

Lecture 21 - catalyst

Tuesday, March 31, 2020 14:27

- today: finish prox example
- catalyst → accelerate
 - non-convex opt.
 - submodular opt.

finish prox example

* to be useful, need $\text{prox}_\gamma^\Omega$ to be efficiently computable

$$\text{prox}_\gamma^{\|\cdot\|_1}(z) = \underset{w}{\text{argmin}} \|w\|_1 + \frac{1}{2\gamma} \|w-z\|^2$$

"soft-thresholding" (component-wise) $\triangleq \begin{cases} \text{sgn}(z_i) [|z_i| - \gamma] & \text{if } |z_i| \geq \gamma \\ 0 & \text{o.w.} \end{cases}$

used esp. for lasso: ℓ_1 -reg. least-square

FISTA → accelerated prox. gradient method

↳ SOTA for batch lasso

* scikit-learn → uses SAGA for lasso $\left\{ \begin{array}{l} \ell_1\text{-reg.} \\ \text{log. reg.} \end{array} \right.$ $\min_w \frac{1}{n} \sum_i f_i(w) + \Omega(w)$

prox SAGA: $w_{t+1} = \text{prox}_\gamma^\Omega \left(w_t - \gamma \left[\nabla_{f_t}(w_t) - \bar{g}_t + \frac{1}{n} \sum_{j=1}^n \bar{g}_j^t \right] \right)$

could accelerate using "Catalyst"

Catalyst algorithm [Lin, Meinel & Hachem NIPS 2015]

"meta-algorithm": outer loop which uses a linearly convergent alg. in inner loop to get overall acceleration (↑)

main idea: use the accelerated proximal point algorithm

with approximation inner loop of prox operator

proximal pt. alg.: is proximal gradient with $f=0$

$$\boxed{w_{t+1} = \text{prox}_\gamma^\Omega(w_t)} \quad (\text{to solve } \min_w \Omega(w))$$

Catalyst alg. (for μ -strongly convex $F(w)$)

let $q \triangleq \frac{\mu}{\mu + \frac{1}{\gamma}}$ (γ is algorithmic parameter)

repeat:

$$w_{t+1} \approx \underset{w}{\operatorname{argmin}} \underbrace{F(w) + \frac{1}{2\gamma} \|w - z_t\|^2}_{\triangleq G_t(w)} \quad \text{s.t. } G_t(w_t) - \min_w G_t(w) \leq \epsilon_t$$

to be specified \downarrow

\uparrow $\operatorname{prox}_\gamma^{F(\cdot)}(z_t)$

using inner loop optimization alg. [e.g. SAGA or AFW]

$$z_{t+1} = w_{t+1} + \beta_{t+1} (w_{t+1} - w_t) \quad \text{[accelerated Nesterov trick piece]}$$

"extrapolation"

β_{t+1} is found using fancy equations so that everything works

• solve for α_{t+1} in eq. $\alpha_{t+1}^2 = (1 - \alpha_{t+1}) \rho_t^2 + q \alpha_{t+1}$

$$\beta_{t+1} \triangleq \frac{\alpha_t (1 - \alpha_t)}{\alpha_t^2 + \alpha_{t+1}} \quad (\text{pick } \alpha_{t+1} \in]0, 1[)$$

Catalyst trick: use γ & ϵ_t
s.t. overall # of inner loop calls
give overall acceleration

with clever analysis of warm starting

acceleration results:

if inner loop alg. has convergence $\exp(-\frac{\tilde{\mu}}{L} t)$ $\tilde{\mu} \geq \mu + \frac{1}{\gamma}$ (strong convexity $G_t(w)$)
then with correct constants

(μ -strongly convex F) linear rate: $\rho = \frac{1}{K}$ $\xrightarrow{\text{becomes}}$ $\tilde{\rho} = \frac{1}{\sqrt{K}}$ for catalyst
(F convex case) $\frac{1}{t}$ on F $\xrightarrow{\text{becomes}}$ $\frac{1}{t^2}$

results can get (theory) accelerated SAGA
" SVRG
" AFW
etc...



15h27

" AFW
etc.,



non-convex optimization

recall: FW with line search on f non-convex

$$\min_{S \subseteq \mathbb{R}^d} g(w_S) \leq O\left(\frac{1}{\sqrt{L}}\right)$$

FW gap

convex: $\mathbb{E} f(w_\epsilon) - f^* \leq \epsilon$

non-convex: $\sqrt{\mathbb{E} \|f(w_\epsilon)\|^2} \leq \epsilon$

gradient method: $f(w) \leq f(w_\epsilon) + \langle \nabla f(w_\epsilon), w - w_\epsilon \rangle + \frac{L}{2} \|w - w_\epsilon\|^2$

$$w_{t+1} = w_\epsilon - \frac{1}{L} \nabla f(w_\epsilon)$$

$$\Rightarrow f(w_{t+1}) \leq f(w_\epsilon) - \frac{1}{2L} \|\nabla f(w_\epsilon)\|^2$$

NIPS 2016 tutorial "Large-Scale Optimization: Beyond Stochastic Gradient Descent and Convexity"
[Suvri Sra slides](#)

Faster nonconvex optimization via VR

(Reddi, Hefny, Sra, Póczos, Smola, 2016; Reddi et al., 2016)

Algorithm	Nonconvex (Lipschitz smooth)
SGD	$O\left(\frac{1}{\epsilon^2}\right)$
GD	$O\left(\frac{n}{\epsilon}\right)$
SVRG	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$
SAGA	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$
MSVRG	$O\left(\min\left(\frac{1}{\epsilon^2}, \frac{n^{2/3}}{\epsilon}\right)\right)$

$$\mathbb{E}[\|\nabla g(\theta_t)\|^2] \leq \epsilon$$

Remarks

New results for convex case too; additional nonconvex results
 For related results, see also (Allen-Zhu, Hazan, 2016)

20

Linear rates for nonconvex problems

$$\min_{\theta \in \mathbb{R}^d} g(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

The Polyak-Łojasiewicz (PL) class of functions

$$g(\theta) - g(\theta^*) \leq \frac{1}{2\mu} \|\nabla g(\theta)\|^2$$

(Polyak, 1963); (Łojasiewicz, 1963)

Linear rates for nonconvex problems

$$g(\theta) - g(\theta^*) \leq \frac{1}{2\mu} \|\nabla g(\theta)\|^2 \quad \Bigg| \quad \mathbb{E}[g(\theta_t) - g^*] \leq \epsilon \quad \text{😄}$$

Algorithm	Nonconvex	Nonconvex-PL
SGD	$O\left(\frac{1}{\epsilon^2}\right)$	$O\left(\frac{1}{\epsilon^2}\right)$
GD	$O\left(\frac{n}{\epsilon}\right)$	$O\left(\frac{n}{2\mu} \log \frac{1}{\epsilon}\right)$
SVRG	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$	$O\left(\left(n + \frac{n^{2/3}}{2\mu}\right) \log \frac{1}{\epsilon}\right)$
SAGA	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$	$O\left(\left(n + \frac{n^{2/3}}{2\mu}\right) \log \frac{1}{\epsilon}\right)$
MSVRG	$O\left(\min\left(\frac{1}{\epsilon^2}, \frac{n^{2/3}}{\epsilon}\right)\right)$	—

Variant of **nc-SVRG** attains this fast convergence!

(Reddi, Hefny, Sra, Póczos, Smola, 2016; Reddi et al., 2016) 22

submodular optimization

submodularity is an analog of convexity for tractability

for set functions (combinatorial opt.)

$$F: 2^V \rightarrow \mathbb{R}$$

convention here: $F(\emptyset) = 0$

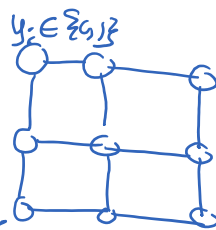
$V = \{1, \dots, d\}$ is "ground set"

$2^V = \{V \rightarrow \{0,1\}\} =$ set of all subsets of V

concrete example:

Ising model
 $y_i \in \{0,1\}$

$$E(y) = \sum_i c_i y_i - \sum_{\substack{i, j \text{ neighbor} \\ \text{of } i}} c_{ij} y_i y_j$$

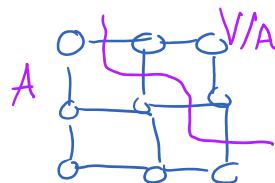


when $c_{ij} > 0$, $E(y)$ is submodular
"attractive potential"

MRF here is called "associative Markov network"
AMN

$F(A_y)$ where $A_y = \{i : y_i = 1\}$

can minimize $E(y)$ by using "graph cut" alg
(or $F(A_y)$)



F is submodular $\Leftrightarrow F(A) + F(B) \geq F(A \cap B) + F(A \cup B) \quad \forall A, B$

\Leftrightarrow function $A \mapsto F(A \cup \{k\}) - F(A)$ is non-increasing for all k

i.e. $F(A \cup \{k\}) - F(A) \leq F(B \cup \{k\}) - F(B)$
 $\forall B \supseteq A$

"diminishing return property"

\Rightarrow intuitively, that greedy alg are not "too bad" for maximization

* $F(A) \triangleq g(|A|)$ if g is concave then F is submodular
continuity

* link with convexity \rightarrow Sovisz extension (cts. fct.)

* embed sets as corners of hypercube in dimension $d \quad v(A) = \mathbb{1}_A \in \{0,1\}^d$

Sovisz extension f extends $F(A)$ from corners to entire hypercube using convex interpolation

(piecewise linear function on $[0,1]^d$)

$f(w) = F(A)$ when $w = v(A)$

F is submodular \Leftrightarrow Sovisz extension is convex

* can write $f(w) = \max_{S \in B(F)} \langle s, w \rangle$
 ρ
 "Base polytope"

\leftarrow this can be computed efficiently using greedy alg.

$\min_{A \subseteq V} F(A) = \min_{w \in [0,1]^d} \left(\max_{S \in B(F)} \langle s, w \rangle \right)$
 $f(w)$

\rightarrow use projected subgradient method

$\partial f(w) = \arg \max_{S \in B(F)} \langle s, w \rangle$

* with l_2 -regularization, use duality to get a smooth alg.

$\min_{S \in B(F)} \frac{1}{2} \|s\|^2$

\rightarrow use "min-norm pt." alg.

\odot SOTA for submodular opt.