**CS208: Applied Privacy for Data Science**
**Machine Learning under DP**

School of Engineering & Applied Sciences
Harvard University

March 8, 2022

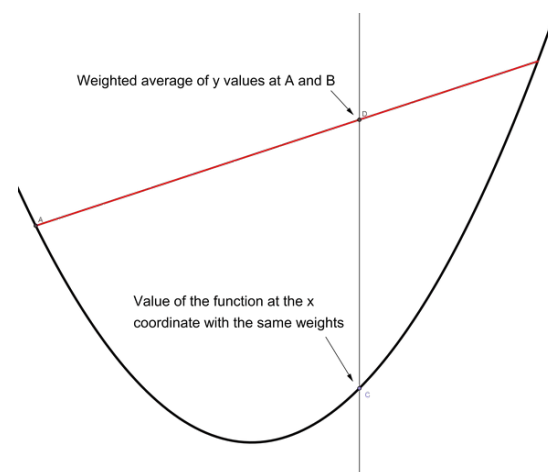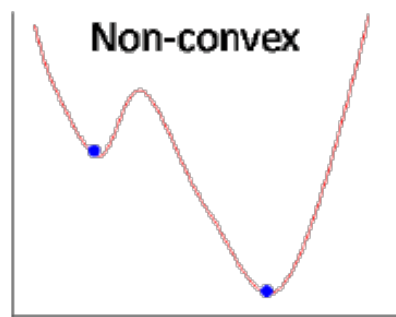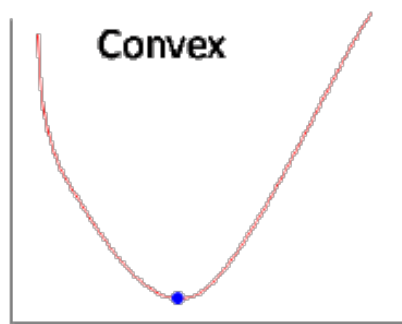# ML Inputs and Loss Functions

- Data: $(x_1, y_1), \ldots, (x_n, y_n) \sim \mathcal{P}$
  - Examples $x_i \in \mathcal{X}$ $d$-dimensional, discrete or continuous
  - Labels $y_i \in \mathcal{Y}$ 1-dimensional, discrete or continuous
  - $\mathcal{P}$ typically unknown
- A loss function:
  - $\ell : \Theta \times \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ $\quad$ $\ell(\theta | x_i, y_i)$ measures ``loss''
  - Define L: $\Theta \to \mathbb{R}$ $\quad$ $L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(\theta | x_i, y_i)$
  - Example: Squared loss: $\ell(\theta | x, y) = |(\beta_1 x + \beta_0) - y|^2$.
- Goal: output $\hat{\theta} \in \Theta$ s.t.
$$L(\hat{\theta}) \approx \min L(\theta)$$

# Convexity

- Def: $L$ is convex if for all points $\vec{a}, \vec{b}$, we have

$$L\left(\frac{\vec{a} + \vec{b}}{2}\right) \leq \frac{L(\vec{a}) + L(\vec{b})}{2}.$$

- Convex functions have no local minima

Weighted average of y values at A and B
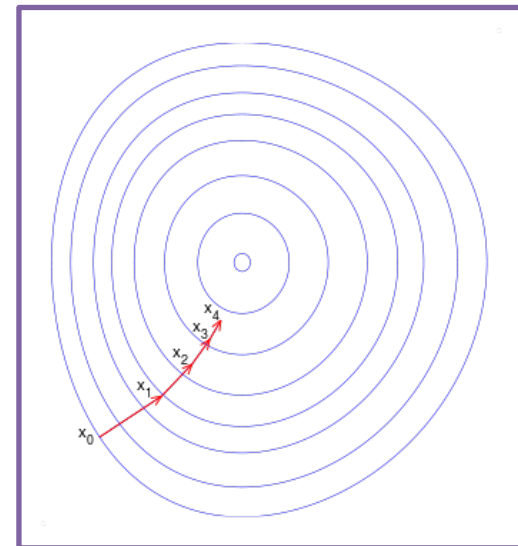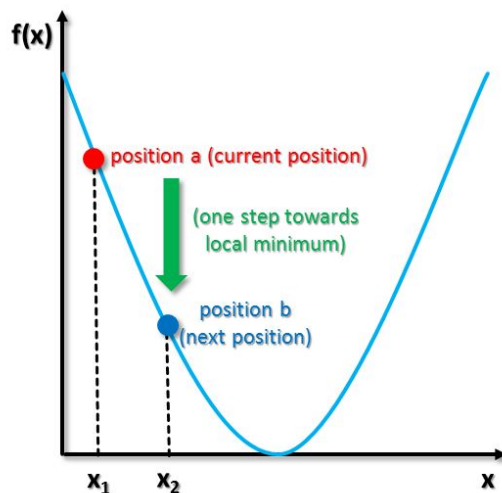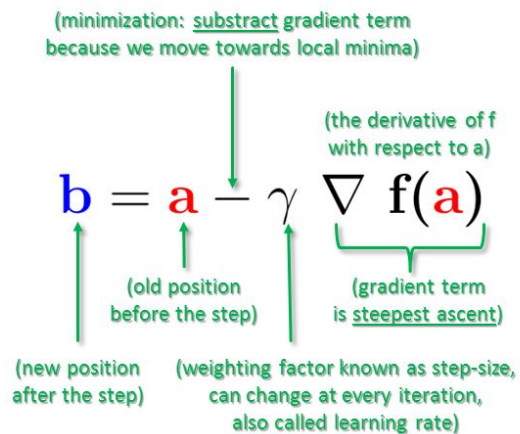
Value of the function at the x coordinate with the same weights

Convex

Non-convex

- Loss function for logistic regression is convex
  - No closed form solution for minimum, but it is easy to find

# Gradient Descent

- Proceed in steps
- Start from (carefully chosen) initial parameters $\hat{\theta}_0$
- At each step, move in direction opposite to the gradient of the loss $\nabla L(\hat{\theta} \mid \vec{x}, \vec{y})$.
- Gradient is the vector of partial derivatives



(minimization: substract gradient term because we move towards local minima)

(the derivative of f with respect to a)

$$\mathbf{b} = \mathbf{a} - \gamma \, \nabla \, \mathbf{f(a)}$$

(old position before the step)

(gradient term is steepest ascent)

(new position after the step)

(weighting factor known as step-size, can change at every iteration, also called learning rate)



f(x)

position a (current position)

(one step towards local minimum)

position b (next position)

$x_1$    $x_2$    x



$x_0$, $x_1$, $x_2$, $x_3$, $x_4$

# Gradient Descent

- Specify
  - Number of steps $T$
  - Learning rate $\eta$
- Pick initial point $\hat{\theta}_0 \in \Theta$
- For $t = 1$ to $T$
  - Compute gradient
  $$g_t = \nabla L\left(\hat{\theta}_{t-1}\right) = \frac{1}{n}\sum_i \nabla \ell\left(\hat{\theta}_{t-1}|x_i, y_i\right)$$
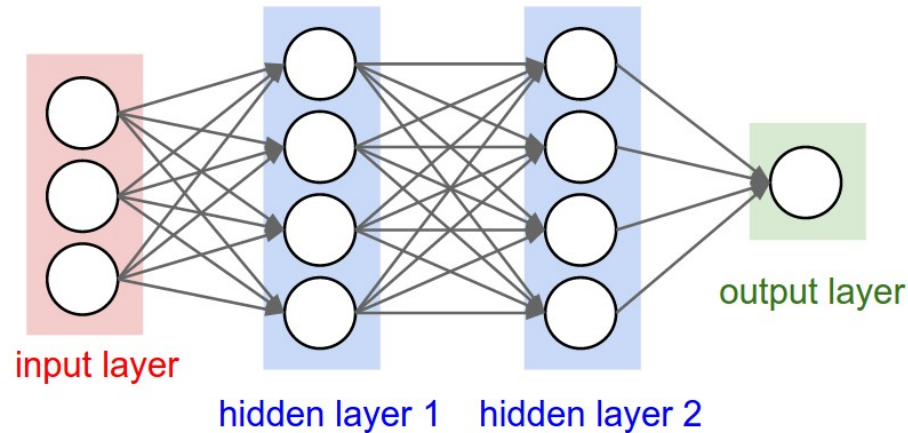  - $\hat{\theta}_t = \hat{\theta}_{t-1} - \eta \cdot g_t$
- Output $\hat{\theta} = \sum_{t=1}^{T} \hat{\theta}_t$   or   $\hat{\theta}_T$
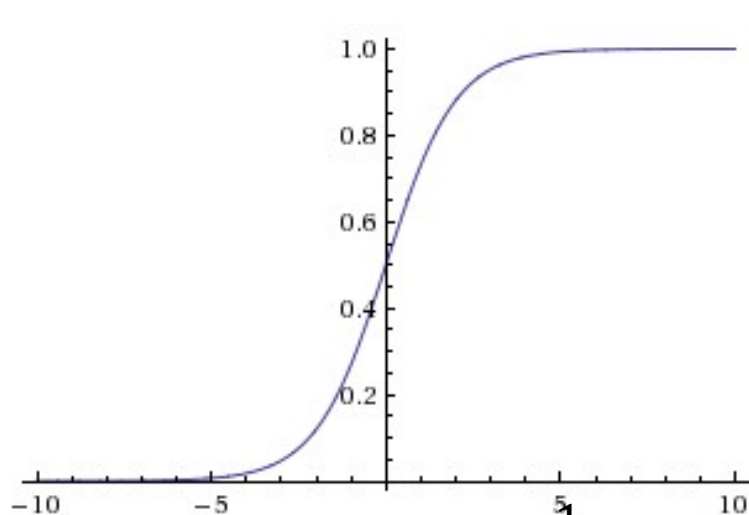
Average iterate

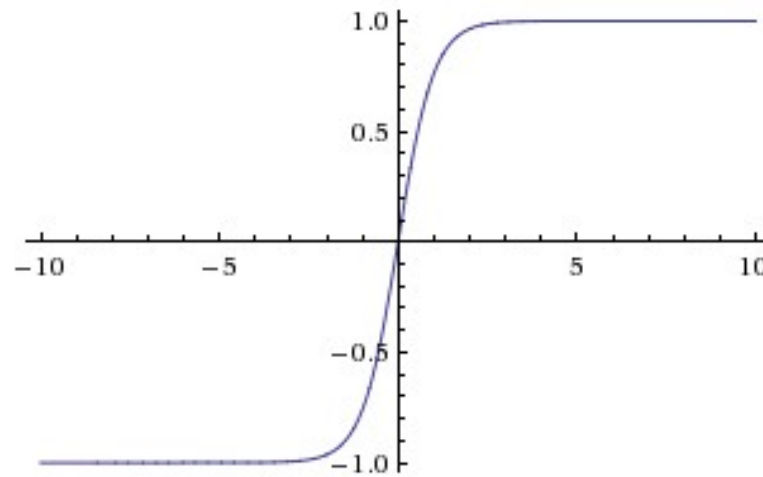Last iterate

# Gradient Descent for Neural Networks



input layer
hidden layer 1    hidden layer 2
output layer

- Each node is a linear function of inputs (specified by $\theta$) composed with a nonlinear "activation" function
- Gradient of Loss function can be computed quickly
  - Using chain rule (technique called "backpropagation")
- But no longer convex, has many local minima
  - Can get stuck in a bad place
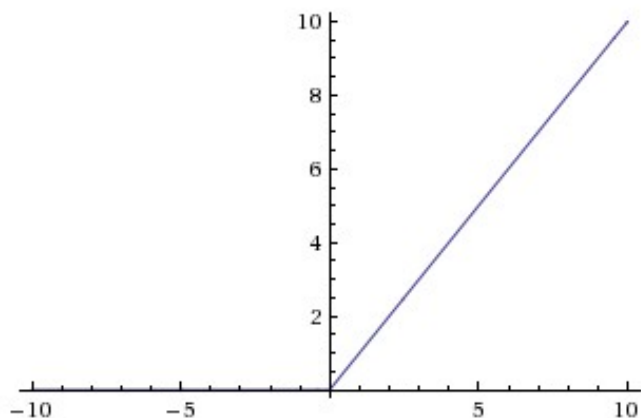  - But works well in practice!
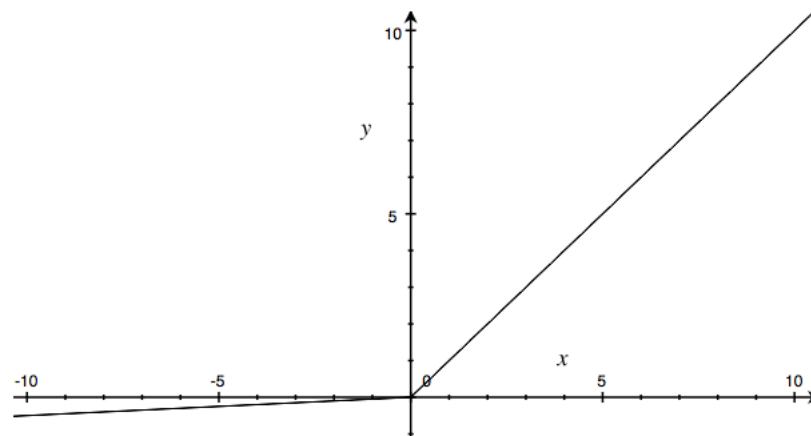
# Common Activation Functions



Sigmoid $\sigma(x) = \dfrac{1}{1+e^{-x}}$



$\tanh(x) = 2\sigma(2x) - 1$



$\text{ReLU}(x) = \max(0, x)$



Leaky $\text{ReLU}(x) = \max(0.05x, x)$

# DP for Vector-Valued Functions

- Let $f : \mathcal{X}^n \to \mathbb{R}^k$, and $M(x) = f(x) + Z$ for noise $Z \in \mathbb{R}^k$.
- global $\ell_2$-sensitivity of $f$ is

$$\text{GS}_{f,\ell_2} \stackrel{\text{def}}{=} \max_{x \sim x'} \|f(x) - f(x')\|_2.$$

$$\|z\|_2 = \left( \sum_j |z_j|^2 \right)^{1/2}$$

- Gaussian Mechanism: $Z \sim \mathcal{N}\left( \vec{0}, 2\left(\frac{\text{GS}_{f,\ell_2}}{\varepsilon}\right)^2 \cdot \ln\frac{1.25}{\delta} \cdot I_k \right)$

  – independent Gaussian noise per coordinate.

# Robustness to Noise in Gradient Estimation

- For efficiency:

  Sample a minibatch $B \in \{1, 2, \ldots, n\}$

  Gradient estimate $\tilde{g}_t = \frac{1}{|B|} \sum_{i \in B} \nabla \ell(\hat{\theta}_{t-1}, x_i, y_i)$

  Stochastic Gradient Descent (SGD)!

- For privacy:

  Add Gaussian Noise $\tilde{g}_t = g_t + \mathcal{N}(0, \sigma^2 I)$

  In both cases, $\tilde{g}_t$ is an unbiased estimate of $g_t$: $E[\tilde{g}_t] = g_t$

# DP Gradient Descent

[Williams-McSherry`10, …]

- Specify
  - Number of steps $T$
  - Learning rate $\eta$
  - Privacy parameters $\varepsilon, \delta$
  - Clipping parameter $C$. Write $[\vec{z}]_\Delta = \vec{z} \cdot \max\left(1, \frac{C}{\|\vec{z}\|_2}\right)$.
  - Noise variance $\sigma^2 = \text{TBD}(T, \varepsilon, \delta, C)$.
- Pick initial point $\hat{\theta}_0$
- For $t = 1$ to $T$
  - Estimate gradient as noisy average of clipped gradients
    $$\hat{g}_t = \frac{1}{n} \sum_i \left[\nabla \ell\left(\hat{\theta}_{t-1} | x_i, y_i\right)\right]_C + \mathcal{N}(0, \sigma^2 I)$$
  - $\hat{\theta}_t = \hat{\theta}_{t-1} - \eta \cdot \hat{g}_t$
- Output $\hat{\theta} = \sum_{t=1}^{T} \hat{\theta}_t$   or   $\hat{\theta}_T$

# Privacy Analysis

- Proof idea: Show releasing $(\hat{g}_1, \hat{g}_2, \ldots, \hat{g}_T)$ satisfies DP
  - Each step (releasing $\hat{g}_t$) satisfies $(\epsilon, \delta)$-DP
  - Adaptive composition across $T$ steps

# Privacy Analysis

- By Gaussian Mechanism, each iteration is $(\varepsilon_0, \delta_0)$-DP if

$$\sigma^2 = 2\left(\frac{C}{\varepsilon_0 n}\right)^2 \cdot \ln\frac{1.25}{\delta_0}$$

- By Advanced Composition for <span style="color:red">adaptive</span> queries, overall algorithm is $(\varepsilon, \delta)$-DP for:

$$\varepsilon = O\left(\varepsilon_0 \cdot \sqrt{T \ln(2/\delta)}\right)$$
$$\delta = 2T \cdot \delta_0$$

- Solving, suffices to use noise variance

$$\sigma^2 = O\left(\left(\frac{C}{\varepsilon n}\right)^2 \cdot T \cdot \ln\frac{T}{\delta} \cdot \ln\frac{1}{\delta}\right)$$

# Improved Analysis with "Concentrated DP"

[Dwork-Rothblum `16, Bun-Steinke `16]

- By Gaussian Mechanism, each iteration is $\varepsilon_0^2$ -zCDP if

$$\sigma^2 = \frac{1}{2}\left(\frac{C}{\varepsilon_0 n}\right)^2 \cdot \cancel{\ln \frac{1.25}{\delta_0}}$$

- By composition of zCDP, overall algorithm is $T \cdot \varepsilon_0^2$-zCDP.

- By properties of zCDP, overall algorithm is $(\varepsilon, \delta)$-DP for:

$$\varepsilon = T \cdot \varepsilon_0^2 + 2\sqrt{T \cdot \varepsilon_0^2 \cdot \ln(1/\delta)}$$

- Solving, suffices to use noise variance

$$\sigma^2 = O\left(\left(\frac{C}{\varepsilon n}\right)^2 \cdot T \cdot \ln \frac{1}{\delta} \cdot \cancel{\ln \frac{T}{\delta}}\right)$$

# DP **Stochastic** Gradient Descent (**S**GD)

[Jain-Kothari-Thakurta `12, Song-Chaudhuri-Sarwate `13, Bassily-Smith-Thakurta `14]

- Specify
  - Number of steps $T$, learning rate $\eta$, privacy parameters $\varepsilon, \delta$, clipping parameter C.
  - Batch size $B \ll n$ (for efficiency)
  - Noise variance $\sigma^2 = \text{TBD}(T, \varepsilon, \delta, C, B)$.
- Pick initial point $\hat{\theta}_0$
- For $t = 1$ to $T$

  - Select a random batch $S_t \subseteq \{1, \dots, n\}$ of size $B$.
  - Estimate gradient as noisy average of clipped gradients
  $$\hat{g}_t = \frac{1}{B} \sum_{i \in S_t} \left[ \nabla \ell \left( \hat{\theta}_{t-1} | x_i, y_i \right) \right]_C + \mathcal{N}(0, \sigma^2 I)$$
  - $\hat{\theta}_t = \hat{\theta}_{t-1} - \eta \cdot \hat{g}_t$
- Output $\hat{\theta} = \sum_{t=1}^T \hat{\theta}_t$   or   $\hat{\theta}_T$

# DP SGD: Improved Privacy Analysis

[Bassily-Smith-Thakurta `14, Abadi-Chu-Goodfellow-McMahan-Mironov-Talwar-Zhang `17]

- Privacy amplification by subsampling:
  If $S : \mathcal{X}^n \to \mathcal{X}^B$ outputs a random subset of $pn$ out of $n$ rows
  and $M : \mathcal{X}^B \to \mathcal{Y}$ is $(\varepsilon, \delta)$-DP, then
  $M'(x) = M(S(x))$ is $(ln(1 + (e^\epsilon - 1)p), p\delta)$-DP.
  - Keep $S_t$ secret; Use their randomness

$p\epsilon$

- Poisson sampling: choosing each point independently with probability p=B/n.
- Choosing B points without replacement
- Choosing B points with replacement

# DP SGD: Improved Privacy Analysis

[Bassily-Smith-Thakurta `14, Abadi-Chu-Goodfellow-McMahan-Mironov-Talwar-Zhang `17]

- Privacy amplification by subsampling:
  If $S : \mathcal{X}^n \to \mathcal{X}^B$ outputs a random subset of $pn$ out of $n$ rows and
  $M : \mathcal{X}^B \to \mathcal{Y}$ is $(\varepsilon, \delta)$-DP, then
  $M'(x) = M(S(x))$ is $(ln(1 + (e^\epsilon - 1)p), p\delta)$-DP.
  - Keep $S_t$ secret; Use their randomness

$p\epsilon$

- We can take $p = B/n$.
  - Unfortunately privacy amplification by subsampling does not hold for zCDP.
  - But similar analysis can be recovered using the "moments accountant" [Abadi et al. `17] or "truncated zCDP" [Bun et al. `18].

# DP SGD: Privacy Analysis

- By Gaussian Mechanism, each iteration is $(\varepsilon_0, \delta_0)$-DP if

$$\sigma^2 = 2\left(\frac{C}{\varepsilon_0 n}\right)^2 \cdot \ln\frac{1.25}{\delta_0}$$

- Subsampling + Gaussian is $(\varepsilon'_0, \delta'_0)$-DP

$$\varepsilon'_0 = \frac{B}{n}\epsilon_0 \qquad \delta'_0 = \frac{B}{n}\delta_0$$

- Advanced Composition over $T$ iterations

$$\varepsilon = O\left(\varepsilon'_0 \cdot \sqrt{T\ln(2/\delta)}\right)$$
$$\delta = 2T \cdot \delta'_0$$

> Moments accountant $O(\frac{B}{n}\epsilon_0 \sqrt{T}, T\delta)$; [WBK19,MTZ19] better analysis; [JUO20] attack

**Main Idea: Privacy amplification by subsampling + Composition**

# Neural Networks & Privacy

- Choice of the model architecture
  - The Gaussian noise proportional to the square root of number of parameters.

- Hyperparameter tuning
  - If we run analyses on the training data with various hyperparameter settings, and choose the best one (any problems?)
  - Doing this privately (with additional cost in privacy)
  - Use public dataset (CIFAR-100 dataset for CIFAR-10 training)

- Still can be improved…
  - MNIST(99.8% baseline)   98.1%   $(2.93, 10^{-5}) - DP$
  - CIFAR-10(99.7% baseline)   66.2%   $(7.53, 10^{-5}) - DP$

Results from [PTS+ 20]

# Differentially Private Empirical Risk Minimization

# Supervised ML Output

Primary Goal (risk minimization):

- Find $\theta \in \Theta$ minimizing $L(\theta) = \mathrm{E}_{(x,y) \sim \mathcal{P}}[\ell(\theta|x, y)]$.
- Difficulty: $\mathcal{P}$ unknown.

Subgoal 1 (empirical risk minimization (ERM)):

- Find $\theta \in \Theta$ minimizing $L(\theta|\vec{x}, \vec{y}) = \frac{1}{n}\sum_{i=1}^{n} \ell(\theta|x_i, y_i)$.
- Turns learning into optimization.
- Difficulty: overfitting*

Subgoal 2 (regularized ERM):

- Find $\theta \in \Theta$ minimizing $L(\theta|\vec{x}, \vec{y}) = \frac{1}{n}\sum_{i=1}^{n} \ell(\theta|x_i, y_i) + R(\theta)$.
- $R(\theta)$ typically penalizes "large" $\theta$, can capture Bayesian prior.

*Fact: DP automatically helps prevent overfitting! [Dwork et al. `15]

# Output Perturbation

[Chaudhuri-Monteleoni-Sarwate `11]

$$M(\vec{x}, \vec{y}) = \text{argmin}_\theta \left( \frac{1}{n} \sum_{i=1}^n \ell(\theta | x_i, y_i) + R(\theta) \right) + \text{Noise}$$

Challenge: bounding sensitivity of $\theta_{opt} = \text{argmin}_\theta(\cdot)$

- Global sensitivity can be infinite (e.g. OLS regression)

- Global sensitivity can be bounded when $\ell$ is strictly convex, has bounded gradient (as a function of $\theta$), and $R$ is strongly convex. Even analyzing local sensitivity seems to require unique global optimum and using an optimizer that is guaranteed to succeed.

# Objective Perturbation

[Chaudhuri-Monteleoni-Sarwate `11]

$$M(\vec{x}, \vec{y}) = \mathrm{argmin}_\theta \left( \frac{1}{n} \sum_{i=1}^{n} \ell(\theta | x_i, y_i) + R(\theta) + R_{\mathrm{priv}}(\theta, \mathrm{noise}) \right)$$

Challenge: how to put noise in the objective function?

- [CMS11] use $R_{\mathrm{priv}}(\theta, v) = \langle \theta, v \rangle + c\|\theta\|^2$ where $v$ is sampled with probability density $\propto \exp(-c'\varepsilon\|v\|)$.

- Privacy proven under similar assumptions on $\ell$ and $R$ as before, plus $\ell$ having bounded Jacobian.

- Has better performance than output perturbation [CMS11].

# Exponential Mechanism for ML

[Kasiwiswanathan-Lee-Nissim-Raskhodnikova-Smith `11]

Use utility function

$$u\left((\vec{x}, \vec{y}), \theta\right) = -L(\theta | \vec{x}, \vec{y}) = -\frac{1}{n}\sum_{i=1}^{n} \ell(\theta | x_i, y_i) - R(\theta).$$

That is,

$$\Pr[M(\vec{x}, \vec{y}) = \theta] \propto e^{-\frac{\varepsilon}{2}\sum_{i=1}^{n} \ell(\theta | x_i, y_i) - \frac{\varepsilon n}{2}R(\theta).}$$
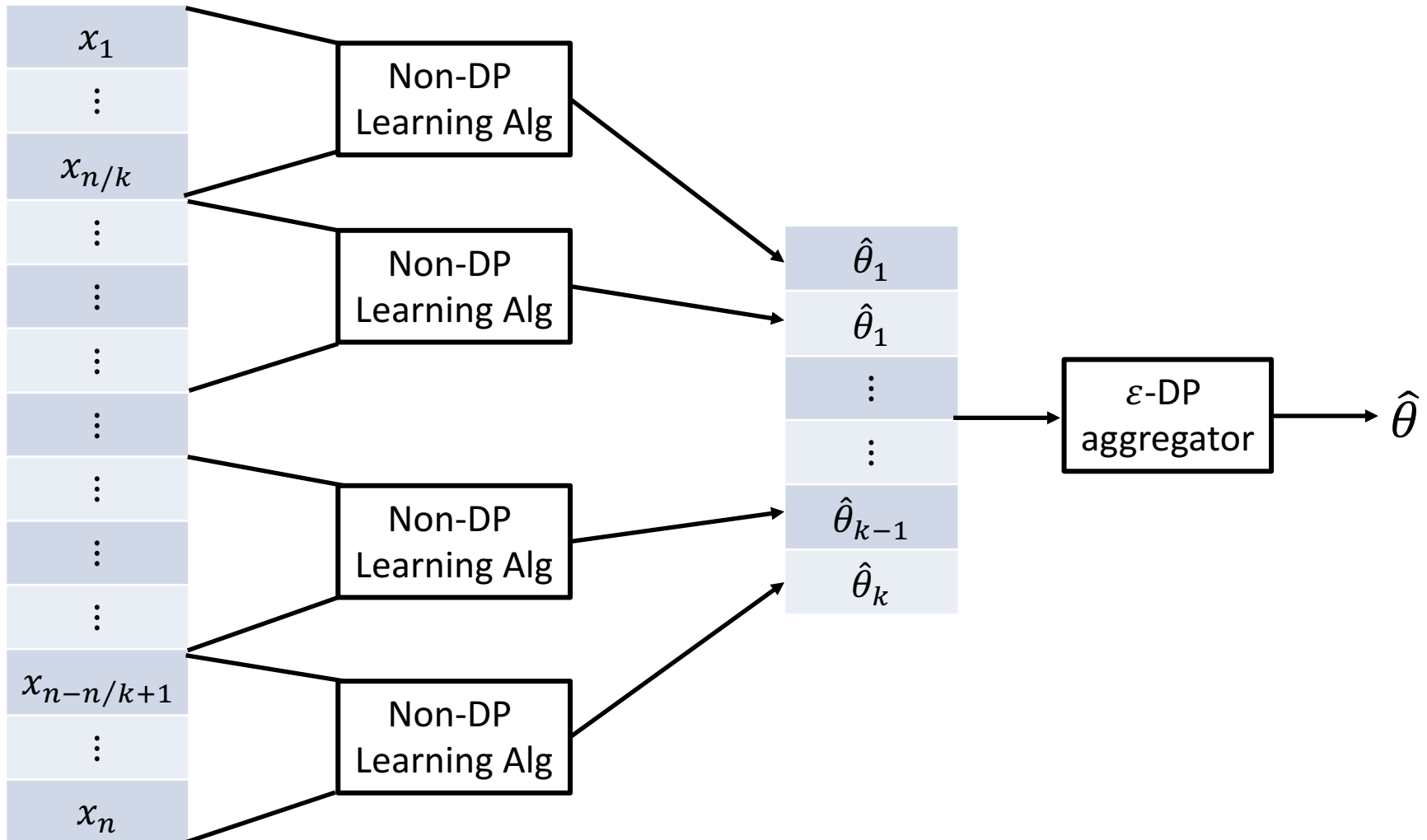
Is $\varepsilon$-DP if the loss functions are clipped to [0,1]. (why?)

Thm [KLNRS `11, informally stated]: anything learnable non-privately on a finite data universe is also learnable with DP (with larger $n$).

Problem: runtime often exponential in dimensionality of $\theta$.

# **Subsample & Aggregate**

[Nissim-Rakhodnikova-Smith `07, Smith `11]



Q: Why is this $\varepsilon$-DP?

# Subsample & Aggregate

[Nissim-Rakhodnikova-Smith `07, Smith `11]

- Typical aggregators: DP (clipped) mean, DP median
- Benefits:
  - Use any non-private estimator as a black box
  - Can give optimal asymptotic convergence rates: for many statistical estimators, variance is asymptotically $c_\theta/$(sample size), so variance of DP mean $\hat{\theta}$ is
  $$(1/k) \cdot (c_\theta \cdot k/n) + O(1/\varepsilon k)^2 = \left(1 + o(1)\right) \cdot c_\theta/n$$
  $$\text{if } k = \omega(\sqrt{n}).$$
- Drawbacks:
  - Dependence on dimension, model parameters, distribution can be bad.
  - Often takes very large sample size to kick in.
- Develop:
  - Private Aggregation of Teacher Ensembles (PATE) [PAE+17, PSM+18]

# Modifying ML Algorithms

- Another approach: decompose existing ML/inference algorithms into steps that can be made DP, like Statistical Queries (estimating means of bounded functions)

- Example: linear regression
  - $S_{xx}/n$, $S_{xy}/n$, $\bar{x}, \bar{y}$ are all statistical queries