

# HW 4B: Stochastic Gradient Descent and Lipschitz Extensions

CS 208 Applied Privacy for Data Science, Spring 2019

**Version 1.2: Due Tuesday, April 30, 11:59pm.**

**Instructions:** Submit a single PDF file containing your solutions, plots, and analyses. Make sure to thoroughly explain your process and results for each problem. Also include your documented code and a link to a public repository with your code (such as GitHub/GitLab). Make sure to list all collaborators and references.

- For each of the following sets  $\mathcal{G}$  of datasets and neighbor relations  $\sim$ , hypotheses  $\mathcal{H} \subseteq \mathcal{G}$ , and functions  $f : \mathcal{G} \rightarrow \mathbb{R}$ , calculate (i) the global sensitivity of  $f$  (denoted  $\text{GS}_f$  or  $\partial f$ ), (ii) the minimum local sensitivity of  $f$ , i.e.  $\min_{x \in \mathcal{G}} \text{LS}_f(x)$ , and (iii) the restricted sensitivity of  $f$  (denoted  $\partial_{\mathcal{H}} f$  or  $\text{RS}_f^{\mathcal{H}}$ ). For Part 1a, also describe an explicit Lipschitz extension of  $f$  from  $\mathcal{H}$  to all of  $\mathcal{G}$ .
  - $\mathcal{G} = \mathbb{R}^n$  where  $x \sim x'$  if  $x$  and  $x'$  differ on one row,  $\mathcal{H} = [a, b]^n$  for real numbers  $a \leq b$ , and  $f(x) = (1/n) \sum_{i=1}^n x_i$ .
  - $\mathcal{G} = \mathbb{R}^n$  where  $x \sim x'$  if  $x$  and  $x'$  differ on one row,  $\mathcal{H} = [a, b]^n$  for real numbers  $a \leq b$ , and  $f(x) = \text{median}(x_1, \dots, x_n)$ .
  - $\mathcal{G}$  = the set of undirected graphs (without self-loops) on vertex set  $\{1, \dots, n\}$  where  $x \sim x'$  if  $x$  and  $x'$  are identical except for the neighborhood of a single vertex (i.e. node privacy),  $\mathcal{H}$  = the set of graphs in  $\mathcal{G}$  in which every vertex has degree at most  $d$  for a parameter  $2 \leq d \leq n - 1$ , and  $f(x)$  = the number of isolated (i.e. degree 0) vertices in  $x$ .
- In our code example,<sup>1</sup> we saw how to release an estimated Logistic regression using differentially private stochastic gradient descent (DP-SGD) to optimize the log-likelihood loss function under the centralized model. Convert this code to once again release the probability of marriage given education level, but using DP-SGD under the *local* model.<sup>2</sup> Recall that local DP does not satisfy privacy amplification by subsampling, but you can achieve a similar effect by rotating through disjoint batches, so that each individual participates in at most  $\lceil T \cdot B/n \rceil$  batches, where  $T$  is the number of iterations and  $B$  is the batch size.<sup>3</sup> Evaluate the performance of your method as a function of  $\epsilon$  (fixing  $\delta = 1 \times 10^{-6}$ ), by showing the classification error over  $\epsilon$ , compared to the RMSE of the coefficients compared to the non-privacy preserving estimates.

<sup>1</sup>See [https://github.com/privacytoolsproject/cs208/blob/master/examples/wk7\\_localmodel/privateSGD.r](https://github.com/privacytoolsproject/cs208/blob/master/examples/wk7_localmodel/privateSGD.r) and `privateSGD.ipynb`.

<sup>2</sup>For some useful guidance, look at the class notes for April 8th at <http://people.seas.harvard.edu/~salil/cs208/spring19/MLwithDP-lecture.pdf>.

<sup>3</sup>Note, in the code example,  $T = \sqrt{n}$ ,  $B = \sqrt{n}$ .