

Lecture 21 - catalyst

Thursday, April 1, 2021 13:28

today : catalyst \rightarrow accelerate
 non-convex opt.
 submodular opt.

catalyst algorithm [Lin, Mairal & Harchaoui NeurIPS 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 point alg. : is proximal gradient with $S=0$

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

catalyst alg. (for μ -strongly $F(w)$)

Let $q \triangleq \frac{\mu}{\mu + L}$ (γ is algorithmic parameter)

repeat:

$$w_{t+1} \approx \arg \min_w F(w) + \frac{1}{2\gamma} \|w - z_t\|_2^2 \quad \text{s.t. } G_\epsilon(w_{t+1}) - \min_w G_\epsilon(w) \leq \epsilon_t$$

$\stackrel{\triangle}{=} G_\epsilon(w_{t+1})$
 $\stackrel{\triangle}{=} \text{prox}_{\gamma}^{F(\cdot)}(z_t)$
 with warm start
 use inner loop optimization [e.g. SAGA or AFU]

$$z_{t+1} = w_{t+1} + \underline{\beta_{t+1}} (w_{t+1} - w_t)$$

[accelerated
 Nesterov trick
 piece]

"extrapolation" / "momentum"

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

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

$$\beta_{t+1} \triangleq \frac{\alpha_t(1-\alpha_t)}{\alpha_t^2 + \alpha_{t+1}}$$

(pick $\alpha_{t+1} \in]0, 1[\Sigma$)

Catalyst trick: use $\gamma \in E_t$
 s.t. overall # of inner loop calls
 gives an overall acceleration

with clever analysis of warm starting

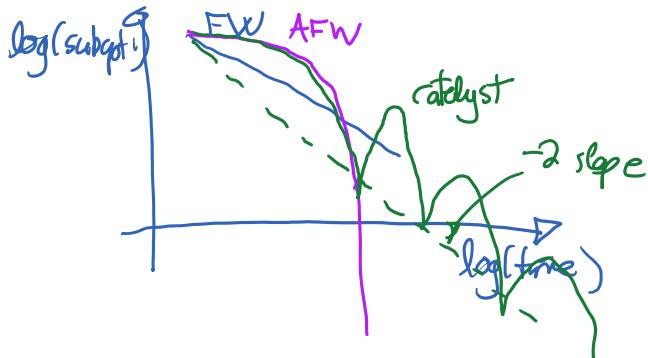
acceleration results:

if inner loop alg. has convergence exp $(-\tilde{\mu}t)$ $\tilde{\mu} \geq \mu + \frac{1}{\gamma}$
 then with correct constants for $\gamma \in E_t$

(μ -strongly conv F) linear rate $\beta = \frac{1}{K}$ becomes $= \frac{1}{\sqrt{K}}$ for catalyst
 (F convex case) $O(\frac{1}{t})$ on F with catalyst $O(\frac{1}{t^2})$

Result: we can get (theory) accelerated SAGA
 " SVRG
 " AFW
 etc.

Issue: catalyst is not adaptive to local strong convexity
 and finicky for choice of γ, E_t, μ etc ...



(4h15)

Non-convex optimization

recall: FW with line search on f non-convex

$$\min_{w \in \mathcal{S}} g(w) \in O\left(\frac{1}{\sqrt{t}}\right)$$

convex: $E f(w_t) - f^* \leq \epsilon$

FW gap

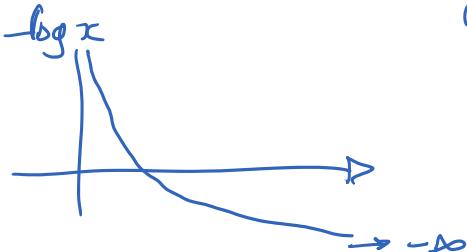
$$(GD) \Rightarrow \frac{1}{\epsilon} \quad \text{Nesterov } \frac{1}{\sqrt{\epsilon}}$$

non-convex: $\mathbb{E} \| \nabla f(w_t) \|_2^2 \leq \epsilon$

blow: if f is μ -strongly convex

$$\Rightarrow f(w_t) - f^* \leq \frac{1}{2\mu} \| \nabla f(w_t) \|_2^2$$

note: $\| \nabla f(w_t) \|$ small $\not\Rightarrow f(w_t) - f^*$ is small
when f is not strongly convex



* can get a $O(\frac{1}{\epsilon})$ rate for gradient descent:

$$f(w) \leq f(w_t) + \langle \nabla f(w_t), w - w_t \rangle + \frac{L}{2} \|w - w_t\|^2 \quad \forall w$$

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

(L -smoothness of f
but no need of
convexity)

$$\Rightarrow f(w_{t+1}) \leq f(w_t) - \frac{1}{2L} \| \nabla f(w_t) \|_2^2$$

if f^* is finite

$$\Rightarrow f(w_t) - f^* \leq f(w_0) - f^* - \frac{1}{2L} \sum_t \| \nabla f(w_t) \|_2^2$$

$$\Rightarrow \sum_t \| \nabla f(w_t) \|_2^2 \leq 2L(f(w_0) - f^*)$$

$$\Rightarrow t \cdot \min_{S.t.} \| \nabla f(w_t) \|_2^2 \leq 2L(f(w_0) - f^*)$$

i.e.
$$\boxed{\min_{S.t.} \| \nabla f(w_t) \|_2^2 \leq \frac{2L}{t} (f(w_0) - f^*)}$$

Faster nonconvex optimization via VR

Faster nonconvex optimization via VR

(Reddi, Hefny, Sra, Poczos, 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

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

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	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$	$O\left((n + \frac{n^{2/3}}{2\mu}) \log \frac{1}{\epsilon}\right)$
SAGA	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$	$O\left((n + \frac{n^{2/3}}{2\mu}) \log \frac{1}{\epsilon}\right)$
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, Poczos, Smola, 2016; Reddi et al., 2016) 22

Submodular optimization

Submodularity is an analog of convexity for tractability of set functions

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

(combinatorial opt.)

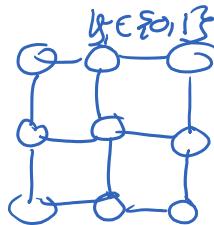
$V = \{1, \dots, d\}$ is "ground set"

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

Concrete example:

Ising model $E(y) = \sum_i \theta_i y_i$

$$-\sum_p \sum_{j \in \text{neighbor}} g_j y_j$$

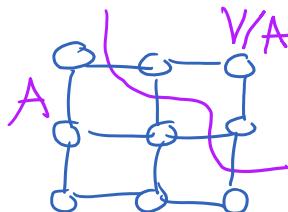


$$A_y = \{ i : y_i = 1 \} \rightarrow$$

$$F(A_4) = \dots$$

when $G_{ij} > 0$, $\Rightarrow E(y)$ is
"attractive potential"

can minimize $E(y)$ by using "graph cut algorithm"
(or $F(Ay)$)



equivalence
with min cost
network flow
problem

$$F \text{ is } \underline{\text{submodular}} \iff F(A) + F(B) \geq F(A \cap B) + F(A \cup B) \quad \forall A, B \subseteq V$$

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

i.e. $F(A \cup \Sigma_{k \geq 1}^c) - F(A) \leq F(B \cup \Sigma_{k \geq 1}^c) - F(B)$ if $B \subseteq A$
 "diminishing return property")

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

$$* F(A) \stackrel{\Delta}{=} g(|A|)$$

↑ cardinality

if g is concave \Rightarrow then F is submodular

* Link with convexity \rightarrow Lovasz extension (cts, fct) ACV

What is the best way to learn English well?

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

Szasz extension f extends $F(\cdot)$ from corners to entire hypercube using convex interpolation

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

$$f(w) = F(A) \text{ when } w = v(A)$$

$$\text{Let's say } w = \sum_i \alpha_i \cdot r_i \downarrow \\ v(A_{r_i}) \Rightarrow f(w) = \sum_i \alpha_i F(A_{r_i})$$

F is submodular \Leftrightarrow Szasz extension f is convex

$$* \text{ Can write } f(w) = \max_{s \in \text{B}(F)} \langle s, w \rangle$$

"Base polytope"

& this can be computed
efficiently using greedy alg

(LMO over $B(F)$ is
efficient)

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

→ use projected subgradient
method

$$\partial f(w) = \operatorname{argmax}_{s \in \text{B}(F)} \langle s, w \rangle$$

* with l_2 -regularization, use duality to get

a smooth obj

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

→ use "min-norm pt." alg
variant of FCFW alg. * SOTA
for submodular
opt.