# Exercise: Minibatch SGD

Let $n, d \geq 1$ be integers, and consider the optimization problem

$$\min_{\theta \in \mathbb{R}^d} F(\theta), \quad \text{where} \quad F(\theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\theta), \quad \forall \theta \in \mathbb{R}^d.$$

In this exercise, we will study a generalization of GD and SGD, called *mini-batch SGD*. The idea is simple:

- Gradient Descent proceeds by computing a full gradient $\frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\theta_k)$ at each step.

- In contrast, SGD proceeds by only computing *one* gradient $\nabla f_{i_k}(\theta_k)$ at each step, for some index $i_k$ picked uniformly at random in $\{1, \ldots, n\}$.

- Mini-batch SGD is a compromise between the two. For some integer $m \in \{1, \ldots, n\}$, it selects uniformly at random (u.a.r.) a subset $v_{k+1} \subseteq \{1, \ldots, n\}$ of cardinality $m$ at each iteration, and computes the partial gradient $\frac{1}{m} \sum_{i \in v_{k+1}} \nabla f_i(\theta_k)$.

The pseudo-code of this algorithm is given below

---

**Algorithm 1:** Mini-batch SGD

---

Start from $\theta_0 \in \mathbb{R}^d$.

Until termination condition, iterate

    Draw a subset $v_{k+1} \subseteq \{1, \ldots, n\}$ of cardinality $m$ u.a.r., independent of the past

    $g_{k+1} = \dfrac{1}{m} \displaystyle\sum_{i \in v_{k+1}} \nabla f_i(\theta_k)$

    $\theta_{k+1} = \theta_k - \gamma_{k+1} g_{k+1}.$

---

**Notation.** For any subset $v \subseteq \{1, \ldots, n\}$ of cardinality $m$, we write for ease

$$f_v = \frac{1}{m} \sum_{i \in v} f_i.$$

In particular, if $m = n$ and $v = \{1, \ldots, n\}$, then $f_v = F$. We also define the noise of a subgradient as

$$\sigma^2 = \mathbb{E}_{v \sim \mathcal{D}}[\|\nabla f_v(\theta^*)\|^2]$$

Here, $\mathcal{D}$ denotes the uniform distribution over subsets of $\{1, \ldots, n\}$ whose cardinality is $m$. The goal of the exercise is to prove a convergence rate for mini-batch SGD, given in the theorem below.

---

**Theorem 1.** *Assume that*

1. *$F$ is $\mu$-strongly convex*

2. *$\mathbb{E}_{\mathcal{D}} \|\nabla f_v(\theta) - \nabla f_v(\theta^*)\|^2 \leq 2L(F(\theta) - F(\theta^*))$ for some constant $L > 0$ (we say that $F$ is "L-smooth in expectation with respect to $\mathcal{D}$")*

3. *The noise $\sigma^2$ is finite.*

---

*Take a precision $\varepsilon > 0$, and choose a step size $\gamma = \min\left\{\frac{1}{2L}, \frac{\varepsilon\mu}{4\sigma^2}\right\}$. Then we have*

$$\mathbb{E}\|\theta_k - \theta^*\|^2 \le \varepsilon \quad \text{as soon as} \quad k \ge \max\left\{\frac{2L}{\mu}, \frac{4\sigma^2}{\varepsilon\mu^2}\right\}\log\left(\frac{2\|\theta_0 - \theta^*\|^2}{\varepsilon}\right).$$

## Questions

1. Using Assumption 2, justify that $\mathbb{E}_{\mathcal{D}}\left[\|\nabla f_v(\theta)\|^2\right] \le 4L(F(\theta) - F(\theta^*)) + 2\sigma^2$.

   (One can use the inequality $\|a + b\|^2 \le 2(\|a\|^2 + \|b\|^2)$.)

   We have
   $$\mathbb{E}_{\mathcal{D}}\left\|\nabla f_v(\theta)\right\|^2 = \mathbb{E}_{\mathcal{D}}\left\|\nabla f_v(\theta) - \nabla f_v(\theta^*) + \nabla f_v(\theta^*)\right\|^2$$
   $$\le 2\mathbb{E}_{\mathcal{D}}\left\|\nabla f_v(\theta) - \nabla f_v(\theta^*)\right\|^2 + 2\mathbb{E}\left\|\nabla f_v(\theta^*)\right\|^2$$
   $$\le 4L\left[F(\theta) - F(\theta^*)\right] + 2\mathbb{E}_{\mathcal{D}}\left\|\nabla f_v(\theta^*)\right\|^2$$
   $$= 4L(F(\theta) - F(\theta^*)) + 2\sigma^2.$$

   The first inequality follows from $\|a + b\|^2 \le 2\|a\|^2 + 2\|b\|^2$, and the second inequality follows from Assumption 2.

2. Justify that

   $$\|\theta_{k+1} - \theta^*\|^2 = \|\theta_k - \theta^*\|^2 - 2\gamma\left\langle\theta_k - \theta^*, \nabla f_{v_{k+1}}(\theta_k)\right\rangle + \gamma^2\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2.$$

   By definition, we have

   $$\|\theta_{k+1} - \theta^*\|^2 = \left\|\theta_k - \theta^* - \gamma\nabla f_{v_{k+1}}(\theta_k)\right\|^2$$
   $$= \|\theta_k - \theta^*\|^2 - 2\gamma\left\langle\theta_k - \theta^*, \nabla f_{v_{k+1}}(\theta_k)\right\rangle + \gamma^2\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2.$$

3. Writing $\mathbb{E}_k$ for the expectation conditional on $\theta_k$, prove that:

   $$\mathbb{E}_k\|\theta_{k+1} - \theta^*\|^2 \le \left(1 - \gamma\mu\right)\|\theta_k - \theta^*\|^2 - 2\gamma\left[F(\theta_k) - F(\theta^*)\right] + \gamma^2\mathbb{E}_k\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2.$$

   (One can use the strong convexity of $F$: $\forall x, y \in \mathbb{R}^d : F(x) - F(y) \ge \left\langle\nabla F(y), x - y\right\rangle + \frac{\mu}{2}\|x - y\|^2$.)

   We use the strong convexity property of $F$ with $x = \theta^*$ and $y = \theta_k$

   $$\mathbb{E}_k\|\theta_{k+1} - \theta^*\|^2 = \|\theta_k - \theta^*\|^2 - 2\gamma\langle\theta_k - \theta^*, \nabla F(\theta_k)\rangle + \gamma^2\mathbb{E}_k\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2$$
   $$= \|\theta_k - \theta^*\|^2 + 2\gamma\langle\nabla F(\theta_k), \theta^* - \theta_k\rangle + \gamma^2\mathbb{E}_k\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2$$
   $$\le \|\theta_k - \theta^*\|^2 + 2\gamma\left[F(\theta^*) - F(\theta_k) - \frac{\mu}{2}\|\theta^* - \theta_k\|^2\right] + \gamma^2\mathbb{E}_k\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2$$
   $$= \left(1 - \gamma\mu\right)\|\theta_k - \theta^*\|^2 - 2\gamma\left[F(\theta_k) - F(\theta^*)\right] + \gamma^2\mathbb{E}_{\mathcal{D}}\left\|\nabla f_{v_{k+1}}(\theta_k)\right\|^2.$$

4. Deduce that, if $\gamma \le \frac{1}{2L}$, we can take the total expectation and write

   $$\mathbb{E}\|\theta_{k+1} - \theta^*\|^2 \le \left(1 - \gamma\mu\right)\mathbb{E}\|\theta_k - \theta^*\|^2 + 2\gamma^2\sigma^2.$$

Taking expectations again and using Question 1, we get

$$\mathbb{E}\big\|\theta_{k+1} - \theta^*\big\|^2 \leq (1 - \gamma\mu)\mathbb{E}\big\|\theta_k - \theta^*\big\|^2 + 2\gamma^2\sigma^2 + 2\gamma(2\gamma L - 1)\mathbb{E}\big[F(\theta_k) - F(\theta^*)\big]$$

$$\leq (1 - \gamma\mu)\mathbb{E}\big\|\theta_k - \theta^*\big\|^2 + 2\gamma^2\sigma^2.$$

In the last inequality, we used the fact that $2\gamma L \leq 1$ since $\gamma \leq \frac{1}{2L}$.

5. Deduce that $\mathbb{E}\big\|\theta_k - \theta^*\big\|^2 \leq (1 - \gamma\mu)^k\big\|\theta_0 - \theta^*\big\|^2 + \frac{2\gamma\sigma^2}{\mu}$.

Recursively applying the above inequality and summing up the resulting geometric series yields

$$\mathbb{E}\big\|\theta_k - \theta^*\big\|^2 \leq (1 - \gamma\mu)^k\big\|\theta_0 - \theta^*\big\|^2 + 2\sum_{j=0}^{k-1}(1 - \gamma\mu)^j\gamma^2\sigma^2$$

$$\leq (1 - \gamma\mu)^k\big\|\theta_0 - \theta^*\big\|^2 + \frac{2\gamma\sigma^2}{\mu}.$$

6. How should one choose $k$ to obtain $\mathbb{E}\big\|\theta_k - \theta^*\big\|^2 \leq \varepsilon$?

By the previous question, we note that

$$\mathbb{E}\big\|\theta_k - \theta^*\big\|^2 \leq (1 - \gamma\mu)^k\big\|\theta_0 - \theta^*\big\|^2 + \frac{2\gamma\sigma^2}{\mu} \leq (1 - \gamma\mu)^k\big\|\theta_0 - \theta^*\big\|^2 + \frac{\varepsilon}{2},$$

since we have chosen $\gamma = \min\left\{\frac{1}{2L}, \frac{\varepsilon\mu}{4\sigma^2}\right\} \leq \frac{\varepsilon\mu}{4\sigma^2}$, so that $\frac{2\gamma\sigma^2}{\mu} \leq \frac{\varepsilon}{2}$. Now, it suffices to take $k$ such that

$$(1 - \gamma\mu)^k\big\|\theta_0 - \theta^*\big\|^2 \leq \frac{\varepsilon}{2}.$$

Using the inequality $(1 - \gamma\mu) \leq e^{-\gamma\mu}$, it is clear that it suffices to take $k$ such that

$$\exp\big(-k\gamma\mu\big)\big\|\theta_0 - \theta^*\big\|^2 \leq \frac{\varepsilon}{2} \quad \text{i.e.} \quad k \geq \frac{1}{\gamma\mu}\log\left(\frac{2\|\theta_0 - \theta^*\|^2}{\varepsilon}\right) = \max\left\{\frac{2L}{\mu}, \frac{4\sigma^2}{\varepsilon\mu^2}\right\}\log\left(\frac{2\|\theta_0 - \theta^*\|^2}{\varepsilon}\right).$$

7. For any $m \in \{1, \ldots, n\}$, we denote by $\mathcal{D}_m$ the uniform distribution over subsets of $\{1, \ldots, n\}$ whose cardinality is $m$, and we define the quantity $\sigma_m^2 = \mathbb{E}_{v\sim\mathcal{D}_m}\big\|\nabla f_v(\theta^*)\big\|^2$. Justify that

$$\sigma_m^2 = \frac{1}{m^2}\sum_{i,j=1}^n \Big\langle\nabla f_i(\theta^*), \nabla f_j(\theta^*)\Big\rangle\mathbb{P}_{v\sim\mathcal{D}_m}\big(i \in v \text{ and } j \in v\big)$$

We have

$$\sigma_m^2 = \mathbb{E}_{v\sim\mathcal{D}_m}\big\|\nabla f_v(\theta^*)\big\|^2 = \mathbb{E}_{v\sim\mathcal{D}_m}\bigg\|\frac{1}{m}\sum_{i\in v}\nabla f_i(\theta^*)\bigg\|^2$$

$$= \mathbb{E}_{v\sim\mathcal{D}_m}\bigg[\frac{1}{m^2}\sum_{i,j\in v}\big\langle\nabla f_i(\theta^*), \nabla f_j(\theta^*)\big\rangle\bigg]$$

$$= \mathbb{E}_{v\sim\mathcal{D}_m}\bigg[\frac{1}{m^2}\sum_{i,j=1}^n\big\langle\nabla f_i(\theta^*), \nabla f_j(\theta^*)\big\rangle\mathbf{1}\{i, j \in v\}\bigg]$$

$$= \frac{1}{m^2}\sum_{i,j=1}^n\Big\langle\nabla f_i(\theta^*), \nabla f_j(\theta^*)\Big\rangle\mathbb{P}_{v\sim\mathcal{D}_m}\big(i \in v \text{ and } j \in v\big).$$

3

8. Justify that $\sigma_n^2 = 0$ and that, for any integers $i, j \in \{1, \ldots, n\}$, we have

$$\mathbb{P}_{v \sim \mathcal{D}_m}\left(i \in v \text{ and } j \in v\right) = \begin{cases} \dfrac{m}{n} & \text{if } i = j \\ \dfrac{m(m-1)}{n(n-1)} & \text{if } i \neq j. \end{cases}$$

We have $\sigma_n^* = \|\nabla F(\theta^*)\|^2 = 0$. Moreover, if $i = j$, we have

$$\mathbb{P}_{v \sim \mathcal{D}_m}\left(i \in v \text{ and } j \in v\right) = \mathbb{P}_{v \sim \mathcal{D}_m}\left(i \in v\right) = \frac{m}{n}.$$

Next, if $i \neq j$, we have

$$\mathbb{P}_{v \sim \mathcal{D}_m}\left(i \in v \text{ and } j \in v\right) = \mathbb{P}_{v \sim \mathcal{D}_m}\left(i \in v \mid j \in v\right) \mathbb{P}_{v \sim \mathcal{D}_m}\left(j \in v\right) = \frac{m-1}{n-1}\frac{m}{n}.$$

9. Deduce that $\sigma_m^2 = \dfrac{\sigma_1^2}{n-1}\left(\dfrac{n}{m} - 1\right)$.

Combining Questions 7 and 8, we obtain

$$\sigma_m^2 = \frac{1}{m^2} \sum_{i,j=1}^{n} \left\langle \nabla f_i(\theta^*), \nabla f_j(\theta^*) \right\rangle \mathbb{P}_{v \sim \mathcal{D}_m}\left(i \in v \text{ and } j \in v\right)$$

$$= \frac{1}{m^2} \sum_{1 \le i \neq j \le n} \left\langle \nabla f_i(\theta^*), \nabla f_j(\theta^*) \right\rangle \frac{m(m-1)}{n(n-1)} + \frac{1}{m^2} \sum_{i=1}^{n} \|\nabla f_i(\theta^*)\|^2 \frac{m}{n}$$

$$= \frac{1}{m^2} \frac{m(m-1)}{n(n-1)} \underbrace{\left[ \sum_{1 \le i,j \le n} \left\langle \nabla f_i(\theta^*), \nabla f_j(\theta^*) \right\rangle + \sum_{i=1}^{n} \|\nabla f_i(\theta^*)\|^2 - \sum_{i=1}^{n} \|\nabla f_i(\theta^*)\|^2 \right]}_{=\|\nabla F(\theta^*)\|^2 = 0}$$

$$+ \frac{1}{m^2} \sum_{i=1}^{n} \|\nabla f_i(\theta^*)\|^2 \frac{m}{n}$$

$$= \frac{n-m}{mn(n-1)} \sum_{i=1}^{n} \|\nabla f_i(\theta^*)\|^2 = \frac{n-m}{m(n-1)} \sigma_1^2.$$

10. What is the runtime of the mini-batch SGD algorithm with mini-batches of size $m$ to reach precision $\varepsilon$, assuming that computing one gradient takes 1 second?

Each iteration requires computing $m$ gradients, and we need to run the algorithm for $k_m$ steps, where

$$k_m = \max\left\{\frac{2L}{\mu}, \frac{4\sigma_m^2}{\varepsilon\mu^2}\right\} \log\left(\frac{2\|\theta_0 - \theta^*\|^2}{\varepsilon}\right).$$

The runtime is proportional to $mk_m$, that is

$$mk_m = \max\left\{\frac{2Lm}{\mu}, \frac{4m\sigma_m^2}{\varepsilon\mu^2}\right\} \log\left(\frac{2\|\theta_0 - \theta^*\|^2}{\varepsilon}\right)$$

$$= \max\left\{\frac{2Lm}{\mu}, \frac{4}{\varepsilon\mu^2}\frac{n-m}{n-1}\sigma_1^2\right\} \log\left(\frac{2\|\theta_0 - \theta^*\|^2}{\varepsilon}\right).$$

11. Prove that the optimal batch size $m^*$ allowing one to reach precision $\varepsilon$ in the shortest possible runtime is given by

$$m^* = \frac{2n\sigma_1^2}{2\sigma_1^2 + L(n-1)\varepsilon\mu},$$

4

assuming for simplicity that this quantity is an integer.

It now suffices to minimize the quantity $\max\left\{\frac{2Lm}{\mu}, \frac{4}{\varepsilon\mu^2}\frac{n-m}{n-1}\sigma_1^2\right\}$ over $m \in \{1, \dots, n\}$. In this maximum, the first term $\frac{2Lm}{\mu}$ is increasing with respect to $m$, while the second term $\frac{4}{\varepsilon\mu^2}\frac{n-m}{n-1}\sigma_1^2$ is decreasing with $m$. The maximum is therefore attained when the two terms are equal, or equivalently

$$\frac{2Lm}{\mu} = \frac{4}{\varepsilon\mu^2}\frac{n-m}{n-1}\sigma_1^2 \quad \Longleftrightarrow \quad m = \frac{2n\sigma_1^2}{2\sigma_1^2 + L(n-1)\varepsilon\mu}.$$

12. Discuss the behavior of $m^*$ when $\varepsilon \to 0$ or $\varepsilon \to 1$, assuming $n$ to be "very large".

When $\varepsilon \to 0$, we get $m \to n$. Therefore, to reach a very high precision (that is, a very small $\varepsilon$), gradient descent is generally preferable.

Conversely, when $\varepsilon = 1$, we get $m \approx \frac{2\sigma_1^2}{L\mu}$, which is a constant independent of $n$. In this case, the mini-batch SGD resembles SGD for which we only compute *one* gradient at a time, except that the constant one is replaced with a constant depending on the problem at hand.

**Remark 1.** *Unfortunately, the quantities $\sigma_1^2, L,$ and $\mu$ are unknown to the practitioner, so $m^*$ cannot be evaluated prior to running the algorithm. Still, this result highlights a trade-off between choosing a small batch size, which allows us to reach low precision (high $\varepsilon$) faster, versus choosing a large batch size, which yields a higher precision $\varepsilon \to 0$ faster.*