

LECTURE-5

Note Title

27-Jun-19

Faster algorithms (better than the black-box bounds)
for structured non-smooth problems

$$f(x) = \underbrace{g(x)} + \underbrace{h(x)}$$

Convex & Smooth

perhaps
Convex & non-smooth

Access to : $\arg \min_x \langle w, x \rangle + h(x) + \frac{1}{2\eta} \|x - y\|^2$
 $\forall w, \eta, y.$

$$f(x) = g(x) + h(x).$$

GD
with
Prox

Algorithm:

$$x_{t+1} = \operatorname{argmin}_x g(x_t) + \langle \nabla g(x_t), x - x_t \rangle + h(x) + \frac{1}{2\eta} \|x - x_t\|^2.$$

$$\text{GD} : x_{t+1} = \operatorname{argmin}_x g(x_t) + \langle \nabla g(x_t), x - x_t \rangle + h(x_t) + \langle \nabla h(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|^2.$$

LASSO: $f(x) = \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$

(Compressed Sensing)

$g(x)$ Linear regression
 $h(x)$ encourage sparsity of x .

$$\begin{aligned}
 x_{t+1} &= \arg \min_x \cancel{g(x_t)} + \langle \nabla g(x_t), x - \cancel{x_t} \rangle + h(x) + \frac{1}{2\eta} \|x - x_t\|^2 \\
 &= \arg \min_x \langle \nabla g(x_t), x \rangle + h(x) + \frac{1}{2\eta} \|x - x_t\|^2 \\
 &\quad \downarrow \\
 &\quad \lambda \|x\|_1
 \end{aligned}$$

$$= \operatorname{argmin}_{\boldsymbol{x}} \sum_{i=1}^d x_i \nabla g(x_t)_i + |x_i| + \frac{1}{2\eta} (x_i - (x_t)_i)^2$$

$$(x_{t+1})_i = \operatorname{argmin}_{x_i} x_i \nabla g(x_t)_i + |x_i| + \frac{1}{2\eta} (x_i - (x_t)_i)^2$$

$$= \operatorname{argmin}_{x_i} \frac{1}{2\eta} (x_i - (x_t)_i + \eta \nabla g(x_t)_i)^2 + |x_i|$$

Exercise:

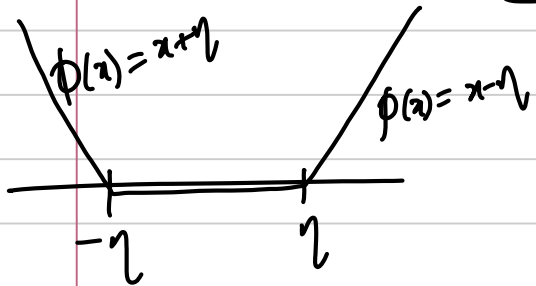
$$\begin{aligned}
 (x_{t+1})_i &= (x_t)_i - \eta \nabla g(x_t)_i - \eta && \text{if } \boxed{(x_t)_i - \eta \nabla g(x_t)_i} \geq \eta \\
 &= \quad \quad \quad \quad \quad \quad + \eta && \text{if } \quad \quad \quad \quad \quad \quad \leq -\eta \\
 &= 0 && \text{otherwise.}
 \end{aligned}$$

$= \Phi(x_t - \eta \nabla g(x_t))$, where Φ is applied to each coordinate.

$$x_{t+1} = \underbrace{x_t - \eta \nabla g(x_t)}_{\hat{x}_{t+1}} + \eta \mathbb{1}_{\{\text{Coordinates} < -\eta\}}$$

$$- \eta \mathbb{1}_{\{\text{Coordinates} > \eta\}}$$

$$- \hat{x}_{t+1} \mathbb{1}_{\{-\eta \leq \text{Coord.} \leq \eta\}}$$



$$\hat{f}_\eta(x_t; x) \triangleq \underbrace{g(x_t)} + \underbrace{\langle \nabla g(x_t), x - x_t \rangle} + \underbrace{h(x)} + \underbrace{\frac{1}{2\eta} \|x - x_t\|^2}$$

$x^* = \underset{x}{\operatorname{argmin}} f(x)$

$$x_{t+1} = \underset{x}{\operatorname{argmin}} \hat{f}_\eta(x_t; x)$$

g is L -Smooth
 $\eta \leq \frac{1}{L}$

$$f(x^*) = g(x^*) + h(x^*)$$

Conv. of $g(\cdot)$

$$\geq g(x_t) + \langle \nabla g(x_t), x^* - x_t \rangle + h(x^*)$$

$$+ \frac{1}{2\eta} \|x^* - x_t\|^2 - \frac{1}{2\eta} \|x^* - x_t\|^2$$

$$= \hat{f}_\eta(x_t; x^*) - \frac{1}{2\eta} \|x^* - x_t\|^2$$

$\left[\begin{array}{l} x_{t+1} = \underset{x}{\operatorname{argmin}} \hat{f}_\eta \\ \hat{f}_\eta - \frac{1}{\eta} \text{str.} \\ \text{cvx.} \end{array} \right] \geq$

$$\hat{f}_\eta(x_t; x_{t+1}) + \frac{1}{2\eta} \|x^* - x_{t+1}\|^2 - \frac{1}{2\eta} \|x^* - x_t\|^2$$

$$\hat{f}_\eta(x_t; x) = \underbrace{g(x_t) + \langle \nabla g(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|^2}_{g(x) + h(x)} + h(x)$$

L -smoothness
 of $g(\cdot)$
 & $\eta \leq \frac{1}{L}$

$$\geq g(x) + h(x) = f(x). \quad \forall x.$$

x_t x
 \uparrow \uparrow

smoothness lemma: $g(y) \leq g(x) + \langle \nabla g(x), y - x \rangle + \frac{L}{2} \|x - y\|^2$

$$f(x^*) \geq f(x_{t+1}) + \frac{1}{2\eta} \|x^* - x_{t+1}\|^2 - \frac{1}{2\eta} \|x^* - x_t\|^2.$$

Rearranging,

$$\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - 2\eta [f(x_{t+1}) - f(x^*)].$$

$$[\text{Telescope}] \leftarrow \leq \|x_0 - x^*\|^2 - 2\eta \sum_{s=1}^{t+1} [f(x_s) - f(x^*)].$$

$$\frac{1}{t+1} \sum_{s=1}^{t+1} [f(x_s) - f(x^*)] \leq \frac{\|x_0 - x^*\|^2 - \|x_{t+1} - x^*\|^2}{2\eta(t+1)}.$$

$\eta = \frac{1}{L}$

$$\leq \frac{L \|x_0 - x^*\|^2}{2(t+1)}.$$

When the non-smooth component is simple, we can
get $O\left(\frac{1}{T}\right)$ convergence rate. [ISTA]
Iterative Shrinkage
Thresholding Algorithm

We can in fact improve this to $O\left(\frac{1}{T^2}\right)$ using
Nesterov's AGD. [FISTA]

STOCHASTIC ALGORITHMS

$$1) \quad f(x) = \mathbb{E}_{(a,b)} \left[(a^T x - b)^2 \right]$$

$$\nabla f(x) = \mathbb{E}_{(a,b)} \left[(a^T x - b) a \right].$$

Oracle : (a_i, b_i)

$$\hat{\nabla} f(x) \triangleq (a_i^T x - b_i) a_i \Rightarrow \mathbb{E}[\hat{\nabla} f(x)] = \underline{\nabla f(x)}.$$

② Empirical Risk Minimization (ERM)

$$f(x) = \frac{1}{2n} \sum_{i=1}^n (a_i^\top x - b_i)^2$$

If $a_i \in \mathbb{R}^d$

$O(nd)$ time $\leftarrow \nabla f(x) = \frac{1}{n} \sum_{i=1}^n (a_i^\top x - b_i) a_i$

$O(d)$ time $\leftarrow \hat{\nabla} f(x) \triangleq (a_i^\top x - b_i) a_i$ where $i \sim \text{Unit}[1, \dots, n]$.

Setting:

$x \rightarrow \boxed{\text{oracle}} \rightarrow \hat{\nabla} f(x)$

$\hat{\nabla} f(x)$ independent of everything else $\left\{ \begin{array}{l} E[\hat{\nabla} f(x)] = \nabla f(x) \\ E[\|\hat{\nabla} f(x) - \nabla f(x)\|^2] \leq \sigma^2 \end{array} \right.$

SGD [Robbins & Monro]

$$x_{t+1} = x_t - \eta \hat{\nabla} f(x_t).$$

$$\begin{aligned} \mathbb{E}[\|x_{t+1} - x^*\|^2] &= \mathbb{E}[\|x_t - x^* - \eta \hat{\nabla} f(x_t)\|^2] \\ &= \mathbb{E}[\|x_t - x^*\|^2 - 2\eta \langle \hat{\nabla} f(x_t), x_t - x^* \rangle + \eta^2 \|\hat{\nabla} f(x_t)\|^2] \\ &= \mathbb{E}[\|x_t - x^*\|^2] - 2\eta \underbrace{\mathbb{E}[\langle \hat{\nabla} f(x_t), x_t - x^* \rangle]}_{=} + \eta^2 \underbrace{\mathbb{E}[\|\hat{\nabla} f(x_t)\|^2]}_{=} \end{aligned}$$

$$\|\nabla f(x)\| \leq G$$

f : Convex & G -Lipschitz
Random noise has bounded variance

$$1) \mathbb{E}[\langle \hat{\nabla} f(x_t), x_t - x^* \rangle | x_t] = \langle \mathbb{E}[\hat{\nabla} f(x_t) | x_t], x_t - x^* \rangle$$

$$= \langle \nabla f(x_t), x_t - x^* \rangle$$

[Convexity] $\Leftarrow \geq f(x_t) - f(x^*) \stackrel{\Delta}{=} \hat{\nabla} f(x_t) - \nabla f(x_t)$

$$2) \mathbb{E}[\langle \hat{\nabla} f(x_t), \hat{\nabla} f(x_t) \rangle | x_t] = \mathbb{E}[\langle \nabla f(x_t) + \eta_t, \nabla f(x_t) + \eta_t \rangle]$$

$$= \mathbb{E}[\|\nabla f(x_t)\|^2] + \mathbb{E}[\|\eta_t\|^2]$$

$$\leq \underbrace{G^2}_{\text{Since } f(\cdot) \text{ is } G\text{-Lipschitz}} + \underbrace{\sigma^2}_{\text{Since noise variance } \leq \sigma^2}$$

$$\text{So, } \mathbb{E}[\|x_{t+1} - x^*\|^2] \leq \mathbb{E}[\|x_t - x^*\|^2] - 2\eta [\mathbb{E}[f(x_t)] - f(x^*)] + \eta^2 (G^2 + \sigma^2).$$

Using the above inequality iteratively

$$\leq \mathbb{E}[\|x_0 - x^*\|^2] - 2\eta \sum_{s=0}^t [\mathbb{E}[f(x_s)] - f(x^*)] + \eta^2 (t+1) (G^2 + \sigma^2).$$

$$\frac{1}{t+1} \sum_{s=0}^t \mathbb{E}[f(x_s)] - f(x^*) \leq \frac{1}{2\eta(t+1)} [\|x_0 - x^*\|^2 - \mathbb{E}[\|x_{t+1} - x^*\|^2]]$$

$$+ \frac{\eta}{2} (G^2 + \sigma^2).$$

Choose $\eta = \sqrt{\frac{\|x_0 - x^*\|^2}{(t+1)(G^2 + \sigma^2)}}$

Av. expected Suboptimality $\leq \frac{\sqrt{(G^2 + \sigma^2)} \cdot \|x_0 - x^*\|}{\sqrt{t+1}}$.

function parameters \rightarrow Stochastic Oracle

For nonSmooth : this is the best possible rate $(\frac{1}{\sqrt{t}})$

For Smooth : Cannot improve on $\frac{1}{\sqrt{t}}$; but can obtain $\frac{L \|x_0 - x^*\|^2}{t^2} + \frac{\sigma \|x_0 - x^*\|}{\sqrt{t}}$.

Exercise: SGD for L -Smooth $f(\cdot)$ gets rate of

$$O\left(\frac{L\|x_0 - x^*\|^2}{t} + \frac{\sigma\|x_0 - x^*\|}{\sqrt{t}}\right).$$

Exercise: If f is G -Lipschitz and μ -strongly convex, then SGD [with diff. η] gets rate of

$$O\left(\frac{G^2 + \sigma^2}{\mu t}\right).$$