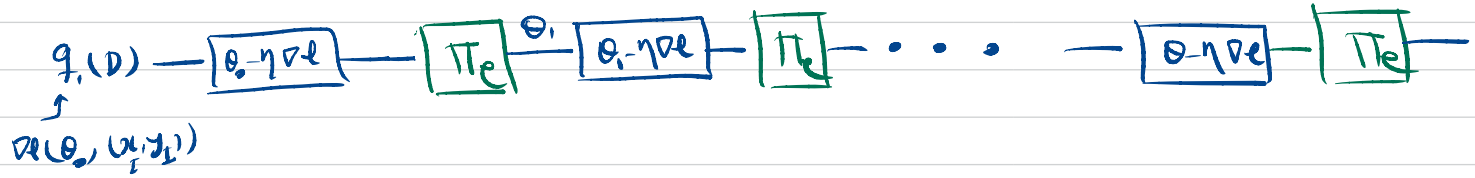
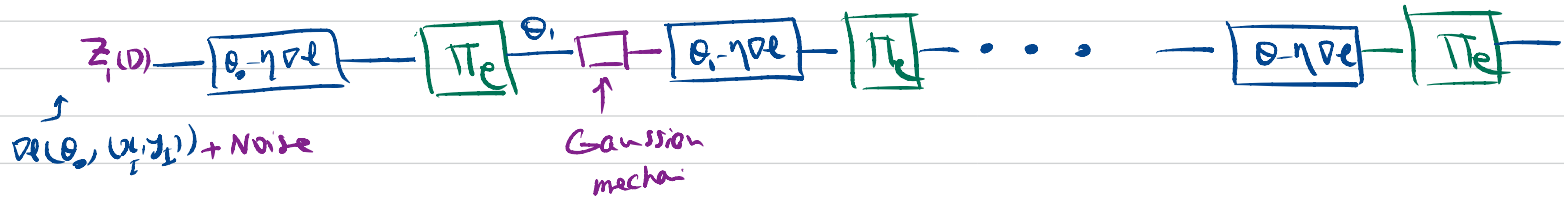


Private Projected SGD (PN-SGD)



To make projected SGD, differentially private, we need to make each iteration private. To do so, we need to pass the query response through a dp mechanism, say Gaussian mechanism.



Thus, each iteration proceeds as follows:

Private Projected GD

Input: Dataset $D = \{(x_i, y_i)\}_{i=1}^n$, loss function ℓ , parameter set \mathcal{C} & η_t

1- Pick $\theta_0 \in \mathcal{C}$ arbitrarily

2- For $t=1$ to T

- Select $I \in \{1, 2, \dots, n\}$ uniformly at random

- set $\theta_t = \Pi_{\mathcal{C}} \left(\theta_{t-1} - \eta_t \left[n \nabla \ell(\theta_{t-1}, (x_I, y_I)) + (N'_{t-1}, \dots, N^d_{t-1}) \right] \right)$

Output $\theta_1, \theta_2, \dots, \theta_T$

$N_t^{i \text{ iid}} \sim N(0, \sigma^2)$

At each iteration, query is $q_i(\theta, x_i, y_i)$ which is an adaptive query as it depends on the previous iteration's output.

Each iteration becomes (ϵ_i, δ_i) -DP with $\epsilon_i = \frac{2nL}{\sigma} \sqrt{2 \log \frac{1}{\delta_i}}$

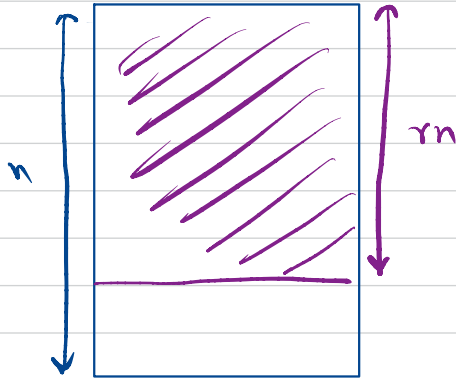
Note that: $\Delta_2^q \leq 2nL$ for the query $q_i(\theta, x_i, y_i)$.

Therefore, advanced composition can be used to obtain the privacy guarantee. But, we can do better! At each iteration, we don't use the whole dataset! We only use One data record & not all the records. Thus, the privacy must be significantly better!

How to quantify this improvement:

* Privacy amplification by Sub-Sampling:

Let M be an (ϵ, δ) -DP mechanism. If it is run on a uniformly selected subset of dataset of size rn ($r \leq 1$), then it will provide $(\log(1 + r(e^\epsilon - 1)), r\delta)$ -DP.



$q(D) \rightarrow \boxed{M} \rightarrow Z_D$
depends only the first part dataset

Z_D provides more privacy compared to the case when $q(D)$ depends on the whole dataset.

In the SGD: $\gamma = n$ (as the size of dataset = 1)

Gaussian mechanism coupled with this sub-sampling is typically called: Sub-sampled Gaussian mechanism.

Remark. we simplify the privacy amplification by sub-sampling as follows:

$$\underbrace{(\log(1 + \gamma \epsilon^2))}_{\leq 2\gamma\epsilon^2}, \gamma\epsilon) - \text{DP} \Rightarrow (2\gamma\epsilon^2, \gamma\epsilon) - \text{DP}$$

There are tighter analysis for "privacy amplification" caused by sub-sampling.

To see the original proof (in terms of HS divergence), see "privacy amplification by sub-sampling: Tight analysis via coupling" by Balle, Barthe, and Gaboardi, ICM2018

Privacy analysis of PN-SGD:

Each iteration is an instance of sub-sampled Gaussian mechanism. To derive its privacy guarantee, we need to find the privacy guarantee of the corresponding Gaussian mechanism & then apply privacy amplification by sub-sampling.

- Gaussian mechanism is $\left(\overset{\text{\textcolor{teal}{\mathcal{L}_2\text{-sensitivity}}}}{\frac{2nL}{\sigma}} \sqrt{2 \log 1/\delta}, \delta \right)$ -DP because the query with noise $N(0, \sigma^2 I)$

$n \nabla \ell(\theta, (x_i, y_i))$ has \mathcal{L}_2 -sensitivity $= 2nL$ with L as Lipschitz constant.

- Applying "privacy amplification by sub-sampling" theorem, each iteration
for $r = \frac{1}{n}$

becomes $\left(\frac{2}{n} \frac{2nL}{\sigma} \sqrt{2 \log \frac{1}{\delta_0}}, \frac{1}{n} \delta_0 \right) - \text{DP}$.

•) Now we can advanced composition:

PN-SGD after $T=n^2$ number of iterations is

$$\left(\frac{4L}{\sigma} \sqrt{2 \log \frac{1}{\delta_0}} \sqrt{2T \log \frac{1}{\delta'}}, \frac{T}{n} \delta_0 + \delta' \right) - \text{DP} \quad \forall \delta' \in (0,1)$$

choosing $\delta' = n \delta_0$, we obtain

$$\frac{T}{n} \delta_0 + n \delta_0 \stackrel{T=n^2}{=} 2n \delta_0 =: \delta$$

$$\text{so } \delta' = \delta/2$$

$$\& \delta_0 = \delta/2n$$

thus, the privacy guarantee becomes

$$\left(\frac{4L}{\sigma} n \sqrt{2 \log \frac{2n}{\delta}} \sqrt{2 \log \frac{2}{\delta}}, \delta \right) - \text{DP}$$

So to have it be (ϵ, δ) -DP, we require

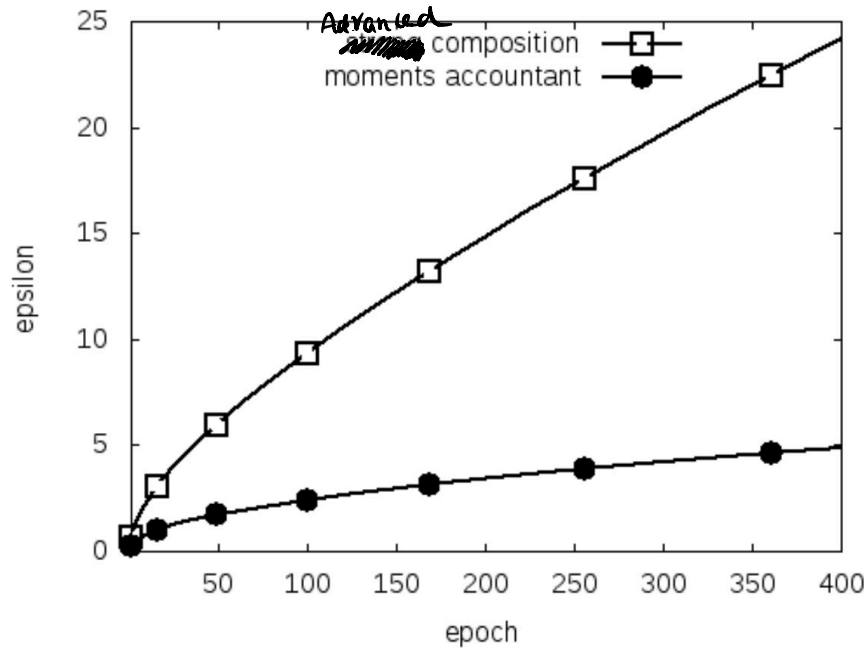
$$\sigma^2 = \frac{64 n^2 L^2}{\epsilon^2} \log \frac{2n}{\delta} \log \frac{2}{\delta}$$

Remarks:

1. A much better approach to study Privacy guarantee of PN-SGD is to invoke Rényi DP, use linear composition of Rényi DP, & then convert back to (ϵ, δ) -DP. A technique known as

moments accountant. The difference is astonishing:

the resulting σ is $\sqrt{\frac{\log n}{\delta}}$ smaller!



[Taken from "Deep learning with differential privacy"]

2- Our privacy analysis relies on the assumption that the loss function is Lipschitz. This is a restrictive assumption for deep learning. For example: MSE, squared hinge loss, or

To make sure sensitivity of \mathcal{L} is bounded even for non-Lipschitz loss, we usually clip the gradient:

$$\text{clip}(\nabla \mathcal{L}(\theta, (x, y)), C) := \begin{cases} C \frac{\nabla \mathcal{L}}{\|\nabla \mathcal{L}\|} & \text{if } \|\nabla \mathcal{L}\|_2 \geq C \\ \nabla \mathcal{L} & \text{if } \|\nabla \mathcal{L}\|_2 < C \end{cases}$$

thus, $\|\text{clip}(\nabla \mathcal{L}, C)\| \leq C \Rightarrow \text{sensitivity of } \text{clip}(\nabla \mathcal{L}, C) = 2C$

Accuracy of PN-SGD:

Recall the folklore result: Given Oracle $G_t(\theta)$ with

$E[G_t(\theta)] = \nabla L(\theta, D)$ & $E[\|G_t(\theta)\|^2] \leq G^2$, we have for

$$\theta_{t+1} = \prod_e (\theta_t - \eta G_t(\theta_t)); \quad E[L(\theta_T, D) - L(\theta^*, D)] \leq \|e\|_2 G \frac{\log T}{\sqrt{T}}.$$

Here, $G_t(\theta) := n \nabla \ell(\theta, (x_I, y_I)) + \underbrace{(N'_t, \dots, N_t^a)}_{N_t}$

Note that $n E[G_t(\theta_t)] = \nabla L(\theta_t, D)$

$$\begin{aligned} \& E[\|G_t(\theta_t)\|^2] &= E[\|n \nabla \ell(\theta_t, (x_I, y_I)) + N_t\|^2] \\ &= n^2 E[\|\nabla \ell(\theta_t, (x_I, y_I))\|^2] + \underbrace{E[n \langle \nabla \ell(\theta_t, (x_I, y_I)), N_t \rangle]}_{=0} \\ &\quad + E[\|N_t\|^2] \end{aligned}$$

$$\leq n^2 L^2 + d \sigma^2$$

so $G := \sqrt{n^2 L^2 + d \sigma^2}$

thus, PN-SGD satisfies:

$$\mathbb{E}[\mathcal{L}(\theta_T, D)] - \mathcal{L}(\omega^*, D) \leq \|\ell\|_2 \sqrt{n^2 L^2 + d \sigma^2} \frac{\log T}{\sqrt{T}}$$

replace
 σ
 by $T = n^2$

$$\lesssim O\left(\frac{\|\ell\|_2 L \sqrt{d \log 48}}{\epsilon}\right)$$

↑
 not decreasing in T

Local Dp

All Privacy problems discussed so far were based on an assumption:

A private dataset is held by a trusted central entity (such as hospital or bank),

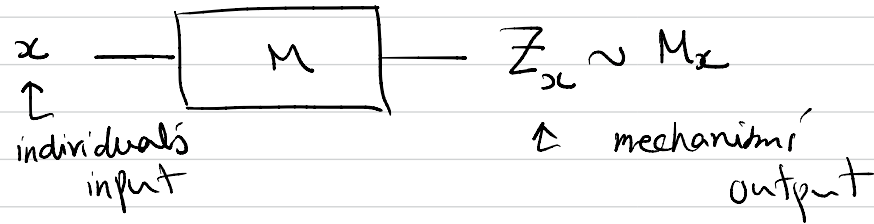
what if there can't be such central entity that's trusted by all individuals?

Example: we wish to find how many women in North Carolina have undergone abortion? We take a sample of women there & ask them. This is what statisticians have been doing

For decades. The caveat here in North Carolina, it is a crime to do abortion. So nobody answers the question truthfully!

To address problem like this, we need to let each individual randomize their data before they release it publicly.

Consider the following the mechanism:



Def. we say mechanism M is ϵ -locally DP (or ϵ -LDP for short)

if

$$\mathbb{E}_z(M_x || M_{x'}) = 0$$

$\forall x \neq x'$ possible inputs

Example. Randomized-response mechanism:



$Z \sim \text{Ber}(r)$ if $x=0$

&

$Z \sim \text{Ber}(1-r)$ if $x=1$

In fact, r indicates the lying probability.

* show that this mechanism is ϵ -LDP if $r = \frac{1}{1+e^\epsilon}$.

Theorem. A mechanism M is Σ -LDP iff & only iff

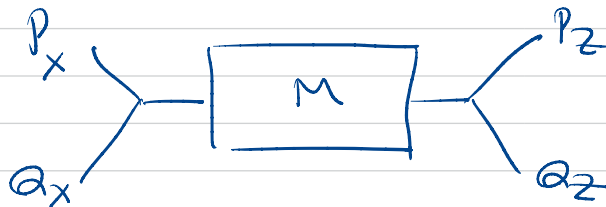
$$\eta_{e\Sigma}(M) = 0$$

Proof: HW4

In other word, an Σ -LDP mechanism always satisfied:

$$E_{e\Sigma}(P_Z \parallel Q_Z) = 0 \quad \forall P_X \& Q_X$$

where $P_Z \& Q_Z$ are the output distributions of M when input distributions are $P_X \& Q_X$.



Example: Prove that the following mechanism is ϵ -LDP:

$$P_{Z|X}(z|x) = \begin{cases} \frac{e^z}{e^z + k - 1} & z = x \in \{1, 2, \dots, k\} \\ \frac{1}{e^z + k - 1} & z \neq x \\ 0 & z \notin \{1, 2, \dots, k\} \end{cases}$$

* In HW1, we already proved that $\eta_{\epsilon}(P_{Z|X}) = 0$, showing that this mechanism is ϵ -LDP.

* This mechanism is known as k -ary randomized response, because both input & output take values in $\{1, 2, \dots, k\}$.