

Balls and Parabolas

Thursday, January 23, 2020 5:32 AM

Yuhua
we have seen convergence rates of $(1 - \frac{1}{n^2})$, $(1 - \frac{1}{n})$
(in function value) for convex f
and $(1 - \frac{\mu}{L})$ for strongly convex f with GD.

Is faster possible? $(1 - \frac{1}{n})$ is best for general convex.

When $\frac{\mu}{L} > \frac{1}{n}$, the rate $(1 - \frac{\mu}{L})$ is better.

Can we go faster?

Part 1. GD as a cutting plane method.

For convex, $f(x) \geq f(x^k) + \langle \nabla f(x^k), x - x^k \rangle$

$$\text{So } \{x: f(x) \leq f(x^k)\} \subseteq \{x: \langle \nabla f(x^k), x - x^k \rangle \leq 0\}$$

Halfspace.

When f is strongly convex,

$$f(x^*) \geq f(x) + \langle \nabla f(x), x - x^* \rangle + \frac{\mu}{2} \|x^* - x\|^2$$

$$f(x^*) - f(x) \geq \frac{\mu}{2} \|x^* - x\|^2 + \frac{\|\nabla f(x)\|^2}{\mu} - \frac{2\|\nabla f(x)\|^2}{\mu}$$

i.e. with $x^+ = x - \frac{\nabla f(x)}{\mu}$

$$\|x^* - x^+\|^2 \leq \frac{\|\nabla f(x)\|^2}{\mu^2} - \frac{2}{\mu} (f(x) - f(x^*))$$

Recall:

$$f(x^*) \leq f(x^+) \leq f(x - \frac{\nabla f(x^*)}{\mu}) \leq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$$

$$f(x^*) \leq f(x^+) \leq f\left(x - \frac{\nabla f(x^*)}{\mu}\right) \leq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$$

$$\begin{aligned} \therefore \|x^* - x^+\|^2 &\leq \frac{\|\nabla f(x)\|^2}{\mu^2} - \frac{2}{\mu} (f(x) - f(x^+)) - \frac{2}{\mu} (f(x^+) - f(x^*)) \\ &\leq \frac{\|\nabla f(x)\|^2}{\mu^2} \left(1 - \frac{\mu}{L}\right) - \frac{2}{\mu} (f(x^+) - f(x^*)). \end{aligned}$$

So the OPT x^* lies in the ball

$$B\left(x^+, \sqrt{1 - \frac{\mu}{L}} \cdot \frac{\|\nabla f(x)\|}{\mu}\right).$$

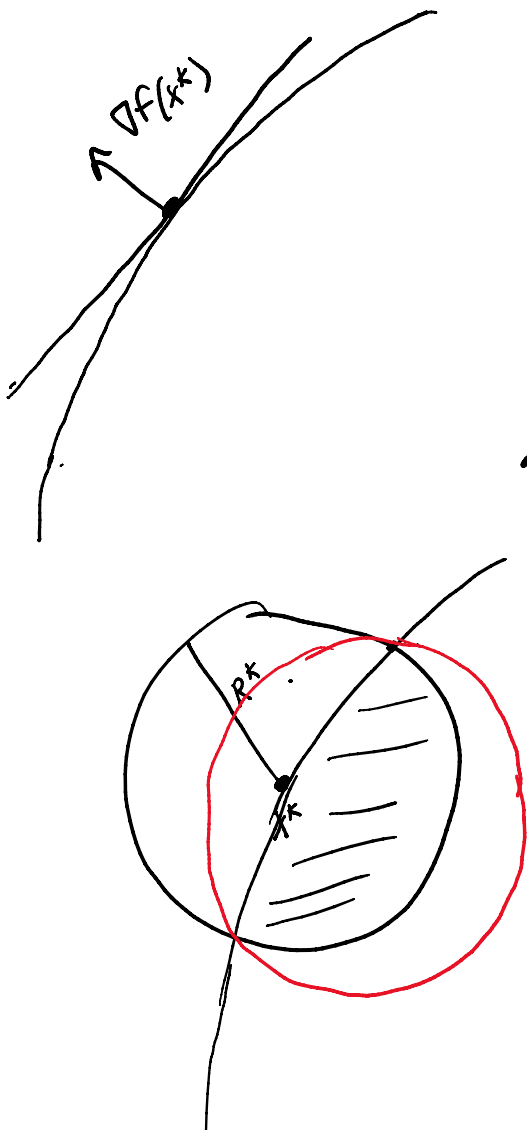
We now use balls instead of Ellipsoids.

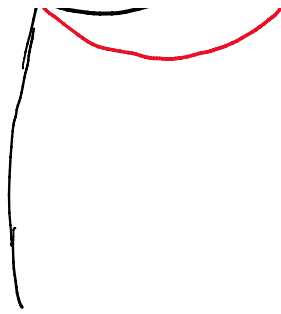
$$E_0 = B\left(0, \frac{\|\nabla f(x^0)\|}{\mu}\right)$$

$$E_k = B(x^k, R_k).$$

$$B^+ = B\left(x^+, \frac{\|\nabla f(x^k)\|}{\mu} \sqrt{1 - \frac{\mu}{L}}\right)$$

$$E_{k+1} = \text{min ball containing } E_k \cap B^+.$$

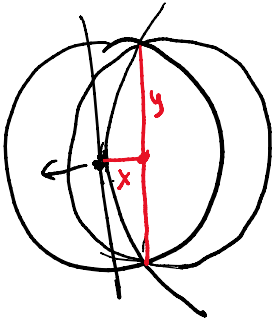




x^+

Lemma. $\exists x: B(0,1) \cap B(g, \|g\|\sqrt{1-\epsilon}) \subseteq B(x, \sqrt{1-\epsilon})$.

Pf.



$$x^2 + y^2 = 1$$

$$(x-g)^2 + y^2 = g^2(1-\epsilon)$$

$$\Rightarrow x^2 + y^2 - 2gx + 1 - x^2 = g^2 - \epsilon g^2$$

$$x = \frac{1 + \epsilon g^2}{2g}$$

$$\text{and } y^2 = 1 - \frac{1}{4g^2} - \frac{\epsilon^2 g^2}{4} - \frac{\epsilon}{2}$$

$$= 1 - \frac{\epsilon}{2} - \frac{1}{2} \left(\frac{1}{2g^2} + \frac{\epsilon^2 g^2}{2} \right)$$

$$> 1 - \frac{\epsilon}{2} - \frac{\epsilon}{2} = 1 - \epsilon.$$

Applying this to $B(x^k, R_k)$, $B(x^k - \frac{\nabla f(x^k)}{L}, \frac{\|\nabla f(x^k)\|}{L} \left(1 - \frac{\mu}{L}\right)^{\frac{k}{2}})$

next ball is $B(x, R_k \sqrt{1 - \frac{\mu}{L}})$.

Lemma. $R_{k+1}^2 \leq \left(1 - \frac{\mu}{L}\right) R_k^2$

$$\Rightarrow R_k^2 \leq \left(1 - \frac{\mu}{L}\right)^k \cdot R_0^2 \leq \left(1 - \frac{\mu}{L}\right)^k.$$

Using $\mathcal{B}_k = R_k$

$$\min_i f(x^i) - f(x^*) \leq \left(\sqrt{1 - \frac{\mu}{L}} \right)^k \cdot (f(x^0) - f(x^*))$$

Recovers GD rate for strongly concave f .

Can we do better?

Observe: our cuts come from

$$\langle \nabla f(x), x^* - x \rangle \leq 0$$

we start with 0

but this improves

as we get better estimates

$$\leq f(x^{\text{new}}) - f(x)$$

i.e. we can "tighten" all previous planes!

accelerated ball method

$$x, \quad x^+ = x - \frac{\nabla f(x^0)}{L}$$

$$x^{++} = x - \frac{\nabla f(x^0)}{\mu} \quad C_0 = x^{0++}$$

We maintain x^k, C^k, R_k .

$$R_0^2 = \frac{\|\nabla f(x^0)\|^2}{\mu^2} \left(1 - \frac{\mu}{L}\right)$$

$$x^{k+1} = \text{line search}(C^k, x^{k+})$$

... .. *termina*

$x = x^{(k)}, \dots$

C^{k+1} is center of min ball containing

$$B\left(C^k, \sqrt{R_k^2 - \frac{\|\nabla f(x^{k+1})\|^2}{\mu^2} \cdot \frac{\mu}{L}}\right) \cap B\left(x^{k+1}, \frac{\|\nabla f(x^{k+1})\|}{\mu} \sqrt{1 - \frac{\mu}{L}}\right)$$

Lemma. $x^* \in B(C^k, R_k)$

$$R_{k+1}^2 \leq \left(1 - \sqrt{\frac{\mu}{L}}\right) R_k^2$$

Pf. We show by induction that

$$x^* \in B\left(C^k, \sqrt{R_k^2 - \frac{2}{\mu} (f(x^{k+1}) - f(x^*))}\right)$$

$k=0$. ✓

$$f(x^{(k+1)}) \leq f(x^{k+1}) - \frac{\|\nabla f(x^{k+1})\|^2}{2L} \leq f(x^{k+1}) - \frac{\|\nabla f(x^{k+1})\|^2}{2L}$$

$$R_{k+1}^2 \leq R_k^2 - \frac{2}{\mu} (f(x^{k+1}) - f(x^*))$$

$$= R_k^2 - \frac{2}{\mu} (f(x^{k+1}) - f(x^{(k+1)+})) - \frac{2}{\mu} (f(x^{(k+1)+}) - f(x^*))$$

$$\leq R_k^2 - \frac{\|\nabla f(x^{k+1})\|^2}{\mu^2} \cdot \frac{\mu}{L} - \frac{2}{\mu} (f(x^{(k+1)+}) - f(x^*))$$

We also have

$$x^* \in B\left(x^{(k+1)+}, \sqrt{\frac{\|\nabla f(x^{k+1})\|^2}{\mu^2} \left(1 - \frac{\mu}{L}\right) - \frac{2}{\mu} (f(x^{(k+1)+}) - f(x^*))}\right)$$

$$x^r \in B\left(x^{(k+1)}, \sqrt{\frac{\|\nabla f(x^{(k+1)})\|^2}{n^2} \left(1 - \frac{\mu}{L}\right) - \frac{2}{\mu} (f(x^{(k+1)}) - f(x^*))}\right)$$

Lemma:

$$B(0, \sqrt{1 - g^2 \epsilon - \delta}) \cap B(a, \sqrt{g^2(1 - \epsilon) - \delta})$$

$$|a| > g \subseteq B(c, \sqrt{1 - \sqrt{\epsilon} - \delta})$$

$$\Rightarrow R_{k+1}^2 \leq R_k^2 \left(1 - \frac{\mu}{L}\right)$$

\therefore using $\mathcal{D}(\cdot) = R_k$

rate is $\left(1 - \frac{\mu}{L}\right)$.

Thm. For any algorithm "using only gradients"

i.e. $x^k \in \text{Span}(x^0, \nabla f(x^0), \nabla f(x^1), \dots, \nabla f(x^{k-1}))$

best possible rate is $\min\left\{n, \sqrt{\frac{L}{\mu}}\right\}$.