June 11 '18

## Lecture 3: Modern $1^{st}$-Order Methods For Structured Problems

▷ GRADIENT DESCENT

$$\boxed{\begin{array}{l} \min_x f(x), \quad \text{assuming } f \text{ is } \mathbb{R}\text{-valued}, \\ \qquad f \in \Gamma_0(\mathbb{R}^n), \; \nabla f \text{ } L\text{-Lipschitz} \end{array}}$$    e.g., $0 \preceq \nabla^2 f(x) \preceq L \cdot I$

$$X_{k+1} = \arg\min_x \; f(x_k) + \langle \nabla f(x_k), x-x_k \rangle + \underbrace{\frac{L}{2} \| x - x_k \|^2}$$

         discuss, vis-a-vis

$$= x_k - \frac{1}{L} \cdot \nabla f(x_k)$$

Cond. Gra / Frank-Wolf

$$= x_k - t \, \nabla f(x_k) \; w, \; \text{stepsize/learning-rate } t = \frac{1}{L}$$

Nesta    See B.bick

... is a
... method.

... like send
... help
... anatically
... worst case)

Convergence Analysis (Vandenberghe's notes, ie., Nesterov's book), for $t = \frac{1}{L}$

let $x^*$ be any optimal sol'n

$$f(x_{k+1}) \le f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \| x_{k+1} - x_k \|^2 \quad \text{"Descent Lemma"}$$

$$= f(x_k) + \langle \nabla f(x_k), -\frac{t}{L} \nabla f(x_k) \rangle + \frac{L}{2} \| \frac{1}{L} \nabla f(x_k) \|^2$$

See picture

$$= f(x_k) - \frac{1}{2L} \| \nabla f(x_k) \|^2 \quad (\Rightarrow \text{descent method}) \; \{ f(x_{k+1}) \le f(x_k) \}$$

$$\le f(x^*) + \langle \nabla f(x_k), x_k - x^* \rangle - \frac{1}{2L} \| \nabla f(x_k) \|^2 \quad \text{via convexity.}$$

$$= f(x^*) + \frac{L}{2} \left( \| x_k - x^* \|^2 - \| x_k - x^* - \frac{1}{L} \nabla f(x_k) \|^2 \right)$$

$$= f(x^*) + \frac{L}{2} \left( \| x_k - x^* \|^2 - \| x_{k+1} - x^* \|^2 \right).$$

Add, for $i = 1, ..., k,$

$_{)-f(\cdot)) \le}$

$$\frac{1}{k} \sum_{i=1}^{k} \left( f(x_i) - f(x^*) \right) \le \frac{1}{k} \frac{L}{2} \sum_{i=1}^{k} \left( \| x_{i-1} - x^* \|^2 - \| x_i - x^* \|^2 \right) \quad \text{telescopes}$$

$$= \frac{1}{k} \frac{L}{2} \left( \| x_0 - x^* \|^2 - \underbrace{\| x_k - x^* \|^2}_{\ge 0} \right)$$

$$\le \frac{1}{k} \frac{L}{2} \| x_0 - x^* \|^2$$

$\heartsuit$

And, since it was a descent method, $f(x_k) \le f(x_i) \quad \forall \, i = 1, ..., k$

So

$$\boxed{ f(x_k) - f^* \le \frac{L}{2k} \| x_0 - x^* \|^2 } \quad < \varepsilon$$

or, if we want $f(x_k) - f^* < \varepsilon$, take $k > \frac{L}{2} \| x_0 - x^* \|^2 \cdot \frac{1}{\varepsilon}$, ie., $\boxed{ O(\frac{1}{\varepsilon}) \text{ iterations} }$

"sublinear"

* ~~Short~~ $\boxed{\text{Asymptotic Worst-Case Result Only !}}$ Discuss...

Is this rate good? Tight?

$x_{k+1} \in \text{span}\{x_0, \nabla f(x_0), ..., \nabla f(x_k)\}$

Thm: Nesterov 1980s, 2003 book:

No 1st order method can ~~beat~~ always guarantee

$$f(x_k) - f(x^*) < \frac{3}{32} \cdot L \cdot \frac{\|x_0 - x^*\|^2}{k^2} \quad \text{for } k \le \frac{1}{2}(n-1).$$

Suppose $f$ is strongly cvx

$\forall x, \quad \forall y, \quad f(y) \ge f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2}\|y-x\|^2 \quad \text{"sc"}$

$f(y) \ge g(y) \implies \min\limits_y f(y) \ge \min\limits_y g(y)$

$y = x - \frac{1}{\mu}\nabla f(x)$

So

$$\min\limits_y f(y) = f(x^*) \ge f(x) - \frac{1}{2\mu}\|\nabla f(x)\|^2 = \min\limits_y g(y)$$

PL Polyak - Łojasiewicz Ineq: $\forall x, \quad \frac{1}{2}\|\nabla f(x)\|^2 \ge \mu \cdot (f(x) - f^*)$

weaker than SC (~~ie~~ unique sol'n, while SC does)

"essential SC", "restricted SC", "error bnd condition" ... $\implies$ PL.

$\begin{bmatrix} \text{Can bound sub-optimality} \\ \text{by gradient} \end{bmatrix}$

Thm (Karimi, Nutini, Schmidt '16)

Grad. Descent w/ $t = \frac{1}{L}$, $\nabla f$ L-Lipschitz, and $\mu$-PL, satisfies

$$f(x_k) - f^* \le \left(1 - \frac{\mu}{L}\right)^k \cdot (f(x_0) - f^*)$$

proof:

as before, $f(x_{k+1}) - f(x_k) \le \frac{-1}{2L}\|\nabla f(x_k)\|^2$

$\le -\frac{\mu}{L}(f(x_k) - f^*)$ via PL

re-arrange, subtract $f^*$,

$$f(x_{k+1}) - f^* \le (1 - \frac{\mu}{L})(f(x_k) - f^*). \quad \blacksquare$$

Rates to reach $\varepsilon$-sol'n

| | | to go from $\varepsilon_0$ to $\varepsilon' = 10^{-2} \cdot \varepsilon_0$, need more it |
|---|---|---|
| subgradient desc. | $O(1/\varepsilon^2)$ | 10,000 × more |
| gradient desc. | $O(1/\varepsilon)$ | 100 × more |
| accel. gra. desc. | $O(1/\sqrt{\varepsilon})$ | 10 × more |
| linear (eg, SC) | $O(-\log(\varepsilon))$ | ~~about 2 more~~ 2·constant more |
| quadratic (eg, Newton) | $O(\log(-\log(\varepsilon)))$ | 1 more |

these depend on $\varepsilon_0$ (eg, $\varepsilon_0 = 1$)

—

▷ NESTEROV'S ACCELERATED METHOD — "optimal"

　　Extends heavy-ball method (analysis for quadratics in, eg, Bertsekas )

　　　"momentum"　$x_{k+1} = x_k - t_k \, \nabla f(x_k) + s_k (x_k - x_{k-1})$.

See

　　distill.pub/2017/momentum

　Nesterov's Method, variant

　　　$x_{k+1} = y_k - t \cdot \nabla f(y_k)$　　　$\boxed{\begin{array}{l} x_k = y_{k-1} - t \, \nabla f(y_{k-1}) \\[4pt] y_k = x_k + \frac{k}{k+3}(x_k - x_{k-1}) \end{array}}$
　　　$y_{k+1} = x_{k+1} +$

　"Thm　For $\nabla f$ L-Lipschitz, if $t = \frac{1}{L}$,　$f(x_k) - f^\# = O(\frac{1}{k^2})$.
　　(see Vandenberghe, or my 4720 notes p 55 ) for proof.


　• Not a descent method !
　• If $f$ is strongly cvx, ought to know the constant $\mu$ in order to exploit.
　• Extends to many variants, eg, proximal "FISTA"　　⟶ discuss slow popularity
　• Tricky to actually get $(x_k)$ to converge, not just $(f(x_k))$　⟶


▷ NON-SMOOTH.

　Suppose $f \in \Gamma(\mathbb{R}^n)$ but $\nabla f$ doesn't exist ( $\partial f$ does )
　Apply smooth method, hope for the best ?
　　(p.40...) No, ex (Shor '98 ), even if you don't hit pts of non-diff,
　　　it messes up. (Wolfe), even if convex


Subgradient method
　• $x_{k+1} = x_k - t_k \cdot \frac{d_k}{}$ for any (deterministic) $d_k \in \partial f(x_k)$, ~~assuming~~
　• $f \in \Gamma_0(\mathbb{R}^n)$, for simplicity assume also $\text{dom}(\partial f) = \mathbb{R}^n$
　• Assume $f$ (not $\nabla f$) is Lipschitz continuous, constant $L_0$, ie, $\boxed{\|d_k\| \le L_0}$


A)　Thm (8.13 in Beck)　p.203　　why? Not a descent method — not even descent direction

　　$\min\limits_{i \in 0,1,\dots,k} f(x_i^*) := f_{best}^k \le \dfrac{L_0 \cdot \text{dist}(x_0, \text{optimal})}{\sqrt{k+1}}$　ie.　$O(\frac{1}{\sqrt{k}})$

　　if $t_k = \dfrac{f(x_k) - f^*}{\|d_k\|^2}$　"Polyak's Stepsize Rule"

(skip mostly)

B) variant: __Thm (8.25 Beck)__  If $t_k$ isn't Polyak, but $\dfrac{\sum_{i=0}^{k} t_i^2}{\sum_{i=0}^{k} t_i} \to 0$ as $k \to \infty$,

then $f_{best}^k - f^* \to 0$ as $k \to \infty$

e.g., $t_k = \dfrac{1}{\sqrt{k+1}}$

c) variant: __Thm (8.28 Beck)__  If $t_k = \dfrac{1}{\|g_k\| \cdot \sqrt{k+1}}$, $f_{best}^k - f^* = O\left(\dfrac{\log(k)}{\sqrt{k}}\right)$

and, ergodic result, if $\bar{x}_k = \dfrac{1}{\sum_{i=1}^{k} t_i} \sum_{i=1}^{k} t_i x_i$ is average, ( note: can compute via a recursive update )

$f(\bar{x}_k) - f^* = O\left(\dfrac{\log(k)}{\sqrt{k}}\right)$

( remove $\log(k)$ if domain is compact )
( see e.g. Bubeck )

D)
p.265 __Thm (3.2 Bubeck)__

Assume we __project__ onto $C$, radius is $R$, $\|g\| \leq L_0$ again ( $\forall g \in \partial f(x)$, $\forall x$ )

Then if $t = \dfrac{R}{L_0} \cdot \dfrac{1}{\sqrt{k}}$, $f\left(\underbrace{\dfrac{1}{k} \sum_{i=1}^{k} x_i}_{\bar{x}_k}\right) - f^* \leq \dfrac{R \cdot L_0}{\sqrt{k}}$

§3.5 Rates are tight.

E) __Thm (8.31 Beck)__  If $f$ is also $\mu$-strongly convex, take $t_k = \dfrac{2}{\mu(k+1)}$

instead of $O\left(\dfrac{1}{\sqrt{k}}\right)$, and $f_{best}^k - f^* = O\left(\dfrac{1}{\mu k}\right)$, instead of $O\left(\dfrac{1}{\sqrt{k}}\right)$.

also
3.9 in Bubeck

▷ __BETTER__ ...

Consider $\left(\min_x f(x) + g(x),\right)$ $f, g \in \Gamma_0(\mathbb{R}^n)$, $\nabla f$ is $L$-Lipschitz ($f = 0$ ok)
                                                                      $g$ isn't.

$\partial(f + g) = \nabla f + \partial g$, so subgradient descent is ⟶ for one

$x_{k+1} = x_k - t_k \cdot (\nabla f(x_k) + \partial g(x_k))$, and as we saw, need $t_k \to 0$, so it's slow.

Instead, follow gra. desc.

$x_{k+1} = \underset{x}{\operatorname{argmin}} \left( f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \dfrac{L}{2} \|x - x_k\|^2 \right) + g(x)$

$= \underset{x}{\operatorname{argmin}} \; g(x) + \dfrac{L}{2} \| x - (x_k - \dfrac{1}{L} \nabla f(x_k)) \|^2 + const.$

$= \operatorname{prox}_{\frac{1}{L} g} (x_k - \dfrac{1}{L} \nabla f(x_k))$.    PROXIMAL GRA. DESC.

Generalizes projected gradient descent,
no penalty from nonsmoothness of $g$ (if you can compute prox).

recall, $x = \text{prox}_g (y)$ means $0 \in \partial g(x) + (x-y)$, ie., $y \in (I + \partial g)(x)$

ie.

$$\boxed{\text{prox}_g (y) = (I + \partial g)^{-1}(y)}$$

Other derivation:

$0 \in \partial f(x) + \partial g(x)$ is optimality

$\Longleftrightarrow \quad 0 \in \frac{1}{L}\partial f(x) + \frac{1}{L}\partial g(x)$

$\Longleftrightarrow \quad x - \frac{1}{L}\partial f(x) \in x + \frac{1}{L}\partial g$

$\Longleftrightarrow \quad (I - \frac{1}{L}\partial f)(x) \in (I + \frac{1}{L}\partial g) x$

$\Longleftrightarrow \quad \boxed{x = (I + \frac{1}{L}\partial g)^{-1}(I - \frac{1}{L}\partial f) x}$    Fixed Pt. Eq'n.    $x = Tx$

iterate                                 iterate $x_{k+1} = Tx_k$

$$x_{k+1} = (I + \frac{1}{L}\partial g)^{-1}(I - \frac{1}{L}\partial f)(x_k)$$

ie.,    $x_{k+1} = x_k - \frac{1}{L}\partial f(x_k) - \frac{1}{L}\partial g(\boxed{x_{k+1}})$    like implicit method, not explicit.

    "SPECIAL CASE":   $f \equiv 0$,   then   $x_{k+1} = (I + \frac{1}{L}\partial g)^{-1} x_k$   is the prox. pt algorithm

$$x_{k+1} = \arg\min_x \; g(x) + \frac{L}{2}\|x - x_k\|^2$$

• General case is prox.gra.desc.
    aka "forward-backward"

○ Extends to acceleration versions ("FISTA")

• No penalty on convergence rates compared to subgra. descent.

• Can you make a Newton version? Yes, but be careful!

Following Bottou, Curtis, Nocedal

▷ STOCHASTIC METHODS

$$\min_x f(x), \qquad f(x) = \mathbb{E}\, F(x; \xi)$$

$$f(x) = \frac{1}{N} \Sigma_i f_i(x), \quad f_i = F(x; \xi_{[i]})$$

SA "Stochastic Approx."
Robbins Munroe,
Polyak

↖ SAA, "Sample Avg Approx"
or "ERM"

SGD

$$x_{k+1} = x_k - t_k d_k, \qquad \mathbb{E}(d_k) = \nabla f(x_k)$$

$$(\text{ex: } d_k = \nabla_i f(x_k) \text{ for } i \sim \text{uniform } \{1, .., N\})$$

Assume

- $\nabla f$ is L-Lipschitz
- $f$ is $\mu$-Polyak-Łojasiewicz (eg., $\mu$-strongly cvx)
- $f$ bdd below (e.g. $f \geq 0$)
- $\mathbb{E}[\|d_k\|^2] \leq M + M_G \|\nabla f(x_k)\|^2, \quad M_G \geq 1$ (ie., $M_G = 1$ is possible assumption)

Thm 1, fixed stepsize (Thm 4.6 Bottou)

let $t_k = t \leq \dfrac{1}{L \cdot M_G}$, then ——— or $\frac{1}{L}$ if $M_G = 1$

$$\mathbb{E}[f(x_k) - f^*] \leq \frac{t L M}{2\mu} + (1 - t\mu)^{k-1} \left( f(x_1) - f^* - \frac{t L M}{2\mu} \right)$$

ie., converge quickly to region of soln.

proof is similar to PL proof, since just use bounds...

Thm 2, diminishing stepsizes (Thm 4-7 Bottou) ——— ie., $\beta = \frac{2}{\mu}$ is a good choice

let $t_k = \dfrac{\beta}{\gamma + k}$ for $\beta > \frac{1}{\mu}$, $\gamma > 0$, and $t_1 \leq \frac{1}{L M_G}$ (or $\frac{1}{L}$ if $M_G = 1$)

Then $\boxed{\mathbb{E}[f(x_k) - f^*] \leq \dfrac{\nu}{\gamma + k}}$, $\nu = \nu(\beta, \gamma)$ is a constant

⊢— Minibatching —⟩

Exploits GPU, and CPU ⌐draw memory hierarchy.

See MATLAB demo

$*$ | if $f$ is $\mu$-strongly cvx, $f(x_k) - f^* \leq \varepsilon \implies \frac{1}{2}\|x - x^*\|^2 \leq \mu^{-1}\varepsilon$ |

$\Rightarrow$ so $\mathbb{E}\left(\|x_k - x^*\|^2\right) \leq O\left(\frac{1}{k}\right)$

▷ VARIANCE REDUCTION ("Gradient Aggregation")

Specific to $f(x) = \frac{1}{N} \sum_{i=1}^{N} f_i(x)$, e.g., $f_i(x) = \varphi(a_i^T x - b_i)$ for GLM

**Algo : SAGA**

Initialize: $x^{(i)} = x_0$ $\forall i = 1, -, N$ (each $x^{(i)}$ is $n$-dimensional)

and store $\{\nabla f_i(x^{(i)})\}_{i=1}^{N}$ in a $n \times N$ table $*$ (see left page: if then can store $j$-st $a_i^T x^{(i)} - b_i$)

For $k = 1, 2, \dots$

$\quad j \sim \text{Uniform}([1, --, N])$

$\quad \bar{z} = \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(x^{(i)})$ via table

$\quad x_{k+1} = x_k - t_k (\nabla f_j(x_k) - \nabla f_j(x^{(j)}) + \bar{z})$

$\quad x^{(j)} \leftarrow x_k$, update table w/ $\nabla_j f(x^{(j)})$. (update $\bar{z}$)

**Thm** for appropriate $t$, this converges linearly!

**Algo : SVRG**

For $k = 1, 2, \dots$ "epoch"

$\quad \bar{z} = \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(x_k)$    Full pass through data

$\quad w_0 \leftarrow x_k$

$\quad$ For $l = 1, 2, --, m-1$

$\quad\quad j \sim \text{Uniform}([1, -, N])$

$\quad\quad w_{l+1} = w_l - t(\nabla f_j(w_l) - \nabla f_j(x_k) + \bar{z})$

$\quad$ Option I $\quad x_{k+1} = w_m$

$\quad$ Option II $\quad x_{k+1} = w_l$ for $l \sim \text{Uni}([0, --, T-1])$

$\quad$ Option III $\quad x_{k+1} = \frac{1}{m} \sum_{l=1}^{m} w_{l+1}$

**Thm** For appropriate $t$, this converges linearly!

▷ ITERATE AVERAGING

First approach, iterate SGD as usual, $x_{k+1} = x_k - t_k \cdot d_k$

but hope $\bar{x}_k := \frac{1}{k} \sum_{j=1}^{k} x_j$ converges faster than $x_j$.

If we use $t_k = O(1/k)$, it doesn't help.

but, if strongly cvx, choose $t_k = O\left(\frac{1}{k^\alpha}\right)$ for $\alpha \in (\frac{1}{2}, 1)$

then

$$\mathbb{E}\left( \| \bar{x}_k - x^* \|^2 \right) = O\left(\frac{1}{k}\right) \text{ while } \mathbb{E}\left( \| x_k - x^* \|^2 \right) = O\left(\frac{1}{k^\alpha}\right)$$

but with right $\alpha$, this has better constants. "optimal"

Helps if ill-conditioned.

see "Robust SA" Nemirovsky

"Primal-Dual Avg" Nesterov 86

▷ STEP-SIZES  ( B.B., Wolfe, Armijo ... )

$c_1 \cong 10^{-4}$

• Sufficient Decrease (Armijo)  $f(x_k + t_k \cdot d_k) \le f(x_k) + c_1 \cdot t_k \cdot \langle d_k, p_k \rangle$

Wolfe Cond.

• Prevent short-steps ($t_k \to 0$) by with

directional deriv.

a) Curvature Conditions: $\langle \nabla f(x_k + t_k d_k), d_k \rangle \ge c_2 \cdot \langle \nabla f(x_k), d_k \rangle$

$c_2 \in (c_1, 1)$  (0.1 to 0.9)

or

Strong Wolfe: Armijo and  $|\langle \ldots \rangle| \le c_2 |\langle \ldots, \ldots \rangle|$

b) backtrack

Goldstein is another possibility, not for quasi-Newton methods