# Optimization for ML: Exercise Session 2

**Exercise 1** (tightness and statistical optimality). The goal of this exercise is to explore whether the bounds proven in the lectures are tight and statistically optimal.

Let $\xi$ be a random variable with some distribution $\mathcal{Q}$ with a finite second moment and define $f : \mathbb{R} \to \mathbb{R}$,

$$f(\theta) = \frac{1}{2}\mathbb{E}\left(\xi - \theta\right)^2 .$$

1. Compute the unique minimizer $\theta_*$ of $f$ and $f(\theta_*)$.

2. Express a stochastic gradient descent on $f$ that does not have a direct access to $\mathcal{Q}$ but only to i.i.d. samples $\xi_1, \xi_2, \cdots \sim \mathcal{Q}$.

3. We consider the stochastic gradient descent with constant stepsize $\gamma$.

   (a) Using a theorem of the lectures, bound $\mathbb{E}(\theta_k - \theta_*)^2$.
   (b) Compute $\mathbb{E}(\theta_k - \theta_*)^2$ exactly and compare with the bound obtained in the previous question.

4. We consider the stochastic gradient descent with variable stepsize $\gamma_k = \beta/(k_0 + k)$.

   (a) Using a theorem of the lectures, bound $\mathbb{E}(\theta_k - \theta_*)^2$.
   (b) Assume $\beta = 1$ and $k_0 = 1$. Express $\theta_k$ as a function of $\xi_1, \ldots, \xi_k$. Compute $\mathbb{E}(\theta_k - \theta_*)^2$.

   *The empirical average is an optimal (minimax) estimator of the mean; stochastic gradient descent is said to be statistically optimal on this problem as its performance differs only by a multiplicative constant.*

   *This exercise motivates the use of decaying stepsizes $\gamma = \Theta(1/k)$ and that it is hopeless to obtain a better rate than $\Theta(1/k)$ in our general setting.*

**Exercise 2** (importance sampling for coordinate gradient descent). Let $F : \mathbb{R}^p \to \mathbb{R}$ be a continuously differentiable, $\mu$-strongly fonction for some $\mu > 0$. As usual, we denote $\theta_*$ the global minimizer of $F$. We assume that $F$ is $(L_1, \ldots, L_p)$-smooth, in the sense that

$$\forall \theta, \theta' \in \mathbb{R}^p, \quad F(\theta') \leqslant F(\theta) + \langle \nabla F(\theta), \theta' - \theta \rangle + \frac{1}{2}\sum_{j=1}^{p} L_j(\theta'(j) - \theta(j))^2 .$$

The notion of $(L_1, \ldots, L_p)$-smoothness refines the notion of $L$-smoothness by allowing different curvatures along the different coordinates.

1. In this question, we consider the coordinate gradient descent algorithm: choose $\theta_0 \in \mathbb{R}^p$ and for all $k \in \mathbb{N}$, sample $j_{k+1} \sim \text{Unif}(\{1, \ldots, p\})$ independently of the past and compute $\theta_{k+1}$ such that

$$\theta_{k+1}(j_{k+1}) = \theta_k(j_{k+1}) - \gamma \partial_{j_{k+1}} F(\theta_k), \tag{1}$$
$$\theta_{k+1}(j) = \theta_k(j), \qquad j \neq j_{k+1}, \tag{2}$$

   where we choose the stepsize $\gamma = \frac{1}{2\max_j L_j}$.

**(a)** Using a result from the lectures, show that the iteration complexity to obtain $\mathbb{E}\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$ is $k = O\left(\left(\log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon}\right) p \max_j \frac{L_j}{\mu}\right)$.

We now show that this upper bound on the iteration complexity is tight. Consider $p \geqslant 2$ and the function $H(\theta) = \frac{1}{2}\sum_{j=1}^{p} L_j \theta(j)^2$.

**(b)** Show that $H$ is $(L_1, \ldots, L_p)$-smooth and $\mu$-strongly convex with $\mu = \min_j L_j$.

**(c)** Denote $j_{\min} = \operatorname{argmin}_j L_j$. When $\theta_0 = e_{j_{\min}}$ is the $j_{\min}$-th element of the canonical basis, show that the coordinate gradient descent (1)–(2) on $F = H$ satisfies

$$\mathbb{E}\|\theta_k - \theta_*\|^2 \geqslant \left(1 - \frac{\mu}{p \max_j L_j}\right)^k \|\theta_0 - \theta_*\|^2 .$$

**(d)** Conclude that in this case, the iteration complexity to obtain $\mathbb{E}\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$ is $k = \Omega\left(\left(\log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon}\right) p \max_j \frac{L_j}{\mu}\right)$.

2. The goal of this question is to show that the iteration complexity of stochastic gradient descent can be improved by an appropriate weighted sampling of the coordinates.

We consider the following weighted generalization of the coordinate gradient descent method. Let $\pi = (\pi_1, \ldots, \pi_p)$ denote a probability distribution on $\{1, \ldots, p\}$. Choose $\theta_0 \in \mathbb{R}^p$ and for all $k \in \mathbb{N}$, sample $j_{k+1} \sim \pi$ independently of the past and compute $\theta_{k+1}$ such that

$$\theta_{k+1}(j_{k+1}) = \theta_k(j_{k+1}) - \gamma_{j_{k+1}} \partial_{j_{k+1}} F(\theta_k) ,$$
$$\theta_{k+1}(j) = \theta_k(j) , \qquad j \neq j_{k+1} ,$$

where $\gamma_1, \ldots, \gamma_p$ are now coordinate-dependent stepsizes.

**(a)** Prove that, if $\gamma_j \propto \pi_j^{-1}$, the weighted coordinate gradient descent is a stochastic gradient descent in the sense of the lectures.

**(b)** Show that for all $\theta, \theta' \in \mathbb{R}^p$,

$$\sum_{j=1}^{p} \frac{1}{L_j} \left(\partial_j F(\theta) - \partial_j F(\theta')\right)^2 \leqslant \langle \theta - \theta', \nabla F(\theta) - \nabla F(\theta')\rangle .$$

**(c)** Consider a weighted coordinate gradient descent with weights $\pi_j = \frac{L_j}{\sum_{j'} L_{j'}}$. Show that, for some appropriate choice of the stepsizes $\gamma_1, \ldots, \gamma_p$ to be determined, the iteration complexity to obtain $\mathbb{E}\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$ is

$$k = O\left(\left(\log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon}\right) \sum_j \frac{L_j}{\mu}\right) .$$

Importance sampling improved the dependence in the worst of the condition numbers $\max_j \frac{L_j}{\mu}$ to the average of the condition numbers $\frac{1}{p}\sum_j \frac{L_j}{\mu}$.

**Exercise 1** (tightness and statistical optimality). The goal of this exercise is to explore whether the bounds proven in the lectures are tight and statistically optimal.

Let $\xi$ be a random variable with some distribution $Q$ with a finite second moment and define $f : \mathbb{R} \to \mathbb{R}$,

$$f(\theta) = \frac{1}{2}\mathbb{E}\left(\xi - \theta\right)^2 .$$

1. Compute the unique minimizer $\theta_*$ of $f$ and $f(\theta_*)$.

---

(1) $\quad f(\theta) = \frac{1}{2} E[\xi^2] + \frac{1}{2}\theta^2 - \theta \cdot E[\xi]$

par la règle de Fermat, $\quad \nabla f(\theta^*) = \theta^* - E[\xi] = 0$

alors $\quad \theta^* = E[\xi]$

---

$$f(\theta^*) = \frac{1}{2}E\left[(\xi - E[\xi])^2\right] = \frac{1}{2}\text{Var}[\xi]$$

---

2. Express a stochastic gradient descent on $f$ that does not have a direct access to $Q$ but only to i.i.d. samples $\xi_1, \xi_2, \cdots \sim Q$.

---

(2) On est dans un cas d'approximation stochastique,

$$E[f(\theta,\xi)] = f(\theta) = E\left[\frac{1}{2}(\xi - \theta)^2\right]$$

$$f(\theta,\xi) = \frac{1}{2}(\xi - \theta)^2$$

$$g(\theta,\xi) = \frac{\partial}{\partial\theta} f(\theta,\xi) = \theta - \xi$$

Descent de gradient stochastique

$\theta_0 \in \mathbb{R}$, pour tout $k \in \mathbb{N}$, on prend $\xi_{k+1} \sim Q$, $\gamma_k \in \mathbb{R}$

$$\theta_{k+1} = \theta_k - \gamma_k \cdot g(\theta_k, \xi_{k+1})$$

$$= \theta_k - \gamma_k \cdot (\theta_k - \xi_{k+1})$$

$$= (1 - \gamma_k) \cdot \theta_k + \gamma_k \cdot \xi_{k+1}$$

**3.** We consider the stochastic gradient descent with constant stepsize $\gamma$.

    **(a)** Using a theorem of the lectures, bound $\mathbb{E}(\theta_k - \theta_*)^2$.

---

(3) (a)      $\dfrac{\delta^2}{\delta\theta^2} f(\theta,\xi) = 1 \leq 1$

donc $f(\theta,\xi)$ est $M$-lisse , avec $M=1$

$\nabla^2 f(\theta) = 1 \geq 1$

donc $f(\theta)$ est $\mu$-fortement convexe , avec $\mu = 1$

et

$$\sigma^2 = \mathbb{E}|g(\theta_*,\xi)|^2 = \mathbb{E}\left[(\xi - \mathbb{E}[\xi])^2\right] = \text{Var}[\xi]$$

D'après le Thm. 4.17, Si $\gamma_k = \gamma < \dfrac{1}{2M} = \dfrac{1}{2}$ , alors

$$\mathbb{E}(\theta_k - \theta_*)^2 \leq (1-\gamma\mu)^k \cdot |\theta_0 - \theta_*|^2 + \dfrac{2\gamma\sigma^2}{\mu}$$

$$= (1-\gamma)^k \cdot |\theta_0 - \theta_*|^2 + 2\gamma \cdot \text{Var}[\xi]$$

    **(b)** Compute $\mathbb{E}(\theta_k - \theta_*)^2$ exactly and compare with the bound obtained in the previous question.

---

(b)    $\theta_k$ est un estimateur stochastique de $\theta_*$

$$\mathbb{E}(\theta_k - \theta_*)^2 = b(\theta_k)^2 + \text{Var}[\theta_k]$$

· $b(\theta_k)^2$

$b(\theta_k) = \mathbb{E}[\theta_k] - \theta_*$

$b(\theta_{k+1}) = \mathbb{E}[\theta_{k+1}] - \theta_* = \mathbb{E}\left[(1-\gamma)\theta_k + \gamma \cdot \xi_{k+1}\right] - \theta_*$

$$= (1-\gamma)\cdot\mathbb{E}[\theta_k] + \gamma \cdot \theta_* - \theta_*$$

$$= (1-\gamma)\cdot(\mathbb{E}[\theta_k] - \theta_*)$$

$$= (1-\gamma)\cdot b(\theta_k)$$

alors $\quad b(\theta_k)^2 = (1-r)^{2k} \cdot b(\theta_0)^2 = (1-r)^{2k} \cdot |\theta_0 - \theta_*|^2$

- $Var[\theta_k]$

$$Var[\theta_{k+1}] = Var[(1-r) \cdot \theta_k + r \cdot \xi_{k+1}]$$
$$= (1-r)^2 \cdot Var[\theta_k] + r^2 \cdot Var[\xi]$$

soit $\quad v = (1-r)^2 \cdot v + r^2 \cdot Var[\xi]$

$$v = \frac{r^2 \cdot Var[\xi]}{1-(1-r)^2} = \frac{r \cdot Var[\xi]}{2-r}$$

et

$$Var[\theta_{k+1}] - v = (1-r)^2 \cdot (Var[\theta_k] - v)$$

alors

$$Var[\theta_k] - v = (1-r)^{2k} \cdot (\underbrace{Var[\theta_0]}_{=0} - v)$$

$$Var[\theta_k] = v - (1-r)^{2k} \cdot v$$
$$= (1-(1-r)^{2k}) \cdot \frac{r \cdot Var[\xi]}{2-r}$$

Donc

$$E(\theta_k - \theta_*)^2 = (1-r)^{2k} \cdot |\theta_0 - \theta_*|^2 + (1-(1-r)^{2k}) \cdot \frac{r \cdot Var[\xi]}{2-r}$$

**4.** We consider the stochastic gradient descent with variable stepsize $\gamma_k = \beta/(k_0 + k)$.

  **(a)** Using a theorem of the lectures, bound $\mathbb{E}(\theta_k - \theta_*)^2$.

---

(4) (a) D'après (3)(a)

  (i) $f(\theta, \xi)$ est $M$-lisse , $M = 1$

  (ii) $f(\theta)$ est $\mu$-fortement convexe , $\mu = 1$

  (iii) $\sigma^2 = \mathbb{E}|g(\theta_*, \xi)|^2 = \text{Var}[\xi]$

Par la théorème 4.19 , si $\beta > \dfrac{1}{\mu}$ et $\gamma_0 = \dfrac{\beta}{k_0} \leqslant \dfrac{1}{2M}$ ,

alors

$$E(\theta_k - \theta_*)^2 \leqslant \frac{\nu}{k_0 + k}$$

avec $\quad \nu = \max\left\{ k_0 \|\theta_0 - \theta_*\|^2 , \dfrac{2\sigma^2\beta^2}{\beta\mu - 1} \right\}$

---

  **(b)** Assume $\beta = 1$ and $k_0 = 1$. Express $\theta_k$ as a function of $\xi_1, \ldots, \xi_k$. Compute $\mathbb{E}(\theta_k - \theta_*)^2$.

  *The empirical average is an optimal (minimax) estimator of the mean; stochastic gradient descent is said to be statistically optimal on this problem as its performance differs only by a multiplicative constant.*

  *This exercise motivates the use of decaying stepsizes $\gamma = \Theta(1/k)$ and that it is hopeless to obtain a better rate than $\Theta(1/k)$ in our general setting.*

---

(b) $\beta = 1$ , $k_0 = 1 \implies \gamma_k = \dfrac{1}{k+1}$

---

$\theta_{k+1} = \theta_k - \gamma_k \cdot g(\theta_k, \xi_k) = \theta_k - \dfrac{\theta_k - \xi_k}{k+1} = \dfrac{k}{k+1} \cdot \theta_k + \dfrac{1}{k+1} \cdot \xi_k$

---

Pour montrer $\theta_k = \dfrac{\xi_1 + \cdots + \xi_k}{k}$

Récurrence

  • $\theta_1 = 0 \cdot \theta_0 + \xi_1 = \xi_1$

- soit $\theta_k = \dfrac{\xi_1 + \cdots + \xi_k}{k}$

$$\theta_{k+1} = \frac{k}{k+1} \cdot \frac{\xi_1 + \cdots + \xi_k}{k} + \frac{\xi_{k+1}}{k+1} = \frac{\xi_1 + \cdots + \xi_{k+1}}{k+1}$$

Donc $\quad \theta_k = \dfrac{\xi_1 + \cdots + \xi_k}{k}$

$$
\begin{aligned}
E(\theta_k - \theta_*)^2 &= E\left[\left(\frac{\xi_1 + \cdots + \xi_k}{k} - E[\xi]\right)^2\right] \\
&= E\left[\left(\frac{\xi_1 + \cdots + \xi_k}{k} - E\left[\frac{\xi_1 + \cdots + \xi_k}{k}\right]\right)^2\right] \\
&= Var\left[\frac{\xi_1 + \cdots + \xi_k}{k}\right] \\
&= \frac{\sigma^2}{k}
\end{aligned}
$$

**Exercise 2** (importance sampling for coordinate gradient descent). Let $F : \mathbb{R}^p \to \mathbb{R}$ be a continuously differentiable, $\mu$-strongly fonction for some $\mu > 0$. As usual, we denote $\theta_*$ the global minimizer of $F$. We assume that $F$ is $(L_1, \ldots, L_p)$-smooth, in the sense that

$$\forall \theta, \theta' \in \mathbb{R}^p, \quad F(\theta') \leqslant F(\theta) + \langle \nabla F(\theta), \theta' - \theta \rangle + \frac{1}{2} \sum_{j=1}^{p} L_j (\theta'(j) - \theta(j))^2.$$

The notion of $(L_1, \ldots, L_p)$-smoothness refines the notion of $L$-smoothness by allowing different curvatures along the different coordinates.

1. In this question, we consider the coordinate gradient descent algorithm: choose $\theta_0 \in \mathbb{R}^p$ and for all $k \in \mathbb{N}$, sample $j_{k+1} \sim \mathrm{Unif}(\{1, \ldots, p\})$ independently of the past and compute $\theta_{k+1}$ such that

$$\theta_{k+1}(j_{k+1}) = \theta_k(j_{k+1}) - \gamma \partial_{j_{k+1}} F(\theta_k), \tag{1}$$
$$\theta_{k+1}(j) = \theta_k(j), \qquad j \neq j_{k+1}, \tag{2}$$

where we choose the stepsize $\gamma = \frac{1}{2 \max_j L_j}$.

(a) Using a result from the lectures, show that the iteration complexity to obtain $\mathbb{E}\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$ is $k = O\left(\left(\log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon}\right) p \max_j \frac{L_j}{\mu}\right)$.

---

(1) (a) Par le corollaire , pour $\mathbb{E}\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$

si $\quad k \geqslant 2 \cdot \left( \log \frac{2\|\theta_0 - \theta_*\|^2}{\varepsilon} \right) \cdot \frac{LP}{\mu}$

$\geqslant 2 \cdot \left( \log \frac{2\|\theta_0 - \theta_*\|^2}{\varepsilon} \right) \cdot p \cdot \max_j \frac{L_j}{\mu}$

donc

$$k = O\left( \left( \log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon} \right) \cdot p \cdot \max_j \frac{L_j}{\mu} \right)$$

---

We now show that this upper bound on the iteration complexity is tight. Consider $p \geqslant 2$ and the function $H(\theta) = \frac{1}{2} \sum_{j=1}^{p} L_j \theta(j)^2$.

(b) Show that $H$ is $(L_1, \ldots, L_p)$-smooth and $\mu$-strongly convex with $\mu = \min_j L_j$.

---

(b) $\forall j = 1, \ldots, p$, soit $H_j(\theta) = \frac{1}{2} L_j \theta(j)^2$

$\forall \theta, \theta' \in \mathbb{R}^p$,

$\left| \partial H_j(\theta) - \partial H_j(\theta') \right| = \left| L_j \theta(j) - L_j \theta'(j) \right| \leqslant L_j \cdot \left| \theta(j) - \theta'(j) \right|$

donc $H_j$ est $L_j$-lisse

Alors on a

$$H_j(\theta') \le H_j(\theta) + \partial H_j(\theta) \cdot (\theta'_{(j)} - \theta_{(j)}) + \frac{L_j}{2}(\theta'_{(j)} - \theta_{(j)})^2$$

donc

$$\sum_{j=1}^{P} H_j(\theta') \le \sum_{j=1}^{P} H_j(\theta) + \sum_{j=1}^{P} \partial H_j(\theta) \cdot (\theta'_{(j)} - \theta_{(j)}) + \frac{1}{2} \sum_{j=1}^{P} L_j \cdot (\theta'_{(j)} - \theta_{(j)})^2$$

$$H(\theta') \le H(\theta) + \langle \nabla H(\theta), \theta' - \theta \rangle + \frac{1}{2} \sum_{j=1}^{P} L_j \cdot (\theta'_{(j)} - \theta_{(j)})^2$$

alors $H$ est $(L_1, \dots, L_p) -$ lisse

D'après la définition, $H_j(\theta)$ est fortement convexe de module $- L_j$, donc on a

$$H_j(\theta) \ge H_j(\theta') + \partial H_j(\theta) \cdot (\theta'_{(j)} - \theta_{(j)}) + \frac{L_j}{2} \cdot (\theta'_{(j)} - \theta_{(j)})^2$$

alors

$$\sum_{j=1}^{P} H_j(\theta) \ge \sum_{j=1}^{P} H_j(\theta') + \sum_{j=1}^{P} \partial H_j(\theta) \cdot (\theta'_{(j)} - \theta_{(j)}) + \sum_{j=1}^{P} \frac{L_j}{2}(\theta'_{(j)} - \theta_{(j)})^2$$

$$H(\theta) \ge H(\theta') + \langle \nabla H(\theta), \theta' - \theta \rangle + \frac{1}{2} \sum_{j=1}^{P} L_j \cdot (\theta'_{(j)} - \theta_{(j)})^2$$

$$H(\theta) \ge H(\theta') + \langle \nabla H(\theta), \theta' - \theta \rangle + \frac{1}{2} \cdot \min_j L_j \sum_{j=1}^{P} (\theta'_{(j)} - \theta_{(j)})^2$$

donc $H$ est $\mu-$ fortement convexe, avec $\mu = \min_j L_j$

**(c)** Denote $j_{\min} = \operatorname{argmin}_j L_j$. When $\theta_0 = e_{j_{\min}}$ is the $j_{\min}$-th element of the canonical basis, show that the coordinate gradient descent (1)–(2) on $F = H$ satisfies

$$\mathbb{E}\|\theta_k - \theta_*\|^2 \ge \left(1 - \frac{\mu}{p \max_j L_j}\right)^k \|\theta_0 - \theta_*\|^2.$$

(c) Soit $\theta_0 = e_{j\min}$, $H(\theta_0) = \frac{1}{2} L_{j\min} = \frac{1}{2} \min_j L_j$

pour $j_{k+1} \sim \operatorname{Unif}(\{1, \dots, p\})$

$$\|\theta_{k+1} - \theta^*\|^2 = \|\theta_k - \gamma \cdot \partial_{j_{k+1}} H(\theta_k) \cdot e_{j_{k+1}} - \theta^*\|^2$$

$$= \|\theta_k - \theta^* - \gamma \cdot \partial_{j_{k+1}} H(\theta_k) \cdot e_{j_{k+1}}\|^2$$

$$= \|\theta_k - \theta^*\|^2 - 2\langle \theta_k - \theta^*, \gamma \cdot L_{j_{k+1}} \cdot \theta_k(j_{k+1}) \cdot e_{j_{k+1}}\rangle + \gamma^2 \cdot L_{j_{k+1}}^2 \cdot \theta^2(j_{k+1})$$

$$\geqslant \|\theta_k - \theta^*\|^2 - \langle \theta_k - \theta^*, 2 \cdot \gamma \cdot L_{j_{k+1}} \cdot \theta_k(j_{k+1}) \cdot e_{j_{k+1}}\rangle$$

donc

$$E\left[\|\theta_{k+1} - \theta^*\|^2\right] \geqslant E\left[\|\theta_k - \theta^*\|^2 - \langle \theta_k - \theta^*, 2 \cdot \gamma \cdot L_{j_{k+1}} \cdot \theta_k(j_{k+1}) \cdot e_{j_{k+1}}\rangle\right]$$

$$= E\left[\|\theta_k - \theta^*\|^2\right] - \frac{2\gamma}{P} \cdot E\left[\langle \theta_k - \theta^*, \nabla H(\theta_k) - \nabla H(\theta^*)\rangle\right]$$

$$\geqslant E\left[\|\theta_k - \theta^*\|^2\right] - \frac{2\gamma}{P} \cdot \mu \cdot E\left[\|\theta_k - \theta^*\|^2\right]$$

$$= \left(1 - \frac{2 \cdot \mu}{P \cdot \max_j L_j}\right) \cdot E\left[\|\theta_k - \theta^*\|^2\right]$$

alors

$$E\|\theta_k - \theta^*\|^2 \geqslant \left(1 - \frac{2 \cdot \mu}{P \cdot \max_j L_j}\right)^k \cdot \|\theta_0 - \theta^*\|^2$$

**(d)** Conclude that in this case, the iteration complexity to obtain $\mathbb{E}\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$ is $k = \Omega\left(\left(\log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon}\right) p \max_j \frac{L_j}{\mu}\right)$.

(d) Pour avoir $E\|\theta_k - \theta^*\|^2 \leqslant \varepsilon$, d'après (c), on a

$$\left(1 - \frac{2 \cdot \mu}{P \cdot \max_j L_j}\right)^k \cdot \|\theta_0 - \theta^*\|^2 \leqslant \varepsilon$$

$$\left(1 - \frac{2 \cdot \mu}{P \cdot \max_j L_j}\right)^k \leqslant \frac{\varepsilon}{\|\theta_0 - \theta^*\|^2}$$

$$k \cdot \log\left(1 - \frac{2 \cdot \mu}{P \cdot \max_j L_j}\right) \leqslant \log \frac{\varepsilon}{\|\theta_0 - \theta^*\|^2}$$

$$k \geqslant \log\left(\frac{\varepsilon}{\|\theta_0 - \theta^*\|^2}\right) \cdot \frac{1}{\log\left(1 - \frac{2 \cdot \mu}{P \cdot \max_j L_j}\right)}$$

$$k \geqslant -\log\left(\frac{\varepsilon}{\|\theta_0 - \theta^*\|^2}\right) \cdot \frac{P \cdot \max_j L_j}{2\mu}$$

$$k \geqslant \log\left(\frac{\|\theta_0 - \theta^*\|^2}{\varepsilon}\right) \cdot \frac{P}{2} \cdot \max_j \frac{L_j}{\mu}$$

2. The goal of this question is to show that the iteration complexity of stochastic gradient descent can be improved by an appropriate weighted sampling of the coordinates.

We consider the following weighted generalization of the coordinate gradient descent method. Let $\pi = (\pi_1, \ldots, \pi_p)$ denote a probability distribution on $\{1, \ldots, p\}$. Choose $\theta_0 \in \mathbb{R}^p$ and for all $k \in \mathbb{N}$, sample $j_{k+1} \sim \pi$ independently of the past and compute $\theta_{k+1}$ such that

$$\theta_{k+1}(j_{k+1}) = \theta_k(j_{k+1}) - \gamma_{j_{k+1}} \partial_{j_{k+1}} F(\theta_k),$$
$$\theta_{k+1}(j) = \theta_k(j), \qquad j \neq j_{k+1},$$

where $\gamma_1, \ldots, \gamma_p$ are now coordinate-dependent stepsizes.

(a) Prove that, if $\gamma_j \propto \pi_j^{-1}$, the weighted coordinate gradient descent is a stochastic gradient descent in the sense of the lectures.

---

(a) Pour $F : \mathbb{R}^p \longrightarrow \mathbb{R}$, $\xi = j \sim D(\pi)$

$$g(\theta, \xi) = g(\theta, j) = \frac{1}{\pi_j} \cdot \partial_j F(\theta) \cdot e_j$$

$$E[g(\theta, \xi)] = \nabla F(\theta)$$

---

· Descente de gradient stochastique

$\theta_0 \in \mathbb{R}^p$ et $\forall k \in \mathbb{N}$, on prend $\xi_{k+1} \sim D(\pi)$

soit $\gamma_k = \gamma = c$,

$$\theta_{k+1} = \theta_k - \gamma \cdot g(\theta_k, \xi_{k+1})$$

$$= \theta_k - c \cdot \frac{1}{\pi_{j_{k+1}}} \cdot \partial_{j_{k+1}} F(\theta) \cdot e_{j_{k+1}}$$

$$= \theta_k - \gamma_{k+1} \cdot \partial_{j_{k+1}} F(\theta) \cdot e_{j_{k+1}}$$

---

(b) Show that for all $\theta, \theta' \in \mathbb{R}^p$,

$$\sum_{j=1}^p \frac{1}{L_j} (\partial_j F(\theta) - \partial_j F(\theta'))^2 \leq \langle \theta - \theta', \nabla F(\theta) - \nabla F(\theta') \rangle.$$

(b) $\forall \theta, \theta', \theta'' \in \mathbb{R}^P$,

$$F(\theta + \theta'') \leq F(\theta) + \langle \nabla F(\theta), \theta'' \rangle + \frac{1}{2} \sum_{j=1}^{P} L_j \cdot \theta''_{(j)}{}^2$$

$$F(\theta + \theta'') \geq F(\theta') + \langle \nabla F(\theta'), \theta + \theta'' - \theta' \rangle$$

alors

$$F(\theta') + \langle \nabla F(\theta'), \theta + \theta'' - \theta' \rangle \leq F(\theta) + \langle \nabla F(\theta), \theta'' \rangle + \frac{1}{2} \sum_{j=1}^{P} L_j \cdot \theta''_{(j)}{}^2$$

$$\langle \nabla F(\theta') - \nabla F(\theta), \theta'' \rangle - \frac{1}{2} \sum_{j=1}^{P} L_j \cdot \theta''_{(j)}{}^2 \leq F(\theta) - F(\theta') - \langle \nabla F(\theta'), \theta - \theta' \rangle$$

par la règle de Fermat, on optim en $\theta''_{(j)}$

$$\forall j = 1, \ldots, P, \qquad \theta''_{(j)} = \frac{1}{L_j} \left( \partial_j F(\theta') - \partial_j F(\theta) \right)$$

donc

$$\frac{1}{2} \sum_{j=1}^{P} \frac{1}{L_j} \cdot \left( \partial_j F(\theta') - \partial_j F(\theta) \right)^2 \leq F(\theta) - F(\theta') - \langle \nabla F(\theta'), \theta - \theta' \rangle$$

on inverse $\theta$ et $\theta'$

$$\frac{1}{2} \sum_{j=1}^{P} \frac{1}{L_j} \cdot \left( \partial_j F(\theta) - \partial_j F(\theta') \right)^2 \leq F(\theta') - F(\theta) - \langle \nabla F(\theta), \theta' - \theta \rangle$$

et on somme

$$\sum_{j=1}^{P} \frac{1}{L_j} \cdot \left( \partial_j F(\theta') - \partial_j F(\theta) \right)^2 \leq \langle \nabla F(\theta') - \nabla F(\theta), \theta' - \theta \rangle$$

(c) Consider a **weighted** coordinate gradient descent with **weights** $\pi_j = \frac{L_j}{\sum_{j'} L_{j'}}$. Show that, for some **appropriate choice** of the stepsizes $\gamma_1, \ldots, \gamma_p$ to be determined, the **iteration complexity** to obtain $\mathbb{E}\|\theta_k - \theta_*\|^2 \leq \varepsilon$ is

$$k = O\left( \left( \log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon} \right) \sum_j \frac{L_j}{\mu} \right).$$

Importance sampling improved the dependence in the worst of the condition numbers $\max_j \frac{L_j}{\mu}$ to the average of the condition numbers $\frac{1}{p} \sum_j \frac{L_j}{\mu}$.

(c) (i) $F$ est $\mu$-fortement convexe.

(ii) $g(\theta, \xi)$ est $M$-lisse

$$E\| g(\theta, j) - g(\theta', j)\|^2 = \sum_{j=1}^{P} \lambda_j \cdot \| \frac{1}{\lambda_j} \partial_j F(\theta) \cdot e_j - \frac{1}{\lambda_j} \partial_j F(\theta') \cdot e_j \|^2$$

$$= \sum_{j=1}^{P} \frac{1}{\lambda_j} \cdot \left( \partial_j F(\theta) - \partial_j F(\theta') \right)^2$$

$$= \left( \sum_{j'=1}^{P} L_{j'} \right) \cdot \sum_{j=1}^{P} \frac{1}{L_j} \left( \partial_j F(\theta) - \partial_j F(\theta') \right)^2$$

$$\leq \left( \sum_{j'=1}^{P} L_{j'} \right) \cdot \left\langle \theta - \theta', \nabla F(\theta) - \nabla F(\theta') \right\rangle$$

donc $M = \sum_{j=1}^{P} L_j$

(iii) $\sigma^2 = E\| g(\theta_*, \xi)\|^2 = \sum_{j=1}^{P} \lambda_j \cdot \| \frac{1}{\lambda_j} \cdot \partial_j F(\theta_*) \cdot e_j \|^2 = 0$

D'après le corollaire du cours, il existe un choix de $\gamma = c$ tel que

pour $k \geqslant O\left( \left( \log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon} \right) \cdot \frac{M}{\mu} \right)$ alors $E\|\theta_k - \theta_*\|^2 \leqslant \varepsilon$

donc

$$k = O\left( \left( \log \frac{\|\theta_0 - \theta_*\|^2}{\varepsilon} \right) \cdot \sum_{j=1}^{P} \frac{L_j}{\mu} \right)$$