

Call a method is a first-order Method if each iteration

"it uses all previous computed gradients $\nabla f(x_1), \nabla f(x_2), \dots, \nabla f(x_{k-1})$ and initial x , to generate a new point x_k .

In short, first-order Method can only access to f and ∇f .

Remark: There are other access model such as f^* and ∇f^* .

Usually, first-order method satisfies

$$x_k \in \text{span}(x_1, \nabla f(x_1), \nabla f(x_2), \dots, \nabla f(x_{k-1}))$$

According to the definition, cutting plane methods are first-order.

However, the common use of the term "first-order method" refers to methods that only use "simple formula"

A typical first-order method:

$$x_{k+1} = x_k - \frac{1}{\beta} \nabla f(x_k) \quad (\text{gradient descent})$$

Common belief: First-order method (those with simple formula) works well for large-scale problems

My belief: There should be something that is better.

Practicality of Cutting Plane Methods.

How is it possible to say compute the center of gravity of 1 billions dimension convex sets if the runtime is $O(n^5)$?

→ I am writing a package with run time closer to linear.
Let see how long does it take.

→ You never need to do so.

Note that all cutting plane method we mentioned satisfies

$$x_k \in \text{span}(x_1, \nabla f(x_1), \nabla f(x_2), \dots, \nabla f(x_{k-1})) \quad (k \text{ th dim})$$

Therefore, wlog, we only need to solve some $k-1$ dim problem at k^{th} step.

Therefore, we only need to solve some k -th dim problem at k th step.

Roughly speaking, if the cutting plane method involves solving a problem that takes $f(k)$ time for $1D^k$.

Then, the total time for the k th iter is

$$T(\nabla f) + nk + f(k)$$

compute gradient iteration solve some crazy
 subspace problem

As I mentioned last time, for problem of the form $\sum_i f_i(a_i^T x - b_i)$,

$T(\nabla f)$ is often dominated by computing Av for some v where $A = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{pmatrix}$

So, a more fair runtime estimate for cutting plane method is

$$\underbrace{\text{nnz}(A)}_{\text{number of nonzeros}} + nk + f(k)$$

As comparison, L-BFGS takes
 $\text{nnz}(A) + nk$.

If it is a large-scale problem, $f(k) \ll nk$.

My belief: Cutting plane methods becomes more competitive as $n \rightarrow \infty$.

If the problem is large scale and we spent lots of time to get ∇f ,
Then it sounds silly to just combine three gradient by simple combination.

Note that IPM, CHOL, Multigrid, AGD, they all takes decades to become recognized. Not sure about cutting plane methods.

Thesis Question: Develop cutting plane method that take advantage of $\sum f_i$ structure and implement it.

Summary of first-order methods

Algorithm	errors	(Assume the function is α strongly convex, β smooth)
Gradient Descent	$\beta R^2/t$ $\beta R^2(1 - \frac{\alpha}{\beta})t$	α strongly convex, β smooth G-Lipschitz, $R = \ x_0 - x^*\ _2$ equivalent
Mirror Descent	$G R / \sqrt{t}$ $G^2 / \alpha t$	tight tight reduction
Armijo-tol	$\beta R^2 / 4\epsilon$	tight ..

Accelerated
Gradient Descent

Certified
Cutting Plane
Methods

$\frac{\beta R^2}{\alpha t}$	tight	\leftarrow
$\beta R^2 (1 - \sqrt{\frac{\alpha}{\beta}}) t$	tight	\downarrow
$(1 - \frac{1}{n})^R(t)$	tight	\leftarrow
$\beta R^2 (1 - \sqrt{\frac{\alpha}{\beta}}) t$	$\tilde{O}(\text{tight})$	

Misconception: linear convergence is better than sublinear convergence
 $(1 - \varepsilon)^t$ vs $\frac{1}{t^c}$

We can always solve convex problem with linear convergence (aka $(1 - \frac{1}{n})^{-2(+)}$)

Misconception: dimension independent result is faster for large-scale problem.

Even for very simple problem like $\sum (x_i - x_{i+1})^2$, we have $\frac{\beta}{\alpha} = n^2$.

Therefore dimension independent result can depends on dimension.

Relations between those bounds

① Suppose we have a way to find x^*

$$f(x) - f(x^*) \leq \frac{G^2}{\alpha t}$$

for any α -strongly convex, G -Lip f

For any G -Lip convex f , we consider $f_\alpha(x) = f(x) + \frac{\alpha}{2} \|x - x_0\|^2$

We have that

$$\begin{aligned} f(x) - f(x^*) &\leq f_\alpha(x) - f_\alpha(x^*) + \frac{\alpha}{2} \|x^* - x_0\|^2 \\ &\leq \frac{G^2}{\alpha t} + \frac{\alpha}{2} R^2 \quad (\text{used the algo above}) \end{aligned}$$

Choosing the best α , we have $f(x) - f(x^*) \leq \frac{G^2 R^2}{5\alpha t}$

Hence, if we have an algorithm achieves $\frac{G^2}{\alpha t}$,

we can turn it to an algorithm achieves $\frac{G^2 R^2}{5\alpha t}$.

② Suppose we know $\frac{BR^2}{t^2}$. If f is α -strongly convex, we have

② Suppose we know $\frac{\beta^2}{t^2}$. If f is α -strongly convex, we have

$$\frac{\alpha}{2} \|x_t - x^*\|^2 \leq f(x_t) - f(x^*) \leq \frac{\beta}{t^2} \|x_0 - x^*\|^2$$

Set $t = 2\sqrt{\frac{\beta}{\alpha}}$, we have $\|x_t - x^*\|^2 \leq \frac{1}{2} \|x_0 - x^*\|^2$

Hence, every $\geq \sqrt{\frac{\beta}{\alpha}}$, the distance to optimum decrease by $\frac{1}{2}$.

By restarting the algorithm, we have $\beta^2(1 - \frac{1}{2}\sqrt{\frac{\alpha}{\beta}})^t$

Lower bounds

Consider the problem:

$$-X + \frac{\beta}{2} \sum_{i=1}^n (x_i - x_{i+1})^2 + \frac{\alpha}{2} \sum_{i=1}^n x_i^2$$

The sol'n is of the form



Any approximate sol'n has support at least $\sqrt{\frac{\beta}{\alpha}}$.

Next, note that if

$$x_k = \text{span}(x_1, \nabla f(x_1), \nabla f(x_2), \dots, \nabla f(x_{k-1}))$$

and if $x_1 = (1, 0, 0, \dots, 0)$

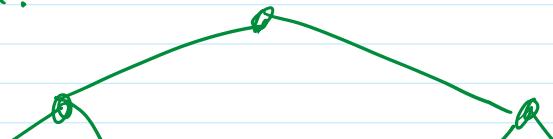
Then x_k has only k non-zeros in the beginning.

Therefore, we need $\sqrt{\frac{\beta}{\alpha}}$ iterations.

All first-order methods including cutting plane methods has this serious issue
 $\# \text{iter} \geq \text{"diameter of the problem"}$.

Open problem: Is there any way to help the connectivity for first order method?

How about



How about

