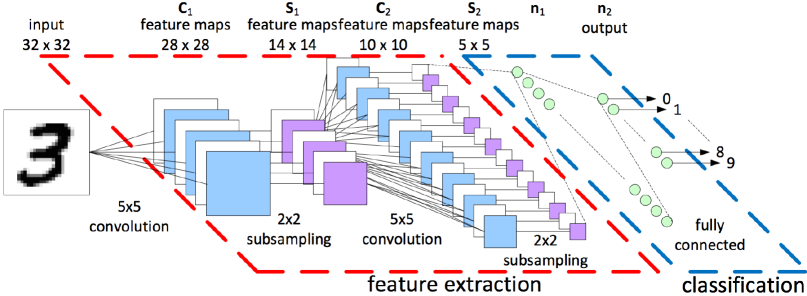


PDE Based Trustworthy Deep Learning

Stan Osher

Department of Mathematics, UCLA

Deep Learning (DL)



DL = Big Data + Deep Nets + SGD + HPC

Deep Learning: Revolution in Technology

Face ID



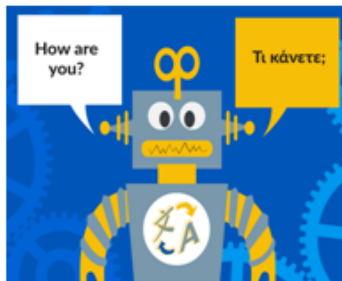
Autonomous Cars



Alpha Go



Machine Translation

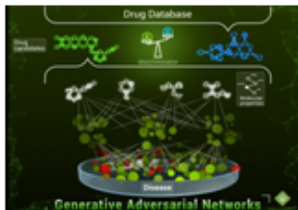


Deep Learning: Revolution in Science

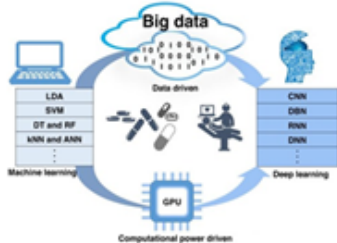
Protein Structure Prediction



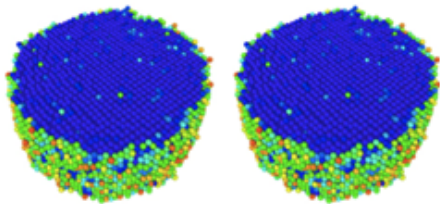
Molecular Generation



Drug Discovery



Material Design



However, Deep Learning is Not Trustworthy!

Trustworthy deep learning:

1. Robust deep learning
2. Accurate deep learning
3. Efficient deep learning
4. Private deep learning

...

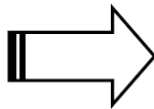
with theoretical guarantees!

Adversarial Vulnerability of Deep Neural Nets

Original Inputs



Modified Inputs



Wrong ML Detection



Deep Learning is Very Expensive



AlphaGO
1202 CPUs, 176 GPUs,
100+ Scientists.

Lee Se-dol
1 Human Brain,
1 Coffee.

Break Privacy of the Face Recognition System



Figure: Recovered (Left), Original (Right)

Membership Attack: determine if a record is in the training set.

Model Inversion Attack: recover the photo of a person given the person's name in face recognition task.

Other Data Abuse: [Netflix Recommendation Competition](#),
[Privacy of the Genome Data](#), ...

Membership Attack – Linear Regression

Suppose three groups of volunteers give us their data to train a linear model:

Group I. $x_1^i, x_2^i, x_3^i \sim 1 + 0.05\mathcal{N}(0, 1)$, and $y^i = 1$ for $i = 1, 2, \dots, 100$

Group II. $x_1^i, x_2^i, x_3^i \sim 1 + 0.05\mathcal{N}(0, 1)$, and $y^i = 2$ for $i = 1, 2, \dots, 100$

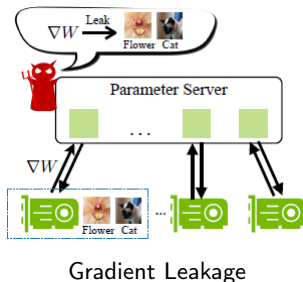
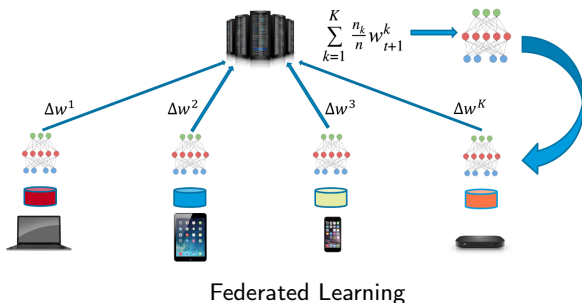
Group III. $x_1^i, x_2^i, x_3^i \sim 1 + 0.05\mathcal{N}(0, 1)$, and $y^i = 3$ for $i = 1, 2, \dots, 100$

If the trained linear model is:

$$y = 0.981x_1 + 1.012x_2 + 1.007x_3. \quad \text{RMSE} = 0.1858.$$

Which group of data is used to train such a linear model?

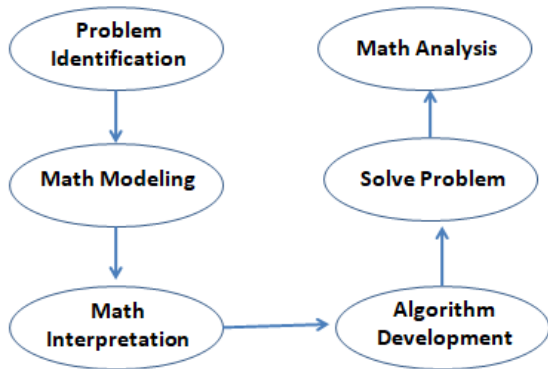
Federated Learning is Not Private



Federated Learning: train a centralized model, w , while training data is distributed over many clients. In each communication-round, clients update their local models with their own private data. The center server then aggregates these local models, and sends the updated model to clients.

Gradient Leakage: update of the local model encodes private data. Gradient is not an encryption of private data.

Our Principle:



A few examples:

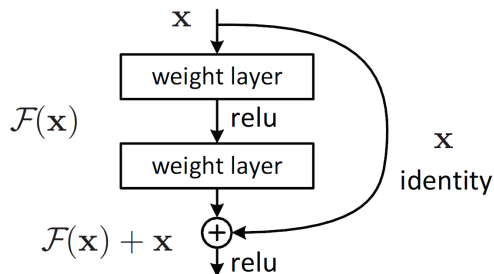
Adversarial Robust Deep Learning I.1

Deep Nets Compression I.2

Privacy-Preserving Machine Learning II.1 & II.2

I. Transport Equation vs. Residual Learning

ResNet vs. Transport Equation



Plain Net: $\mathbf{x}_{l+1} = \mathcal{G}(\mathbf{x}_l)$

ResNet: $\mathbf{x}_{l+1} = \mathbf{x}_l + \mathcal{F}(\mathbf{x}_l)$

Forward propagation (FP) of ResNet for any data-label pair $(\hat{\mathbf{x}}, y)$

$$\begin{cases} \mathbf{x}(0) = \hat{\mathbf{x}}, \\ \mathbf{x}(t_{k+1}) = \mathbf{x}(t_k) + \Delta t \cdot \bar{F}(\mathbf{x}(t_k), \mathbf{w}(t_k)), k = 1, 2, \dots, L-1 \text{ with } \bar{F} \doteq \frac{1}{\Delta t} \mathcal{F} \\ \hat{y} \doteq f(\mathbf{x}(1)) = \text{softmax}(\mathbf{w}_{\text{FC}} \cdot \mathbf{x}). \end{cases}$$

Continuum limit: $\frac{d\mathbf{x}(t)}{dt} = \bar{F}(\mathbf{x}(t), \mathbf{w}(t))$.

Transport equation (TE): $\frac{\partial u}{\partial t}(\mathbf{x}, t) + \bar{F}(\mathbf{x}, \mathbf{w}(t)) \cdot \nabla u(\mathbf{x}, t) = 0, \quad \mathbf{x} \in \mathbb{R}^d$.

ResNet vs. Transport Equation

Forward and backward propagation

1. Let $u(\mathbf{x}, 1) = f(\mathbf{x})$, note $u(\hat{\mathbf{x}}, 0) = u(\mathbf{x}(1), 1) = f(\mathbf{x}(1))$. Therefore, we model FP as computing $u(\hat{\mathbf{x}}, 0)$ along the characteristics of the following TE

$$\begin{cases} \frac{\partial u}{\partial t}(\mathbf{x}, t) + \bar{F}(\mathbf{x}, \mathbf{w}(t)) \cdot \nabla u(\mathbf{x}, t) = 0, & \mathbf{x} \in \mathbb{R}^d, \\ u(\mathbf{x}, 1) = f(\mathbf{x}). \end{cases}$$

2. Backpropagation (BP): find $\mathbf{w}(t)$ for the following control problem

$$\begin{cases} \frac{\partial u}{\partial t}(\mathbf{x}, t) + \bar{F}(\mathbf{x}, \mathbf{w}(t)) \cdot \nabla u(\mathbf{x}, t) = 0, & \mathbf{x} \in \mathbb{R}^d, \\ u(\mathbf{x}, 1) = f(\mathbf{x}), \\ u(\mathbf{x}_i, 0) = y_i, & \mathbf{x}_i \in T, \text{ with } T \text{ being the training data.} \end{cases}$$

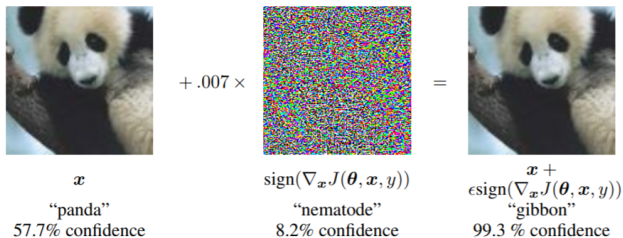
$\mathbf{x}(1)$ is the transport of $\hat{\mathbf{x}}$ along the the characteristics.

I.1 Feynman-Kac Formalism Principled Adversarial Defense

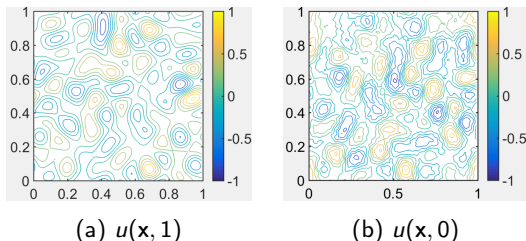
B. Wang, B. Yuan, Z. Shi, and S. Osher, ResNets Ensemble via the Feynman-Kac Formalism to Improve Natural and Robust Accuracies, NeurIPS, 2019.

Code: <https://github.com/BaoWangMath/EnResNet>

How Adversarial Attacks Happen? – A PDE Interpretation



In the TE model, $u(x, 0)$ serves as the decision function for classification.

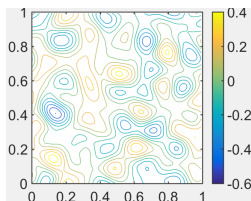


The decision boundary is highly erratic, exposed to adversarial attacks!

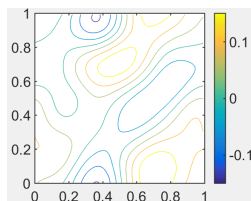
Improving Robustness via Diffusion

We use diffusion to regularize the decision function $u(x, 0)$, which resulting in

$$\begin{cases} \frac{\partial u}{\partial t} + \bar{F}(\mathbf{x}, \mathbf{w}(t)) \cdot \nabla u + \frac{1}{2}\sigma^2 \Delta u = 0, & \mathbf{x} \in \mathbb{R}^d, t \in [0, 1), \\ u(\mathbf{x}, 1) = f(\mathbf{x}). \end{cases}$$



(a) $u(\mathbf{x}, 0), \sigma = 0.01$



(b) $u(\mathbf{x}, 0), \sigma = 0.1$

Theorem (Stability) Let $\bar{F}(\mathbf{x}, t)$ be Lipschitz in both \mathbf{x} and t , and $f(\mathbf{x})$ is bounded. For the above terminal value problem of the convection-diffusion equation, $\sigma \neq 0$, we have

$$|u(\mathbf{x} + \delta, 0) - u(\mathbf{x}, 0)| \leq C \left(\frac{\|\delta\|_2}{\sigma} \right)^\alpha$$

for some constant $\alpha > 0$ if $\sigma \leq 1$. $C := C(d, \|f\|_\infty, \|\bar{F}\|_{L^\infty_{\mathbf{x},t}})$ is a constant.

Feynman-Kac Formula and Deep Nets Design

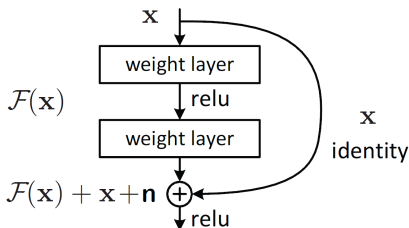
By Feynman-Kac formula, we have

$$u(\hat{\mathbf{x}}, 0) = \mathbb{E} [f(\mathbf{x}(1)) | \mathbf{x}(0) = \hat{\mathbf{x}}],$$

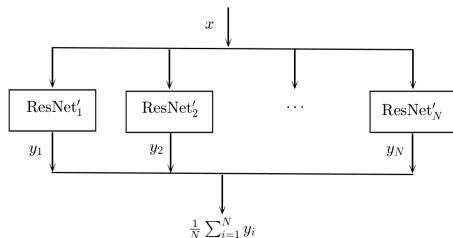
where $\mathbf{x}(t)$ is an Itô process,

$$d\mathbf{x}(t) = \bar{F}(\mathbf{x}(t), \mathbf{w}(t))dt + \sigma dB_t.$$

Deep Nets Design!



Residual mapping + Gaussian noise



Average multiple jointly trained ResNets

Empirical Adversarial Risk Minimization (Robust Training)

Adversarial training: $\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\max_{\delta \in S} \mathcal{L}(f(\mathbf{w}, \mathbf{x} + \delta), y)]$

Adversarial attacks:

FGSM

$$\mathbf{x}' = \mathbf{x} + \epsilon \operatorname{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y)).$$

IFGSM

$$\mathbf{x}^{(m)} = \mathbf{x}^{(m-1)} + \alpha \cdot \operatorname{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^{(m-1)}, y)), \quad m = 1, 2, \dots, M.$$

C&W

$$\min_{\delta} \|\delta\|_2^2, \quad \text{subject to } f(\mathbf{x} + \delta, \mathbf{w}) = t, \quad \mathbf{x} + \delta \in [0, 1]^n.$$

Table: Natural and robust acc of EnResNets on the CIFAR10. Unit: %.

Model	\mathcal{A}_{nat}	\mathcal{A}_{rob} (FGSM)	\mathcal{A}_{rob} (IFGSM ²⁰)	\mathcal{A}_{rob} (C&W)
ResNet20	75.11	50.89	46.03	58.73
En ₁ ResNet20	77.21	55.35	49.06	65.69
En ₅ ResNet20	82.52	58.92	51.48	67.73

I.2 Deep Neural Nets Compression

Channel-Pruning for Adversarial Robust Deep Nets

Deep Nets Compression

Common approaches to improve inference efficiency of deep learning:

Sparse weights

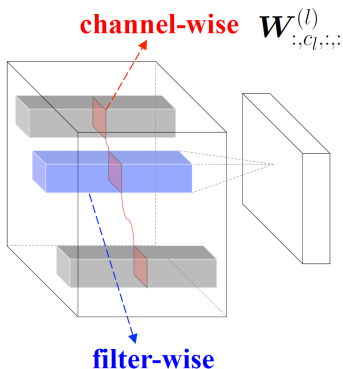
Quantized weights

We focus on sparsifying deep nets (structured & unstructured)!

Neural architecture redesign!

+

Structured & Unstructured weights pruning!

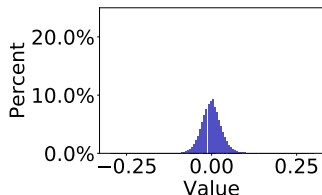


$$n_f \times n_c \times n_w \times n_h$$

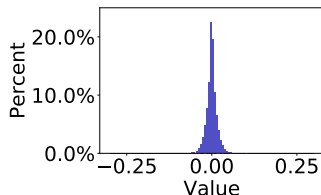
Remark. Structured sparsity can remarkably speed up inference.

Neural Architecture Redesign

Sparsity meets Robustness: ResNet20 vs. En₅ResNet20



(a) ResNet20 (AT) (3.64% ($\leq 10^{-3}$))



(b) En₅ResNet20 (AT) (11.57% ($\leq 10^{-3}$))

Table: Natural and robust acc of EnResNet on the CIFAR10. Unit: %.

Model	\mathcal{A}_{nat}	\mathcal{A}_{rob} (FGSM)	\mathcal{A}_{rob} (IFGSM ²⁰)	\mathcal{A}_{rob} (C&W)
ResNet20	75.11	50.89	46.03	58.73
En ₅ ResNet20	82.52	58.92	51.48	67.73

Ensemble of the noise injected ResNets can improve weights sparsity, accuracy and robustness of the adversarially trained models!

Maximize Sparsity: Structured & Unstructured Sparsity

Adversarial training:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) := \min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\delta \in S} \mathcal{L}(f(\mathbf{w}, \mathbf{x} + \delta), y) \right]$$

Augmented Lagrangian:

$$\mathcal{L}_{\beta}(\mathbf{w}, \mathbf{u}, \mathbf{z}) = \mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{u}\|_1 + \langle \mathbf{z}, \mathbf{w} - \mathbf{u} \rangle + \frac{\beta}{2} \|\mathbf{w} - \mathbf{u}\|^2, \quad \lambda, \beta \geq 0$$

Unstructured sparsity: ℓ_1 -penalty; Structured sparsity: group ℓ_1 -penalty.

ADMM or Split Bregman:

$$\begin{cases} \mathbf{w}^{t+1} \leftarrow \arg \min_{\mathbf{w}} \mathcal{L}_{\beta}(\mathbf{w}, \mathbf{u}^t, \mathbf{z}^t) \\ \mathbf{u}^{t+1} \leftarrow \arg \min_{\mathbf{u}} \mathcal{L}_{\beta}(\mathbf{w}^{t+1}, \mathbf{u}, \mathbf{z}^t) \\ \mathbf{z}^{t+1} \leftarrow \mathbf{z}^t + \beta(\mathbf{w}^{t+1} - \mathbf{u}^{t+1}) \end{cases}$$

Remark 1. One can improve the sparsity of the final learned weights by replacing $\|\mathbf{u}\|_1$ with $\|\mathbf{u}\|_0$; but $\|\cdot\|_0$ is not differentiable.

Remark 2. The Lagrange multiplier term, $\langle \mathbf{z}, \mathbf{w} - \mathbf{u} \rangle$, seeks to close the gap between \mathbf{w}^t and \mathbf{u}^t , and this in turn reduces sparsity of \mathbf{w}^t .

Group the weights into: $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_G\}$, group Lasso: $\sum_{g=1}^G \|\mathbf{w}_g\|_2$; group ℓ_0 : $\sum_{g=1}^G \mathbf{1}_{\|\mathbf{w}_g\|_2 \neq 0}$.

Relaxed Augmented Lagrangian

Relaxed Augmented Lagrangian:

$$\mathcal{L}_\beta(\mathbf{w}, \mathbf{u}) = \mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{u}\|_0 + \frac{\beta}{2} \|\mathbf{w} - \mathbf{u}\|^2.$$

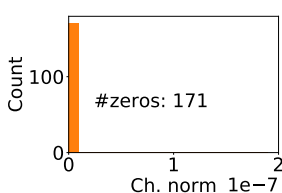
Remark 1. For a fixed \mathbf{w}^t , we have

$$\mathbf{u}^t = H_{\sqrt{2\lambda/\beta}}(\mathbf{w}^t) = (\mathbf{w}_1^t \chi_{\{|\mathbf{w}_1| > \sqrt{2\lambda/\beta}\}}, \dots, \mathbf{w}_d^t \chi_{\{|\mathbf{w}_d| > \sqrt{2\lambda/\beta}\}}),$$

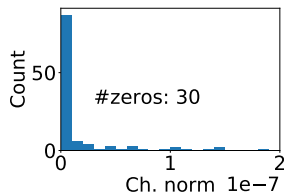
where $H_\alpha(\cdot)$ is the hard-thresholding operator with parameter α .

Remark 2. Fixed \mathbf{u}^t , \mathbf{w}^t can be updated by gradient descent.

Remark 3. \mathbf{w} here is sparser than that in the augmented Lagrangian.



(a) Relaxed (Zoom in)



(b) Original (Zoom in)

Figure: Channel norms of the adversarially trained ResNet20.

Convergence of the Relaxed Augmented Lagrangian

Theorem. Assume \mathcal{L}_β is L -smooth in \mathbf{w} , then the relaxed augmented Lagrangian $\mathcal{L}_\beta(\mathbf{w}^t, \mathbf{u}^t)$ decreases monotonically and converges sub-sequentially to a limit point $(\bar{\mathbf{w}}, \bar{\mathbf{u}})$ provided the stepsize η such that $\eta < 2/(\beta + L)$

Proof.

Step 1. $\mathcal{L}_\beta(\mathbf{w}^{t+1}, \mathbf{u}^{t+1}) \leq \mathcal{L}_\beta(\mathbf{w}^{t+1}, \mathbf{u}^t)$ by the arg min update of \mathbf{u}^t .

Step 2. Note $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla f(\mathbf{w}^t) - \eta \beta (\mathbf{w}^t - \mathbf{u}^t)$, and by the definition of \mathcal{L} :

$$\mathcal{L}_\beta(\mathbf{w}^{t+1}, \mathbf{u}) - \mathcal{L}_\beta(\mathbf{w}^t, \mathbf{u}) \leq \left(\frac{L}{2} + \frac{\beta}{2} - \frac{1}{\eta} \right) \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2,$$

thus when $\eta \leq 2/(\beta + L)$, we have $\mathcal{L}_\beta(\mathbf{w}^{t+1}, \mathbf{u}) \leq \mathcal{L}_\beta(\mathbf{w}^t, \mathbf{u})$.

Step 3. $\mathcal{L}_\beta(\mathbf{w}^t, \mathbf{u}^t)$ decreases monotonically and bounded below by 0, hence converges.

Sparsity vs. Accuracy & Robustness

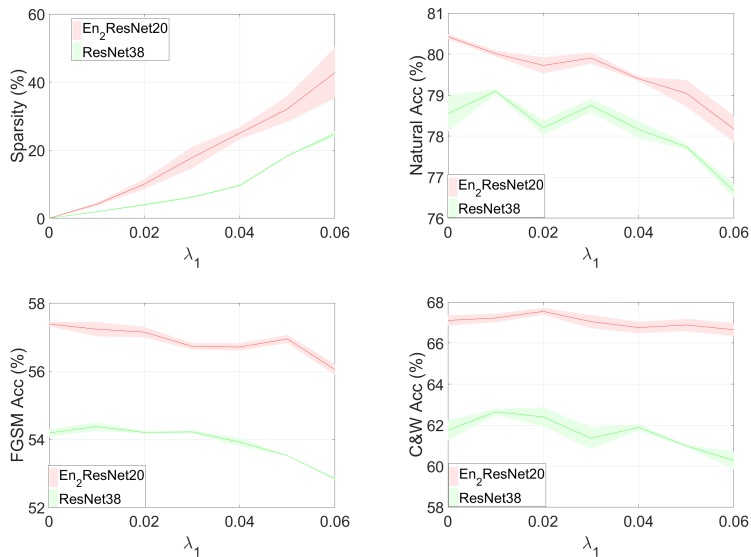


Figure: En₂ResNet20 vs. ResNet38 under different λ_1 . (5 runs, $\beta = 1$).

II. Privacy-Preserving Machine Learning with Laplacian Smoothing

II.1 Privacy-Preserving Empirical Risk Minimization (ERM)

B. Wang, Q. Gu, M. Boedihardjo, F. Barekat, and S. Osher. DP-LSSGD: A Stochastic Optimization Method to Lift the Utility in Privacy-Preserving ERM, ArXiv:1906.12056, 2019

Code: <https://github.com/BaoWangMath/DP-LSSGD>

Differential Privacy



Figure: Recovered (Left), Original (Right)

Differential privacy (DP) is a successful countermeasure to adversaries that try to break the privacy of machine learning.

Add differential privacy constraint in training machine learning models!

F. McSherry and I. Mironov, Differentially Private Recommender Systems: Building Privacy into the Netflix Prize Contenders, KDD, 2009.

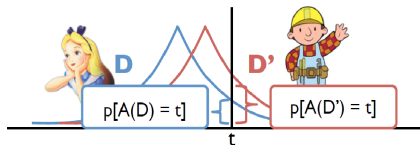
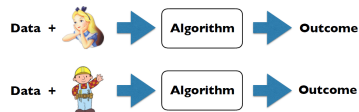
M. Fredrikson, S. Jha, T. Ristenpart, Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures, CCS, 2015.

Differential Privacy

Definition. A randomized algorithm \mathcal{A} is (ϵ, δ) -differentially private if for any two neighboring datasets D, D' that differ in only one entry and for all events S in the output space of \mathcal{A} , we have

$$\Pr(\mathcal{A}(D) \in S) \leq e^\epsilon \Pr(\mathcal{A}(D') \in S) + \delta.$$

DP promises to protect individuals from any additional harm that they might face due to their data being in the private database x that they would not have faced had their data not been part of x .



For all D, D' that differ in one person, if \mathcal{A} is (ϵ, δ) -DP, then:

$$\Pr \left[\left| \ln \left(\frac{\Pr[\mathcal{A}(D) \in S]}{\Pr[\mathcal{A}(D') \in S]} \right) \right| \geq \epsilon \right] \leq \delta$$

Privacy-Preserving Empirical Risk Minimization

Empirical risk minimization (ERM):

$$\min F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{w}, \mathbf{x}_i, y_i).$$

Differentially private SGD (DP-SGD)

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta_k \left(\frac{1}{m} \sum_{k=1}^m \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n} \right), \quad \mathbf{n} \sim \mathcal{N}(0, \nu^2 I_{d \times d}), \quad \{i_k\}_{k=1}^m \subset [n]$$

How to quantify n to guarantee (ϵ, δ) -DP?

Major difficulty: quantifying privacy loss aggregation during SGD.

Theorem (Privacy Budget) Suppose that each f_i is L -Lipschitz. Given the number of iterations T , for any $(\epsilon, \delta > 0)$, DP-SGD, with injected Gaussian noise $\mathcal{N}(0, \nu^2 I)$, satisfies (ϵ, δ) -DP with $\nu^2 = 20T\alpha G^2/(\mu n^2 \epsilon)$, where $\alpha = \log(1/\delta)/((1-\mu)\epsilon) + 1$, if $\exists \mu \in (0, 1)$ s.t. $\alpha \leq \log(\mu n^3 \epsilon / (5b^3 T \alpha + \mu b n^2 \epsilon))$ and $5b^2 T \alpha / (\mu n^2 \epsilon) \geq 1.5$.

Proof.

Step 1. Construct the query $\mathbf{q}_k = (n/b)\nabla F(\mathbf{w}^k)$, which has the ℓ_2 -sensitivity $\Delta(\mathbf{q}_k) = \|\nabla f_{i_k}(\mathbf{w}^k) - \nabla f_{i'_k}(\mathbf{w}^k)\|_2 \leq 2G/b$.

Step 2. A randomized mechanism $\mathcal{M} : \mathcal{S}^n \rightarrow \mathcal{R}$ satisfies (α, ρ) -Rényi DP (RDP), if for all adjacent datasets $S, S' \in \mathcal{S}^n$, we have

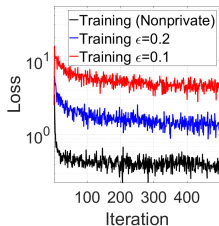
$$D_\alpha(\mathcal{M}(S) \parallel \mathcal{M}(S')) := \frac{1}{\alpha - 1} \log \mathbb{E}_{\mathcal{M}(S')} \left(\frac{\mathcal{M}(S)}{\mathcal{M}(S')} \right)^\alpha \leq \rho.$$

Step 3. Gaussian mechanism $\mathcal{M} = q(S) + \mathbf{n}$ where $\mathbf{n} \sim \mathcal{N}(0, \nu^2 I)$, satisfies $(\alpha, \alpha \Delta^2(\mathbf{q}) / (2\nu^2))$ -RDP. In addition, \mathcal{M} satisfies $(\alpha, 5\tau^2 \Delta^2(\mathbf{q}) \alpha / \nu^2)$ -RDP given $\nu'^2 = \nu^2 / \Delta^2(\mathbf{q}) \geq 1.5$, $\alpha \leq \log(1/\tau(1 + \nu'^2))$, if we apply \mathcal{M} to a subset by uniform sampling without replacement with τ being the subsample rate.

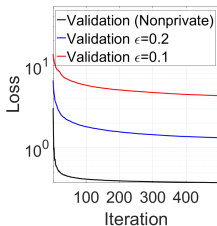
Step 4. If $\mathcal{M}_i : \mathcal{S}^n \rightarrow \mathcal{R}$, for $i \in [k]$, satisfy (α, ρ_i) -RDP, then their composition $(\mathcal{M}_1(S), \dots, \mathcal{M}_k(S))$ satisfies $(\alpha, \sum_{i=1}^k \rho_i)$ -RDP. Moreover, the input of the i -th mechanism can be based on outputs of the previous $(i-1)$ mechanisms.

Step 5. (α, ρ) -RDP $\Leftrightarrow (\rho + \log(1/\delta)/(\alpha - 1), \delta)$ -DP for all $\delta \in (0, 1)$.

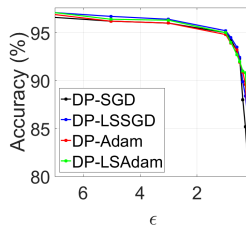
SGD vs. DP-SGD



Training Loss



Validation Loss



Validation Acc

Figure: Logistic regression on the MNIST trained by DP-SGD with $(\epsilon, 10^{-5})$ -DP guarantee (left & middle). LeNet on the MNIST trained by DP-SGD with $(\epsilon, 10^{-5})$ -DP guarantee (right).

DP-SGD reduces the utility of the trained model severely.

Question: Can we do better than DP-SGD with negligible extra computation and memory costs?

DP-SGD with Laplacian Smoothing (DP-LSSGD)

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta_k A_\sigma^{-1} \left(\frac{1}{m} \sum_{k=1}^m \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n} \right).$$

where

$$A_\sigma = (I - \sigma L) = \begin{bmatrix} 1 + 2\sigma & -\sigma & 0 & \dots & 0 & -\sigma \\ -\sigma & 1 + 2\sigma & -\sigma & \dots & 0 & 0 \\ 0 & -\sigma & 1 + 2\sigma & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ -\sigma & 0 & 0 & \dots & -\sigma & 1 + 2\sigma \end{bmatrix}_{d \times d}$$

A_σ^{-1} p.s.d with condition number $1 + 4\sigma$; FFT Implementation

DP-LSSGD has the same privacy budget as DP-SGD!

Proposition (Post-processing) Let $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R$ be a randomized algorithm that is (ϵ, δ) -DP. Let $f : R \rightarrow R'$ be an arbitrary mapping. Then $f \circ \mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R'$ is (ϵ, δ) -DP.

S. Osher, B. Wang, P. Yin, X. Luo, F. Baretat, M. Pham, and A. Lin, arXiv:1806.06317, 2018

Code: <https://github.com/BaoWangMath/LaplacianSmoothing-GradientDescent>

For any pair of $\|x - y\|_1 \leq 1$, and any $S \subset R'$, let $T = \{r \in R : f(r) \in S\}$, we have

$$Pr[f(\mathcal{M}(x)) \in S] = Pr[\mathcal{M}(x) \in T] \leq \exp(\epsilon) Pr[\mathcal{M}(y) \in T] + \delta = \exp(\epsilon) Pr[f(\mathcal{M}(x)) \in S] + \delta$$

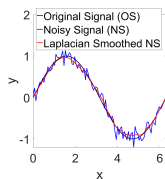
Laplacian Smoothing as a Denoiser

Consider the following diffusion equation with the Neumann BC

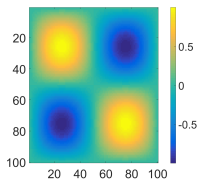
$$\begin{cases} \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}, & (x, t) \in [0, 1] \times [0, +\infty), \\ \frac{\partial u(0, t)}{\partial x} = \frac{\partial u(1, t)}{\partial x} = 0, & t \in [0, +\infty) \\ u(x, 0) = f(x), & x \in [0, 1] \end{cases}$$

Backward Euler in time and central finite difference in space with v^0 being a discretization of $f(x)$. **Unconditionally stable!**

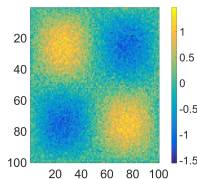
$$\mathbf{v}^{\Delta t} - \mathbf{v}^0 = \Delta t \mathbf{L} \mathbf{v}^{\Delta t} \Rightarrow \mathbf{v}^{\Delta t} = (\mathbf{I} - \Delta t \mathbf{L})^{-1} \mathbf{v}^0 \quad (\sigma = \Delta t)$$



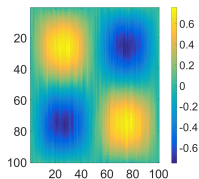
(a)



(b)



(c)



(d)

Figure: Illustration of LS ($\sigma = 10$ for \mathbf{v}_1 and $\sigma = 100$ for \mathbf{v}_2). (a): 1D signal sampled uniformly from $\sin(x)$ for $x \in [0, 2\pi]$. (b), (c), (d): 2D original, noisy, and denoised signals sampled from $\sin(x)\sin(y)$ for $(x, y) \in [0, 2\pi] \times [0, 2\pi]$.

DP-LSSGD Improves Utility of Logistic Regression over DP-SGD (MNIST)

$$\min_{\mathbf{w}} F(\mathbf{w}) = \min_{\mathbf{w}} \left\{ \frac{1}{n} \sum_{i=1}^n -\log \left(\frac{\exp \langle \mathbf{w}, \mathbf{x}_i \rangle_{y_i}}{\sum_j \exp \langle \mathbf{w}, \mathbf{x}_i \rangle_j} \right) + \lambda \|\mathbf{w}\|_2 \right\}, \quad \lambda = 1 \times 10^{-4}.$$

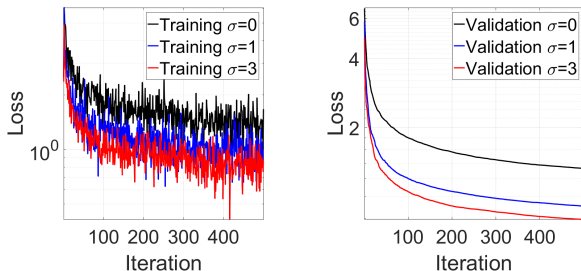


Figure: $(0.2, 10^{-5})$ -DP guarantee. Step size: $1/t$.

Table: Acc of logistic regression with $(\epsilon, \delta = 10^{-5})$ -DP guarantee.

ϵ	0.25	0.20	0.15	0.10
$\sigma = 0$	81.45 ± 1.59	78.92 ± 1.14	77.03 ± 0.69	73.49 ± 1.60
$\sigma = 1$	83.27 ± 0.35	81.56 ± 0.79	79.46 ± 1.33	76.29 ± 0.53
$\sigma = 2$	83.65 ± 0.76	82.15 ± 0.59	80.77 ± 1.26	76.31 ± 0.93

Mitigate Membership Attack – Linear Regression

Suppose three groups of volunteers give us their data to train a linear model:

Group I. $x_1^i, x_2^i, x_3^i \sim 1 + 0.05\mathcal{N}(0, 1)$, and $y^i = 1$ for $i = 1, 2, \dots, 100$

Group II. $x_1^i, x_2^i, x_3^i \sim 1 + 0.05\mathcal{N}(0, 1)$, and $y^i = 2$ for $i = 1, 2, \dots, 100$

Group III. $x_1^i, x_2^i, x_3^i \sim 1 + 0.05\mathcal{N}(0, 1)$, and $y^i = 3$ for $i = 1, 2, \dots, 100$

Linear model trained by DP-SGD with $(0.1, 10^{-5})$ -DP guarantee:

$$y = 1.441x_1 + 3.074x_2 - 1.069x_3. \quad \text{RMSE} = 4.4795.$$

Linear model trained by DP-LSSGD with $(0.1, 10^{-5})$ -DP guarantee:

$$y = 0.827x_1 + 0.532x_2 + 1.728x_3. \quad \text{RMSE} = 0.8839.$$

Utility Guarantees

Algorithm	Privacy	Assumption	Utility	Measurement
DP-SGD	(ϵ, δ)	convex	$\tilde{\mathcal{O}}\left(\sqrt{(D_0 + G^2)d}/(\epsilon n)\right)$	optimality gap
DP-SGD	(ϵ, δ)	nonconvex	$\tilde{\mathcal{O}}\left(\sqrt{d}/(\epsilon n)\right)$	ℓ_2 -norm of gradient
DP-LSSGD	(ϵ, δ)	convex	$\tilde{\mathcal{O}}\left(\sqrt{\gamma(D_\sigma + G^2)d}/(\epsilon n)\right)$	optimality gap
DP-LSSGD	(ϵ, δ)	nonconvex	$\tilde{\mathcal{O}}\left(\sqrt{\beta d}/(\epsilon n)\right)^1$	ℓ_2 -norm of gradient

¹ Measured in the norm induced by \mathbf{A}_σ^{-1} .

$D_\sigma = \|\mathbf{w}^0 - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2$ and \mathbf{w}^* is the global minimizer.

$$\gamma = \frac{1}{d} \sum_{i=1}^d \frac{1}{1 + 2\sigma - 2\sigma \cos\left(\frac{2\pi}{d}\right)} = \frac{1 + \alpha^d}{(1 - \alpha^d)\sqrt{4\sigma + 1}}, \quad \text{with } \alpha = \frac{2\sigma + 1 - \sqrt{4\sigma + 1}}{2\sigma}.$$

$$\beta = \frac{2\alpha^{2d+1} - \xi\alpha^{2d} + 2\xi d\alpha^d - 2\alpha + \xi}{\sigma^2 \xi^3 (1 - \alpha^d)^2},$$

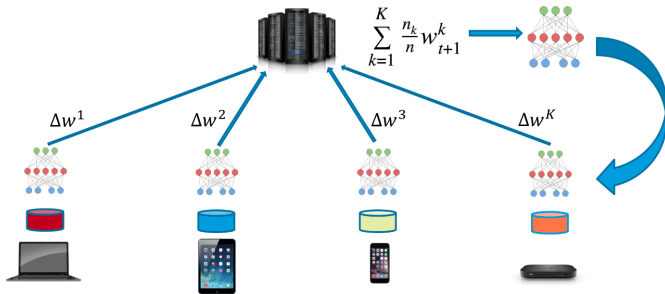
where

$$\alpha = \frac{2\sigma + 1 - \sqrt{4\sigma + 1}}{2\sigma}, \quad \text{and } \xi = -\frac{\sqrt{1 + 4\sigma}}{\sigma}.$$

Poisson summation formula + Residue theorem.

II.2 Differentially Private Federated Learning

Federated Learning (FL)



Train a centralized model (\mathbf{w}) while training data is distributed over many clients. In communication-round t , the server distributes the current model \mathbf{w}_t to a subset M_t of m clients. These clients update the model based on their local data. Let the updated local models be $\mathbf{w}_1^t, \mathbf{w}_2^t, \dots, \mathbf{w}_{n_t}^t$, so the update is

$$H_i^t \doteq \mathbf{w}_i^t - \mathbf{w}^t, \text{ for } i \in M_t.$$

These updates could be a single gradient computed on the client. Then the server collects these updates to update the global model

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta^t H^t, \quad H^t \doteq \frac{1}{m} \sum_{i \in M_t} H_i^t.$$

How to protect the privacy of clients' data?

Differentially Private Federated Learning with Laplacian Smoothing

Algorithm Differentially-Private Federated Learning with Laplacian Smoothing (DP-Fed-LS)

Server executes:

initialize \mathbf{w}^0

for each round $t = 1, 2, \dots, T$ **do**

$m \leftarrow \max(\tau \cdot K, 1)$ where $0 < C \leq 1$

$M_t \leftarrow$ (random set of m clients)

for each client $j \in M_t$ **in parallel do**

$\mathbf{w}_j^t \leftarrow \mathbf{w}^{t-1}$

$\mathcal{B} \leftarrow$ (split local data set into batches of size B)

for each local epoch $i = 1, 2, \dots, E$ **do**

for batch $b \in \mathcal{B}$ **do**

$\mathbf{w}_j^t \leftarrow \mathbf{w}_j^t - \eta_t \cdot \frac{1}{B} \sum_{i \in b} \nabla \ell(\mathbf{w}_j^t; b_i)$

$\mathbf{w}_j^t \leftarrow \mathbf{w}^{t-1} + \text{clip}(\mathbf{w}_j^t - \mathbf{w}^{t-1})$, where $\text{clip}(\mathbf{v}) \leftarrow \mathbf{v} / \max(1, \|\mathbf{v}\|_2 / G)$

$\Delta_j^t \leftarrow \mathbf{w}_j^t - \mathbf{w}^{t-1}$

$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} + \frac{1}{m} \mathbf{A}_\sigma^{-1} \left(\sum_{j=1}^m \Delta_j^t + \mathbf{n} \right)$, where $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \nu^2 \mathbf{I})$

Output \mathbf{w}^T

Privacy Budget for DP-Fed-LS

Theorem (Privacy Budget for DP-Fed-LS) For any $\delta \in (0, 1)$, and ϵ satisfying

$$(2 \log(1/\delta) + (1 + \tau)\epsilon)^2 \leq \frac{3(1 - \tau)\epsilon^3}{8\tau^2 T} \quad \text{and} \quad \epsilon \leq \tau \sqrt{\frac{8T}{3} \log\left(\frac{1}{\delta}\right)}$$

the DP-Fed-LS algorithm (with or without Laplacian smoothing), satisfies (ϵ, δ) -DP if its injected Gaussian noise $\mathcal{N}(0, \nu^2 I)$ is chosen to be

$$\nu \geq (4\tau G)/\epsilon$$

where G is the ℓ_2 -bound of clipped gradient, $\tau := m/K$ is the subsampling ratio of active clients, T is the total number of communication rounds.

Proof.

Step 1. Construct the query $\mathbf{q}_t = \frac{1}{m} \sum_{j=1}^K \mathbf{w}_j^t - \mathbf{w}^t$ with the Gaussian mechanism $\mathcal{M}_t = \frac{1}{m} \sum_{j=1}^K \mathbf{w}_j^t - \mathbf{w}^t + \frac{1}{m} \mathbf{n}$. The ℓ_2 -sensitivity of \mathbf{q}_t is $\Delta \mathbf{q}_t = \|\mathbf{w}_j^t - \mathbf{w}_j^{t'}\|_2 / m \leq 2G/m$.

Step 2.–Step 5. The same as that in proving privacy for DP-SGD.

DP-Fed-LS Improves Utility of Logistic Regression over DP-Fed (MNIST)

Table: Testing accuracy of logistic regression trained by DP-Fed ($\sigma = 0$) and DP-Fed-LS ($\sigma = 1, 2, 3$) on MNIST with $(\epsilon, 1/K^{1.1})$ -DP guarantee with $K = 2000$ be the number of clients.

ϵ	2	3	4	5
$\sigma = 0$	60.23 \pm 2.7	73.50 \pm 1.0	80.72 \pm 0.49	82.24 \pm 0.27
$\sigma = 1$	66.11 \pm 2.7	76.98 \pm 0.68	82.85 \pm 0.26	84.09 \pm 1.1
$\sigma = 2$	67.84 \pm 2.1	79.57 \pm 1.1	82.88 \pm 0.19	84.85 \pm 0.80
$\sigma = 3$	68.52 \pm 1.1	80.60 \pm 0.84	82.54 \pm 0.12	84.51 \pm 0.44

Thank You

I. TE modeling of DNN

I.1 Feynman-Kac formalism for robust and efficient DL

I.2 Channel-pruning for the Feynman-Kac formula principled deep nets

II. Laplacian smoothing

II.1 Differentially-private ERM

II.2 Differentially-private federated learning

References:

1. B. Wang, X. Luo, W. Zhu, Z. Li, Z. Shi, and S. Osher, NeurIPS, 2018.
2. B. Wang, B. Yuan, Z. Shi, and S. Osher, NeurIPS, 2019.
3. T Dinh*, B. Wang*, A. Bertozzi, S. Osher, and J. Xin, Preprint, 2019.
4. B. Wang, Q. Gu, M. Boedihardjo, and S. Osher, arXiv:1906.12056, 2019.
5. Z. Liang, B. Wang, Q. Gu, S. Osher. and Y. Yao, Preprint, 2019.

Code: <https://github.com/BaoWangMath>