## Exam Stochastic Optimization – part I 31 January 2023

duration 1.5 hours, paper documents allowed

**Exercise 1.** We consider the following optimization problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N f_i(x) + g(x)$$

where for all  $i \in \{1, ..., N\}$ ,  $f_i$  is a differentiable convex function whose gradient is *L*-Lipschitz and g is a convex, lower-semi-continuous function whose proximal operator is easy to compute. We shall denote  $f(x) = \frac{1}{N} \sum_{i=1}^{N} f_i(x)$  and F(x) = f(x) + g(x).

Notation:  $U(\{1, \ldots, N\})$  is the uniform distribution over the set  $\{1, \ldots, N\}$  and B(p) is the Bernoulli distribution with mean p. We denote the conditional expection with  $\mathbb{E}_k[X_k] = \mathbb{E}[X_k|i_0, b_0, \ldots, i_k, b_k].$ 

The goal of this exercise is to study the convergence of the following prox-SVRG algorithm

$$x_{0} \in \mathbb{R}^{d}, w_{0} = x_{0}, \gamma > 0, p \in ]0, 1[$$
  

$$\forall k \in \mathbb{N} :$$
  

$$i_{k+1} \sim U(\{1, \dots, N\})$$
  

$$G_{k+1} = \nabla f(w_{k}) + \nabla f_{i_{k+1}}(x_{k}) - \nabla f_{i_{k+1}}(w_{k})$$
  

$$x_{k+1} = \operatorname{prox}_{\gamma g} (x_{k} - \gamma G_{k+1})$$
  

$$b_{k+1} \sim B(p)$$
  

$$w_{k+1} = (1 - b_{k+1})w_{k} + b_{k+1}x_{k}$$

We shall not assume the strong convexity of F, only its convexity. We will proceed by proving the following points.

1. Show that

$$f(x_{k+1}) \le f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

2. Show that for any  $x_* \in \arg \min_x F(x)$ 

$$g(x_{k+1}) + \langle G_{k+1}, x_{k+1} - x_k \rangle + \frac{1}{2\gamma} \|x_{k+1} - x_k\|^2 \le g(x_*) + \langle G_{k+1}, x_* - x_k \rangle + \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_{k+1} - x_*\|^2 \le g(x_*) + \langle G_{k+1}, x_* - x_k \rangle + \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_{k+1} - x_*\|^2 \le g(x_*) + \langle G_{k+1}, x_* - x_k \rangle + \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_{k+1} - x_*\|^2 \le g(x_*) + \langle G_{k+1}, x_* - x_k \rangle + \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_{k+1} - x_*\|^2 \le g(x_*) + \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_k - x_*\|^2 \le g(x_*) + \frac{1}{2\gamma} \|x_k - x_*\|^2 - \frac{1}{2\gamma} \|x_k - x_*\|^2 \le g(x_*) + \frac{1}{2\gamma$$

3. Show that

$$\mathbb{E}_k[\langle G_{k+1}, x_* - x_k \rangle] = \langle \nabla f(x_k), x_* - x_k \rangle$$

4. Show that for all  $\alpha > 0$ 

$$\langle G_{k+1} - \nabla f(x_k), x_k - x_{k+1} \rangle \le \frac{\alpha}{2} \|G_{k+1} - \nabla f(x_k)\|^2 + \frac{1}{2\alpha} \|x_k - x_{k+1}\|^2$$

5. Show that

$$||G_{k+1} - \nabla f(x_k)||^2 \le 2||\nabla f_{i_{k+1}}(x_k) - \nabla f_{i_{k+1}}(x_*)||^2 + 2||\nabla f_{i_{k+1}}(w_k) - \nabla f_{i_{k+1}}(x_*)||^2$$

6. Denote  $\mathcal{D}_k = \frac{1}{N} \sum_{i=1}^N \|\nabla f_i(w_k) - \nabla f_i(x_*)\|^2$ . Show that

$$\mathbb{E}[\mathcal{D}_{k+1}|i_0, b_0, \dots, i_k] = p\mathbb{E}_k \Big[ \|\nabla f_{i_{k+1}}(x_k) - \nabla f_{i_{k+1}}(x_*)\|^2 \Big] + (1-p)\mathcal{D}_k$$

7. Show that

$$\mathbb{E}_{k}\Big[\|\nabla f_{i_{k+1}}(x_{k}) - \nabla f_{i_{k+1}}(x_{*})\|^{2}\Big] \le 2L\big(f(x_{*}) - f(x_{k}) - \langle \nabla f(x_{k}), x_{*} - x_{k}\rangle\big)^{2}$$

8. Chose values for the proof constants  $\alpha$  and  $\beta$  that ensure

$$\mathbb{E}_{k}\left[F(x_{k+1}) + \beta \mathcal{D}_{k+1} + \frac{1}{2\gamma} \|x_{k+1} - x_{k}\|^{2}\right] \le F(x_{*}) + \beta \mathcal{D}_{k} + \frac{1}{2\gamma} \|x_{k} - x_{*}\|^{2} + \left(\frac{L}{2} - \frac{1}{2\gamma} + \frac{1}{2\alpha}\right) \|x_{k+1} - x_{k}\|^{2}$$

- 9. What range of step sizes ensures that  $\left(\frac{L}{2} \frac{1}{2\gamma} + \frac{1}{2\alpha}\right) \|x_{k+1} x_k\|^2 \le 0$ ?
- 10. Suppose  $\gamma$  satisfies the condition of the previous question and denote  $\bar{x}_K = \frac{1}{K} \sum_{k=1}^{K}$ . Show that

$$\mathbb{E}[F(\bar{x}_K) - F(x_*)] \le \frac{\beta \mathcal{D}_0 + \frac{1}{2\gamma} \|x_0 - x_*\|^2}{K}$$

11. Compare with the convergence rate of other algorithms we have seen during the course.