

- SAG for CRF
- proximal method, acceleration & non-convexity

SAG for CRF:

$$-\nabla f_i(w) = \sum_c \sum_{\tilde{y}_c} p(\tilde{y}_c|x^{(i)}, w) \Psi_i(\tilde{y}_c)$$

$$\Rightarrow \text{store } (p(\tilde{y}_c|x^{(i)}, w))_{\tilde{y}_c, c} \text{ in memory for SAG/SAGA correction}$$

uses other tricks:

- NUS
 - line search on step-size \Rightarrow this requires marginal inference "oracle call" for every step-size tried

SAG for CRF

OEG

SDCA
a stochastic dual coordinate ascent

$$w^{(t+1)} = (1-\gamma_t) w^{(t)} - \gamma_t [\nabla f_i(w^{(t)}) - \alpha_i^t + \frac{1}{n} \sum_j g_j^t]$$

$$\alpha_i^{(t)}(\tilde{y})^{(t+1)} \propto \alpha_i^{(t)}(\tilde{y})^{(t)} \exp(-\gamma_t \sum_j D(\alpha_j^{(t)}))$$

$$\text{do } \alpha_i^{(t+1)} = (1-\gamma_t) \alpha_i^{(t)}(\tilde{y})^{(t)} + \gamma_t \tilde{s}_i(\tilde{y})$$

$$w^{(t)} \xrightarrow{\text{KKT relationship}} \begin{bmatrix} \dots \\ \frac{1}{n} \Psi_i(\tilde{y}) \end{bmatrix}$$

get this using
"marginal inference"

[roughly: line search on $H(\alpha)$ + 2d line search]

Proximal gradient method:

- generalization of Armijo's gradient method for other non-smooth functions

→ generalization of projected gradient method to other non-smooth functions

composite framework : $F(w) = f(w) + \Omega(w)$ where f is convex and L -smooth

indicator of M but Ω is convex, not necessarily smooth

constrained opt. setting : $\Omega(w) = S_M(w) \triangleq \begin{cases} 0 & \text{if } w \in M \\ +\infty & \text{otherwise} \end{cases}$

l_1 -regularization : $\Omega(w) = \alpha \|w\|_1$

proximal gradient algorithm:

$$w_{t+1} = \arg \min_w f(w_t) + \langle \nabla f(w_t), w - w_t \rangle + \frac{L}{2\gamma_t} \|w - w_t\|^2 + \Omega(w)$$

$\underbrace{\Omega(w)}_{B_t(w)}$

if $\gamma_t \leq \frac{1}{L}$; then $f(w) \leq B_t(w) \quad \forall w$

we can rewrite $B_t(w) = \frac{1}{2\gamma_t} \|w - [w_t - \gamma_t \nabla f(w_t)]\|^2 + \text{cst.}$ (by completing the square)

aside:

optimality condition on $F(w)$:

$$w^* = \text{prox}_{\gamma}^{\Omega}(w^* - \gamma \nabla f(w^*))$$

$$w_{t+1} = \text{prox}_{\gamma_t}^{\Omega}(w_t - \gamma_t \nabla f(w_t))$$

inf condition

$$\left(\Omega \rightleftharpoons \inf_{w \in \mathbb{R}^d} \|w\|^2) \Rightarrow \inf_w \Omega(w) + \|w - z\|^2$$

"proximal operator": $\text{prox}_{\gamma}^{\Omega}(z) \triangleq \arg \min_w \left\{ \Omega(w) + \frac{1}{2\gamma} \|w - z\|^2 \right\}$

strongly convex
⇒ unique minimum

like for project, prox is non-expansive i.e. 1-Lipschitz

$$\text{i.e. } \|\text{prox}_{\gamma}^{\Omega}(w) - \text{prox}_{\gamma}^{\Omega}(w')\|_2 \leq \|w - w'\|_2$$

⇒ rates for unconstrained gradient method

transfer to proximal analog

* to be useful, we need prox_F^L to be efficiently computable

$$\text{prox}_F^L(z) = \underset{w}{\operatorname{arg\,min}} \|w\|_2 + \frac{1}{2\gamma} \|w - z\|^2$$

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

Catalyst algorithm [Lin, Marcin & Harchaoi NIPS 2015]

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

main idea: use the accelerated proximal point algorithm
with approximation in inner loop of prox operator

proximal point algorithm: is proximal gradient with $f \equiv 0$

$$w_{t+1} = \text{prox}_F^L(w_t)$$

Catalyst alg.: (for μ -strongly convex F)

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

repeat:

any w_{t+1}

$$w_{t+1} \approx \underset{w}{\operatorname{arg\,min}} F(w) + \frac{1}{2\gamma} \|w - z_t\|^2$$

s.t. $G_t(w_{t+1}) - \min_w G_t(w) \leq \varepsilon_t$

to be specified

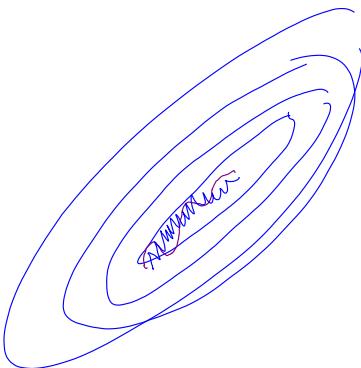
$\nabla G_t(\omega)$
 (using our inner alg. [e.g. SAGA FW])

(accelerated Nesterov trick piece) $Z_{t+1} = \omega_{t+1} + \beta_{t+1} (\omega_{t+1} - \omega_t)$

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

solve for α_{t+1} in equation $\alpha_{t+1}^2 = (1-\alpha_{t+1})\alpha_t^2 + q\alpha_{t+1}$
 (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 is: use γ and ε_b .

S.t. overall # of inner loop calls
 give acceleration

with clever analysis of warm starting
 result is if inner loop alg. has convergence $\exp(-\frac{1}{L} t)$ strong convexity of inner loop problem
 then with correct constants

go from $\frac{1}{\varepsilon} \rightarrow \frac{1}{\varepsilon^2}$ for outer loop convex minimization

$\frac{1}{K} \rightarrow \alpha \frac{1}{\sqrt{K}}$ for strongly-convex case

Non-convex optimization:

convex: $\mathbb{E} f(\mathbf{z}_t) - f^* \leq \varepsilon_b$

non-convex: $\mathbb{E} \|\nabla f(\mathbf{z}_t)\|^2 \leq \varepsilon_b$

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

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

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

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

NIPS 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

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

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!

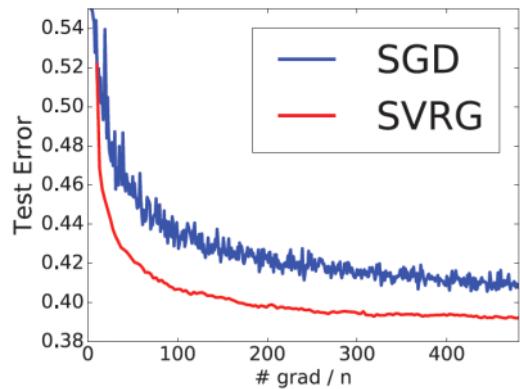
(Reddi, Hefny, Sra, Poczos, Smola, 2016; Reddi et al., 2016) 22

Suvrit Sra (mi.mit.edu)

Beyond stochastic gradients and convexity: Part 2



Empirical results



CIFAR10 dataset; 2-layer NN

25

Suvrit Sra (mi.mit.edu)

Beyond stochastic gradients and convexity: Part 2



Pointers

- SAG for CRF paper:
 - Non-Uniform Stochastic Average Gradient Method for Training Conditional Random Fields
M. Schmidt, R. Babanezhad, M.O. Ahmed, A. Defazio, A. Clifton, A. Sarkar
[AISTATS 2015](#) | [code](#)
- Proximal gradient method
 - see [slides](#) of this [great optimization class by L. Vandenberghe](#)
- Catalyst-- meta-algorithm for acceleration:
 - A Universal Catalyst for First-Order Optimization
Hongzhou Lin, Julien Mairal, Zaid Harchaoui
[NIPS 2015](#)
- non-convex optimization:
 - see [slides](#) from Suvrit Sra at the NIPS 2016 tutorial on "Large-Scale Optimization: Beyond Stochastic Gradient Descent and Convexity" (e.g. table of rates on p. 20 and 22)

Other good optimization pointers:

- great coverage of convex optimization by Mark Schmidt at the Machine Learning Summer School in 2015) - [slides](#) | [video](#)
- two great classes on optimization:
 - EE236C- Optimization Methods for Large-Scale Systems (Spring 2016) - Prof. L. Vandenberghe, UCLA - [link](#)
 - Convex optimization class by Rayn Tibshirani at CMU - Fall 2015 - [link](#)
 - [slides](#) on the Frank-Wolfe lecture
 - [nice summary table](#)
- tutorial by Francis Bach on SAG et al. at NIPS 2016 -- [slides](#)