

Stochastic gradient descent proof

Shuvomoy Das Gupta

December 25, 2022

We study a simple proof of stochastic gradient descent.

Algorithm description

Problem setup. We are interested in solving the problem

$$p^* = \left(\begin{array}{ll} \underset{x \in \mathbb{R}^d}{\text{minimize}} & f(x) \\ \text{subject to} & x \in C, \end{array} \right) \quad (\mathcal{P})$$

Notation

Π_C : projection onto the set C

where we have the following assumptions regarding the nature of the problem.

Assumption 1 (Structure of \mathcal{P}). We assume:

- $f : \mathbb{R}^d \rightarrow (-\infty, \infty]$ is a closed, proper, and convex function,
- C is a nonempty, closed, convex set, with $C \subseteq \text{intdom } f$, and
- $\text{argmin} f(C) = X^* \neq \emptyset$.

Stochastic gradient descent. The stochastic gradient descent (SGD) algorithm to solve (\mathcal{P}) is described by Algorithm 1, where we make the following assumption regarding the nature of the oracle.

Algorithm 1 SGD to solve (\mathcal{P})

input: f, C , iteration limit K

algorithm:

1. initialization:

pick $x_0 \in C$ arbitrarily

2. main iteration:

for $k = 0, 1, 2, \dots, K - 1$

 pick stepsizes $\alpha_k > 0$ and random $g_k \in \mathbb{R}^d$ satisfying Assumption 2

$x_{k+1} \leftarrow \Pi_C(x_k - \alpha_k g_k)$

end for

3. return x_K

Assumption 2 (Stochastic oracle). We assume that given an iterate x_k , the stochastic oracle is capable of producing a random vector g_k with the following properties:

- (unbiased) $\forall_{k \geq 0} \mathbf{E}[g_k \mid x_k] \in \partial f(x_k)$, and
- (bounded variance) $\exists_{G > 0} \forall_{k \geq 0} \mathbf{E}[\|g_k\|^2 \mid x_k] \leq G^2$.

Convergence analysis. First, note that, for all $k \geq 0$:

$$\begin{aligned}
& \mathbf{E} \left[\left\| \underbrace{\Pi_C(x_k - \alpha_k g_k)}_{x_{k+1}} - \underbrace{\Pi_C(x_*)}_{x_*} \right\|^2 \mid x_k \right] \\
&= \mathbf{E} \left[\underbrace{\left\| \Pi_C(x_k - \alpha_k g_k) - \Pi_C(x_*) \right\|^2}_{\leq \|x_k - \alpha_k g_k - x_*\|^2} \mid x_k \right] \\
&\leq \mathbf{E} \left[\underbrace{\|x_k - \alpha_k g_k - x_*\|^2}_{= \|x_k - x_*\|^2 + \alpha_k^2 \|g_k\|^2 - 2\alpha_k \langle x_k - x_*, g_k \rangle} \mid x_k \right] \\
&= \mathbf{E} \left[\|x_k - x_*\|^2 + \alpha_k^2 \|g_k\|^2 - 2\alpha_k \langle x_k - x_*, g_k \rangle \mid x_k \right] \\
&= \mathbf{E} \left[\|x_k - x_*\|^2 \mid x_k \right] + \alpha_k^2 \mathbf{E} \left[\|g_k\|^2 \mid x_k \right] - 2\alpha_k \mathbf{E} \left[\langle x_k - x_*, g_k \rangle \mid x_k \right] \quad \triangleright \text{using linearity of expectation} \\
&= \|x_k - x_*\|^2 + \alpha_k^2 \mathbf{E} \left[\|g_k\|^2 \mid x_k \right] - 2\alpha_k \langle x_k - x_*, \mathbf{E} [g_k \mid x_k] \rangle \quad \triangleright \text{using "taking out what's known" rule} \\
&\leq \|x_k - x_*\|^2 + \alpha_k^2 G^2 - 2\alpha_k \langle x_k - x_*, \mathbf{E} [g_k \mid x_k] \rangle \quad \mathbf{E} [h(X)Y \mid X] = h(X)\mathbf{E} [Y \mid X] \\
&\quad / * \\
&\quad \text{we have } \mathbf{E} \left[\|g_k\|^2 \mid x_k \right] \leq G^2 \\
&\quad \Leftrightarrow \forall y, f(y) \geq f(x_k) + \langle \mathbf{E} [g_k \mid x_k]; y - x_k \rangle \\
&\quad \xRightarrow{y \leftarrow x_*} f(x_*) \geq f(x_k) - \langle \mathbf{E} [g_k \mid x_k]; x_k - x_* \rangle \\
&\quad \Rightarrow -\langle \mathbf{E} [g_k \mid x_k]; x_k - x_* \rangle \leq f(x_*) - f(x_k) \\
&\quad */ \\
&= \|x_k - x_*\|^2 + \alpha_k^2 G^2 - 2\alpha_k (f(x_k) - f(x_*)),
\end{aligned}$$

So, we have proved

$$\mathbf{E} \left[\|x_{k+1} - x_*\|^2 \mid x_k \right] \leq \|x_k - x_*\|^2 + \alpha_k^2 G^2 - 2\alpha_k (f(x_k) - f(x_*)),$$

so taking expectation with respect to x_k on both sides, we get:

$$\begin{aligned}
& \mathbf{E} \left[\mathbf{E} \left[\|x_{k+1} - x_*\|^2 \mid x_k \right] \right] \\
&= \mathbf{E} \left[\|x_{k+1} - x_*\|^2 \right] \quad \triangleright \text{using Adam's law } \mathbf{E} [\mathbf{E} [Y \mid X]] = \mathbf{E} [Y] \\
&\leq \mathbf{E} \left[\|x_k - x_*\|^2 + \alpha_k^2 G^2 - 2\alpha_k (f(x_k) - f(x_*)) \right] \\
&= \mathbf{E} \left[\|x_k - x_*\|^2 \right] - 2\alpha_k \mathbf{E} [f(x_k) - f(x_*)] + \alpha_k^2 G^2,
\end{aligned}$$

so

$$\mathbf{E} \left[\|x_{k+1} - x_*\|^2 \right] - \mathbf{E} \left[\|x_k - x_*\|^2 \right] \leq -2\alpha_k \mathbf{E} [f(x_k) - f(x_*)] + \alpha_k^2 G^2.$$

Now, let us do a telescoping sum:

$$\begin{aligned}
\mathbf{E} \left[\|x_{k+1} - x_\star\|^2 \right] - \mathbf{E} \left[\|x_k - x_\star\|^2 \right] &\leq -2\alpha_k \mathbf{E} [f(x_k) - f(x_\star)] + \alpha_k^2 G^2 \\
\mathbf{E} \left[\|x_k - x_\star\|^2 \right] - \mathbf{E} \left[\|x_{k-1} - x_\star\|^2 \right] &\leq -2\alpha_k \mathbf{E} [f(x_{k-1}) - f(x_\star)] + \alpha_{k-1}^2 G^2 \\
&\vdots \\
\mathbf{E} \left[\|x_{m+1} - x_\star\|^2 \right] - \mathbf{E} \left[\|x_m - x_\star\|^2 \right] &\leq -2\alpha_m \mathbf{E} [f(x_m) - f(x_\star)] + \alpha_m^2 G^2,
\end{aligned}$$

and adding the equations above, we get:

$$\begin{aligned}
\mathbf{E} \left[\|x_{k+1} - x_\star\|^2 \right] - \mathbf{E} \left[\|x_m - x_\star\|^2 \right] &\leq -2 \sum_{i=m}^k \alpha_i \mathbf{E} [f(x_i) - f(x_\star)] + G^2 \sum_{i=m}^k \alpha_i^2 \\
\Leftrightarrow 0 &\leq \mathbf{E} \left[\|x_{k+1} - x_\star\|^2 \right] \leq \mathbf{E} \left[\|x_m - x_\star\|^2 \right] - 2 \sum_{i=m}^k \alpha_i \mathbf{E} [f(x_i) - f(x_\star)] + G^2 \sum_{i=m}^k \alpha_i^2 \\
\Rightarrow 0 &\leq \mathbf{E} \left[\|x_m - x_\star\|^2 \right] - 2 \sum_{i=m}^k \alpha_i \mathbf{E} [f(x_i) - f(x_\star)] + G^2 \sum_{i=1}^m \alpha_i^2 \\
\Leftrightarrow \sum_{i=m}^k \alpha_i \mathbf{E} [f(x_i) - f(x_\star)] &\leq \frac{1}{2} \left(\mathbf{E} \left[\|x_m - x_\star\|^2 \right] + G^2 \sum_{i=m}^k \alpha_i^2 \right) \\
\Rightarrow \left(\sum_{i=m}^k \alpha_i \right) \left(\min_{i \in \{m, \dots, k\}} \mathbf{E} [f(x_i) - f(x_\star)] \right) &\leq \frac{1}{2} \left(\mathbf{E} \left[\|x_m - x_\star\|^2 \right] + G^2 \sum_{i=m}^k \alpha_i^2 \right) \\
&\quad \triangleright \text{for } b_k \geq 0, \text{ we have } (\min_k a_k) \sum_k b_k \leq \sum_k a_k b_k \\
\Rightarrow \mathbf{E} \left[\min_{i \in \{m, \dots, k\}} \{f(x_i) - f(x_\star)\} \right] &\leq \min_{i \in \{m, \dots, k\}} \mathbf{E} [f(x_i) - f(x_\star)] \\
&\leq \frac{\mathbf{E} \left[\|x_m - x_\star\|^2 \right] + G^2 \sum_{i=m}^k \alpha_i^2}{2 \sum_{i=m}^k \alpha_i}. \\
&\quad \triangleright \text{using } \mathbf{E} [\min_i X_i] \leq \min_i \mathbf{E} [X_i]
\end{aligned}$$

In the last inequality, m is arbitrary, so set $m \leftarrow 0$, which leads to:

$$\mathbf{E} \left[\min_{i \in \{0, \dots, k\}} f(x_i) \right] - f(x_\star) \leq \frac{\mathbf{E} \left[\|x_0 - x_\star\|^2 \right] + G^2 \sum_{i=0}^k \alpha_i^2}{2 \sum_{i=0}^k \alpha_i},$$

so if we have $\sum_{i=0}^k \alpha_i^2 < \infty$ and $\sum_{i=0}^k \alpha_i = \infty$, then we have

$$\mathbf{E} \left[\min_{i \in \{0, \dots, k\}} f(x_i) \right] \rightarrow f(x_\star).$$