

## Exercise: Convergence rate of Adagrad

For any three vectors  $a, a', b \in \mathbb{R}^d$ , we define

$$a \odot b = \begin{pmatrix} a_1 b_1 \\ \vdots \\ a_d b_d \end{pmatrix}, \quad a \oslash b \text{ or } \frac{a}{b} = \begin{pmatrix} a_1/b_1 \\ \vdots \\ a_d/b_d \end{pmatrix},$$

$$\|a\|_b^2 = \sum_{i=1}^d a_i^2 b_i, \quad \langle a, a' \rangle_b = \sum_{i=1}^d a_i a'_i b_i$$

$$a^{\odot 2} = a \odot a,$$

### Algorithm 1: Adagrad

1. Pick  $x_0 \in \mathbb{R}^d$ .
2. Until termination condition, iterate:
  - Generate  $\xi_{k+1}$  independent of the past

$$g_{k+1} = \begin{pmatrix} g_{k+1}(1) \\ \vdots \\ g_{k+1}(d) \end{pmatrix} = \nabla f_x(x_k, \xi_{k+1})$$

$$\gamma_{k+1}(j) = \frac{\alpha}{\sqrt{\sum_{s=0}^k g_{s+1}^2(j)}}, \quad \forall j = 1, \dots, d$$

$$x_{k+1} = x_k - \gamma_{k+1} \odot g_{k+1}$$

The convergence guarantees of Adagrad are given in the theorem below.

**Theorem 1.** *Suppose that*

1.  $f(\cdot, \xi)$  is convex for all  $\xi$
2. There exists  $x^* \in \arg \min F$ , where  $F(x) = \mathbb{E}[f(x, \xi)]$
3. There exists  $D > 0$  s.t. for all  $k \geq 0$ , for all  $i \in \{1, \dots, d\} : |x_{k,i} - x_i^*| \leq D$
4. For all  $x, \xi : \|\nabla f(x, \xi)\| \leq G$ .

*Then the iterates of Adagrad satisfy*

$$\mathbb{E} [F(\bar{x}_K) - F(x^*)] \leq \frac{dG}{\sqrt{K}} \left( \frac{D^2}{\alpha} + 2\alpha \right)$$

where  $\bar{x}_K = \frac{1}{K} \sum_{k=0}^{K-1} x_k$ .

1. Justify that  $f(x_k, \xi_{k+1}) - f(x^*, \xi_{k+1}) \leq \langle g_{k+1}, x_k - x^* \rangle$ . One can use the convexity of  $f(\cdot, \xi_{k+1})$  and the fact that, for any convex function  $\varphi$ , we have  $\varphi(y) - \varphi(x) \geq \langle \nabla \varphi(x), y - x \rangle$ .

2. Using the relation  $x_{k+1} = x_k - \gamma_{k+1} \odot g_{k+1}$ , show that

$$\frac{1}{2} \|x_k - x^*\|_{\gamma_{k+1}^{-1}}^2 - \frac{1}{2} \|x_{k+1} - x^*\|_{\gamma_{k+1}^{-1}}^2 + \frac{1}{2} \|g_{k+1}\|_{\gamma_{k+1}}^2 = \langle g_{k+1}, x_k - x^* \rangle.$$

3. Deduce that

$$\sum_{k=0}^{K-1} f(x_k, \xi_{k+1}) - f(x^*, \xi_{k+1}) \leq \sum_{k=0}^{K-1} \left( \frac{1}{2} \|x_k - x^*\|_{\gamma_{k+1}^{-1}}^2 - \frac{1}{2} \|x_{k+1} - x^*\|_{\gamma_{k+1}^{-1}}^2 \right) + \sum_{k=0}^{K-1} \frac{1}{2} \|g_{k+1}\|_{\gamma_{k+1}}^2. \quad (1)$$

From now on, we will control the two terms in equation (1) separately.

**Control of the first term.**

4. Prove that for any  $z, a, b \in \mathbb{R}^d$  we have  $\|z\|_a^2 - \|z\|_b^2 = \|z\|_{a-b}^2$ .

5. Deduce that the first term in (1) can be rewritten as

$$\sum_{k=0}^{K-1} \left( \|x_k - x^*\|_{\gamma_{k+1}^{-1}}^2 - \|x_{k+1} - x^*\|_{\gamma_{k+1}^{-1}}^2 \right) = \|x_0 - x^*\|_{\gamma_1^{-1}}^2 - \|x_K - x^*\|_{\gamma_K^{-1}}^2 + \sum_{k=1}^{K-1} \|x_k - x^*\|_{(\gamma_{k+1}^{-1} - \gamma_k^{-1})}^2$$

6. Show that the Adagrad learning rates  $\gamma_k$ 's have coordinates that are non-increasing over time:

$$\forall j \in \{1, \dots, d\} : \quad \gamma_{k+1}(j) \leq \gamma_k(j).$$

7. Using Assumption 3 in the theorem, deduce that the first term in (1) can be controlled as

$$\sum_{k=0}^{K-1} \left( \|x_k - x^*\|_{\gamma_{k+1}^{-1}}^2 - \|x_{k+1} - x^*\|_{\gamma_{k+1}^{-1}}^2 \right) \leq \sum_{j=1}^d \left[ \frac{D^2}{\gamma_1(j)} + \sum_{k=1}^{K-1} D^2 \left( \frac{1}{\gamma_{k+1}(j)} - \frac{1}{\gamma_k(j)} \right) \right].$$

8. Using Assumption 4 from the theorem, prove that  $\gamma_K(j) \geq \frac{\alpha}{G\sqrt{K}}$ .

9. Deduce the final upper bound on the first term

$$\sum_{k=0}^{K-1} \left( \|x_k - x^*\|_{\gamma_{k+1}^{-1}}^2 - \|x_{k+1} - x^*\|_{\gamma_{k+1}^{-1}}^2 \right) \leq \frac{D^2 d G \sqrt{K}}{\alpha}.$$

**Control of the second part of (1):**  $\sum_{k=0}^{K-1} \|g_{k+1}\|_{\gamma_{k+1}}^2$ .

10. For any  $k \in \{0, K-1\}$ , define  $a_k^{(i)} = g_{k+1,i}^2 \geq 0$ . Justify that

$$\sum_{k=0}^{K-1} \|g_{k+1}\|_{\gamma_{k+1}}^2 = \alpha \sum_{k=0}^{K-1} \sum_{i=1}^d \frac{a_k^{(i)}}{\sqrt{\sum_{s=0}^k a_s^{(i)}}}.$$

We will use the following Lemma

**Lemma 1.** *Let  $(a_k)$  be a sequence of nonnegative numbers. Then*

$$\sum_{k=0}^{K-1} \frac{a_k}{\sqrt{\sum_{s=0}^k a_s}} \leq 2 \sqrt{\sum_{s=0}^{K-1} a_s}$$

The proof of this lemma is a bonus question at the end of the exercise. For now, we will simply use this result without proof.

11. Using the lemma above, show that

$$\sum_{k=0}^{K-1} \|g_{k+1}\|_{\gamma_{k+1}}^2 \leq 2\alpha dG\sqrt{K}.$$

12. Deduce that

$$\sum_{k=0}^{K-1} f(x_k, \xi_{k+1}) - f(x^*, \xi_{k+1}) \leq \frac{dD^2G\sqrt{K}}{\alpha} + 2\alpha dG\sqrt{K}.$$

13. Applying the expectation on both sides, and deduce that

$$\frac{dD^2G}{\alpha\sqrt{K}} + 2\frac{\alpha dG}{\sqrt{K}} \geq \frac{1}{K} \mathbb{E} \left[ \sum_{k=0}^{K-1} F(x_k) - F(x^*) \right].$$

14. Using the convexity of  $F$ , conclude the proof of the theorem.

**Bonus question:** Prove the lemma.