

1 A Deterministic Convergence Analysis

2 Here we analyze the AdaRem method in full batch setting and show that when the objective function
3 is L -lipschitz convex, AdaRem converges with rate $O(1/k)$.

4 **Notation** Given two vectors $u, v \in \mathbb{R}^d$, we use $\langle u, v \rangle$ for inner product, $u \odot v$ for element-wise
5 product, u/v to denote element-wise division. Given a vector $x \in \mathbb{R}^d$ we denote its i -th coordinate by
6 x_i and its ℓ_2 -norm by $\|x\|_2$. For a vector x_t in the t -th iteration, the i -th coordinate of x_t is denoted
7 as $x_{t,i}$ by adding a subscript i . We use g_i to denote $\nabla f(x)_i$.

8 **Claim A.1.** *If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -lipschitz convex, then for all x, y ,*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|_2^2$$

9 **Theorem A.2.** *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a L -lipschitz convex function and $x^* = \arg \min_x f(x)$. For
10 AdaRem, we set adjustment coefficient a_t bounded: $c \leq a_t \leq 2$, and weight decay factor $\lambda = 0$.
11 Then, AdaRem with learning rate $\eta \leq \frac{1}{2L}$ satisfies the following:*

$$f(x_k) \leq f(x^*) + \frac{\|x_0 - x^*\|_2^2}{2c\eta k}$$

12 *Proof.* For $i = 0, \dots$, define

$$x_{i+1} = x_i - t_i \odot g_i$$

13

$$t_i = \eta a_i$$

14 First, by convexity of f , we have:

$$f(x_i) \leq f(x^*) + \langle g_i, x_i - x^* \rangle. \quad (1)$$

15 Further, as f is L -lipschitz, by the previous lemma,

$$\begin{aligned} f(x_{i+1}) &\leq f(x_i) + \langle g_i, x_{i+1} - x_i \rangle + \frac{L}{2} \|x_{i+1} - x_i\|_2^2 \\ &= f(x_i) + \langle g_i, -t_i \odot g_i \rangle + \frac{L}{2} \|t_i \odot g_i\|_2^2 \\ &= f(x_i) - \sum_{j=1}^d g_{i,j}^2 t_{i,j} + \frac{L}{2} \sum_{j=1}^d g_{i,j}^2 t_{i,j}^2 \\ &= f(x_i) - \sum_{j=1}^d t_{i,j} \left(1 - \frac{t_{i,j} L}{2}\right) g_{i,j}^2 \\ &\leq f(x_i) - \sum_{j=1}^d \frac{t_{i,j}}{2} g_{i,j}^2 \end{aligned} \quad (2)$$

16 where the last inequality follows as $t_{i,j} L = \eta a_{i,j} L \leq 1$. In particular, the above shows that AdaRem
17 is monotonic: the objective value is non-decreasing. Combining the above two equations we get,

$$\begin{aligned}
f(x_{i+1}) &\leq f(x^*) + \langle g_i, x_i - x^* \rangle - \sum_{j=1}^d \frac{t_{i,j}}{2} g_{i,j}^2 \\
&= f(x^*) + \sum_{j=1}^d \frac{(x_i - x^*)_j^2}{2t_{i,j}} - \left[\sum_{j=1}^d \frac{(x_i - x^*)_j^2}{2t_{i,j}} - \langle g_i, x_i - x^* \rangle + \sum_{j=1}^d \frac{t_{i,j}}{2} g_{i,j}^2 \right] \\
&= f(x^*) + \sum_{j=1}^d \frac{(x_i - x^*)_j^2}{2t_{i,j}} - \sum_{j=1}^d \frac{1}{2t_{i,j}} \left[(x_i - x^*)_j^2 - 2t_{i,j}g_{i,j}(x_i - x^*)_j + t_{i,j}^2 g_{i,j}^2 \right] \\
&= f(x^*) + \sum_{j=1}^d \frac{(x_i - x^*)_j^2}{2t_{i,j}} - \sum_{j=1}^d \frac{((x_i - x^*)_j - t_{i,j}g_{i,j})^2}{2t_{i,j}} \\
&= f(x^*) + \sum_{j=1}^d \frac{(x_i - x^*)_j^2 - (x_{i+1} - x^*)_j^2}{2t_{i,j}} \\
&\leq f(x^*) + \frac{1}{2c\eta} \left(\|x_i - x^*\|_2^2 - \|x_{i+1} - x^*\|_2^2 \right). \tag{3}
\end{aligned}$$

18 where the last inequality follows as $t_{i,j} \geq c\eta$. Summing the above equations for $i = 0, \dots, k-1$, we
19 get

$$\sum_{i=0}^{k-1} (f(x_{i+1}) - f(x^*)) \leq \frac{1}{2c\eta} \left(\|x_0 - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right) \leq \frac{\|x_0 - x^*\|_2^2}{2c\eta}$$

20 Finally, by Equation (2), $f(x_0), \dots, f(x_k)$ is non-increasing. Therefore, $f(x_k) - f(x^*) \leq f(x_i) -$
21 $f(x^*)$ for all $i < k$. Thus

$$k \cdot (f(x_k) - f(x^*)) \leq \frac{\|x_0 - x^*\|_2^2}{2c\eta}.$$

22 The theorem now follows. □