Non-convex optimization when the solution is not unique: A kaleidoscope of favorable conditions

February 2023 Nicolas Boumal, with Quentin Rebjock OPTIM, Institute of Mathematics, EPFL



The paper, arXiv:2303.00096, treats optimization on manifolds. This talk specialized to \mathbb{R}^n . Funds: SEFRI/ERC GEOSYM

This talk is about fast local convergence

 $\min_{x\in\mathbf{R}^n}f(x)$

We use algorithms to compute $x \in \mathbf{R}^n$ such that $f(x) \in \mathbf{R}$ is small.

Ideally, we want a *global* minimum, but that's hard. *Local* will do. *x* is a local minimum if $f(x) \le f(y)$ for all *y* in a neighborhood of *x*.

Say f is C^2 (continuous Hessian). Focus on local convergence rates.

A simple look at the one-dimensional case

Say $f(x) = x^4$. Minimizer is $x^* = 0$.

Then $\nabla f(x) = 4x^3$ and $\nabla^2 f(x) = 12x^2$.

Gradient descent: $x_{k+1} = x_k - \alpha \nabla f(x_k) = (1 - 4\alpha x_k^2) x_k$. Only sublinear convergence (if $0 < \alpha < 1/2x_0^2$). Newton's method: $x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} [\nabla f(x_k)] = \frac{2}{3} x_k$. Only linear convergence.

The culprit: $\nabla^2 f(x^*) \ge 0$ can kill local rates

Indeed, if we assume a positive definite Hessian at x^* , then typical algorithms enjoy their "normal" fast local convergence rates.

This is what we find in classic optimization textbooks.

The issue is: for $f \in C^2$, at a critical point x^* ,

 $\nabla^2 f(x^*) > 0 \implies x^*$ is an isolated local minimum.

But quite often, minima are not isolated

And therefore, there is no way $\nabla^2 f(x^*) > 0$ for such applications.

Overparameterized regression / neural network (e.g., $\min_{x} ||F(x) - b||^2$) Redundant parameterization (e.g., $(L, R) \mapsto LR^{\mathsf{T}}$)

Symmetry (e.g., f(x) invariant to rotation)

Yet, we often still see fast convergence.

Why?



Insights from a simple 2-D example

$$f(x,y) = \frac{\mu}{2}x^2 \qquad \nabla f(x,y) = \begin{bmatrix} \mu x \\ 0 \end{bmatrix} \qquad \nabla^2 f(x,y) = \begin{bmatrix} \mu \\ 0 \end{bmatrix} \qquad 0$$

The set of minimizers is a line $S = \{(0, y) : y \in \mathbf{R}\}$, where $f^* = 0$. The gradient and Hessian "ignore" the tangent direction. As a result, typical algorithms only "see" the direction that matters.

$$\ker \nabla^2 f(0, y) = \{(0, u) : u \in \mathbf{R}\}$$

$$\text{All other eigenvalues} = \mu$$

$$f(x, y) - f^* = \frac{1}{2\mu} \|\nabla f(x, y)\|^2$$

$$\|\nabla f(x, y)\| = \mu \operatorname{dist}((x, y), \mathcal{S})$$

Take 1: Morse–Bott

$$\ker \nabla^2 f(0, y) = \{(0, u) : u \in \mathbf{R}\}$$

All other eigenvalues = μ
$$f(x, y) - f^* = \frac{1}{2\mu} \|\nabla f(x, y)\|^2$$
$$\|\nabla f(x, y)\| = \mu \operatorname{dist}((x, y), \mathcal{S})$$

If the minima are not singletons, maybe they are still "nice" sets. Say the set of minima S is an embedded submanifold of \mathbf{R}^{n} .

Surely, $\nabla f(x) = 0$ and $\nabla^2 f(x) \ge 0$ for $x \in S$. Also, $\nabla^2 f(x)[v] = 0$ if $v \in T_x S$.

Def.: *f* is MB if *S* is smooth and ker $\nabla^2 f(x) = T_x S$. Can then show good rates for various methods. We did not find many refs; see e.g. Fehrman, Gess & Jentzen 2020.

To avoid topological curiosities, here S is the set of global minimizers of f. The paper treats local minima.

 $\ker \nabla^2 f(0, y) = \{(0, u) : u \in \mathbf{R}\}$ All other eigenvalues = μ

$$f(x,y) - f^* = \frac{1}{2\mu} \|\nabla f(x,y)\|^2$$

Take 2: Polyak–Łojasiewicz $\int f(x,y) - f^* = \frac{\mu}{2} \operatorname{dist}((x,y), S)^2$

 $\|\nabla f(x,y)\| = \mu \operatorname{dist}((x,y),\mathcal{S})$

In 1963, Polyak studied f where gradient norm² ~ optimality gap. This is also called gradient dominance and is a particular case of Kurdyka–Łojasiewicz.

Def.:
$$f$$
 is PŁ at x^* if $f(x) - f(x^*) \le \frac{1}{2\mu} \|\nabla f(x)\|^2$ for x around x^* .

This includes strongly convex functions, and much more.

Linear cvgce for GD, and superlinear cvgce for cubic regularization. Polyak 1963; Nesterov & Polyak 2006. Many, many, many recent analyses of algorithms under PŁ.

Take 3: Quadratic Growth

$$\ker \nabla^2 f(0, y) = \{(0, u) : u \in \mathbf{R}\}$$

All other eigenvalues = μ
$$f(x, y) - f^* = \frac{1}{2\mu} \|\nabla f(x, y)\|^2$$

$$f(x, y) - f^* = \frac{\mu}{2} \operatorname{dist}((x, y), \mathcal{S})^2 \qquad \|\nabla f(x, y)\| = \mu \operatorname{dist}((x, y), \mathcal{S})$$

We could also assume that f grows fast as we move away from S. Already in Bonnans & Ioffe 1995, but likely much older.

Def.: f has QG at x^* if $f(x) - f(x^*) \ge \frac{\mu}{2} \operatorname{dist}(x, \mathcal{S})^2$ for x around x^* .

Interestingly, this one is well defined even if f is nonsmooth. This has been used to show fast convergence in that setting. Drusvyatskiy & Lewis 2016; Davis & Jiang 2022; Lewis & Tian 2022.

Take 4: Error Bound

$$\ker \nabla^2 f(0, y) = \{(0, u) : u \in \mathbf{R}\}$$

All other eigenvalues = μ
$$f(x, y) - f^* = \frac{1}{2\mu} \|\nabla f(x, y)\|^2$$
$$\|\nabla f(x, y)\| = \mu \operatorname{dist}((x, y), \mathcal{S})$$

We could assume the gradient grows fast as we move away from S. This seems to have originated with Luo & Tseng 1993.

Def.: f has **EB** at x^* if $||\nabla f(x)|| \ge \mu$ dist(x, S) for x around x^* .

Luo & Tseng showed linear convergence for several methods with EB.

Yue, Zhou & Man-Cho So 2019 showed quadratic cvgce for cubic regularization.

Some conditions imply others for $f \in C^2$



Morse–Bott is explicitly strong, partly because it requires a smooth solution set S. With a few simple Taylor expansion arguments, we can see MB \Rightarrow QG, EB, PŁ. 11

Some conditions imply others for $f \in C^2$



* $PL \Rightarrow QG$ has a rich history: see Ioffe 2000, Bolte et al. 2017 (arXiv 2015); Drusvyatskiy et al. 2015, Karimi et al. 2016, Zhang 2017.

Actually, all conditions coincide for $f \in C^2$! $[f \in C^1 \text{ and PL}] \Rightarrow \text{smooth } S$ $f(x, y) = \frac{x^2 y^2}{x^2 + y^2}$ C^2 Morse Polyak Łojasiewicz Bott C^1 with some constant, C^2 with better constant. C² Quadratic Error Growth Bound C^2 $[f \in C^1 \text{ and } QG] \Rightarrow EB, PŁ$ 13 $f(x) = 2x^2 + x^2 \sin(1/\sqrt{|x|})$

Technical details of note

 $\ker \nabla^2 f(x) = T_x S$ All other eigenvalues $\ge \mu$ $f(x) - f^* \le \frac{\mu}{2} \operatorname{dist}(x, S)^2$ $\|\nabla f(x)\| \ge \mu \operatorname{dist}(x, S)$

PŁ, QG and EB hold *in a neighborhood*. It may shrink along " \Rightarrow ".

If one holds with μ , all hold for all $\mu' < \mu$ (trade-off with ngbhd).

For $f \in C^k$ with $k \ge 2$, we have $P_{L} \Rightarrow MB$ with S of class C^{k-1} .

$PŁ \Rightarrow MB$: Elements of proof

The most interesting bit is to show: $PL \Rightarrow S$ smooth

Pick $\bar{x} \in S$. Let $P(\bar{x})$ be the projector to the image of $\nabla^2 f(\bar{x})$, and:

$$\mathcal{Z} = \{x \text{ close to } \bar{x} : P(\bar{x}) \nabla f(x) = 0\}$$

Clearly, \mathcal{Z} contains \mathcal{S} (locally).

Also, Z is a smooth submanifold: study rank of $P(\bar{x})\nabla^2 f(x)$. And we can show that PŁ implies Z = S (locally).

arXiv:2303.00096

Take away for non-isolated minima

 $\ker \nabla^2 f(x) = T_x S$ All other eigenvalues $\ge \mu$ $f(x) - f^* \le \frac{\mu}{2} \operatorname{dist}(x, S)^2$ $\|\nabla f(x)\| \ge \mu \operatorname{dist}(x, S)$

If f is C², those four conditions are equivalent.
Thus, assuming one of MB, PŁ, QG or EB, we can use all.
This helps analysis. In our paper: cubic regularization, trust regions.

E.g.: Nesterov & Polyak '06 compared to Yue et al. '19 for cubic regularization.