

CSE 597: LARGE-SCALE MACHINE LEARNING  
MATHEMATICAL FOUNDATIONS AND APPLICATIONS

FINAL EXAM

FALL 2021

Name:.....

PSU ID:..... @psu.edu

There are four problems with total score of 60 with 10 bonus points (exam will be graded for 50 points).  
Please write all the details in your derivations and be as rigorous as possible.

Q1: ----- / 10

Q2: ----- / 20

Q3: ----- / 20

Q4: ----- / 10

Total ----- / 50

**Result:**

**Question 1** (10 pts). Consider the binary classification problem, where we are given a training set  $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^n$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{+1, -1\}$ . In SVM for binary classification we would like to find the minimizer of the following problem

$$\min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{(\mathbf{x}, y) \in \mathcal{S}} \ell(\mathbf{w}; (\mathbf{x}, y))$$

where

$$\ell(\mathbf{w}; (\mathbf{x}, y)) = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$$

is the Hinge loss.

By examining the dual form of the SVM problem and using the strong duality theorem (i.e., if  $\mathbf{w}_*$  and  $\alpha_*$  are primal and dual solutions, respectively, then it holds that  $\mathbf{w}_* = \sum_{i=1}^n \alpha_{*,i} y_i \mathbf{x}_i$ ), show that  $\|\mathbf{w}_*\| \leq 1/\sqrt{\lambda}$ .

**Question 2** (20 pts). Consider a distributed optimization setting where  $p$  machines aim to jointly optimize the following ERM objective

$$f(\mathbf{w}) = \frac{1}{p} \sum_{i=1}^p f_i(\mathbf{w}),$$

under the orchestration of a central server where  $f_i(\mathbf{w})$  is the training (empirical) loss over data at  $i$ th device.

A simple distributed optimization algorithm (synchronous SGD) to optimize above objective is as follows. Starting at an initial solution  $\mathbf{w}_1$ , at each iteration  $t$ , the server sends the global model  $\mathbf{w}_t$  to all  $p$  devices, each device computes local unbiased stochastic gradient  $\mathbf{g}_{i,t}$  at  $\mathbf{w}_t$  ( $\mathbb{E}[\mathbf{g}_{i,t}] = \nabla f(\mathbf{w}_t)$ ) and sends it back to the server. Then server aggregates the stochastic gradients from all devices by averaging  $\mathbf{g}_t = (1/p) \sum_{i=1}^p \mathbf{g}_{i,t}$  and updates the global model by  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t$ . And this proceeds for  $T$  iterations. Since at each step,  $p$  stochastic gradients are computed in parallel, we can show that for  $O(T)$  communication rounds, we can achieve an  $O(\frac{1}{pT})$  convergence rate for smooth and strongly convex functions (follows from the analysis of mini-batch SGD in lecture with batch size  $B = p$ ).

In this question we consider a variant of above distributed optimization algorithm and would like to establish its convergence rate. In the modified algorithm, instead of using a single sample (i.e., fixed batch size) to compute local stochastic gradients at all iterations, the devices are allowed to utilize growing mini-batches with growing coefficient  $\rho > 1$  at different steps. Specifically, consider the following distributed optimization algorithm with growing (dynamic) mini-batches:

**inputs:**  $p, T, \mathbf{w}_1$ , initial mini-batch size  $B_1$ , mini-batch growing coefficient  $\rho > 1$   
**initialize:**  $t = 1$   
1: **while**  $\sum_{s=1}^t B_s \leq T$  **do**  
2:     Each device  $i$  samples  $B_t$  samples IID and computes unbiased stochastic gradient  $\mathbf{g}_{i,t}$  at  $\mathbf{w}_t$   
3:     Server aggregates gradients from all devices via averaging:  $\mathbf{g}_t = \frac{1}{p} \sum_{i=1}^p \mathbf{g}_{i,t}$   
4:     Server updates the model by  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t$  and broadcast to all devices  
5:     Set  $B_{t+1} = \lfloor \rho^t B_1 \rfloor$  (growing the batch size)  
6:     Update  $t \leftarrow t + 1$

For above algorithm,

- (a) Assume the global objective function  $f(\cdot)$  is  $\alpha$ -strongly convex and  $\beta$ -smooth. Show that if we choose  $\eta < \frac{1}{\beta}$  in above algorithm, then for all  $t \in \{1, 2, \dots\}$ , we have

$$\mathbb{E}[f(\mathbf{w}_{t+1}) - f^*] \leq (1 - \nu) \mathbb{E}[f(\mathbf{w}_t) - f^*] + \frac{\eta(2 - \beta\eta)}{2pB_t} \sigma^2$$

where  $f^*$  is the global minimum,  $\sigma$  is the variance of stochastic gradients, and  $\nu = \frac{1}{2}\eta\alpha(1 - \beta\eta)$  satisfies  $0 < \nu < 1$

- (b) Use above result and show that with  $t = O(\log T)$  communication rounds we can achieve an  $O(\frac{1}{pT})$  convergence rate (i.e., using dynamic mini-batches we can reduce the number of communications from  $O(T)$  in vanilla distributed SGD with  $B = 1$  to  $O(\log T)$  while achieving the same convergence rate).

**Question 3** (20 pts). This problem is mostly a reading exercise. Consider a dataset  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{-1, +1\}$ . Assume that the data is linearly separable by margin  $\gamma \in (0, 1]$  defined as  $\gamma = \max_i y_i \frac{\langle \mathbf{w}_*, \mathbf{x}_i \rangle}{\|\mathbf{w}_*\| \|\mathbf{x}_i\|}$  where  $\mathbf{w}_*$  is the optimal classifier.

Let  $\mathbf{M} \in \mathbb{R}^{m \times d}$  be a random Gaussian matrix with  $M_{ij} = \frac{1}{\sqrt{m}} \mathcal{N}(0, 1)$ . Show that for any  $\delta, \varepsilon \in (0, 1)$  if  $m > \mathcal{O}\left(\frac{1}{\varepsilon^2} \ln \frac{n}{\delta}\right)$ , the dataset  $S' = \{(\mathbf{M}\mathbf{x}_1, y_1), \dots, (\mathbf{M}\mathbf{x}_n, y_n)\}$  is linearly separable with margin  $\gamma - \frac{2\varepsilon}{1-\varepsilon}$ . (What you need to do is to read this paper “Is margin preserved after random projection?”, and simplify/rewrite the proof.)

**Question 4** (10 pts). Consider the function  $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$ . Let  $\nabla f_S(\mathbf{x}) = \frac{1}{|S|} \sum_{i \in S} \nabla f_i(\mathbf{x})$  be the stochastic gradient of function computed at a subset  $S \subseteq \{1, 2, \dots, n\}$  sampled uniformly at random. By using concentration inequalities, show that what would be the size of  $S$  (batch size) to guarantee that

$$\mathbb{P} [\|\nabla f_S(\mathbf{x}) - \nabla f(\mathbf{x})\|_2^2 \leq \varepsilon] \geq 1 - \delta$$

for any  $\varepsilon, \delta \in (0, 1)$  and  $\mathbf{x} \in \mathbb{R}^d$ .

Good luck.