataset D, subset of data to be deleted  $D_u$ , mini-batch size m'.

**Output**: unlearned model parameters  $\theta^{\mathcal{U}}$ .

- 1: procedure INFLUENCEUNLEARN $(\theta, D, D_u; m')$
- 2: assign the number of batches  $s \leftarrow \lceil \frac{m}{m'} \rceil$  (*m* is the number of samples in  $D_u$ )
- 3: split  $D_u$  into s mini-batches  $D_u^1, D_u^2, \dots, D_u^s$ , each of size m'
- 4: initialise  $D' \leftarrow D$
- 5: initialise  $\theta^{\mathcal{U}} \leftarrow \theta$
- 6: **for**  $i = 1; i \le s; i + t$ **do**
- 7:  $D' \leftarrow D' \setminus D_u^i$
- 8:  $\Delta_{m'} \leftarrow \nabla L\left(\theta^{\mathcal{U}}, D_u^i\right)$
- 9:  $H \leftarrow \nabla^2 L\left(\theta^{\mathcal{U}}, D'\right)$ 10:  $\theta^{\mathcal{U}} \leftarrow \theta^{\mathcal{U}} + H^{-1}\Delta_{m'}$
- 11: end for
- 12: return  $\theta^{\mathcal{U}}$

13: end procedure

# DeltaGrad (ICML 2020)

- At each step of training, t, the model parameters {θ<sub>0</sub>, θ<sub>1</sub>, ..., θ<sub>t</sub>} and the gradients of the loss function {∇L(θ<sub>0</sub>), ∇L(θ<sub>1</sub>), ..., ∇L(θ<sub>t</sub>)} are saved.
- Suppose  $D_u = \{z_i | i \in M\} \subseteq D, m = |M|$ .





# DeltaGrad (ICML 2020)



#### Algorithm 7 DeltaGrad removal mechanism, (Wu et al., 2020)

**Input:** model training parameters  $\theta := \{\theta_0, \theta_1, \dots, \theta_t\}$ , training data D, indices of removed training samples M, stored training gradients  $\nabla L(\theta) := \{\nabla L(\theta_0), \nabla L(\theta_1), ... \nabla L(\theta_t)\}$ , period  $T_0$ , total iteration number T, history size k, burn-in iteration number  $j_0$ , learning rate  $\eta_t$ . **Output:** unlearned model parameters  $\theta^{\mathcal{U}} = \theta_{\mathcal{T}}^{\mathcal{U}}$ . 1: procedure DELTAGRADUNLEARN( $\boldsymbol{\theta}, D, M; \nabla L(\boldsymbol{\theta}), T_0, T, k, j_0, \eta_t$ ) initialise  $\theta_0^{\mathcal{U}} \leftarrow \theta_0$ 2: initialise an array  $\Delta G \leftarrow []$ 3: initialise an array  $\Delta \Theta \leftarrow []$ 4:  $\ell \leftarrow 0$ 5: for t = 0; t < T; ++ do 6: if  $[(t_0 - j_0) \pmod{T_0} == 0]$  or  $t \le j_0$  then 7: compute  $\nabla L(\theta_t^{\mathcal{U}})$  exactly 8: compute  $\nabla L(\theta_t^{\mathcal{U}}) - \nabla L(\theta_t)$ , using the cached gradient  $\nabla L(\theta_t)$ 9:  $\Delta G[\ell] \leftarrow \nabla L(\theta_t^{\mathcal{U}}) - \nabla L(\theta_t)$ 10: $\Delta \Theta[\ell] \leftarrow \theta_t^{\mathcal{U}} - \theta_t$ 11:  $\ell \leftarrow \ell + 1$ 12:compute  $\theta_{t+1}^{\mathcal{U}}$  by using exact GD update 13:else 14: $B_{i_k} \leftarrow \text{L-BFGS}(\Delta G[-k:], \Delta \Theta[-k:])$ 15: $\nabla L(\theta_t^{\mathcal{U}}) \leftarrow \nabla L(\theta_t^{\mathcal{U}}) + B_{j_k}(\theta_t^{\mathcal{U}} - \theta_t)$  approximate the gradient 16:compute  $\theta_{t+1}^{\mathcal{U}}$  via the modified gradient formula, Eq. (5.9), using approximated  $\nabla L(\theta_t^{\mathcal{U}})$ 17:end if 18:end for 19:return  $\theta_t^{\mathcal{U}}$ 20:21: end procedure

# DeepObliviate

南京航空航天大学 Nanjing University of Aeronautics and Astronautics

The method further improves on the slicing component of SISA by using the so-called **temporal residual memory** to identify which intermediate models need to be retrained, adding approximation into the process.



$$h^{\mathcal{U}} = \operatorname{TRAIN}(\{D'_d, D_{d+1}, \dots, D_{d+t}\} | h_{d-1}) \oplus (h_B \ominus h_{d+t})$$
$$D'_d = D_d \setminus \{\mathbf{z}\}$$



# DeepObliviate

The value of t is determined by the temporal residual memory. The authors define the temporal residual memory as the  $\ell^1$  distance (or Manhattan distance) between the influence of the deleted data z on successive models when z is included in the training and when it is not. Formally, the temporal residual memory  $\Delta t$  at step t is:

$$\Delta t \coloneqq \| I(D_{d+t} | h_{d+t-1}) - I(D_{d+t} | h_{d+t-1}^U) \|_{1},$$

where  $I(D_i|h_{i-1}) \coloneqq h_i \ominus h_{i-1}$ .

The authors use detrended fluctuation analysis (DFA) to eliminate noise in  $\Delta$  and systematically determine whether  $\Delta$  has stabilized. DFA is used to determine the statistical self-affinity of a time-series signal by fitting  $\Delta t$  with a decaying power-law function, whose derivative is easily-computable and can be used to determine stationarity.

# DeepObliviate



Algorithm 9 Unlearning with DeepObliviate, (He et al., 2021). **Input:** parameters for intermediary trained models  $\theta = \{\theta_1, \ldots, \theta_B\}$ , training data  $D = D_1 \cup \cdots \cup D_B$ , data point to be removed  $\mathbf{z} \in D_d$ , stationarity hyperparameter  $\varepsilon$ . **Output:** unlearned model  $h^{\mathcal{U}}$ . 1: procedure DEEPOBLIVIATEUNLEARN( $\boldsymbol{\theta}, D, \mathbf{z}; \varepsilon$ ) initialise  $h^{\mathcal{U}} \leftarrow h_{d-1}$ 2: initialise  $\theta_{d-1}^{\mathcal{U}} \leftarrow \theta_{d-1}$ 3: 4:  $D_d \leftarrow D_d \setminus \{\mathbf{z}\}$ for  $t = 0; t \le B - d; t + do$ 5:  $h^{\mathcal{U}} \leftarrow \text{TRAIN}(D_{d+t} \mid h^{\mathcal{U}})$ 6:  $\theta_{d+t}^{\mathcal{U}} \leftarrow h^{\mathcal{U}}$  get parameters from  $h^{\mathcal{U}}$ 7: compute temporal influence  $V_{d+t} \leftarrow \theta_{d+t} - \theta_{d+t-1}$  of  $D_{d+t}$  on  $\theta_{d+t-1}$ 8: compute temporal influence  $V_{d+t}^{\mathcal{U}} \leftarrow \theta_{d+t}^{\mathcal{U}} - \theta_{d+t-1}^{\mathcal{U}}$  of  $D_{d+t}$  on  $\theta_{d+t-1}^{\mathcal{U}}$ 9: compute temporal residual memory  $\Delta_{d+t} \leftarrow ||V_{d+t} - V_{d+t}^{\mathcal{U}}||_1$ 10: compute power-law exponent using DFA  $\alpha \leftarrow \text{DFA}(\{\Delta_d, \Delta_{d+1}, \dots, \Delta_{d+t}\})$ 11:  $Y(x) := ax^{-\alpha} + b$ 12: $f(x) := \partial Y(x) / \partial x = a \cdot (-\alpha) \cdot x^{-\alpha - 1}$ 13: $a, b \leftarrow \arg\min_{a,b} (Y(x) - \Delta_x)$  ising least squares,  $x \in \{d, \dots, d+t\}$ 14:15: $q \leftarrow f(d+t)$ if  $|g| < \varepsilon$  then 16:17:break end if 18:end for 19: $h^{\mathcal{U}} \leftarrow h^{\mathcal{U}} \oplus (h_B \ominus h_{d+t})$ 20: $\{\theta_d, \ldots, \theta_{d+t}\} \leftarrow \{\theta_d^{\mathcal{U}}, \ldots, \theta_{d+t}^{\mathcal{U}}\}$ 21:**Return**  $h^{\mathcal{U}}$ 22:23: end procedure

# Comparison







Unlearning Type									
Method	Applicability	Properties	Certificate of Unlearning						
SISA	Incrementally trained models <sup>a</sup>	Exact Weak <sup>b</sup>	Exact						
DaRE	DaRE trees DaRE RF	Exact	Exact						
Fisher	LQ loss function	Approximate	3.3, 33.4 kNATs <sup>c</sup>	$\begin{bmatrix} 0.0\%^{\mathbf{d}} & (m' = \lfloor m/8 \rfloor) \\ 0.26\% & (m' = m) \end{bmatrix}$					
Influence	SC, BG, Lipschitz Hessian	Approximate Weak <sup>e</sup>	$(\epsilon, \delta)$ -certified	$\begin{array}{c c} 0.1\%^{\mathbf{d}} & (m' = \lfloor m/8 \rfloor) \\ 1.1\% & (m' = m) \end{array}$					
DeltaGrad	SC, smooth loss, BG, Lipschitz Hessian, SI	Approximate Strong <sup>f</sup>	N/A <sup>g</sup>						
Descent-to- Delete	SC, smooth (Bounded, Lipschitz Hessian) <sup>h</sup>	Approximate Strong <sup>i</sup>	$(\epsilon, \delta)$ -certified						
DeepObliviate	Any deep-learning model	Approximate Strong	$egin{array}{llllllllllllllllllllllllllllllllllll$						

# Comparison



Performance											
Method		Efficiency		Effectiveness	Consistency						
		min	g. mean	max	Effectiveness	min	g. mean	max			
SIGAa	Bourtoule	$ 1.36\times$	$2.49 \times$	$4.63 \times$	< 2% (18.76%)		Exact				
DIDA	et al (2021)			Exact							
	He et al.	$ 14.0\times$	39.6  imes	$75.0 \times$	< 5% (15.1%)						
	(2021)										
$D_{a}BE_{p}$	$\min d_{\max}$	$10\times$	$366 \times$	$9735 \times$	$\leq 0.5\%$	Exact					
Dant	$\max d_{\max}$	$145 \times$	$1272\times$	$35856 \times$	$\leq 2.5\%$						
Fishor <sup>c</sup>	$m' = \lfloor m/8 \rfloor$	$1.0 \times$	3.3  imes	$36.8 \times$	pprox 0.0%		N/A				
1 151101	m' = m	$1.5 \times$	13.0  imes	$282.4 \times$	< 0.2%		$\Pi/\Pi$				
Influencec	$m' = \lfloor m/8 \rfloor$	0.3  imes	5.9  imes	$75.7 \times$	< 0.55%	N/A					
minuence	m' = m	$2.5 \times$	$38.2 \times$	$215.1 \times$	< 0.57%		N/A				
$DeltaGrad^d$		$1.6 \times$	$2.7 \times$	6.5  imes	pprox 0.0%	$1.7 \times 10^{-6}$	$1.1 \times 10^{-5}$	$1.4\times 10^{-4}$			
	PGD	$1 + \mathcal{I}^{-1}\log(\epsilon n/\sqrt{p})$			$p/(e^{\mathcal{I}}\epsilon^2 n^2)$						
Descent-to-	sRPGD	$\mathcal{I}^{-\frac{3}{5}}(\epsilon n/\sqrt{p})^{\frac{2}{5}}$		N/A	$\frac{\left(\frac{\sqrt{p}}{\epsilon n\mathcal{I}}\right)^{\frac{2}{5}}}{\sqrt{\frac{\sqrt{p}}{\epsilon p\sqrt{\mathcal{I}}}}}$						
Delete <sup>e</sup>	wRPGD	$\mathcal{I}^{-rac{3}{4}}\sqrt{\epsilon n/\sqrt{p}}$									
	DPGD	$\frac{\min\{\mathcal{I}^{-1}\log n, n^{\frac{4-3\xi}{2}} + \mathcal{I}^{-1}\log(\epsilon n/\sqrt{p})\}}{n^{\frac{4-3\xi}{2}} + \mathcal{I}^{-1}\log(\epsilon n/\sqrt{p})}$		$\log n$ ,		$rac{2\exp\left(-\mathcal{I}n^{rac{4-3\xi}{2}} ight)}{\epsilon^2 n^2}+rac{1}{n^{\xi}}$					
				$(\epsilon n/\sqrt{p})\}$							
$DeepObliviate^{f}$	$\min \varepsilon$	$6.17 \times$	$13.97 \times$	$45.45 \times$	< 0.18%	97.25	98.82	99.95			
	arepsilon=0.1	$9.26 \times$	$22.81\times$	$66.67 \times$	< 0.35%	95.78	97.61	99.85			