

- Last time: SGD
 - #iter(ε) worse than GD
 - cost/iter $\approx n \times$ faster than GD

Today's Lecture:

Q: Can we have the best of both worlds, i.e.,
 convergence as fast as GD
and $O(1)$ grad. computation/iter like SGD

A: SVRG.

We would like to understand what causes SGD to have slower rates than GD

We will revisit the finite sums setup:

$$f(\omega) = \frac{1}{n} \sum_i f_i(\omega)$$

Quick comparison:

- SUD on γ -str. cvx:

$$\mathbb{E} \|\omega_{k+1} - \omega^*\|^2 \leq (1-\gamma\gamma) \mathbb{E} \|\omega_k - \omega^*\|^2 + \gamma^2 \mathbb{E} \|\nabla f_{sk}(\omega_k)\|^2$$

$$\left[\begin{array}{l} \text{Due to bound} \\ \text{from last} \\ \text{lecture} \end{array} \right] \leq (1-\gamma\gamma)^{k+1} \|\omega_0 - \omega^*\|^2 + \frac{\gamma}{\eta} M^2$$

- UD on γ -str CVX + B-smooth

Due to smoothness we have:

$$\begin{aligned} \|\nabla f(\omega_k)\|^2 &= \|\nabla f(\omega_k) - \nabla f(\omega^*)\|^2 \\ &\leq \beta^2 \|\omega_k - \omega^*\|^2 \end{aligned}$$

Hence a simple conv. rate bound gives:

$$\begin{aligned} \|\omega_{k+1} - \omega^*\|^2 &\leq (1-\gamma\gamma) \|\omega_k - \omega^*\|^2 + \gamma^2 \beta^2 \|\omega_k - \omega^*\|^2 \\ &\leq (1-\gamma\gamma + \gamma^2 \beta^2)^{k+1} \cdot \|\omega_k - \omega^*\|^2 \end{aligned}$$

Observe that SGD rates look like:

$$C_1 \cdot \|w_0 - w^*\|^2 + V$$

And for GD : $C_2^k \|w_0 - w^*\|^2$

The "V"-term causes worse rates in SGD

Q: Why? Can we fix it?

Remark: We can't take advantage of smoothness in SGD since

$$\|\nabla f_{S_k}(w_k)\|^2 \leq B_{S_k} \|w_k - w_{S_k}^*\|^2$$

That is smoothness gives us an upper bound, but only with respect to the global opt of a single function and in general

$$\arg \min_w f_i(w) \neq \arg \min_w S f_i(w)$$

However, we can do a trick:

$$\begin{aligned}\|\nabla f_{SK}(\omega_k)\|^2 &\leq \|\nabla f_{SK}(\omega_k) - \nabla f_{SK}(\omega^*) + \nabla f_{SK}(\omega^*)\|^2 \\ &\leq 2\|\nabla f_{SK}(\omega_k) - \nabla f_{SK}(\omega^*)\|^2 + 2\|\nabla f_{SK}(\omega^*)\|^2 \\ &\leq 2 \cdot \beta_{SK} \|\omega_k - \omega^*\| + 2 \cdot \sigma_{SK}\end{aligned}$$

$\underbrace{\phantom{2 \cdot \beta_{SK}}}_{A}$ $\underbrace{\phantom{2 \cdot \sigma_{SK}}}_{B}$

A looks like the term in GD and
B measures how large the grad. of
 f_{SK} is at the global min of $\sum_i f_i(\omega)$.

Remark:

When $A \geq B \Rightarrow$ SGD is in the
linear rate regime
i.e. variance decays with # iter.

What we want: A variant of SGD, e.g.

$$\omega_{k+1} = \omega_k - \gamma \cdot g_k(\omega_k)$$

- such that:
- $E g_k = \nabla f$
 - g_k is "cheap" on average
 - $A \geq B$ always

This is possible!

SVRG:

Stochastic variance reduced gradient method.

In SVRG we choose g_k as follows:

$$g_k(\omega) = \nabla f_{S_k(\omega)} - \nabla f_{S_k(\omega_0)} + \nabla f(\omega_0)$$

full grad.

This term will allow the $A \geq B$ property

Let's bound the variance of g_k

$$\begin{aligned}
 \mathbb{E} \|g_k(\omega)\|^2 &= \mathbb{E} \|g_{k(\omega)} \pm \nabla f_{sk}(\omega^*)\|^2 \\
 &= \mathbb{E} \|\nabla f_{sk}(\omega^*) - \nabla f_{sk}(\omega_0) + \nabla f(\omega_0) \pm \nabla f_{sk}(\omega^*)\|^2 \\
 &\leq 2 \mathbb{E} \|\nabla f_{sk}(\omega_0) - \nabla f_{sk}(\omega^*)\|^2 + \\
 &\quad 2 \mathbb{E} \|\nabla f_{sk}(\omega_0) - \nabla f_{sk}(\omega^*) - \nabla f(\omega_0)\|^2 \\
 &\leq 2 \beta \mathbb{E} \|\omega_k - \omega^*\|^2 + (\dots)
 \end{aligned}$$

Observe now that

$$\begin{aligned}
 &\mathbb{E} \|\nabla f_{sk}(\omega_0) - \nabla f_{sk}(\omega^*) - \nabla f(\omega_0)\|^2 \\
 &= \mathbb{E} \|\nabla f_{sk}(\omega_0) - \nabla f_{sk}(\omega^*) - \underbrace{\nabla f(\omega_0) + \nabla f(\omega^*)}_{\mathbb{E} X}\| \\
 &\quad \underbrace{\qquad\qquad\qquad}_{X} \qquad \underbrace{\qquad\qquad\qquad}_{\mathbb{E} X} \\
 &= \mathbb{E} \|X - \mathbb{E} X\|^2 \leq \mathbb{E} \|X\|^2 \\
 &= \mathbb{E} \|\nabla f_{sk}(\omega_0) - \nabla f_{sk}(\omega^*)\|^2 \\
 &\leq \beta \mathbb{E} \|\omega_0 - \omega^*\|^2
 \end{aligned}$$

With the above "update rule" we then get:

If f is γ -str. conv and each f_i is β -smooth

$$\begin{aligned} \mathbb{E}\|\omega_{k+1} - \omega^*\|^2 &\leq \mathbb{E}\|\omega_k - x^*\|^2 - (2\gamma\lambda + 2\gamma^2\beta)\mathbb{E}\|\omega_k - \omega^*\|^2 \\ &\quad + 2\gamma^2\beta^2\mathbb{E}\|\omega_0 - \omega^*\|^2 \\ &\leq (1 - 2\gamma\lambda + 2\gamma^2\beta)^{k+1}\|\omega_0 - \omega^*\|^2 + 2(k+1)\gamma^2\beta^2\|\omega_0 - \omega^*\|^2 \end{aligned}$$

$\underbrace{\phantom{(1 - 2\gamma\lambda + 2\gamma^2\beta)^{k+1}}}_{1/4}$ $\underbrace{}_{1/4}$

Set $2(k+1)\gamma^2\beta^2 = 1/4$

$$\Rightarrow \gamma = O(1) \cdot \frac{1}{\beta \cdot k}$$

Set $(1 - 2\gamma\lambda + 2\gamma^2\beta)^{k+1} = 1/4$

$$\Rightarrow \left(1 - \frac{C\gamma}{\beta \cdot k} + \frac{C'\cdot 1}{\beta \cdot k^2}\right)^{k+1} = 1/4$$

Setting $k = O(n) \cdot \beta^2/\gamma^2$ gives the above

Hence if $\gamma = O(n) \frac{\gamma}{\beta^2}$, $k = O(n) \cdot \beta^2/\gamma^2$

$$\Rightarrow \mathbb{E}\|\omega_T - \omega^*\|^2 \leq 0.5 \|\omega_0 - \omega^*\|^2$$

- Remarks:
- The above only decreases dist. to opt by a constant factor
 - The step

$$g_k(\omega) = \nabla f_{sk}(\omega) - \nabla f_{sk}(\omega_0) + \nabla f(\omega)$$
 costs 1 full grad, i.e., its complexity is proportional to $\mathcal{O}(D)$

Solutions to above:

Repeat in "epochs"

SVRG:

```

y = w₀
t = 0
for epoch = 1 : E
    q = ∇f(y)
    for s = 1 : S
        s_t ~ unif {1, ..., n}
        w_{t+1} = w_t - γ (∇f_{s_t}(w_t) - ∇f_{s_t}(y) + q)
        t = t + 1
    y = w_{K-1}
  
```

Observe that the cost of $g = \nabla f(y)$ becomes amortized

Also overall rate

$$E \|\omega_E - \omega^*\|^2 \leq 0.5^E \|\omega_0 - \omega^*\|$$

similar to GD

Hence,

epochs to ε accuracy: $O(\log(\frac{1}{\varepsilon}))$

iterations/epoch: $K = O(\frac{\beta^2}{\alpha\varepsilon})$

computational cost/epoch: 1 full grad
+ "small" grads

Overall complexity of SVRG

$$O\left(\log\left(\frac{1}{\varepsilon}\right) \cdot \text{cost}(\nabla f) + \frac{\beta^2}{\alpha^2} \log\left(\frac{1}{\varepsilon}\right) \cdot \frac{\text{cost}(\nabla f)}{n}\right)$$

Compare to GD: $O\left(\frac{\beta^2}{\alpha^2} \log\left(\frac{1}{\varepsilon}\right) \cdot \text{cost}(\nabla f)\right)$.

Remarks:

- SVRG has linear rate of convergence like GD and small amortized cost/iter like SGD
- However we have to tune more hyperparam i.e., stepsize + length of each epoch.

Open Problems:

- What happens if we first run vanilla SGD and then full GD or SVRG
- Adaptively change estimate $g_k(w)$ can be coarser in the beginning and finer towards the end
- What is the best way to choose g_k to optimize rates?
- Performance on non-convex problems is questionable. Why?