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

Proof. Since ∇f has Lipschitz constant L, if x_* is the minimizer of f, then

$$
\frac{1}{2L} \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}_\star)\|_2^2 = \frac{1}{2L} \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}_\star)\|_2^2 + \langle \boldsymbol{x} - \boldsymbol{x}_\star, \nabla f(\boldsymbol{x}_\star) \rangle \le f(\boldsymbol{x}) - f(\boldsymbol{x}_\star); \tag{A.1}
$$

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

$$
G(z) = f(z) - \langle \nabla f(x), z \rangle
$$
, and $H(z) = f(z) - \langle \nabla f(y), z \rangle$,

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

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

By their definitions, this implies that

$$
f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{x} \rangle \le f(\mathbf{y}) - \langle \nabla f(\mathbf{x}), \mathbf{y} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2^2
$$

$$
f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{y} \rangle \le f(\mathbf{x}) - \langle \nabla f(\mathbf{y}), \mathbf{x} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2.
$$

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

 \Box

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

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

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 $F(\bm{x}) = \mathbb{E} f_i(\bm{x})$ and $\sigma^2 = \mathbb{E} \|\nabla f_i(\bm{x}_\star)\|^2$. Then $\mathbb{E} \nabla f_i(\bm{x}) = \nabla F(\bm{x})$, and we obtain:

$$
\mathbb{E} \|x_{k+1} - x_{\star}\|_2^2 \le \|x_k - x_{\star}\|_2^2 - 2\gamma \langle x_k - x_{\star}, \nabla F(x_k) \rangle \n+ 2\gamma^2 \mathbb{E} [L_i \langle x_k - x_{\star}, \nabla f_i(x_k) - \nabla f_i(x_{\star}) \rangle] + 2\gamma^2 \mathbb{E} \|\nabla f_i(x_{\star})\|_2^2 \n\le \|x_k - x_{\star}\|_2^2 - 2\gamma \langle x_k - x_{\star}, \nabla F(x_k) \rangle \n+ 2\gamma^2 \sup_i L_i \mathbb{E} \langle x_k - x_{\star}, \nabla f_i(x_k) - \nabla f_i(x_{\star}) \rangle + 2\gamma^2 \mathbb{E} \|\nabla f_i(x_{\star})\|_2^2 \n= \|x_k - x_{\star}\|_2^2 - 2\gamma \langle x_k - x_{\star}, \nabla F(x_k) \rangle \n+ 2\gamma^2 \sup L \langle x_k - x_{\star}, \nabla F(x_k) - \nabla F(x_{\star}) \rangle + 2\gamma^2 \sigma^2
$$

We now utilize the strong convexity of $F(x)$ and obtain that

$$
\leq \|\boldsymbol{x}_k - \boldsymbol{x}_\star\|_2^2 - 2\gamma\mu(1-\gamma\sup L)\|\boldsymbol{x}_k - \boldsymbol{x}_\star\|_2^2 + 2\gamma^2\sigma^2
$$

=
$$
(1 - 2\gamma\mu(1-\gamma\sup L))\|\boldsymbol{x}_k - \boldsymbol{x}_\star\|_2^2 + 2\gamma^2\sigma^2
$$

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

$$
\mathbb{E} \|\bm{x}_{k} - \bm{x}_{\star}\|_{2}^{2} \leq \left(1 - 2\gamma\mu(1 - \gamma \sup L)\right)\right)^{k} \|\bm{x}_{0} - \bm{x}_{\star}\|_{2}^{2} + 2\sum_{j=0}^{k-1} \left(1 - 2\gamma\mu(1 - \gamma \sup L)\right)\right)^{j} \gamma^{2} \sigma^{2}
$$

$$
\leq \left(1 - 2\gamma\mu(1 - \gamma \sup L)\right)\right)^{k} \|\bm{x}_{0} - \bm{x}_{\star}\|_{2}^{2} + \frac{\gamma\sigma^{2}}{\mu(1 - \gamma \sup L)}.
$$

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,

$$
\mathbb{E} \|x_{k+1} - x_{\star}\|_{2}^{2} \leq (1 - 2\mu\gamma(1 - \gamma \sup L)) \|x_{k} - x_{\star}\|_{2}^{2} + 2\gamma^{2}\sigma^{2},
$$
 (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

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

and

$$
F_{\xi}(\gamma_{\xi}^{*}) = \left(1 - \frac{\mu^{2}\xi}{2\mu \sup L\xi + 2\sigma^{2}}\right)\xi
$$
 (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 $\gamma^* = \gamma_{\epsilon}^*$,

$$
\mathbb{E} \|x_{k+1} - x_{\star}\|_{2}^{2} \le F_{\|x_{k} - x_{\star}\|_{2}^{2}}(\gamma^{*})
$$
\n(A.5)

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

$$
\leq \left(1 - \frac{\mu^2 \varepsilon}{2\mu\varepsilon \sup L + 2\sigma^2}\right) \|x_k - x_\star\|_2^2. \tag{A.7}
$$

Iterating the expectation,

$$
\mathbb{E}||x_{k+1}-x_{\star}||_2^2 \le \left(1-\frac{\mu^2 \varepsilon}{2\mu\varepsilon \sup L+2\sigma^2}\right)^k \varepsilon_0.
$$
 (A.9)

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

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

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

or, equivalently

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

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

 \Box