

A Proofs

Our main results utilize an elementary fact about smooth functions with Lipschitz continuous gradient, called the co-coercivity of the gradient. We state the lemma and recall its proof for completeness.

A.1 The Co-coercivity Lemma

Lemma A.1 (Co-coercivity) *For a smooth function f whose gradient has Lipschitz constant L ,*

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

Proof. Since ∇f has Lipschitz constant L , if \mathbf{x}_* is the minimizer of f , then

$$\frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}_*)\|_2^2 = \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}_*)\|_2^2 + \langle \mathbf{x} - \mathbf{x}_*, \nabla f(\mathbf{x}_*) \rangle \leq f(\mathbf{x}) - f(\mathbf{x}_*); \quad (\text{A.1})$$

see, for instance, [[13], page 26]. Now define the convex functions

$$G(\mathbf{z}) = f(\mathbf{z}) - \langle \nabla f(\mathbf{x}), \mathbf{z} \rangle, \quad \text{and} \quad H(\mathbf{z}) = f(\mathbf{z}) - \langle \nabla f(\mathbf{y}), \mathbf{z} \rangle,$$

and observe that both have Lipschitz constants L and minimizers \mathbf{x} and \mathbf{y} , respectively. Applying (A.1) to these functions therefore gives that

$$G(\mathbf{x}) \leq G(\mathbf{y}) - \frac{1}{2L} \|\nabla G(\mathbf{y})\|_2^2, \quad \text{and} \quad H(\mathbf{y}) \leq H(\mathbf{x}) - \frac{1}{2L} \|\nabla H(\mathbf{y})\|_2^2.$$

By their definitions, this implies that

$$\begin{aligned} f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{x} \rangle &\leq f(\mathbf{y}) - \langle \nabla f(\mathbf{x}), \mathbf{y} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2^2 \\ f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{y} \rangle &\leq f(\mathbf{x}) - \langle \nabla f(\mathbf{y}), \mathbf{x} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2. \end{aligned}$$

Adding these two inequalities and canceling terms yields the desired result. \square

A.2 Proof of Theorem 2.1

With the notation of Theorem 2.1, and where i is the random index chosen at iteration k , we have

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 &= \|\mathbf{x}_k - \mathbf{x}_* - \gamma \nabla f_i(\mathbf{x}_k)\|_2^2 \\ &= \|(\mathbf{x}_k - \mathbf{x}_*) - \gamma (\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)) - \gamma \nabla f_i(\mathbf{x}_*)\|_2^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma \langle \mathbf{x}_k - \mathbf{x}_*, \nabla f_i(\mathbf{x}_k) \rangle + \\ &\quad \gamma^2 \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*) + \nabla f_i(\mathbf{x}_*)\|_2^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma \langle \mathbf{x}_k - \mathbf{x}_*, \nabla f_i(\mathbf{x}_k) \rangle + \\ &\quad 2\gamma^2 \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)\|_2^2 + 2\gamma^2 \|\nabla f_i(\mathbf{x}_*)\|_2^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma \langle \mathbf{x}_k - \mathbf{x}_*, \nabla f_i(\mathbf{x}_k) \rangle \\ &\quad + 2\gamma^2 L_i \langle \mathbf{x}_k - \mathbf{x}_*, \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*) \rangle + 2\gamma^2 \|\nabla f_i(\mathbf{x}_*)\|_2^2, \end{aligned}$$

where we have employed Jensen's inequality in the first inequality and the co-coercivity Lemma A.1 in the final line. We next take an expectation with respect to the choice of i . By assumption, $i \sim \mathcal{D}$ such that $F(\mathbf{x}) = \mathbb{E} f_i(\mathbf{x})$ and $\sigma^2 = \mathbb{E} \|\nabla f_i(\mathbf{x}_*)\|^2$. Then $\mathbb{E} \nabla f_i(\mathbf{x}) = \nabla F(\mathbf{x})$, and we obtain:

$$\begin{aligned} \mathbb{E} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 &\leq \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma \langle \mathbf{x}_k - \mathbf{x}_*, \nabla F(\mathbf{x}_k) \rangle \\ &\quad + 2\gamma^2 \mathbb{E} [L_i \langle \mathbf{x}_k - \mathbf{x}_*, \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*) \rangle] + 2\gamma^2 \mathbb{E} \|\nabla f_i(\mathbf{x}_*)\|_2^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma \langle \mathbf{x}_k - \mathbf{x}_*, \nabla F(\mathbf{x}_k) \rangle \\ &\quad + 2\gamma^2 \sup_i L_i \mathbb{E} \langle \mathbf{x}_k - \mathbf{x}_*, \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*) \rangle + 2\gamma^2 \mathbb{E} \|\nabla f_i(\mathbf{x}_*)\|_2^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma \langle \mathbf{x}_k - \mathbf{x}_*, \nabla F(\mathbf{x}_k) \rangle \\ &\quad + 2\gamma^2 \sup L \langle \mathbf{x}_k - \mathbf{x}_*, \nabla F(\mathbf{x}_k) - \nabla F(\mathbf{x}_*) \rangle + 2\gamma^2 \sigma^2 \end{aligned}$$

We now utilize the strong convexity of $F(\mathbf{x})$ and obtain that

$$\begin{aligned} &\leq \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 - 2\gamma\mu(1 - \gamma \sup L)\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 + 2\gamma^2\sigma^2 \\ &= (1 - 2\gamma\mu(1 - \gamma \sup L))\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 + 2\gamma^2\sigma^2 \end{aligned}$$

when $\gamma\mu \leq 1$. Recursively applying this bound over the first k iterations yields the desired result,

$$\begin{aligned} \mathbb{E}\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 &\leq \left(1 - 2\gamma\mu(1 - \gamma \sup L)\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + 2 \sum_{j=0}^{k-1} \left(1 - 2\gamma\mu(1 - \gamma \sup L)\right)^j \gamma^2\sigma^2 \\ &\leq \left(1 - 2\gamma\mu(1 - \gamma \sup L)\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + \frac{\gamma^2\sigma^2}{\mu(1 - \gamma \sup L)}. \end{aligned}$$

We next turn to the second part of the theorem, where we optimize the step size γ for a fixed tolerance ε . Recall the main recursive step in the previous proof,

$$\mathbb{E}\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 \leq (1 - 2\mu\gamma(1 - \gamma \sup L))\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 + 2\gamma^2\sigma^2, \quad (\text{A.2})$$

which is valid as long as $\mu\gamma \leq 1$. The minimal value of the quadratic

$$F_\xi(\gamma) = (1 - 2\gamma\mu(1 - \gamma \sup L))\xi + 2\sigma^2\gamma^2$$

is achieved at

$$\gamma_\xi^* = \frac{\mu\xi}{2\xi\mu \sup L + 2\sigma^2}, \quad (\text{A.3})$$

and

$$F_\xi(\gamma_\xi^*) = \left(1 - \frac{\mu^2\xi}{2\mu \sup L\xi + 2\sigma^2}\right)\xi \quad (\text{A.4})$$

Note that because $\sup L/\mu \geq 1$, it follows that $\mu\gamma_\xi^* \leq 1/2$. Thus if we choose step-size $\gamma^* = \gamma_\xi^*$,

$$\mathbb{E}\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 \leq F_{\|\mathbf{x}_k - \mathbf{x}_*\|_2^2}(\gamma^*) \quad (\text{A.5})$$

$$= \left(F_{\|\mathbf{x}_k - \mathbf{x}_*\|_2^2}(\gamma^*) - F_\varepsilon(\gamma^*)\right) + F_\varepsilon(\gamma^*) \quad (\text{A.6})$$

$$\leq \left(1 - \frac{\mu^2\varepsilon}{2\mu\varepsilon \sup L + 2\sigma^2}\right)\|\mathbf{x}_k - \mathbf{x}_*\|_2^2. \quad (\text{A.7})$$

(A.8)

Iterating the expectation,

$$\mathbb{E}\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2 \leq \left(1 - \frac{\mu^2\varepsilon}{2\mu\varepsilon \sup L + 2\sigma^2}\right)^k \varepsilon_0. \quad (\text{A.9})$$

It follows that if $\varepsilon \leq \mathbb{E}\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2^2$, then

$$\log(\varepsilon/\varepsilon_0) \leq k \log\left(1 - \frac{\mu^2\varepsilon}{2\mu\varepsilon \sup L + 2\sigma^2}\right) \quad (\text{A.10})$$

$$\leq -k\left(\frac{\mu^2\varepsilon}{2\mu\varepsilon \sup L + 2\sigma^2}\right) \quad (\text{A.11})$$

or, equivalently

$$k \leq \log(\varepsilon_0/\varepsilon)\left(\frac{2\mu\varepsilon \sup L + 2\sigma^2}{\mu^2\varepsilon}\right) \quad (\text{A.12})$$

$$= \log(\varepsilon_0/\varepsilon)\left(\frac{2\sup L}{\mu} + \frac{2\sigma^2}{\mu^2\varepsilon}\right). \quad (\text{A.13})$$

□