

# Lecture 18: Differentially Private Machine Learning

Foundations of Privacy  
Carnegie Mellon University

# Announcement

- HW 3 released. Due Nov. 14th

① Written Component (pdf)  
② Programming component (ipynb)

} zip file  
for submission

# Support Jeremy

<https://gofund.me/d922f9f9>

🔍 Search

How it works ▾

Start a GoFundMe



## Support the Lacomis family after brain surgery



Afsoon Afzal is organizing this fundraiser on behalf of Jessica Lacomis.

Created 2 days ago



Medical, Illness & Healing

# Model Training with DP

Given private data  $x_1, \dots, x_n$ , solve

$$\min_{w \in \mathbb{R}^d} L(w) \equiv \frac{1}{n} \sum_{i=1}^n \ell(w; \cancel{w_i}^{x_i})$$

*Empirical Risk.*

subject to differential privacy

# DP-SGD (in Theory)

## Differentially Private SGD [BST14, SCS13]

- At each iteration  $t$ ,
- Gradient estimate on a mini-batch  $B_t$  :

$$g_t = \left( \frac{1}{|B_t|} \sum_{i \in B_t} \nabla \ell(w_t; x_i) \right)$$

- Noisy gradient update :

$$w_{t+1} = w_t - \eta (g_t + Z_t),$$

$$Z_t \sim \mathcal{N}(0, \sigma^2 I_d)$$

Privacy Proof  
assumes  $\ell$  is  $L$ -Lipschitz  
for all  $x$ , for  $w$



$$\|\nabla \ell(\overset{w_t; x_i}{x_t}; s_i)\|_2 \leq L$$

Set  $\sigma$  to scale with  $L$

# DP-SGD (in Practice)

## Differentially Private SGD [ACGMMTZ16]

- At each iteration  $t$ ,
- Average *clipped* gradient estimate:

$$g_t = \left( \frac{1}{|B_t|} \sum_{i \in B_t} \text{Clip}(\nabla \ell(w_t; x_i), G) \right)$$

- Noisy gradient update :

$$w_{t+1} = w_t - \eta (g_t + Z_t), \quad Z_t \sim \mathcal{N}(0, \sigma^2 I_d)$$

If  $\|g\|_2 > G$ ,  
what is  $\text{Clip}(g, G)$ ?  
 $\left( g \cdot \frac{1}{\|g\|_2} \right) \cdot G$

Gradient Clipping:

$$\text{Clip}(g, G) = g \min \left\{ 1, \frac{G}{\|g\|_2} \right\}$$

$\uparrow$   
 $\nabla \ell(w; x_i)$



Set  $\sigma$  to scale with  $G$

# Privacy Guarantee for DP-SGD (with Clipping)

[BST14, ACGMMTZ16]

$T = \#$  iterations

Theorem: DP-SGD with gradient clipping of threshold  $G$  satisfies  $(\epsilon, \delta)$ -differential privacy, if the noise rate

$\sqrt{T}$  growth  $\rightarrow$   
advanced  
Composition

$$\sigma \geq a \frac{Gq \sqrt{T \ln(1/\delta)}}{\epsilon}$$

for some constant  $a$  and  $q = \frac{|B_t|}{n}$ .

sub-sampling  
rate.

privacy amplification  
via sub-sampling.

*How about convergence and optimality?*

# Gradient clipping can create bias

- *Xiangyi Chen, Z. S.W., Mingyi Hong*  
“Understanding Gradient Clipping in Private SGD: A Geometric Perspective”  
In NeurIPS 2020 (Spotlight)



# Bad Example I

$$\text{Loss: } L(x) = \frac{1}{3} \sum_{i=1}^3 \frac{1}{2} (w - x_i)^2$$

$w, x_i \in \mathbb{R}$

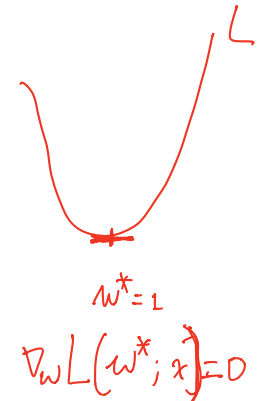
where  $x_1 = x_2 = -3$  and  $x_3 = 9$ .

$\Rightarrow$  Optimum  $w^* = 1$

Clipped gradient at  $w^*$

$$\mathbb{E}[\text{Clip}(\nabla_x \ell(w^*; x_i), 1)] = 1/3$$

$\Rightarrow$  push iterates away from opt



# Bad Example 2

$$\text{Loss: } L(w; x) = \frac{1}{2} \sum_{i=1}^2 \frac{1}{2} (w - x_i)^2$$

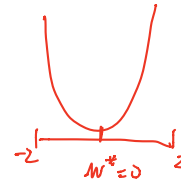
$$\text{where } x_1 = 3, x_2 = -3$$

$$\Rightarrow \text{Optimum } w^* = 0$$

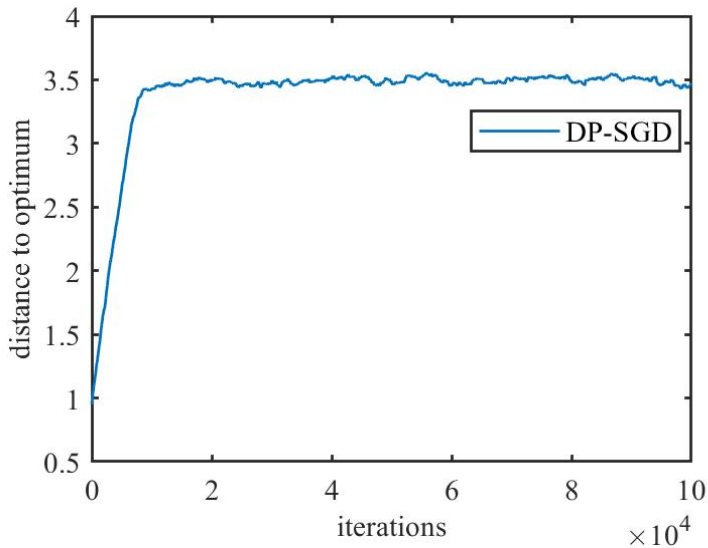
Clipped gradient for any  $w \in [-2, 2]$

$$\mathbb{E}[\text{Clip}(\nabla_x \ell(w; x_i), 1)] = 0$$

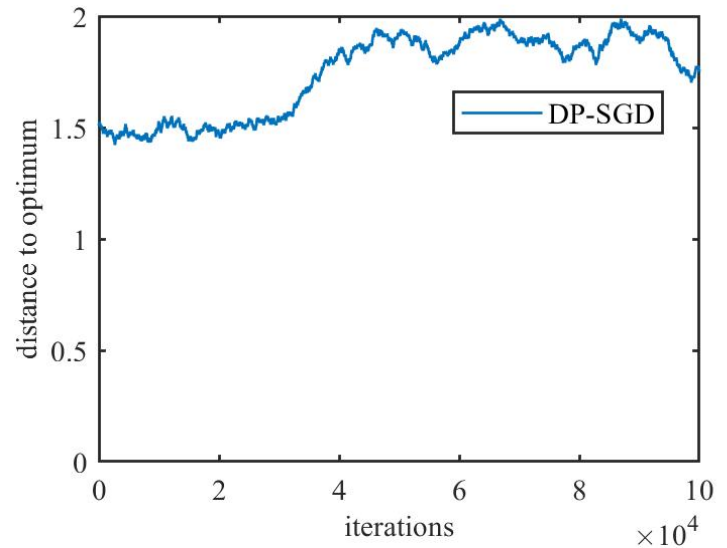
$\Rightarrow$  does not converge to opt



# Adversarial Effects of Clipping



Example 1



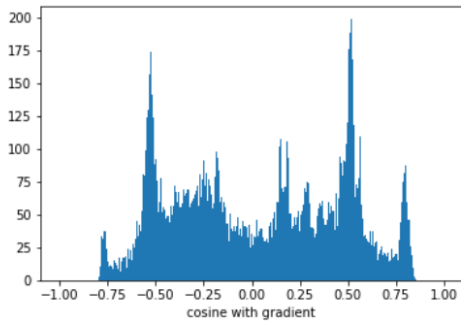
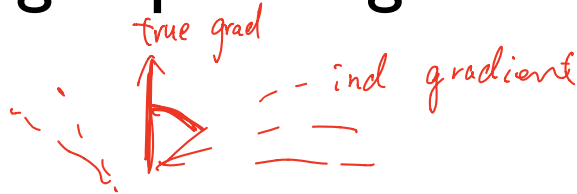
Example 2

*Do these occur in practical instances?*

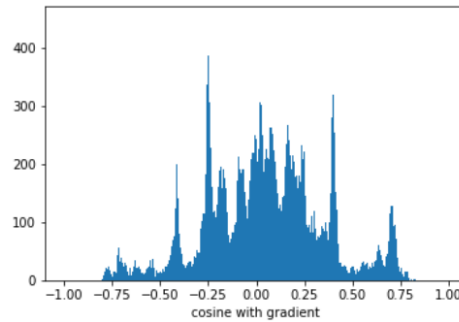
# DP-SGD on MNIST

- DP-SGD with Clip norm  $G = 1$   
60 epochs,  $\epsilon \approx 3$ , test accuracy  $\approx 96.5\%$
- DP-SGD with Clip norm  $G = 0.1$   
60 epochs,  $\epsilon \approx 3$ , test accuracy  $\approx 92\%$

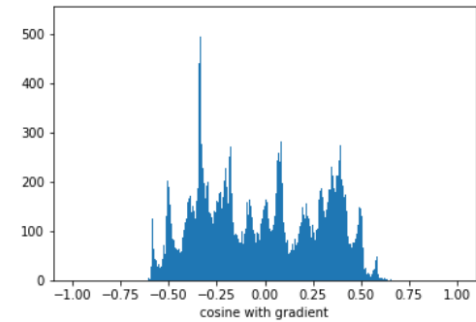
# A glimpse of gradient distribution



(a) Epoch 4



(b) Epoch 10



(c) Epoch 59

Histogram of cosine between stochastic gradients and true gradient

*Symmetric structures in gradients still lead to convergence under clipping.*

# Gradient Distribution of NN

## Visualization with random projection

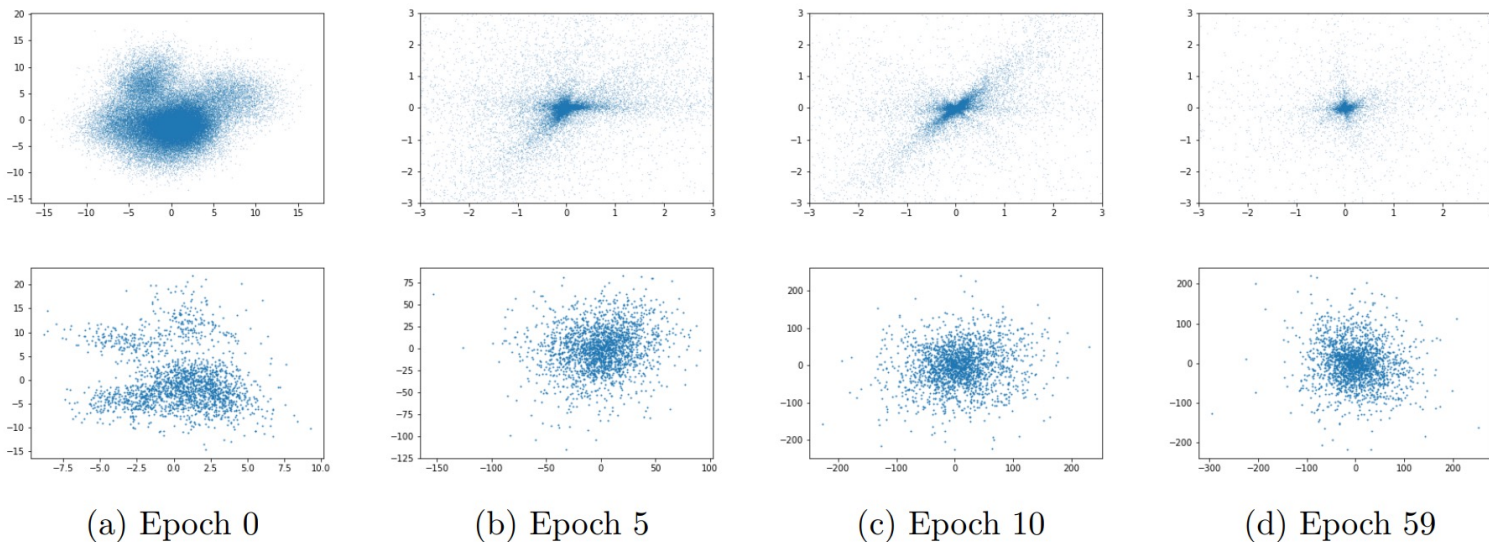
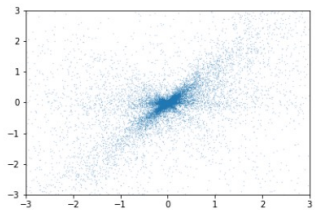


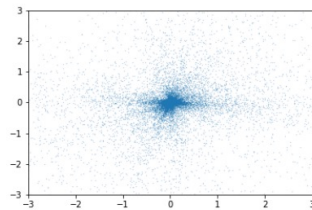
Figure 1: Gradient distributions on MNIST (top row) and CIFAR10 (bottom row) at the end of different epochs (indexed by columns). The gradients for epoch 0 are computed at initialization (before training).

# Gradient Distribution of NN

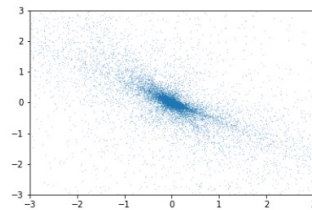
## Multiple random projections



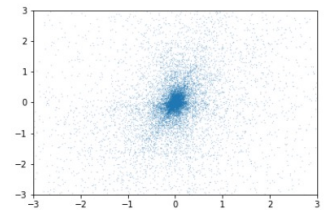
(a) Repeat 1



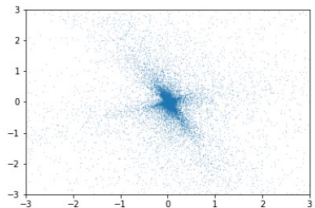
(b) Repeat 2



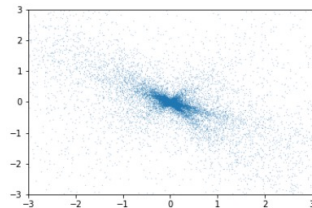
(c) Repeat 3



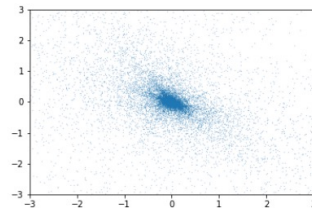
(d) Repeat 4



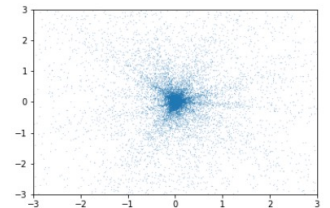
(e) Repeat 5



(f) Repeat 6



(g) Repeat 7



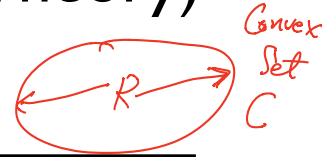
(h) Repeat 8

Figure 2: Gradient distributions on MNIST at the end of epoch 9 projected using different random matrices.



# Convergence Guarantee for DP-SGD (in Theory)

Consider DP-SGD with Projection



Theorem: Let  $L: C \rightarrow \mathbb{R}$  be convex and  $L$ -Lipschitz.

Suppose  $C \subseteq \mathbb{R}^d$  is a convex set with diameter  $R$ .

Let  $w^*$  be the minimizer of  $L$  in the set  $C$ .

• For regular SGD (w/ projection)

$$L(\hat{w}) - L(w^*) \leq \frac{RL}{\sqrt{T}}$$

• For DP-SGD (w/ projection),

$$\mathbb{E} [L(\hat{w}) - L(w^*)] \leq \mathcal{O} \left( \frac{RL \sqrt{d \ln(1/\delta)}}{n \epsilon} \right)$$

# Leveraging low-dimensional structure in gradients

- *Yingxue Zhou, Z. S.W., Arindam Banerjee*  
“Bypassing the Ambient Dimension: Private SGD with Gradient Subspace”  
In ICLR 2021

# Dimensionality

$$O\left(\frac{C\sqrt{d \ln(1/\delta)}}{n\epsilon}\right) + \frac{1}{T} \sum_{t=1}^T W_{\nabla f(x_t), C}(\tilde{p}_t, p_t)$$



DP-SGD without clipping  
Depends on ambient  
dimension  $d$

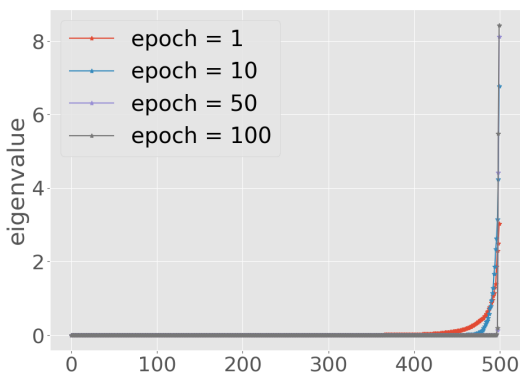


Clipping bias

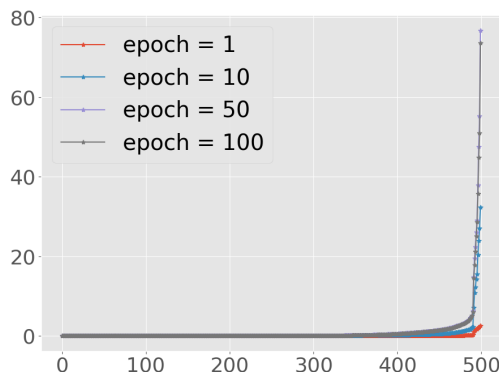
# Spectrum of Gradient Second Moments

Eigenvalues of  
$$M_t = \mathbb{E}[\nabla \ell(x_t, s_i) \nabla \ell(x_t, s_i)^\top]$$

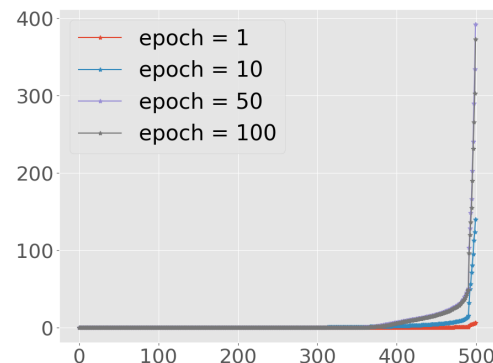
SGD



DP-SGD  $\sigma = 1$



DP-SGD  $\sigma = 2$



Order of eigenvalues from largest to smallest  
Ambient dimension  $d \approx 130,000$

[ZWB21]

# Projected DP-SGD (PDP-SGD)

Assume small amount of public data (no privacy concern)

## PDP-SGD [ZWB21]

- For  $t = 1, \dots, T$ 
  - Gradient estimate on a mini-batch  $B_t$  :  
 $\tilde{g}_t \leftarrow$  noisy gradient estimate with Gaussian noise
  - Use public data to compute projection  $\Pi_k$  onto the top- $k$  eigenspace of  $M_t$
  - Update :  
$$x_{t+1} = x_t - \eta \Pi_k \tilde{g}_t$$

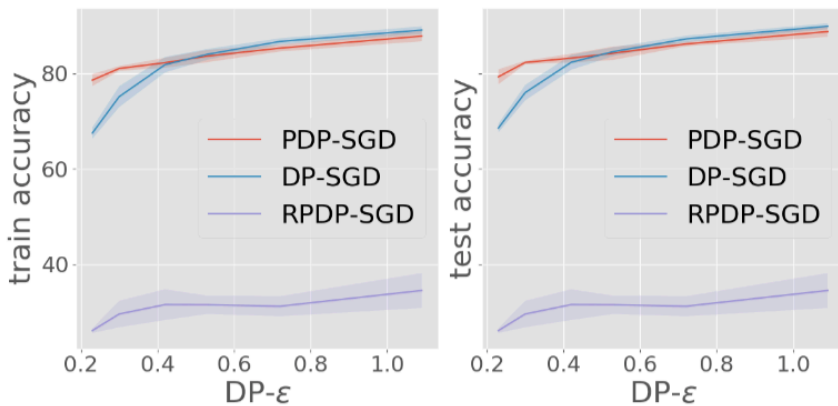
# Balancing two sources of error

- Error due to projection

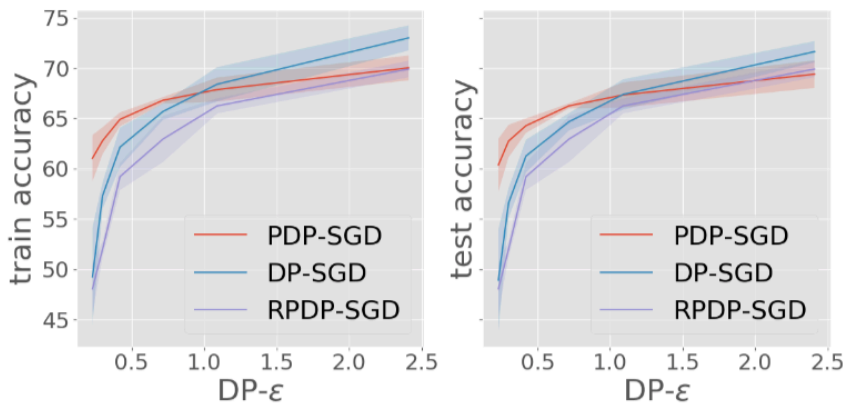
$$\|\Pi_k \nabla \ell(x; s_i) - \nabla \ell(x; s_i)\|$$

- Gradient perturbation in the subspace  $\approx \frac{\sqrt{k}}{n\epsilon}$   
(from  $\sqrt{d}$  to  $\sqrt{k}$ )

*PDP-SGD  $\rightarrow$  DP-SGD  
with gradient  
projection  
onto low-dim  
subspace.*

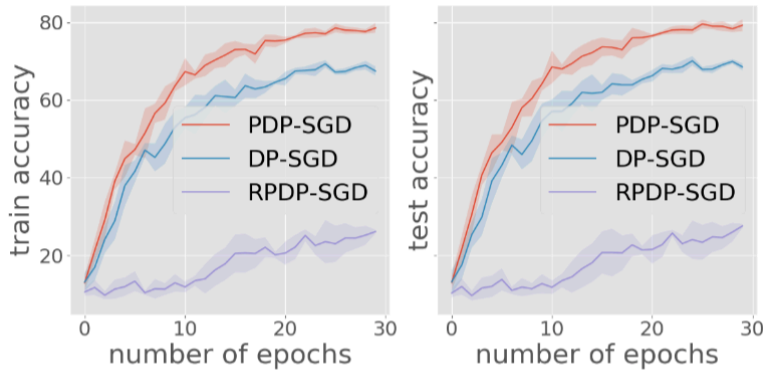


(a) MNIST

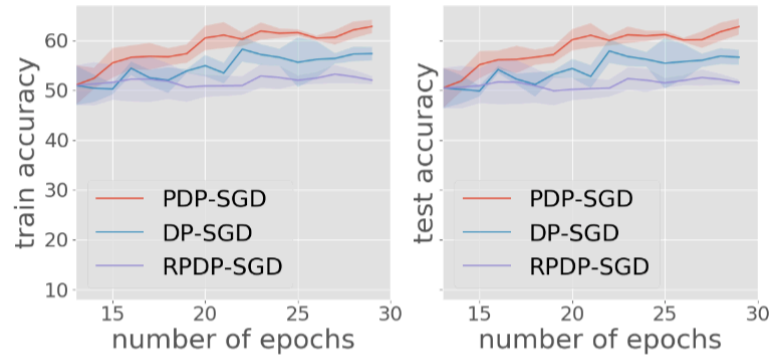


(b) Fashion MNIST

# Training Dynamics



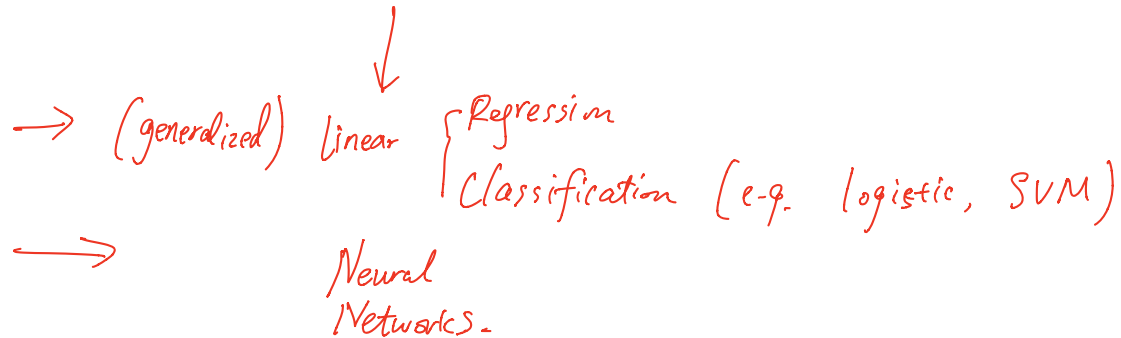
(a) MNIST,  $\epsilon = 0.23$



(b) Fashion MNIST,  $\epsilon = 0.30$

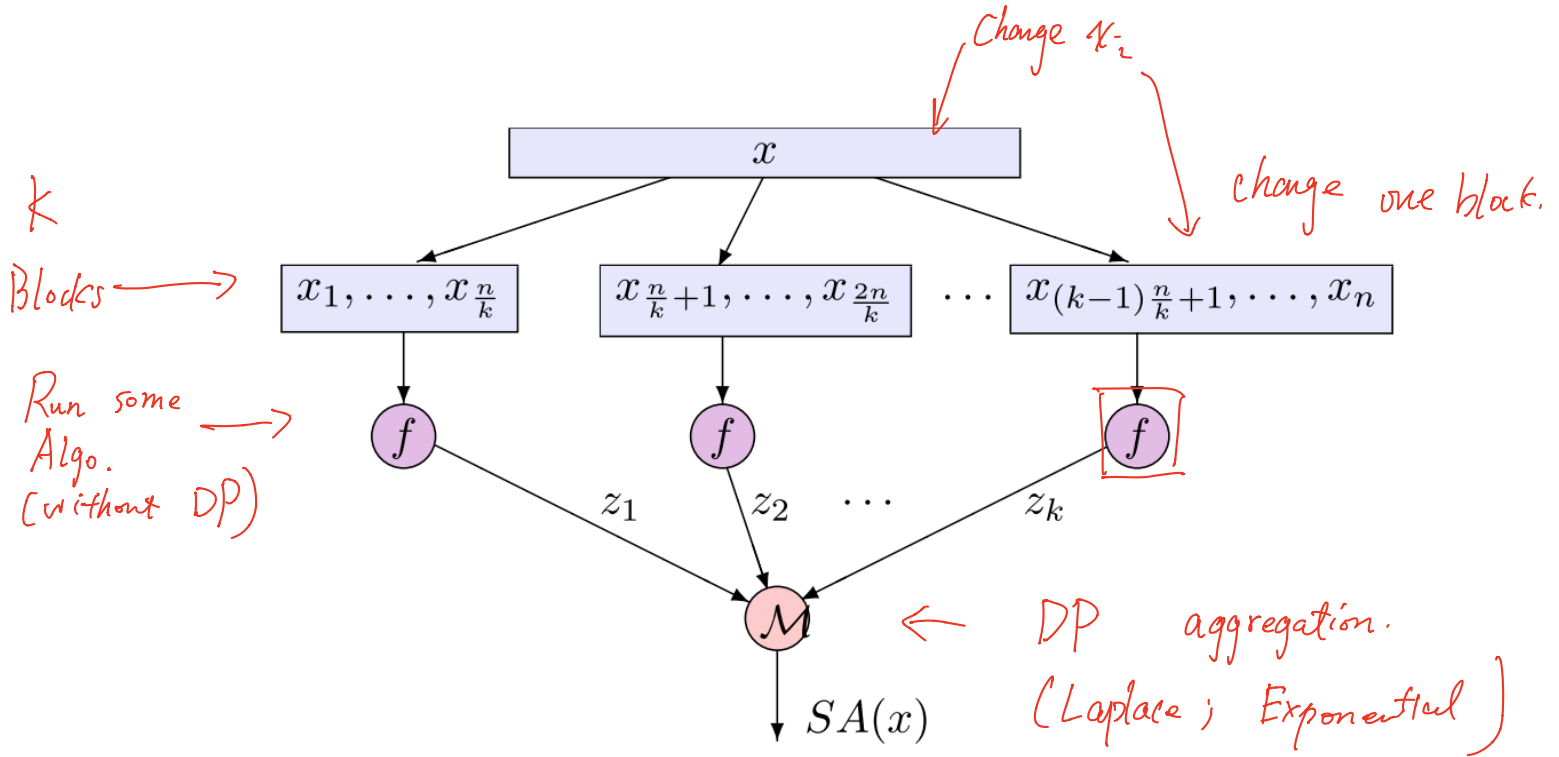


# What if DPSGD is not applicable?



Reduce the problem to  
Non-private ML.

# Subsample and Aggregate



# Private Aggregation of Teacher Ensembles (PATE)

