# Privacy in Statistics and Machine Learning       Spring 2021
## In-class Exercises for Lecture 12 (Gradient Descent)
## March 9, 2021

**Adam Smith and Jonathan Ullman**

*Problems with marked with an asterisk (*) are more challenging or open-ended.*

1. The performance of gradient descent can vary a lot depending on how we choose the step size, even for convex, one-dimensional problems. Suppose we run gradient descent with the loss function $L(w) = w^2$ and no constraints (that is, $C = \mathbb{R}$).

   - If we start at $w_0 = 1$ and use step size $\eta = 2$, how will the algorithm behave? Will it converge?
   - How would your answer to the previous part change if we used $\eta = 1$? What about $\eta = \frac{1}{2}$?
   - How would you answer to the first part change if we imposed the constraint $w \, in \, C = [-10, 10]$? What about $w \in [10, 20]$?

2. Prove the following variant of the Amplification by Subsampling Lemma.

   Suppose that given an algorithm $A$ whose input can be a data set of any size, we build a new algorithm $A'_p$ as follows: on input $\mathbf{x}$, construct a smaller data set $\mathbf{x}'$ by including each data record from $\mathbf{x}$ with probability $p$, independently of other data records. Finally, return $A(\mathbf{x}')$.

   If $A$ is $(\varepsilon, 0)$-DP under insertion/removal, then show that $A'_p$ is $(\varepsilon', 0)$-DP under insertion/removal, where $\varepsilon' = \ln\left(1 + p(e^\varepsilon - 1)\right)$.

3. Let's generalize the analysis of private SGD to the version where at each step, we use a uniformly random batch $B_t$ of $m$ records to estimate the gradient, so

   $$\tilde{g}_t = \left(\frac{1}{m} \sum_{i \in B_t} \nabla \ell(w_{t-1}; x_i)\right) + N(0, \sigma^2).$$

   Given $\delta$, we want to understand for which $\varepsilon$ this step is $(\varepsilon, \delta)$-DP. Show that, as long as

   $$\sigma \geq 2G\sqrt{\ln(1/\delta)} \cdot \frac{1}{m}$$

   the privacy cost of one step of gradient descent with subsampling is at most $e$ times higher than it would be if we had used the entire data set to estimate the gradient. In other words, subsampling has virtually no effect on privacy as long the noise level is sufficiently high.

4. We analyzed gradient descent for the setting where the diameter $R$ of $C$ is bounded. But suppose $C$ is not bounded—say $C = \mathbb{R}^d$. We could still hope to get a good bound if our initial point $w_0$ is not too far from a true optimum $w^*$. A friend conjectures that if one has a good idea of $\|w_0 - w^*\|$, one should be able to set $\eta$ to get a bound of the form

   $$L(\hat{w}) - L(w^*) \leq \frac{G \times \left(\text{some function of } \|w_0 - w^*\|\right)}{\sqrt{T}}.$$

Are they correct? What function belongs there?

5. It is common for the learning rate $\eta$ to decrease over the course of gradient descent. Suppose we set $\eta_t = \frac{1}{\sqrt{t}}$, and update the estimate as $u_t = w_{t-1} + \eta_t \nabla L(w_{t-1})$. This way of doing things has the benefit that we don't need to set the number of iterations $T$ ahead of time.

Show that, when $G = R = 1$, we can get the same asymptotic risk bound of $O(1/\sqrt{T})$ for gradient descent.