

Perturbed Proximal-Gradient Algorithms

Gersende Fort

LTCI, CNRS, Telecom ParisTech, Université Paris-Saclay, 75013 Paris France

Joint work with Y. Atchadé (Univ. Michigan, USA) and Eric Moulines (Ecole Polytechnique, France)



Problem

How to minimize / find the minimum

- on a convex subset Θ of some finite dimensional Euclidean space with norm $\|\cdot\|$
- of a convex function $f : \theta \mapsto f(\theta)$ from Θ to \mathbb{R} , which is smooth enough:

$$\exists L > 0 \text{ s.t. } \forall \theta, \theta' \in \Theta, \quad \|\nabla f(\theta) - \nabla f(\theta')\| \leq L \|\theta - \theta'\|$$

- under non-smooth convex constraints $g : \theta \mapsto g(\theta)$ from Θ to $(-\infty, +\infty]$

when ∇f is not explicit ?

Problem 1: $\min_{\theta \in \Theta} (f(\theta) + g(\theta))$ **Problem 1':** $\operatorname{argmin}_{\theta \in \Theta} (f(\theta) + g(\theta))$

Example: Penalized Maximum Likelihood Inference in Latent Variable Models

The function f is "the log-likelihood of the observations Y ": $f(\theta) = -\log \int p(Y|x; \theta) \phi(x) \mu(dx)$

The feasible set: $\theta \in \Theta \subseteq \mathbb{R}^d$

The penalty term is a sparsity constraint: $g(\theta) = \lambda \sum_{i=1}^d |\theta_i|$ which is not a differentiable function

Non explicit gradient of f , but can be approximated:

$$-\nabla f(\theta) = \int \nabla_{\theta} (\log p(Y|x; \theta)) \pi_{\theta}(x|Y) \mu(dx) \approx \frac{1}{m} \sum_{k=1}^m \nabla_{\theta} \log p(Y|X_k; \theta)$$

where $(X_k)_k$ is from a MCMC with target $\pi_{\theta}(\cdot|Y) d\mu$, the cond. dist. of the latent variables X given Y .

The Proximal-Gradient Algorithm

Iterative algorithm: see [1] for convergence results

$$\begin{aligned} \theta_{n+1} &= \operatorname{argmin}_{\theta \in \Theta} \left(\gamma_{n+1} g(\theta) + \frac{1}{2} \|\theta - \{\theta_n - \gamma_{n+1} \nabla f(\theta_n)\}\|^2 \right) \\ &= \operatorname{Prox}_{\gamma_{n+1} g} (\theta_n - \gamma_{n+1} \nabla f(\theta_n)) = T_{\gamma_{n+1}, g}(\theta_n) \end{aligned}$$

Examples:

Projection on a closed convex set $\mathcal{K} \subseteq \Theta$

$$g(\theta) = \begin{cases} +\infty & \theta \notin \mathcal{K} \\ 0 & \theta \in \mathcal{K} \end{cases}$$

$$\theta_{n+1} = \operatorname{Proj}_{\mathcal{K}} (\theta_n - \gamma_{n+1} \nabla f(\theta_n))$$

Elastic net penalty

$$g(\theta) \propto \alpha \sum_{i=1}^d |\theta_i| + \frac{1-\alpha}{2} \|\theta\|^2$$

$$\theta_{n+1, i} = \text{shrinkage/thresholding of } (\theta_n - \gamma_{n+1} \nabla f(\theta_n))_i$$

↔ Unapplicable since $\nabla f(\theta_n)$ is not explicit in our framework

Questions: Can we replace $\nabla f(\theta_n)$ with an approximation while keeping the same asymptotic behavior ? How to choose the step-size γ_n ? In the Monte Carlo case, how to choose the (possibly) time-dependent batch-size m_n ?

REFERENCES

- [1] A. Beck, and M. Teboulle. *Gradient-based algorithms with applications to signal-recovery problems*. Convex Optimization in Signal Processing and Communications, 2009.
 [2] Y. Atchadé, G. Fort and E. Moulines. On Stochastic Proximal Gradient Algorithms *arXiv:1402:2365*, revised in Dec 2015.

The Perturbed Proximal Gradient Algorithm

Iterative algorithm:

$$\theta_{n+1} = \operatorname{Proj}_{\mathcal{K}} (\operatorname{Prox}_{\gamma_{n+1} g} (\theta_n - \gamma_{n+1} H_{n+1}))$$

where \mathcal{K} is a convex closed subset of Θ and H_{n+1} is a (possibly deterministic) approximation of $\nabla f(\theta_n)$.

Monte Carlo case: when $\nabla f(\theta) = \mathbb{E}_{\theta} [H_{\theta}(X)]$ with $X \sim \pi_{\theta}$

$$H_{n+1} = \frac{1}{m_{n+1}} \sum_{k=1}^{m_{n+1}} H_{\theta_n}(X_{n+1, k}) \quad \text{with Markov (or i.i.d) samples with inv. dist. } \pi_{\theta_n}$$

A General Convergence Result FROM [2, SECTION 3]

Set $\eta_{n+1} = H_{n+1} - \nabla f(\theta_n)$ $\mathcal{L} = \{\text{minimizers of } f + g\}$

Theorem. If $\gamma_n \in (0, 1/L]$, $\sum_n \gamma_n = +\infty$ and the following series converge

$$\sum_n \gamma_{n+1} \eta_{n+1}, \quad \sum_n \gamma_{n+1} \langle T_{\gamma_{n+1}, g}(\theta_n); \eta_{n+1} \rangle, \quad \sum_n \gamma_{n+1}^2 \|\eta_{n+1}\|^2$$

then there exists $\theta_{\infty} \in \mathcal{L}$ such that $\lim_n \tilde{\theta}_n = \theta_{\infty}$.

Other results. Explicit expression for U_n s.t.

$$(f + g)(\tilde{\theta}_n) - \min(f + g) \leq \sum_{k=1}^n \frac{a_k}{\sum_{t=1}^n a_t} (f + g)(\tilde{\theta}_k) - \min(f + g) \leq U_n \quad \text{with } \tilde{\theta}_n = \sum_{k=1}^n \frac{a_k}{\sum_{t=1}^n a_t} \tilde{\theta}_k$$

where a_1, \dots, a_n are non-negative real numbers

When applied to the Monte Carlo Proximal-Gradient Algorithm FROM [2, SECTION 4]

Under conditions on the Monte Carlo samples (geometric ergodicity, containment condition, ...):

$$\mathbb{E} [\|\eta_{n+1}\|^2 | \mathcal{F}_n] = O_{L^1} \left(\frac{1}{m_{n+1}} \right) \quad \|\mathbb{E} [\eta_{n+1} | \mathcal{F}_n]\| = O_{L^1} \left(\frac{1}{m_{n+1}} \right)$$

with fixed batch-size ($m_n = m$) but decreasing step-size γ_n s.t. $\sum_n \gamma_n = +\infty$ and $\sum_n \gamma_n^2 < \infty$,

the above convergence result

U_n is $O(1/\sqrt{n})$ for different choices of (a_k, γ_k) .

with increasing batch-size ($m_n \leq m_{n+1}$) at a linear rate $m_n \sim n$, and with constant step-size $\gamma_n = \gamma$

the above convergence result

U_n is $O(\ln/n)$ with a uniform weight $a_k = 1$; rate after $O(n^2)$ MC samples.

Conclusions: We provided sufficient conditions for

- (a) the same asymptotic behavior and the same rate of convergence as the exact algorithm,
 (b) which hold for both the cases of a **biased and unbiased approximation** H_{n+1}

