

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

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

Induction on this ensures

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

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

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

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

Pick $\alpha_k = \frac{1}{L\sqrt{k+1}}$, $\Rightarrow 1 - \frac{L\alpha_k}{2} \geq \frac{1}{2}$, then RHS $\leq \frac{f(x_0) - \min f + \frac{\sigma^2}{2L}}{\frac{1}{2}(T+1) \frac{1}{L\sqrt{T+1}}} = \frac{L(f(x_0) - \min f) + \frac{\sigma^2}{2}}{\frac{1}{2}\sqrt{T+1}}$

$$\Rightarrow \mathbb{E} \left[\min_{i \leq K} \|\nabla f(x_i)\| \right] \leq O\left(\frac{1}{K^{1/4}}\right).$$

Hope for better guarantees under convexity.

3. Convex Guarantees

For ease, let's look at $\alpha_k = \alpha \leq 1/L$.

Theorem In addition to the previous theorem, suppose f is convex. Then

$$\mathbb{E} \left[\min_{i \leq K+1} \{f(x_i) - f(x^*)\} \right] \leq \frac{\|x_0 - x^*\|^2}{2\alpha(K+1)} + \alpha\sigma^2$$

In particular, $\alpha = 1/\sqrt{K+1}$, for $K \geq L^2$

$$\Rightarrow \text{RHS} = \frac{\|x_0 - x^*\|^2 + 2\sigma^2}{2\sqrt{K+1}}.$$

Proof. By previous proof, we have a "stochastic" descent lemma

$$\mathbb{E}[f(x_{k+1}) | x_k] \leq f(x_k) - \frac{\alpha}{2} \underbrace{\|\nabla f(x_k)\|^2}_{\substack{\leq f(x^*) - \nabla f(x_k)^T (x^* - x_k) \\ = \mathbb{E}[\|g(x_k)\|^2 | x_k] - \sigma^2}} + \frac{\alpha\sigma^2}{2}.$$

$$\leq \underbrace{f(x^*) - \nabla f(x_k)^T (x^* - x_k)}_{\substack{\leq f(x^*) - \nabla f(x_k)^T (x^* - x_k) \\ = \mathbb{E}[\|g(x_k)\|^2 | x_k] - \sigma^2}} - \frac{\alpha}{2} \mathbb{E}[\|g(x_k)\|^2 | x_k] + \alpha\sigma^2.$$

by unbiased
and using
Linearity.

$$= f(x^*) - \mathbb{E} \left[\alpha g(x_k)^T (x^* - x_k) + \frac{\alpha}{2} \|g(x_k)\|^2 \mid x_k \right] + \alpha \sigma^2$$

by def
of x_{k+1}

$$= f(x^*) - \mathbb{E} \left[g(x_k)^T (x^* - x_k) + \frac{1}{2\alpha} \|x_{k+1} - x_k\|^2 \mid x_k \right] + \alpha \sigma^2$$

$$= f(x^*) - \mathbb{E} \left[\frac{1}{2\alpha} \left(\|x_{k+1} - x^*\|^2 - \|x_k - x^*\|^2 \right) \mid x_k \right] + \alpha \sigma^2$$

" $\|x_k - x^*\|^2 - \alpha \|g(x_k)\|^2$

Law of total expectation:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \mathbb{E} \left[\frac{1}{2\alpha} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) \right] + \alpha \sigma^2$$

Summing these up and dividing by K gives

$$\frac{1}{K} \sum \mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \frac{\|x_0 - x^*\|^2}{2\alpha \cdot K} + \alpha \sigma^2. \quad \square$$

Note nonsmooth (but deterministic) also $\frac{1}{\sqrt{K}}$ rate.

Show nonsmooth stochastic opt has $\frac{1}{\sqrt{K}}$ rate.

HW 4 Q1.

(there went $g(x)$ s.t. $\mathbb{E}[g(x)] \in \partial f(x)$)

4. Improvements

Acceleration? Not really.

$$O\left(\frac{L\|x_0 - x^*\|^2}{k^2} + \frac{\sigma^2}{\sqrt{k}}\right) \quad \text{best possible.}$$



Randomized Coordinate Descent

Fix $g(x) = \underbrace{d}_{\sim \text{Uniformly}} \cdot \frac{\partial f}{\partial x_i}(x) \cdot e_i$, $i \sim \{1, \dots, d\}$

$$x_{k+1} = x_k - \frac{1}{L_d} g(x_k).$$

This is a descent method...

~~$$\begin{aligned} \mathbb{E}[f(x_{k+1})] &\leq \mathbb{E}\left[f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2\right] \\ &= f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \mathbb{E}[\|x_{k+1} - x_k\|^2] \end{aligned}$$~~

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2 \\ &= f(x_k) - \frac{1}{L_d} \nabla f(x_k)^T \left(d \frac{\partial f}{\partial x_i}(x_k) e_i\right) \\ &\quad + \frac{1}{2Ld^2} \left(d \frac{\partial f}{\partial x_i}(x_k)\right)^2 \\ &= f(x_k) - \frac{1}{Ld} \left(\frac{\partial f}{\partial x_i}(x_k)\right)^2 + \frac{1}{2L} \left(\frac{\partial f}{\partial x_i}(x_k)\right)^2 \\ &= f(x_k) - \frac{1}{2L} \left(\frac{\partial f}{\partial x_i}(x_k)\right)^2. \end{aligned}$$

$$\begin{aligned} \mathbb{E}[f(x_{k+1})] &\leq \mathbb{E}[f(x_k)] - \frac{1}{2L} \mathbb{E}\left[\left(\frac{\partial f}{\partial x_i}(x_k)\right)^2\right] \\ &= \mathbb{E}[f(x_k)] - \frac{1}{2L} \cdot \frac{1}{d} \mathbb{E}[\|\nabla f(x_k)\|^2] \end{aligned}$$

Iteratively apply that

$$\Rightarrow \mathbb{E}[f(x_{k+1})] \leq \mathbb{E}[f(x_0)] - \frac{1}{2Ld} \sum \mathbb{E}[\|\nabla f(x_k)\|^2]$$

$$\Rightarrow \mathbb{E}\left[\min_{k \leq T} \|\nabla f(x_k)\|^2\right] \leq \frac{2L \cdot d \cdot (f(x_0) - \min f)}{T}$$

Same order of magnitude as our classic
nonconvex, smooth GD rate.

Select i with largest $\frac{\partial f}{\partial x_i}(x_k)$ (Gauss-Southwell Rule)

[ICML'15, Nutini et al.]

Cyclic Coordinate Descent.

(similar guarantees)

[Beck, Tetrushvili
2015, SIOPF]

Stochastic Variance Reduced Gradient Method (SVRG)

Family: $\begin{cases} \text{SAG} \\ \text{SAGA} \\ \text{SDCA} \\ \vdots \end{cases}$

$$\tilde{x}_0 \leftarrow x_0$$

Compute full gradient $\nabla f(\tilde{x})$

outer loop
for $i = 0, \dots,$
 $y_0 = \tilde{x}_0$
 inner loop
 for $j = 0, \dots, 2d$
 $y_{j+1} = y_j - \alpha_k g(y_j)$
 end for
 $\tilde{x}_{i+1} = \text{average}(y_0, \dots, y_{2d})$
 Compute full gradient $\nabla f(\tilde{x}_{i+1})$
end for

SVRG
 $\nabla f(\tilde{x}_i) + \nabla f_i(y_j) - \nabla f_i(\tilde{x}_i)$
is Uniformly!

Storage: $\begin{cases} \nabla f(\tilde{x}_i) \\ \nabla f_i(y_j) \\ \nabla f_i(\tilde{x}_i) \end{cases}$
Each step $\begin{cases} y_j \\ \tilde{x}_i \end{cases}$

[Johnson
Zhang
2013]

Theorem If f is L -smooth μ -strongly convex,
with SVRG converges linearly.

[For comparison, stochastic GD, $O(1/k)$ rate]