

# Lecture 14

## Modern Machine Learning Neural Networks

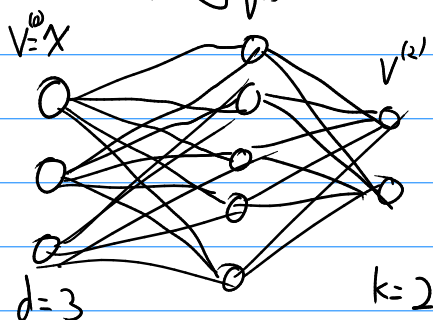
MLP

Dataset  $D, (x, y)$

Loss fn  $l$

$$\mathcal{L}(\theta, D) = \sum_{i \in [n]} l(\theta, x_i, y_i)$$

$x \in \mathbb{R}^d, y \in [k]$

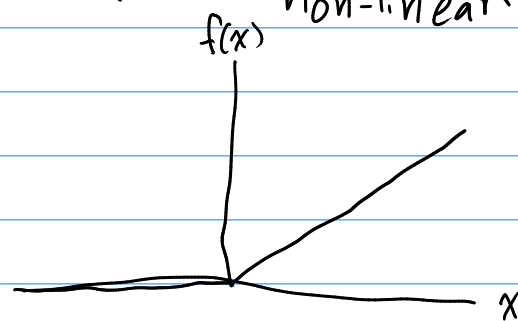


$$v^{(i)} = f^{(i)}(W^{(i)} v^{(i-1)} + b^{(i)})$$

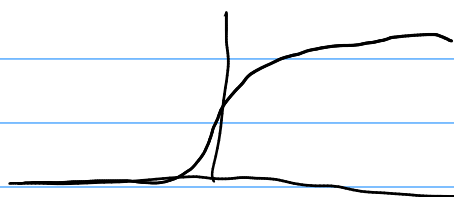
vector  $f$  applied entrywise

Params  $\begin{cases} W^{(i)} \in \mathbb{R}^{m_i \times m_{i-1}} \\ b^{(i)} \in \mathbb{R}^{m_i} \end{cases}$  activation fn  $f^{(i)}: \mathbb{R} \rightarrow \mathbb{R}$  Non-linear

ReLU Rectified linear unit  $f(x) = \max(0, x)$



Sigmoid:  $f(x) = \frac{1}{1+e^{-x}}$   
tanh:  $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



$$\dim(o) = k$$

$$\hat{y} = \text{softmax}(o)$$

$$\hat{y}_i = \frac{\exp(o_i)}{\sum_{j \in [k]} \exp(o_j)}$$

Not convex

$$l(\theta, x, y) = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

$\hat{y}_j = 1$  in coord of true label  $l$   
 $0$  elsewhere

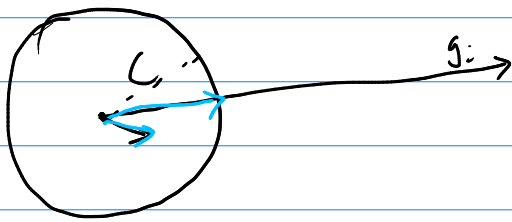
# Differentially Private Stochastic Gradient Descent

Algorithm

BST '14

ACGM MT2 '16

- each p.t. selected w.p.  $\frac{L}{n}$
- |  |   |   |
|--|---|---|
| 1. Select $(x_i, y_i)$ randomly                      | ← | 1. "Lot" of points, size $L$ (in exp.)                                      |
| 2. Compute $n \nabla \ell(\theta_t, x_i, y_i) + G_t$ |   | 2. Compute $g_i = \nabla \ell(\theta_t, x_i, y_i) \forall i \in \text{lot}$ |
| 3. Step in - direction                               |   | 3. Clip $g_i$ to $l_2$ ball of rad. $C$ , average them                      |
|  |   | 4. Add noise  |
|  |   | 5. Step in - direction  |



Clipping

## Sampling Schemes

Select 1 rnd pt.



Choose  $L$  each p.t. w.p.  $\frac{L}{n}$ .

- Poisson Sampling
- Sampling w/o replacement
- Sampling w/ replacement
- Permute dataset, iterate over it, selecting  $B$  pts each time

Privacy Loss  $q = \frac{L}{n}$

BST'14: Gauss + subsamp + adv. comp.

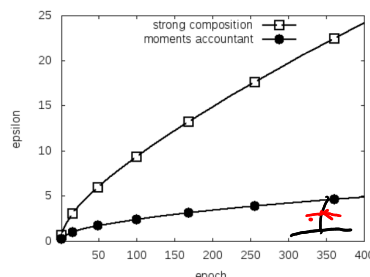
$(\epsilon, \delta)$ -DP  $\rightarrow (O(q \epsilon \sqrt{T \log(1/\delta)}), \delta)$ -DP

$(O(q \epsilon \sqrt{T}), \delta)$ -DP

Jagielski, Ulman, Oprea

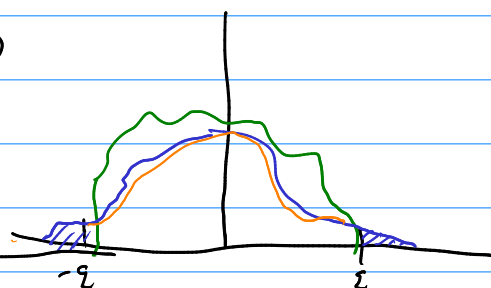
$$L_{Y||Z} = \ln \left( \frac{\Pr[Y=t]}{\Pr[Z=t]} \right)$$

# of iters



$\epsilon$ -DP:  $|L_{M(x)||M(x')}| \leq \epsilon$  w.p.1 ( $\forall x, x'$ )

$(\epsilon, \delta)$ -DP:  $|L_{M(x)||M(x')}| \leq \epsilon$  w.p.  $1 - \delta$



$$\ln E_{t \sim M(x)} \left[ \left( \frac{\Pr[M(x)=t]}{\Pr[M(x')=t]} \right)^\lambda \right] \leq \gamma \Rightarrow (\epsilon, \delta)\text{-DP}$$

$$\delta = \Pr_{t \sim M(x)} \left[ \ln \left( \frac{\Pr[M(x)=t]}{\Pr[M(x')=t]} \right) \geq \epsilon \right]$$

$$= \Pr \left[ \exp \left( \lambda \ln \left( \frac{\Pr[M(x)=t]}{\Pr[M(x')=t]} \right) \right) \geq \exp(\lambda \epsilon) \right]$$

$$\leq \frac{\exp(\gamma)}{\exp(\lambda \epsilon)} = \exp(\gamma - \lambda \epsilon)$$

Rényi-DP [Mironov'17] RDP

## Moments Accountant

1.  $\forall \lambda$ , compute  $\gamma$  for:  
Poisson subsample + Gaussian Mech.)
  2. Apply comp. for  $T$  iterations
  3. RDP  $\rightarrow (\epsilon, \delta)$ -DP (use best  $\lambda$ )
- Not data dependent

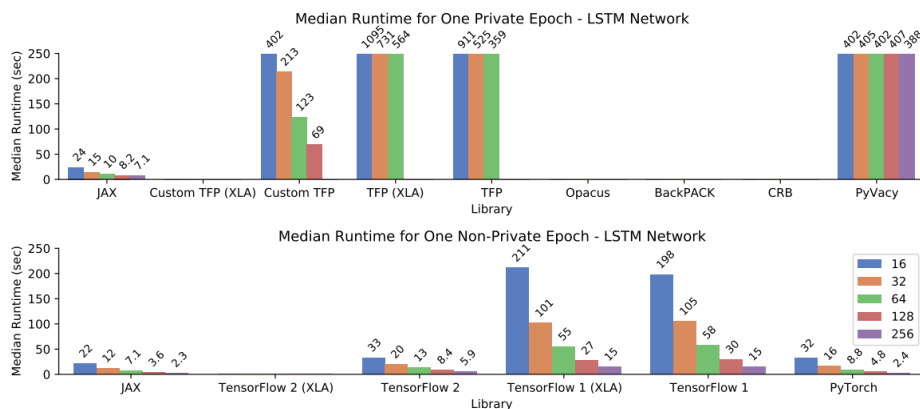
Accuracy? 0 1 3 5

MNIST 99.8% Non private  
98.1% w/privacy.  $(2.93, 10^{-5})$ -DP

CIFAR-10 99.7% Non-privately  
66.2%  $(7.53, 10^{-5})$ -DP

## Speed of DPSGD

Clipping = Slow!



## Changes in Architecture

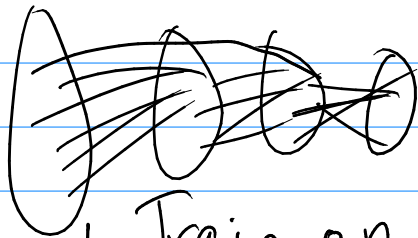
- Activation fn's  
Non-private: ReLU  
Private: tanh!
- Size vs. Noise

## Exploiting Public Data

- Hyperparameter tuning

Liu Talwar '19

Tune on public data



1. Train on a public dataset
2. Freeze most layers, privately train rest

# Private Aggregation of Teacher Ensembles

PATE

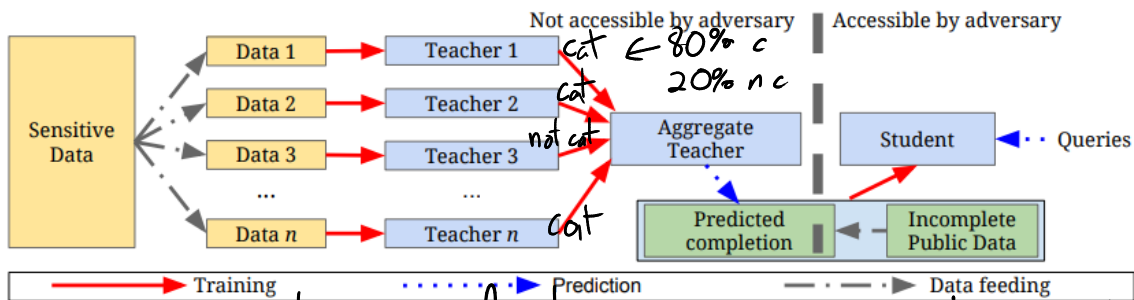
PAEGT '17  
BTT '18

Papernot et al '18

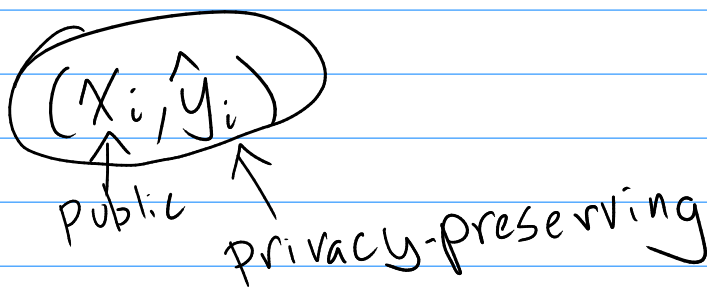
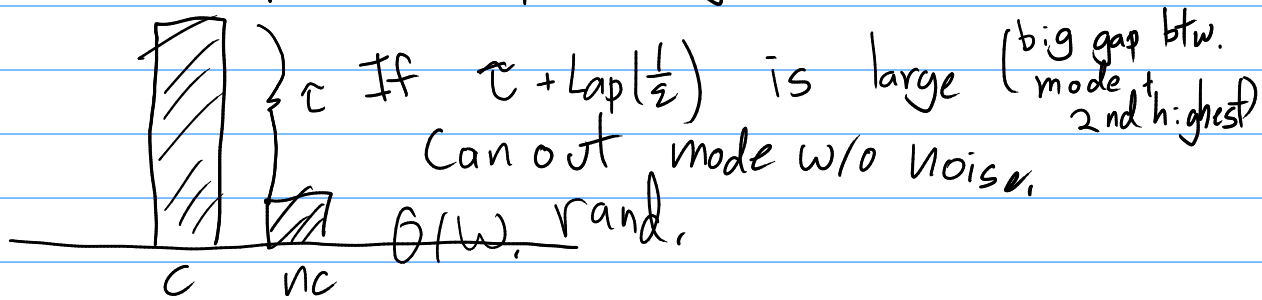
- Black Box - Non Private learners at core
- Private labeled training data + Public unlabeled training data

Stability - Sample and aggregate Nissim Raskhodnikova Smith '07

"Sample" 1. Run  $n_p$  algo on disjoint datasets  
 "Aggregate" 2. Combine results in a DP way.  
 ( $n_p$  complexity)  $\cdot \frac{1}{\epsilon}$



1. One private prediction (feature vector public)  
 Given  $x$ , output  $\hat{y}$  privately TS'13





$(x_1, \hat{y}_1), (x_2, \hat{y}_2), \dots$

synthetic private dataset

2. Many  $(T)$  classification tasks.

- Advanced comp  $\rightarrow \sqrt{T}$  cost

- Sparse vector  $\rightarrow \sqrt{c}$  cost

$\uparrow$  # of  $\rightarrow$  thresh.

which  $q$ 's are unstable

3. A classifier??

trained private classifier

N.P. learning algo