Pinning Down the Łojasiewicz Exponent for Structured Non-Convex Optimization

Anthony Man-Cho So (苏文藻) Department of Systems Engineering & Engineering Management The Chinese University of Hong Kong

(Joint Work w/ Huikang Liu (刘慧康) & Weijie Wu (吴玮 婕))

> MLA'15 7 November 2015

Non-Convex Opt: Something to Avoid?

Non-convex optimization

- typically considered difficult
- algorithms may get stuck in stationary points
- convergence rates of algorithms seem difficult to establish

Non-Convex Opt: Something to Avoid?

Non-convex optimization

- typically considered difficult
- algorithms may get stuck in stationary points
- convergence rates of algorithms seem difficult to establish

Structured non-convex optimization models

 arise in many applications (examples follow shortly)
 structure could be exploited for algorithm design and possibly also for analysis

 sometimes admit faster/cheaper computation than the corresponding convex approximations Structured Non-Convex Opt Models Low-Rank Matrix Recovery $\min_{X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r}} \{ \|A(XY^T) - b\|_F^2 + r(X, Y) \}$ Can be tackled by alternating optimization Compare with the convex approximation: $\min_{Z \in \mathbb{R}^m \times n} \{ \|A(Z) - b\|_F^2 + \|Z\|_* \}$ Standard algorithms (e.g., proximal gradient) require an SVD per iteration Wide range of applications

Structured Non-Convex Opt Models Non-negative Tensor Factorization $\min_{X_1,\dots,X_n > 0} \{ \|T - X_1 \circ \dots \circ X_n\|_F^2 + r(X_1,\dots,X_n) \}$ - Here, T is a tensor (multidimensional array) and \circ represents a particular tensor decomposition (e.g., CANDECOMP/PARAFAC or Tucker) Can be tackled by block coordinate descent methods Applications in signal and image processing

Structured Non-Convex Opt Models

Principal Component Analysis

 $\max_{X \in \mathbb{R}^{m \times n} : X^T X = I} \|A^T X\|_F^2$

Here, m >> n and we aim at recovering the top n left singular vectors of A

A fundamental matrix computation problem with long history and numerous fast algorithms

Analysis of Non-Convex Opt Problems

In view of

- prevalence of structured non-convex optimization models
- availability of fast algorithms for solving them
- a natural question is to understand the convergence behavior of those algorithms convergence (rate) properties of the limit point

 Find conditions that guarantee fast convergence of certain algorithms to global optimum
 assumptions on data (e.g., RIP, incoherence, etc.)

- and initialization of the algorithms typically required
- need to understand the behavior of the non-convex objective function around the global optima

Find conditions that guarantee fast convergence of certain algorithms to global optimum

- assumptions on data (e.g., RIP, incoherence, etc.) and initialization of the algorithms typically required
- need to understand the behavior of the non-convex objective function around the global optima
- focus of many recent efforts
 - (matrix completion) Jain et al.'13, Hardt'14, Sun-Luo'15
 - (dictionary learning) Agarwal et al.'14
 - (PCA/SVD) Shamir'15
 - **-** ...

Use the Łojasiewicz inequality to study the convergence behavior of descent methods
 apply to a wide range of objective functions (convex or non-convex) and descent methods
 no assumptions on the data or initialization of the algorithms required

 typically can only establish convergence to stationary points

- Use the Łojasiewicz inequality to study the convergence behavior of descent methods
 - apply to a wide range of objective functions (convex or non-convex) and descent methods
 - no assumptions on the data or initialization of the algorithms required
 - typically can only establish convergence to stationary points
 - focus of this talk

Łojasiewicz Inequality

• Let f be a real analytic function and x^* be one of its stationary points (i.e., $\nabla f(x^*) = 0$). Then, there exist $\delta, \eta > 0, \theta \in (0, 1/2]$ such that for all $y \in B(x^*, \delta)$,

 $\overline{|f(y) - f(x^*)|^{1-\theta}} \le \eta \|\nabla f(y)\|_2$

Thus, the Łojasiewicz inequality gives local growth information around stationary points *f* of .

 The growth rate is determined by the Łojasiewicz exponent

Łojasiewicz Inequality: An Example

Let $f: \mathbb{R}^n \to \mathbb{R}$ be given by $f(x) = ||x||_2^2$ and $x^* = 0$. Then, $\nabla f(x^*) = 0$ and

 $\|y\|_{2} = |f(y) - f(x^{*})|^{1 - (1/2)} \le \frac{1}{2} \|\nabla f(y)\|_{2} = \|y\|_{2}$

for all $y \in \mathbb{R}^n$

Łojasiewicz Inequality: Implications

(Absil et al.'05) Suppose that a bounded sequence of iterates $\{x_k\}$ satisfies Primary Descent: There exists $\sigma > 0$ such that for sufficiently large k, $f(x_{k+1}) - f(x_k) \le -\sigma \|\nabla f(x_k)\|_2 \|x_{k+1} - x_k\|_2$ • Stationarity: For sufficiently large k, $[f(x_{k+1}) = f(x_k)] \Rightarrow [x_{k+1} = x_k]$ Safeguard: There exists $\kappa > 0$ such that for sufficiently large k, $||x_{k+1} - x_k||_2 \ge \kappa ||\nabla f(x_k)||_2$

Łojasiewicz Inequality: Implications

Then, together with the Łojasiewicz inequality, the sequence {x_k} converges to a stationary point x* of f.

Łojasiewicz Inequality: Implications

Then, together with the Łojasiewicz inequality, the sequence {x_k} converges to a stationary point x* of f.
Moreover, the convergence rate can be estimated as follows:
(Sublinear convergence) If θ ∈ (0,1/2), then

 $\|x_{k} - x^{*}\|_{2} = O(k^{-\theta/(1-2\theta)})$ • (Linear convergence) If $\theta = 1/2$, then $\|x_{k} - x^{*}\|_{2} = e^{-O(k)}$

Many first-order methods for the unconstrained minimization of *f* generate iterates that satisfy the three properties.
 Real analytic functions include, e.g., all

polynomials, regardless of their convexity.

Many first-order methods for the unconstrained minimization of f generate iterates that satisfy the three properties. Real analytic functions include, e.g., all polynomials, regardless of their convexity. However, the Łojasiewicz exponent is difficult to estimate, even for some simple f ! In fact, the exponent is not known in most cases.

• (D'Acunto-Kurdyka'05) For a real polynomial $f: \mathbb{R}^n \to \mathbb{R}$ of degree $d \ge 2$, we have

 $\theta = \frac{1}{d(3d-3)^{n-1}}$

This leads to very weak convergence rate result.

• (D'Acunto-Kurdyka'05) For a real polynomial $f: \mathbb{R}^n \to \mathbb{R}$ of degree $d \ge 2$, we have

$$\theta = \frac{1}{d(3d-3)^{n-1}}$$

This leads to very weak convergence rate result.

 <u>Question</u>: For structured problems, could we get a sharp estimate of the Łojasiewicz exponent?
 This will lead to meaningful convergence rate results for a host of first-order methods. **Orthogonality-Constrained QPs** Consider the following problem (denoted OCQP): $\min_{X \in \mathbb{R}^{m \times n} : X^T X = I} \operatorname{Tr}(X^T A X B)$ non-convex objective function and constraint includes the PCA formulation as special case, as $||A^T X||_F^2 = \operatorname{Tr}(X^T A A^T X)$

Orthogonality-Constrained QPs Consider the following problem (denoted OCQP): $\min_{X \in \mathbb{R}^{m \times n} : X^T X = I} \operatorname{Tr}(X^T A X B)$ non-convex objective function and constraint includes the PCA formulation as special case, as $||A^T X||_F^2 = \operatorname{Tr}(X^T A A^T X)$ The feasible set is known as the Stiefel manifold, denoted by St(m, n). The manifold structure allows us to utilize techniques from manifold optimization.

 Instead of the usual gradient, we use the projected gradient, which is the usual gradient projected onto the tangent space to the manifold.

 For the Stiefel manifold, the projected gradient of f(X) = Tr(X^TAXB) is grad f(X) = 2AXB - XX^TAXB - XBX^TAX

Instead of the usual gradient, we use the projected gradient, which is the usual gradient projected onto the tangent space to the manifold.

For the Stiefel manifold, the projected gradient of f(X) = Tr(X^TAXB) is grad f(X) = 2AXB - XX^TAXB - XBX^TAX
 The set of stationary points is simply X = {X ∈ St(m, n):grad f(X) = 0}

The projected gradient allows us to treat (OCQP) like an unconstrained problem.

 (Schneider-Uschmajew'15) The convergence theorem of Absil et al.'05 carries over to the manifold setting if we replace the gradient in the Łojasiewicz inequality by the projected gradient.

The projected gradient allows us to treat (OCQP) like an unconstrained problem.

 (Schneider-Uschmajew'15) The convergence theorem of Absil et al.'05 carries over to the manifold setting if we replace the gradient in the Łojasiewicz inequality by the projected gradient.

Can be used to establish convergence rate results for a host of retraction-based line-search methods.

Lojasiewicz Inequality for (OCQP)

• (Liu-Wu-S.'15) There exist $\delta, \eta > 0$ such that for all $Y \in St(m,n), X^* \in X$ with $Y \in B(X^*, \delta)$, $|f(Y) - f(X^*)|^{1/2} \le \eta \|\text{grad } f(Y)\|_F$

Łojasiewicz Inequality for (OCQP)

(Liu-Wu-S.'15) There exist δ, η > 0 such that for all Y ∈ St(m, n), X* ∈ X with Y ∈ B(X*, δ), |f(Y) - f(X*)|^{1/2} ≤ η||grad f(Y)||_F
The starting point of the proof is to show that everyX* ∈ X can be factorized asX* = PQ, where columns of P (resp.Q) are eigenvectors of A (resp.B).

Łojasiewicz Inequality for (OCQP)

• (Liu-Wu-S.'15) There exist $\delta, \eta > 0$ such that for all $Y \in St(m, n), X^* \in X$ with $Y \in B(X^*, \delta),$ $|f(Y) - f(X^*)|^{1/2} \le \eta \|\text{grad } f(Y)\|_F$

The starting point of the proof is to show that everyX* ∈ X can be factorized asX* = PQ, where columns of P (resp.Q) are eigenvectors of A (resp.B).

The above inequality implies the linear convergence of any algorithm that satisfies primary descent, stationarity, and safeguard.

A Stiefel-SVRG for PCA

The PCA formulation can be rewritten as

$$\min_{X \in \mathbb{R}^{m \times n} : X^T X = I} - \operatorname{Tr} \left[X^T \left(\frac{1}{N} \sum_{i=1}^N a_i a_i^T \right) X \right]$$

where a₁,..., a_nre the columns of .A
The structure of the problem motivates the design of an SVRG-type algorithm that works on the Stiefel manifold.

A Stiefel-SVRG for PCA

• Let $f_i(X) = -\operatorname{Tr}(X^T a_i a_i^T X)$. Initialize $\tilde{X}^0 \in \operatorname{St}(m, n)$. • Outer loop s = 1, 2, ...• set $\tilde{X} = \tilde{X}^{s}^{-1}$, $\tilde{\mu} = \nabla f(\tilde{X}), X_0 = \tilde{X}$ **Inner loop** $k = 1, \dots, T$ ■ sample $i_k \in \{1, ..., N\}$ uniformly at random • compute $g_k = \nabla f_{i_k}(X_{k-1}) - \nabla f_{i_k}(\tilde{X}) + \tilde{\mu}$ • project g_k onto tangent space at X_{k-1} • apply retraction to get X_k • set $\tilde{X}^s = X_k$, where $k \in \{1, ..., T\}$ is chosen uniformly at random

A Stiefel-SVRG for PCA

 Using the Łojasiewicz inequality, we can prove (Liu-Wu-S.'15) Stiefel-SVRG converges linearly in expectation.

Conclusions

Our work represents a very preliminary step towards understanding the Łojasiewicz exponents associated with structured nonconvex optimization problems.

Many intriguing questions remain

- other structured non-convex (and non-smooth) models?
- implications on algorithm design?
- convergence to global optimum under the Łojasiewicz inequality-based framework?

Thank You!