

Stochastic Gradient Methods

My "Moreau"
dog
 $f_{\mathcal{R}}(x) = \min f(x) + \frac{1}{2\alpha} \|x - x^*\|^2$
↑ Moreau Envelope

Previously gradient oracle

$$x \mapsto \nabla f(x)$$

Now we have stochastic gradient oracle

$$x \mapsto \underbrace{g(x)}_{\text{random variable/vector}} \quad \text{"randomized independently each time we call"}$$

such that

$$\mathbb{E}[g(x)] = \nabla f(x) \quad \text{"unbiased estimator"}$$

$$\text{var}[g(x)] \leq \sigma^2$$

$$\begin{aligned} &:= \mathbb{E}[\|g(x) - \nabla f(x)\|_2^2] \\ &= \mathbb{E}[\|g(x)\|_2^2] - \|\mathbb{E}g(x)\|_2^2 = \sigma^2 \end{aligned}$$

Brief Probability Review:

Linearity of Expectation: For random variables X_1, \dots, X_n and constants $\lambda_1, \dots, \lambda_n$

$$\mathbb{E}[\lambda_1 X_1 + \dots + \lambda_n X_n] = \lambda_1 \mathbb{E}[X_1] + \dots + \lambda_n \mathbb{E}[X_n]$$

Law of Total Expectation:

$$\mathbb{E}_{x_k}[\mathbb{E}[f(x_{k+1}) | x_k]] = \mathbb{E}[f(x_{k+1})]$$

Outline

1. Example of Stochastic Gradient Oracles
2. Nonconvex Guarantees: $\mathbb{E} \|\nabla f(x_k)\|^2$ small
3. Convex Guarantees: $\mathbb{E} f(x_k) - f(x^*)$ small
4. Improvements: Coordinate Methods,
Acceleration,
~~Variance~~ Variance Reduction.

1. Examples

Example 1 "Coordinate Approach", $\min f(x)$, $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Pick $i \in \{1, \dots, d\}$ uniformly at random

Set $g(x) = d \cdot \frac{\partial f}{\partial x_i}(x) \cdot e_i$
 \uparrow i -th basis vector

(Check unbiased oracle)

$$\begin{aligned} \mathbb{E}[g(x)] &= \frac{1}{d} \sum_i d \cdot \frac{\partial f}{\partial x_i}(x) e_i \\ &= \sum_i \frac{\partial f}{\partial x_i}(x) e_i = \nabla f(x). \end{aligned}$$

Example 2 "Finite Sum"

$$\min f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$$

Pick $i \in \{1, \dots, n\}$ uniformly

Set $g(x) = \nabla f_i(x)$

(Check unbiased oracle).

$$\min \frac{1}{n} \|Ax - b\|_2^2 = \frac{1}{n} \sum_{i=1}^n (a_i^T x - b_i)^2$$

Each (a_i, b_i) is a data point

a_i feature vector pictures of dogs or not

b_i measurement.

Example 3 (or 2.1). Infinite Sum / Expectation

There is a distribution of data points

$$(a_i, b_i) \sim \mathcal{D}$$

$$\min \mathbb{E}_{(a_i, b_i)} [(a_i^T x - b_i)^2]$$

$$\begin{aligned} (\text{Linearity Expectations}) &\Rightarrow \mathbb{E} \nabla f(x, a_i, b_i) \\ &= \nabla \mathbb{E} f(x, a_i, b_i) \end{aligned}$$

Example 4 (or 2.2) Better (Lower Variance) Oracles for finite sums.

Idea 1: Look at batches/minibatches of samples at each step

Pick $S \subseteq \{1, \dots, n\}$ with $|S| = k$
uniformly at random
with or without replacement

$$g(x) = \frac{1}{k} \sum_{i \in S} \nabla f_i(x).$$

Idea 2: Variance Reduction, given \tilde{x}

Compute the full gradient $\nabla f(\tilde{x})$
 $= \frac{1}{n} \sum \nabla f_i(\tilde{x})$

Pick $i \in \{1, \dots, n\}$ uniformly

$$g(x) = \nabla f(\tilde{x}) + \underbrace{\nabla f_i(x) - \nabla f_i(\tilde{x})}_{\text{small when } \underline{x - \tilde{x}} \text{ small}}$$

(Check unbiased:

$$\begin{aligned} \mathbb{E}[g(x)] &= \nabla f(\tilde{x}) + \mathbb{E}[\nabla f_i(x)] - \mathbb{E}[\nabla f_i(\tilde{x})] \\ &\quad \underbrace{\qquad \qquad \qquad \nabla f(x) \qquad \qquad \nabla f(\tilde{x}) \qquad \qquad}_{\text{Cancel}} \\ &= \nabla f(x). \end{aligned}$$

SVRG [Johnson, Zhang, 2013].

2. Nonconvex Stochastic Gradient Method Analysis

Consider $x_{k+1} = x_k - \alpha_k \underline{g(x_k)}$
 \uparrow independent at each iteration.

Aside: Not a descent method. For example,

$$g(x) = \begin{cases} 3 \nabla f(x) & \text{with prob } 1/2 \\ -\nabla f(x) & \text{with prob } 1/2 \end{cases}$$



Theorem Suppose f has L -Lipschitz gradient and $g(x)$ is unbiased estimator of $\nabla f(x)$ with variance σ^2 . Then for any $0 \leq \alpha_k \leq \frac{2}{L}$,

$$\mathbb{E} \left[\min_{i \leq k} \|\nabla f(x_i)\|^2 \right] \leq \frac{f(x_0) - \min f + \frac{\sigma^2 L}{2} \sum_{i=1}^k \alpha_i^2}{\sum_{i=1}^k \alpha_i \left(1 - \frac{L \alpha_i}{2} \right)}.$$

Proof. Our Taylor Approximation Theorem ensures

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2 \\ &= f(x_k) - \alpha_k \nabla f(x_k)^T g(x_k) + \frac{L \alpha_k^2}{2} \|g(x_k)\|^2 \end{aligned}$$

Fixing x_k (just consider randomness in step k),

$$\begin{aligned} \mathbb{E}[f(x_{k+1}) | x_k] &\stackrel{\text{Linearity}}{\leq} f(x_k) - \alpha_k \mathbb{E}[\nabla f(x_k)^T g(x_k)] + \frac{L \alpha_k^2}{2} \mathbb{E}[\|g(x_k)\|^2 | x_k] \\ &\stackrel{\text{Linearity}}{=} f(x_k) - \alpha_k \nabla f(x_k)^T \mathbb{E}[g(x_k) | x_k] \\ &\quad + \frac{L \alpha_k^2}{2} \mathbb{E}[\|g(x_k)\|^2 | x_k] \\ &\stackrel{\text{by unbiased and variance bound}}{\leq} f(x_k) - \alpha_k \|\nabla f(x_k)\|_2^2 \\ &\quad + \frac{L \alpha_k^2}{2} (\sigma^2 + \|\nabla f(x_k)\|^2) \\ &= f(x_k) - \left(\alpha_k - \frac{L \alpha_k^2}{2} \right) \|\nabla f(x_k)\|^2 + \frac{L \alpha_k^2 \sigma^2}{2} \end{aligned}$$

By law of total expectation, we have on the overall result of this stochastic process.

$$\mathbb{E}[f(x_{k+1})] \leq \mathbb{E}[f(x_k)] - (\alpha_k - \frac{L\alpha_k^2}{2}) \mathbb{E}[\|\nabla f(x_k)\|^2] + \frac{L\alpha_k^2\sigma^2}{2}$$

Induction on this ensures

$$\min f(x) \leq \mathbb{E}[f(x_{k+1})] \leq \underbrace{\mathbb{E}[f(x_0)]}_{f(x_0)} - \sum_{i=0}^k (\alpha_i - \frac{L\alpha_i^2}{2}) \mathbb{E}[\|\nabla f(x_i)\|^2] + \sum_{i=0}^k \frac{L\alpha_i^2\sigma^2}{2}$$

$$\stackrel{\text{Linearity}}{\Rightarrow} \mathbb{E}\left[\underbrace{\sum_i \frac{\alpha_i (1 - \frac{L\alpha_i}{2}) \|\nabla f(x_i)\|^2}{\sum_j \alpha_j (1 - \frac{L\alpha_j}{2})}}_{\text{a weighted average of the gradient norm squared}}\right] \leq \frac{f(x_0) - \min f + \sum \frac{L\alpha_i^2\sigma^2}{2}}{\sum \alpha_i (1 - \frac{L\alpha_i}{2})}$$

$$\Rightarrow \mathbb{E}\left[\min_{i \leq k} \|\nabla f(x_i)\|^2\right] \leq \frac{f(x_0) - \min f + \sum \frac{L\alpha_i^2\sigma^2}{2}}{\sum \alpha_i (1 - \frac{L\alpha_i}{2})} \quad \square$$

Next time this is a $O(\frac{1}{\sqrt{k}})$ rate.