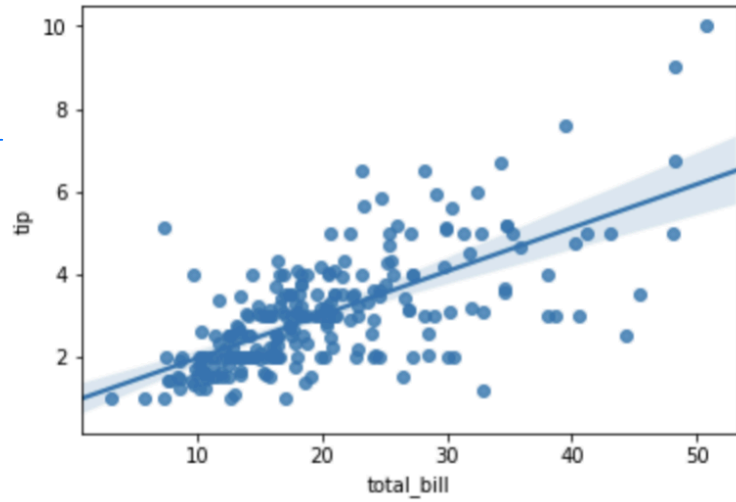


Lecture 16

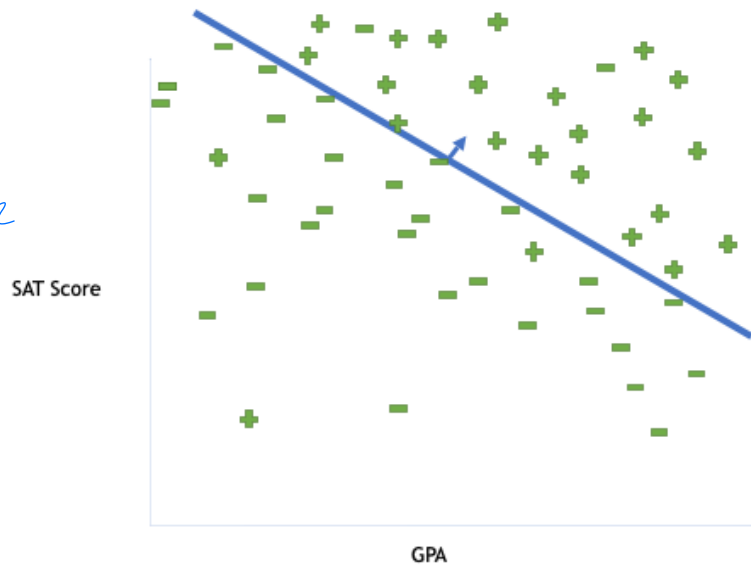
- Private Machine Learning
 - DP Gradient Descent
 - Privacy Analysis.

Optimization in ML

Linear Regression



Linear Classification



(Private) Optimization.

Given a data set $\mathcal{X} = (x_1, \dots, x_n)$

loss function: l

feasible set of parameters: $C \subseteq \mathbb{R}^d$
(weights)

Empirical Risk Minimization (ERM):

$$\min_{w \in C} \underbrace{L(w; \mathcal{X})}_{\text{Empirical Risk}} = \frac{1}{n} \sum_{i=1}^n l(w; x_i)$$

+ Optional
Regularization

e.g., $\lambda \|w\|_1$

Often: $C = \mathbb{R}^d$

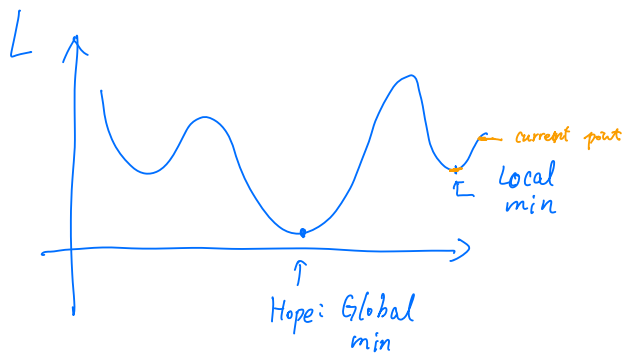
$$l: C \times X \mapsto \mathbb{R}$$

$l(w, x)$ measures "loss"

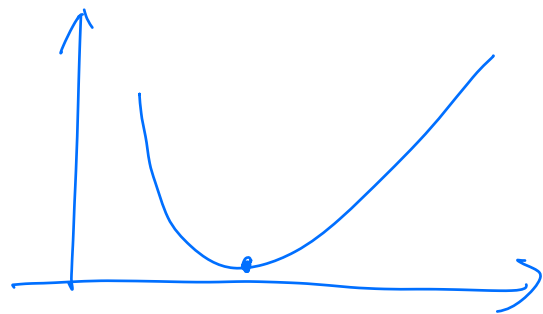
$$L: C \mapsto \mathbb{R}$$

$$L(w) = \frac{1}{n} \sum_{i=1}^n l(w, x_i)$$

Non-convex loss



Convex loss



Projected Gradient Descent (PGD)

$$\text{PGD} \left(\underset{\substack{\text{Loss} \\ \downarrow}}{L}, \underset{\substack{\text{Feasible} \\ \text{Set} \\ \downarrow}}{C}, \underset{\substack{\text{Learning} \\ \text{Rate} \\ \downarrow}}{\eta}, \underset{\substack{\text{\# rounds} \\ \downarrow}}{T} \right):$$

$$\text{Init: } w_0 \in C \quad (\text{any point})$$

$$\text{For } t = 1, \dots, T:$$

$$\text{gradient: } g_t = \nabla L(w_{t-1})$$

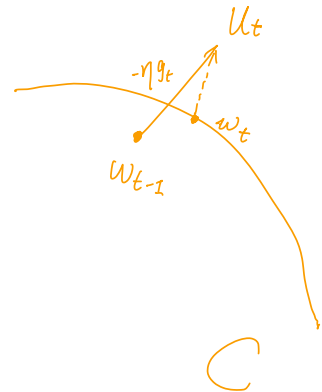
$$u_t \leftarrow w_{t-1} - \eta g_t$$

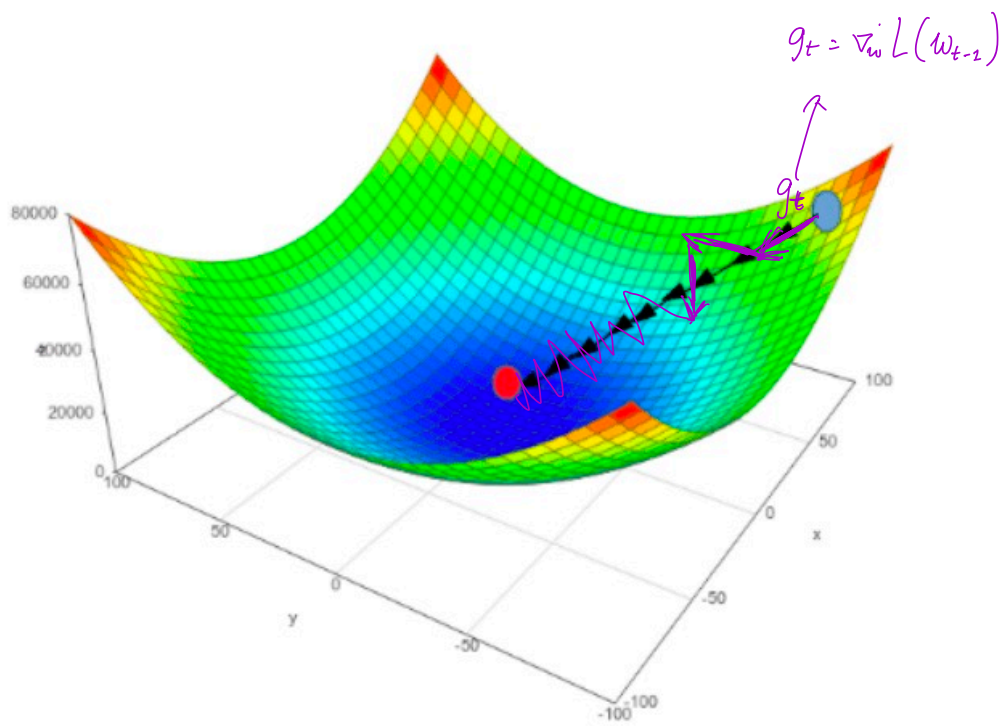
$$\text{projection: } w_t \leftarrow \underset{w \in C}{\text{argmin}} \|w - u_t\|_2$$

$$\text{Output} = w_T \quad \leftarrow \text{Last iterate}$$

or

$$\frac{1}{T} \sum_{t=1}^T w_t \quad \leftarrow \text{Average iterate.}$$





Robustness to noise in gradient estimation. (g_t)

Two sources of noise:

→ For efficiency:

Sample a minibatch $B \subseteq \{1, 2, \dots, n\}$
gradient estimate $\tilde{g}_t = \frac{1}{|B|} \sum_{i \in B} \nabla_w L(w_{t-1}, x_i)$

50 ↓ *50 million* ↓

→ For privacy: Add Gaussian Noise

$$\tilde{g}_t = g_t + N(0, \beta^2 I_d)$$

↑ from Gaussian mech.

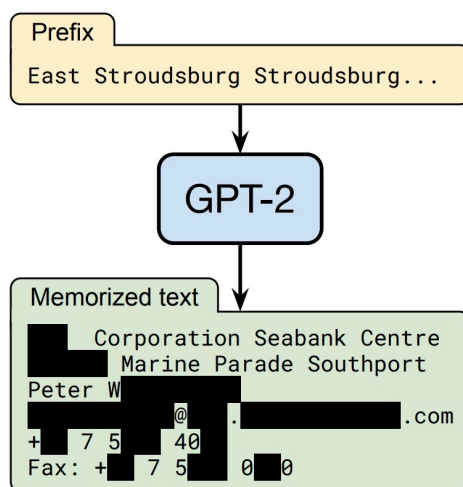
In both cases,

\tilde{g}_t is an unbiased estimate of g_t

$$\mathbb{E}[\tilde{g}_t] = g_t$$

Stochastic Gradient Descent (SGD)

"Memorization" Attack.



Extracting Training Data from Large Language Models

Nicholas Carlini¹ Florian Tramèr² Eric Wallace³ Matthew Jagielski⁴
Ariel Herbert-Voss^{5,6} Katherine Lee¹ Adam Roberts¹ Tom Brown⁵
Dawn Song³ Úlfar Erlingsson⁷ Alina Oprea⁴ Colin Raffel¹
¹Google ²Stanford ³UC Berkeley ⁴Northeastern University ⁵OpenAI ⁶Harvard ⁷Apple

Private SGD. (DP-SGD)

$$\text{Private SGD} \left(\overset{\text{Loss}}{\downarrow} L(\cdot) = \frac{1}{n} \sum_{i=1}^n \ell(\cdot; x_i), \overset{\text{feasible set}}{\downarrow} C, \overset{\text{learning rate}}{\downarrow} \eta, \overset{\text{Noise rate}}{\swarrow} \beta \right) =$$

$$\text{Init} = w_0 \in C$$

For $t=1, \dots, T$:

Random subsample $B_t \subseteq \{1, \dots, n\}$
"mini-batch"

$$g_t = \frac{1}{|B_t|} \sum_{i \in B_t} \nabla_w \ell(w_{t-1}; x_i)$$

$$\tilde{g}_t = g_t + N(0, \beta^2 I_d)$$

\approx Gaussian Mechanism.

$$u_t = w_{t-1} - \eta \cdot \tilde{g}_t$$

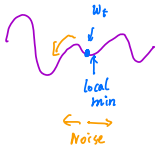
$$w_t = \underset{m \in C}{\text{argmin}} \|w - u_t\|_2$$

$$\text{Output} = \frac{1}{T} \sum_{t=1}^T w_t$$

or

$$w_T$$

Stochastic
Gradient
Langevin
Dynamics (SGLD)



Privacy Proof

Proof idea: • Think of releasing w_1, w_2, \dots, w_T .

- Suffices to release the update between iterates

$$w_0 \left[\xleftarrow{\text{update}} \rightarrow \right] w_1 \left[\xleftarrow{\text{update}} \rightarrow \right] w_2 \dots \dots w_T$$

- Suffices to release the sequence of gradient estimates

$$\tilde{g}_1, \tilde{g}_2, \dots, \tilde{g}_T$$

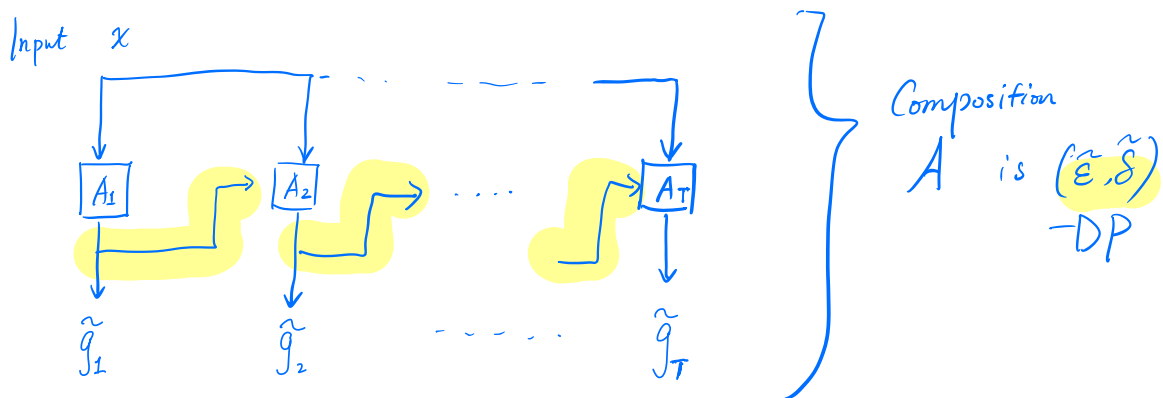
The output is a post-processing

Show releasing $(\tilde{g}_1, \tilde{g}_2, \dots, \tilde{g}_T)$ satisfies DP.

① Each step (releasing \tilde{g}_t) satisfies (ϵ, δ) -DP

② Adaptive composition across T steps

Adaptive Composition.



Suppose each step A_1, \dots, A_T is (ϵ, δ) -DP.

What are the values of $\tilde{\epsilon}$ and $\tilde{\delta}$?

- (Basic) Composition: $\tilde{\epsilon} = T\epsilon$, $\tilde{\delta} = T\delta$.
 - Advanced Composition: $\tilde{\epsilon} = \epsilon \cdot \sqrt{2T \ln(\frac{1}{\delta'})} + T \cdot \epsilon \frac{e^\epsilon - 1}{e^\epsilon + 1}$
 $\tilde{\delta} = T\delta + \underbrace{\delta'}_{>0}$ | for $\epsilon \leq 1$, $e^\epsilon \approx 1 + \epsilon \Rightarrow \frac{e^\epsilon - 1}{e^\epsilon + 1} \approx \frac{\epsilon}{2}$
 $\Rightarrow T \cdot \epsilon \frac{e^\epsilon - 1}{e^\epsilon + 1} \approx \frac{T}{2} \cdot \epsilon^2$
- If $\epsilon < \frac{1}{\sqrt{T}}$
- $(\epsilon \sqrt{T})^2$ is "smaller" than $\epsilon \sqrt{T}$
- Then $\tilde{\epsilon}$ is in the order of $\epsilon \cdot \sqrt{2T \ln(\frac{1}{\delta'})}$

Numeric Example.

$$\varepsilon = \frac{1}{1000}, \quad \delta$$

$$T = 500,$$

$$\text{Basic Composition: } \tilde{\varepsilon} = 0.5, \quad \tilde{\delta} = T\delta$$

$$\text{Advanced Composition: } \tilde{\varepsilon} \leq 0.1, \quad \tilde{\delta} = 10^{-6} + T\delta$$