

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

catalyst algorithm [Lin, Marcial & Huangzou NeurIPS 2015]

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

main idea: use the accelerated proximal point alg.

with approximation inner loop of prox operator

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

$$w_{t+1} = \text{prox}_{\frac{\lambda}{2}}^F(w_t)$$

[to solve $\min_w Q(w)$]

Catalyst alg.: (for μ -strongly $F(w)$)
repeat:
 $w_{t+1} \approx \underset{w}{\text{argmin}} F(w) + \frac{1}{2\gamma} \|w - z_t\|_2^2$ s.t. $G_t(w_{t+1}) - \min_w G_t(w) \leq \epsilon_t$
 $\in \text{prox}_\gamma^F(z_t)$ $\rightarrow (\gamma \text{ is algorithmic parameter})$
use inner loop optimization with warm start [e.g. SAGA, AFW, etc.]

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

[accelerated Nesterov trick piece]
"extrapolation" / "momentum"

Let $q \triangleq \frac{\mu}{\mu + 1}$; β_{t+1} is found using fancy equations so that everything works
• solve for α_{t+1} in eq.: $\alpha_{t+1}^2 = (1 - \alpha_{t+1})\alpha_t^2 + q\alpha_{t+1}$
(pick $\alpha_{t+1} \in]0, 1[$)

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

catalyst trick: use $\gamma \leq \epsilon_t$ s.t. overall # of inner loop calls gives an overall acceleration

with clever analysis of warm starting.

acceleration results

Strong convexity of $G_t(w)$

• if inner loop alg. has convergence $\exp(-\frac{\tilde{\mu}}{L}t)$ $\tilde{\mu} \geq \mu + \frac{1}{\gamma}$

then with correct constants for $\gamma \in \mathcal{E}_k$

(μ -strongly convex F) linear rate $P = \frac{1}{\lambda}$ $\xrightarrow[\text{with catalyst}]{} \approx \frac{1}{\sqrt{k}}$ for catalyst

(F convex case) $O(\frac{1}{t})$ on F $\xrightarrow{\text{becomes}} O(\frac{1}{\epsilon^2})$

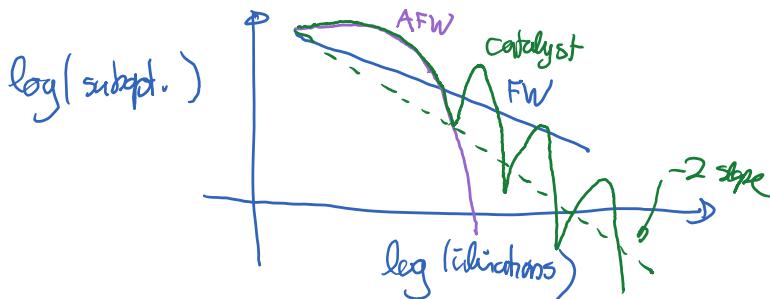
result: we can get (theory) accelerated SAGA

\downarrow SVRG

AFW

etc...

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



non-convex optimization

recall: FW with line search on f non-convex $\min_{w \in \mathcal{S}} g(w) \leq O(\frac{1}{\sqrt{t}})$
FW gap

$$\text{convex } \mathbb{E} f(w_t) - f^* \leq \epsilon$$

$$GD \rightarrow \frac{1}{\epsilon} \quad \text{Nesterov} \rightarrow \frac{1}{\epsilon}$$

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

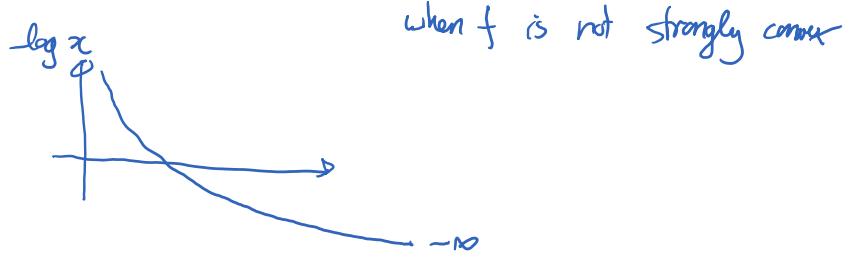
btw: if f is μ -strongly convex

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

note: $\|Df(w_t)\| \text{ small} \Rightarrow f(w_t) - f^* \text{ is small}$

$\log \chi_{\phi_i}$

when f is not strongly convex



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

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

$$w_{t+1} = w_t - \frac{1}{L} \nabla f(w_t) \quad (\text{L-smoothness of } f \text{ but no need of concavity})$$

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

If f^* is finite

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

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

$$\leq f(w_0) - f^* - \frac{1}{2L} \left(\sum_{s=0}^{t-1} \|\nabla f(w_s)\|_2^2 \right)$$

$$\Rightarrow \underbrace{\sum_{s=0}^{t-1} \|\nabla f(w_s)\|_2^2}_{\leq t} \leq 2L(f(w_0) - f^*)$$

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

$$\Rightarrow \boxed{\min_{s \leq t} \|\nabla f(w_s)\|_2^2 \leq \frac{2L}{t} (f(w_0) - f^*)}$$

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

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)

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

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

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

15h35

Submodular optimization

Submodularity is an analog of convexity/concavity for tractability of set functions
(combinatorial opt.)

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

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

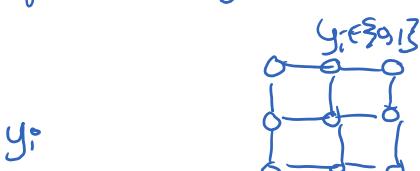
$2^V = \{V \subseteq \{1, \dots, d\}\}$ = set of all subsets of V

concrete example:

$$\text{Ising model} \quad E(y) = \sum_i \Theta_i y_i - \sum_{i,j \in \text{Neigh}(i)} \Theta_{ij} y_i y_j$$

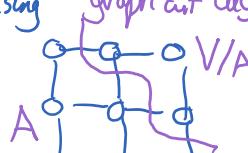
$$A_y = \{i : y_i = 1\}$$

$$F(A_y) = E(y)$$

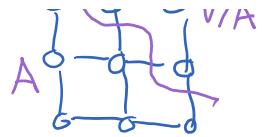


when $\Theta_{ij} \geq 0$ "attractive potential" $\Rightarrow E(y)$ is submodular
MRF here is "associative Markov network" (AMN)

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



equivalent with min cost network flow problem



network flow problem

F is submodular $\Leftrightarrow F(A) + F(B) \geq F(A \cap B) + F(A \cup B)$ for all $A, B \subseteq V$

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

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

"diminishing return property"

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

(a bit like concavity)

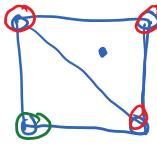
* $F(A) \stackrel{\text{def}}{=} g(|A|)$ if g is concave \Rightarrow then F is Submodular

* link with capacity \rightarrow Lovasz extension (cf. sat. ext.) $A \subseteq V$

\rightarrow embed sets as corners of hypercube in dimension d $V(A) = \mathbf{1}_A \in \{0, 1\}^d$

Lovasz extension f extends $F(\cdot)$ from corners to entire hypercube using Convex Interpolation

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



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

let's say $w = \sum_i \alpha_i v_i \downarrow \begin{matrix} \\ V(A_i) \end{matrix} \Rightarrow f(w) = \sum_i \alpha_i F(A_i) / f(v(A_i))$

F is submodular \Leftrightarrow Lovasz extension f is convex

(it turns out)

* can write $f(w) = \max_{S \in B(F)} \langle s, w \rangle$ \leftarrow this can be computed efficiently using sorting + greedy alg.
"base polytope" (LMO over $B(F)$ is efficient)

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

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

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

\rightarrow use "min-norm pt." alg. \oplus SOTA for variant of FCFW alg. \hookrightarrow submodular qf.

variant of FCFW alg.

submodular
gr.?