Accelerated gradient descent for large-scale optimization

On Gradient descent solving quadratic problems

Andersen Ang	
ECS, Uni. Southampton, UK andersen.ang@soton.ac.uk	Content
Homepage angms.science	Introduction & preliminaries Gradient Descent as an iterative algorithm Gradient Descent is slow when level set is elliptic Accelerated gradient method
Version: March 13, 2025	Nonnegative least squaress Adaptive restarts
First draft: November 12, 2018 Guest lecture of MARO 201 - Advanced Optimization Faculte Polytechnique de Mons	Convergence rate of Gradient Descent on (\mathcal{P}) Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality Convergence rate of accelerated gradient method on (\mathcal{P}) Convergence of accelerated gradient method explained using ODE

Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Introduction

► Theme: solve

$$(\mathcal{P})$$
: given $A \in \mathbb{R}^{m imes n}, b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ by solving
 $\operatorname*{argmin}_{x} f(x) := rac{1}{2} \|Ax - b\|_2^2.$

▶ 1st-order optimality condition $\nabla f(x) = 0$ gives

$$oldsymbol{A}^ opoldsymbol{A} oldsymbol{A} - oldsymbol{A}^ opoldsymbol{b} \ = \ oldsymbol{A} oldsymbol{A} - oldsymbol{b} \ = \ oldsymbol{0} \in \mathbb{R}^m$$

we denote $\nabla f(\overline{x})$ the gradient vector of f in the standard basis in \mathbb{R}^n with respect to the variable x at the point \overline{x}

• Solving (\mathcal{P}) is the same as solving linear inverse problem: given

ven
$$egin{array}{c|c} oldsymbol{A} \in \mathbb{R}^{m imes n} \ oldsymbol{b} \in \mathbb{R}^m \end{array}$$
 , solve $oldsymbol{A} oldsymbol{x} = oldsymbol{b}$.

• Linear algebra solution: $x = A^{-1}b$ if A^{-1} exists.

This lecture: constrained convex quadratic problem

$$(\mathcal{P})$$
: given $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$, find $x \in \mathcal{C} \subset \mathbb{R}^n$ by solving
 $\operatorname*{argmin}_{x} f(x) := \frac{1}{2} \|Ax - b\|_2^2.$

- \blacktriangleright We focus on the following simple convex set ${\mathcal C}$
 - $C = \mathbb{R}^n$ (no constraint), (P) = least squares (LS).
 - $C = \mathbb{R}^n_+$ (nonnegative constraint), (\mathcal{P}) = Nonnegative LS (NNLS).
- ► Some other C (not the focus) :

$$\blacktriangleright \ \mathcal{C} = \left\{ \boldsymbol{x} \in \mathbb{R}^n \ \big| \ \| \boldsymbol{x} \|_2 \leq \epsilon \right\} \ (l_2 \text{-norm constraint}).$$

$$\blacktriangleright \ \mathcal{C} = \left\{ \boldsymbol{x} \in \mathbb{R}^n \ \big| \ \| \boldsymbol{x} \|_1 \leq \epsilon \right\} \ (l_1 \text{-norm constraint}).$$

$$\blacktriangleright \ \mathcal{C} = \left\{ \boldsymbol{x} \in \mathbb{R}^n \ \big| \ \| \boldsymbol{x} \|_0 \leq \epsilon \right\} \ (l_0 \text{-norm constraint}).$$

$$\blacktriangleright \ \mathcal{C} = \left\{ \boldsymbol{x} \in \mathbb{R}^n \ \big| \ \| \boldsymbol{x} \|_{\mathrm{TV}} \leq \epsilon \right\} \text{ (Total variation constraint).}$$

$$\blacktriangleright \ \mathcal{C} = \left\{ \boldsymbol{x} \in \mathbb{R}^n \ \big| \ \boldsymbol{C} \boldsymbol{x} \leq \boldsymbol{d} \right\} \text{ (polytope constraint)}$$

Why study (\mathcal{P})

$$\begin{array}{ll} (\mathcal{P}): \text{ given } \boldsymbol{A} \in \mathbb{R}^{m \times n}, \boldsymbol{b} \in \mathbb{R}^m, \text{ find } \boldsymbol{x} \in \mathcal{C} \subset \mathbb{R}^n \text{ by solving} \\ & \operatorname*{argmin}_{\boldsymbol{x}} f(\boldsymbol{x}) \ \coloneqq \ \frac{1}{2} \| \boldsymbol{A} \boldsymbol{x} - \boldsymbol{b} \|_2^2. \end{array}$$

1. (\mathcal{P}) is easy to understand.

We will look at gradient descent, and the accelerated variants through the lens of (\mathcal{P}) .

- (P) is applicable to wide range of situations.
 We see (P) when we take 2nd-order Taylor's approximation of a model.
- 3. (\mathcal{P}) is easy to generalize/modify.

Facts: basic properties of quadratic problem (\mathcal{P})

- ▶ (\mathcal{P}) is a convex problem: local minima → global minima
 - C is a convex set [assumption]
 - cost function f is convex

$$\bullet f(\boldsymbol{x}) = \frac{1}{2}\boldsymbol{x}^{\top}\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p}^{\top}\boldsymbol{x} + c, \text{ where } \boldsymbol{Q} = \boldsymbol{A}^{\top}\boldsymbol{A}, \ \boldsymbol{p} = \boldsymbol{A}^{\top}\boldsymbol{b}, \ c = \frac{1}{2}\|\boldsymbol{b}\|_{2}^{2}.$$

(𝒫) has a unique global minimum x^{*}
 if A full rank [assumption] ⇒ Q is positive definite ⇔ λ_{min}(Q) > 0.

▶ Gradient
$$\nabla f(\boldsymbol{x}) = \boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p} = \boldsymbol{A}^{\top}\boldsymbol{A}\boldsymbol{x} - \boldsymbol{A}^{\top}\boldsymbol{b}$$

- f is L-smooth with $L = \|\boldsymbol{Q}\|_2 = \sqrt{\lambda_{\max}(\boldsymbol{Q}^{\top}\boldsymbol{Q})} = \lambda_{\max}(\boldsymbol{Q}) = \sigma_{\max}(\boldsymbol{A}).$
 - ► A function g(x) is L-smooth if the gradient ∇g is L-Lipschitz : $\|\nabla g(a) \nabla g(b)\| \le L \|a b\| \ \forall a, b \in \text{dom}g$

$$\begin{split} \bullet & \|\boldsymbol{Q}\|_2 = \sqrt{\lambda_{\max}((\boldsymbol{A}^{\top}\boldsymbol{A})^2)} = \sqrt{\lambda_{\max}^2(\boldsymbol{A}^{\top}\boldsymbol{A})} = \lambda_{\max}(\boldsymbol{A}^{\top}\boldsymbol{A}) = \sigma_{\max}(\boldsymbol{A}) \\ & \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\|_2 = \|(\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p}) - (\boldsymbol{Q}\boldsymbol{y} - \boldsymbol{p})\|_2 \\ & = \|\boldsymbol{Q}(\boldsymbol{x} - \boldsymbol{y})\|_2 \\ & \leq \|\boldsymbol{Q}\|_2 \|\boldsymbol{x} - \boldsymbol{y}\|_2. \end{split}$$
oni: operator norm inequality $\|\boldsymbol{A}\boldsymbol{x}\| \leq \|\boldsymbol{A}\| \|\boldsymbol{x}\|. \end{split}$

Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Iterative algorithm to solve (\mathcal{P})

▶ Iterative solver: produce a sequence $\{x_k\}_{k\in\mathbb{N}} = \{x_1, x_2, \dots\}$ such that

 $f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k)$ for all k. (Descent condition)

• How: starting with an initial guess x_0 , perform the iteration/update



until a stopping condition is met.

- Two remarks
 - 1. Δ is called "direction" but it can have magnitude larger than 1. The true "stepsize" is $t \|\Delta\|$, the true "direction" is $\Delta \|\Delta\|^{-1}$.
 - 2. We focus on algo., but choosing x_0 also affect the performance of algo. ightarrow a topic on its own.

Stopping condition (not the focus here)

Туре	order	Stop if
Variable change	_	$\operatorname{dist}({oldsymbol x}_k,{oldsymbol x}_{k-1})\ \leq\ \epsilon$
Functional value	0th	$f(oldsymbol{x}_k)~\leq~\epsilon$
Successive functional value	0.5th	$ig f(oldsymbol{x}_k) - f(oldsymbol{x}_{k-1})ig \ \le \ \epsilon$
Gradient value	1st	$\ abla f(oldsymbol{x}_k)\ _2 ~\leq~ \epsilon$
Hessian value	2nd	$\lambda_{\min} \Big(abla^2 f(oldsymbol{x}_k) \Big) \; \geq \; \epsilon > 0$

Which one to use and what are their pros & cons will take another hour to explain.

I call successive functional value as "0.5th" order because $|f(\boldsymbol{x}_k) - f(\boldsymbol{x}_{k-1})|$ resembles $||\nabla f||$: by the fact that f is convex and L-smooth, we have that for $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)$, its holds

$$f(\boldsymbol{x}_{k+1}) \ge f(\boldsymbol{x}_k) + \langle \nabla f(\boldsymbol{x}_k), t_k \nabla f(\boldsymbol{x}_k) \rangle + \frac{L}{2} \| t_k \nabla f(\boldsymbol{x}_k) \|_2^2$$

Hence

$$\left|f(\boldsymbol{x}_{k+1}) - f(\boldsymbol{x}_{k})\right| \geq \left(t_{k} + \frac{L}{2}t_{k}^{2}\right) \|\nabla f(\boldsymbol{x}_{k})\|_{2}^{2}$$

Gradient descent (GD)

- GD uses $\Delta = -\nabla f(\boldsymbol{x})$ so we have $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k t\nabla f(\boldsymbol{x}_k)$.
- Choices of stepsize *t*
 - 1. Optimal stepsize / exact line search: $t_{\text{Exact}} = \underset{t \ge 0}{\operatorname{argmin}} f(\boldsymbol{x} + t\Delta)$. Practically the best but can be expensive to compute.
 - 2. Line search stepsize: backtracking / sufficient descent / Armijo Rule / Wolfe condition. Practically okay and commonly used.
 - Constant stepsize: t is fix.
 Practically not the best, but theoretically easier to analyze.
 - 4. * Constant stepsize using Lipschitz constant: $t_L = \frac{1}{L(f)}$. A common choice of stepsize for theorist.
 - 5. Diminishing stepsize: $t_{k+1} = t_k \cdot \theta$, $\theta < 1$, or $t_k = \frac{1}{k}$. Good for stochastic gradient.

Our focus: 1 and 4.

GD with optimal stepsize on (\mathcal{P}) : $\underset{\boldsymbol{x}}{\operatorname{argmin}} f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^{\top} \boldsymbol{Q} \boldsymbol{x} - \boldsymbol{p}^{\top} \boldsymbol{x}$

$$t_{\mathsf{Exact}} = \operatorname*{argmin}_{t \ge 0} f(\boldsymbol{x} - t\nabla f(\boldsymbol{x})) = \frac{1}{\nabla f(\boldsymbol{x})^\top \boldsymbol{Q} \nabla f(\boldsymbol{x})} = \frac{1}{\|\nabla f(\boldsymbol{x})\|_{\boldsymbol{Q}}^2}$$

Derivation: reading exercise

$$f(\boldsymbol{x}) = \frac{1}{2}\boldsymbol{x}^{\top}\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p}^{\top}\boldsymbol{x}. \text{ To get } t_{\text{Exact}} = \underset{t \ge 0}{\operatorname{argmin}} g(t) = f(\boldsymbol{x} - t\nabla f(\boldsymbol{x})), \text{ set } \frac{\partial g(t)}{\partial \boldsymbol{x}} = 0.$$

$$g(t) = \frac{1}{2} (\boldsymbol{x} - t\nabla f(\boldsymbol{x}))^{\top}\boldsymbol{Q} (\boldsymbol{x} - t\nabla f(\boldsymbol{x})) - \boldsymbol{p}^{\top} (\boldsymbol{x} - t\nabla f(\boldsymbol{x})), \text{ ignore terms without } t:$$

$$g(t) = \frac{1}{2} (-\boldsymbol{x}^{\top}\boldsymbol{Q}t\nabla f(\boldsymbol{x}) - t\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\boldsymbol{x} + t^{2}\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x})) + t\boldsymbol{p}^{\top}\nabla f(\boldsymbol{x})$$

$$= -\boldsymbol{x}^{\top}\boldsymbol{Q}t\nabla f(\boldsymbol{x}) + \frac{1}{2}t^{2}\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x}) + t\boldsymbol{p}^{\top}\nabla f(\boldsymbol{x})$$

$$\frac{\partial g(t)}{\partial t} = -\boldsymbol{x}^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x}) + t\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x}) + \boldsymbol{p}^{\top}\nabla f(\boldsymbol{x})$$

$$= -(\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p})^{\top}\nabla f(\boldsymbol{x}) + t\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x})$$

$$= -\nabla f(\boldsymbol{x})^{\top}\nabla f(\boldsymbol{x}) + t\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x}).$$
Finally, $\frac{\partial g(t)}{\partial \boldsymbol{x}} = 0 \implies \nabla f(\boldsymbol{x})^{\top}\nabla f(\boldsymbol{x}) = t\nabla f(\boldsymbol{x})^{\top}\boldsymbol{Q}\nabla f(\boldsymbol{x}).$

GD algo. with t_{Exact} on (\mathcal{P}) : $\underset{\boldsymbol{x}}{\operatorname{argmin}} f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^{\top} \boldsymbol{Q} \boldsymbol{x} - \boldsymbol{p}^{\top} \boldsymbol{x}, \ \boldsymbol{Q} = \boldsymbol{A}^{\top} \boldsymbol{A}, \ \boldsymbol{p} = \boldsymbol{A}^{\top} \boldsymbol{b}$

$$\begin{aligned} \mathsf{GD} \ \mathsf{update}: & \boldsymbol{x}_{k+1} &= \boldsymbol{x}_k - t_{\mathsf{Exact}} \nabla f(\boldsymbol{x}_k) \\ &= \boldsymbol{x}_k - \frac{\|\nabla f(\boldsymbol{x}_k)\|_2^2}{\|\nabla f(\boldsymbol{x}_k)\|_{\boldsymbol{Q}}^2} \nabla f(\boldsymbol{x}_k) \\ &= \boldsymbol{x}_k - \frac{\|\boldsymbol{Q}\boldsymbol{x}_k - \boldsymbol{p}\|_2^2}{\|\boldsymbol{Q}\boldsymbol{x}_k - \boldsymbol{p}\|_{\boldsymbol{Q}}^2} (\boldsymbol{Q}\boldsymbol{x}_k - \boldsymbol{p}) \end{aligned}$$

Algorithm 1: GD (with t_{Exact}) for (\mathcal{P})

Result: A sol. x that approximately solves (\mathcal{P}) 1 Initialization Set $x_0 \in \mathbb{R}^n$, $p = A^{\top}b$, $Q = A^{\top}A$

2 while stopping condition is not met do

$$egin{array}{ccc} \mathbf{3} & egin{array}{c} egin{array}{c} egin{array}{c} \mathbf{g} & egin{array}{c} \mathbf{g} & \mathbf{g} \end{array} \end{bmatrix} & \mathbf{g} & \mathbf{g} & \mathbf{g} \end{bmatrix} & \mathbf{g} & \mathbf{g} \end{bmatrix} & \mathbf{g} & \mathbf{g} \end{bmatrix} = \mathbf{g} + \mathbf{g} +$$

5 end



13/92

Observations

- If the level set of f(x) is circular, GD goes to x^* very fast. (In fact, in 1 step)
- When the level set of f(x) is elliptic, GD zigzags (and slow).



Questions

- ► Why zigzags? Where does this zigzag come from?
- ► How to deal with it: how to improve GD?

Remarks on GD with t_{Exact}

• GD update with t_{Exact} :

$$\boldsymbol{x} = \boldsymbol{x} - \frac{\|\nabla f(\boldsymbol{x})\|_2^2}{\|\nabla f(\boldsymbol{x})\|_{\boldsymbol{Q}}^2} \nabla f(\boldsymbol{x}) = \boldsymbol{x} - \frac{(\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p})^\top (\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p})}{(\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p})^\top \boldsymbol{Q} (\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p})} (\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{p}).$$

- The step size is not constant, it depends on $x \rightarrow$ hard to analyze.
- The step size contains many vector/matrix-vector/matrix products \rightarrow not suitable on problem with big m, n.
- We now consider GD with fixed step size t_L .

GD alg. with fixed $t_L = \frac{1}{L}$ on (\mathcal{P}) : $\operatorname*{argmin}_{\boldsymbol{x}} f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^\top \boldsymbol{Q} \boldsymbol{x} - \boldsymbol{p}^\top \boldsymbol{x}, \ \boldsymbol{Q} = \boldsymbol{A}^\top \boldsymbol{A}, \ \boldsymbol{p} = \boldsymbol{A}^\top \boldsymbol{b}$

GD update:
$$x^+ = x - \frac{1}{\|Q\|_2} \nabla f(x) = x - \frac{1}{\|Q\|_2} (Qx - p) = \cdots$$

we can further simplify this expression, come back to this later.

Algorithm 2: GD for (\mathcal{P})

Result: A solution x that approximately solves (\mathcal{P}) 1 Initialization Set $x_0 \in \mathbb{R}^n$, $p = A^{\top}b$, $Q = A^{\top}A$ 2 Pre-compute $L = ||Q||_2$ 3 while stopping condition is not met **do**

4
$$x = x - \frac{1}{L}(Qx - p).$$

5 end

Same example in \mathbb{R}^2 , but with stepsize t_L

Using t_L has slower convergence, but more applicable to big problem.



Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE



ine?

Deeper understanding of GD: Why GD moves along the red line not the blue line?

- Two more questions:
 - Where does the expression $\boldsymbol{x} = \boldsymbol{x} t \nabla f(\boldsymbol{x})$ "come from"?
 - What does the expression $\boldsymbol{x} = \boldsymbol{x} t \nabla f(\boldsymbol{x})$ "actually do"?
- Answers :
 - It comes from a local quadratic model (denote as F) of f.
 - ► It minimizes that local quadratic model *F*.
- ► GD is stupid because GD is "local"
 - ► GD uses local info. (which is F) to makes local decision so it follows the red path.
 - Blue path is a global decision that needs global info.
 - Being a local decision maker, there is no guarantee the decision made by GD will be as good as the global optimal one.
 - ► GD is a greedy algorithm.

Question: then why GD made such a good move here?

Answer: it just happens when f is "nice", the local decision made by GD is coincidentally as good as the global one.



The local model F that GD minimizes

• GD step takes local info x_k to minimize a local model F of f

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t \nabla f(\boldsymbol{x}_k) = \operatorname*{argmin}_{\boldsymbol{x}} F(\boldsymbol{x}; \boldsymbol{x}_k).$$

- ► Such local model *F* :
 - Takes x_k as parameter
 - Is a 2nd-order function of x
 - Expression of $F(\boldsymbol{x}; \boldsymbol{x}_k) = f(\boldsymbol{x}_k) + (\boldsymbol{x} \boldsymbol{x}_k)^\top \nabla f(\boldsymbol{x}_k) + \frac{1}{2t_k} \|\boldsymbol{x} \boldsymbol{x}_k\|_2^2$. How come: $\frac{\partial F}{\partial \boldsymbol{x}} = 0$ yields $\boldsymbol{x}_k - t \nabla f(\boldsymbol{x}_k)$.
 - Equivalent expression of F

$$F(\boldsymbol{x};\boldsymbol{x}_k) = \frac{1}{2t_k} \left\| \boldsymbol{x} - \left(\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k) \right) \right\|_2^2.$$

► Note that *F* is **Spherical**.

GD makes progress based on the local spherical model

For t_{Exact} :



GD zigzaging theorem

Question: why GD zigzags with t_{Exact} ? Answer: in fact this always occurs.

Theorem (reading exercise) Consecutive gradient directions ³ with t_{Exact} are orthogonal to each other: $\nabla f(\boldsymbol{x}_{k+1}) \perp \nabla f(\boldsymbol{x}_k)$.

 $\textbf{Proof: } \nabla f(\boldsymbol{x}_{k+1}) \perp \nabla f(\boldsymbol{x}_k) \text{ means } \langle \nabla f(\boldsymbol{x}_{k+1}), \nabla f(\boldsymbol{x}_k) \rangle = 0.$

To show this, recall t_{Exact} minimizes $g(t) = f(\boldsymbol{x}_k - t \nabla f(\boldsymbol{x}_k)).$

Consider
$$\frac{\partial g}{\partial t} = 0$$
, we have

$$\begin{array}{ll} \frac{\partial g}{\partial t} & \stackrel{\text{chain rule}}{=} & \left\langle -\nabla f(\boldsymbol{x}_k), \nabla f\left(\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)\right) \right\rangle \\ & = & -\left\langle \nabla f(\boldsymbol{x}_k), \nabla f\left(\boldsymbol{x}_{k+1}\right) \right\rangle \\ & = & 0. \quad \Box \end{array}$$



The problem of gradient descent



To improve GD:

- 1. On dynamic step size t_{Exact} : reduce the zigzag.
- 2. On fixed step size t_L : "move more" when appropriate.

Summary so far

$$(\mathcal{P})$$
: given $A \in \mathbb{R}^{m imes n}$, $b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ by solving
 $x := \operatorname*{argmin}_{x} f(x) = rac{1}{2} \|Ax - b\|_2^2 = rac{1}{2} x^\top Q x - p^\top x$,
where $Q = A^\top A$, $p = A^\top b$, $L = \|Q\|_2$.

Gradient Descent algorithm $\boldsymbol{x} = \boldsymbol{x} - t\nabla f(\boldsymbol{x})$ $\blacktriangleright t_{\text{Exact}} = \frac{\|\nabla f(\boldsymbol{x})\|_2^2}{\|\nabla f(\boldsymbol{x})\|_2^2}$ or $t_L = \frac{1}{L}$.

- ► GD is a local decision maker.
- GD = minimizing a local quadratic model F of f at point x_k .
- When level sets of f is elliptic:
 - GD zigzags with t_{Exact} .
 - GD is slow with t_L .

Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Acceleration by extrapolation

- We will look at two extrapolation schemes:
 - Polyak's Heavy Ball Method (HBM)
 - Nesterov's acceleration
- ► Idea: add "momentum" to the current iterate.
- ▶ Momentum = the "previous history".

$$\begin{split} \text{Push } \boldsymbol{x}_{\text{updated along }} - \nabla f(\boldsymbol{x}_{k-1}) \text{ for } * \text{ amount } \\ * = \beta_k t_{k-1} \| \nabla f(\boldsymbol{x}_{k-1}) \|. \end{split}$$

- ► Other momentum (not the focus)
 - Adam (a variable metric gradient method)
 - Donald G. Anderson acceleration



Three update schemes

Normal gradient

$$\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)$$

Move the point x_k in the direction $-\nabla f(x_k)$ for $t_k \|\nabla f(x)\|$ amount.

► Polyak's Heavy Ball Method

$$oldsymbol{x}_k - t_k
abla f(oldsymbol{x}_k) + eta_k (oldsymbol{x}_k - oldsymbol{x}_{k-1})$$

Perform a GD, move the updated-x in the direction of the previous step for $\beta_k \|x_k - x_{k-1}\|$ amount.

Nesterov's acceleration

$$\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k + \beta_k(\boldsymbol{x}_k - \boldsymbol{x}_{k-1})) + \beta_k(\boldsymbol{x}_k - \boldsymbol{x}_{k-1})$$

Move the not-yet-updated-x in the direction of the previous step for $\beta_k ||x_k - x_{k-1}||$ amount, perform a GD on the shifted-x, then move the updated-x in the direction of the previous step for $\beta_k ||x_k - x_{k-1}||$ amount.

Algorithm 3: HBM for (\mathcal{P}) , with step size t_{Exact}

Result: A solution x that approximately solves (\mathcal{P}) 1 Initialization Set $x_0 \in \mathbb{R}^n$, $p = A^{\top}b$, $Q = A^{\top}A$, $\beta \ge 0$ 2 while stopping condition is not met **do**

3
$$egin{array}{c|c} egin{array}{c|c} egin{array}{c|c} egin{array}{c|c} egin{array}{c|c} egin{array}{c|c} egin{array}{c|c} egin{array}{c|c} egin{array}{c} egin{arra$$

5 end

Algorithm 4: HBM for (\mathcal{P}) , with step size t_L

Result: A solution x that approximately solves (\mathcal{P}) 1 **Initialization** Set $x_0 \in \mathbb{R}^n$, $p = A^{\top}b$, $Q = A^{\top}A$, $\beta \ge 0$ 2 while stopping condition is not met **do**

3 g = Qx - p4 $x = x - \frac{1}{\|Q\|_2}g + \beta(x - x^-)$

5 end

- Fixed β is used here.
- ▶ When $\beta = 0$, HBM reduces to GD

HBM with $t_{\rm Exact}$

Same example set up, 4 iterations.

Red: normal gradient. Blue: HBM with different fixed β .



Observation: for nice f (with spherical level sets), GD is already good enough and HBM adds a little effect. However, for bad f (with elliptic level sets), HBM is better in some cases.

HBM with t_L

Same example set up, 4 iterations.

Red: normal gradient. Blue: HBM with different fixed β .



If f is nice, GD doesn't need acceleration. If f is nice, GD doesn't need acceleration. If f is nice, GD doesn't need acceleration.

This is so important so I repeated 3 times.

The message: do not use acceleration blindly, for some problems GD don't need acceleration. Using acceleration blindly doesn't make you look cool. Knowing when to use it makes you look cool.

Effect of different β on HBM (on elliptic f)

In HBM, we need to guess a good β . A bad β gives bad effect.



33 / 92

Effect of different β on HBM, more iterations

Question: the smaller β , the better HBM?



Is there a way to find best β ? What about dynamic β ?

Nesterov's acceleration

 $\begin{array}{ll} \mathsf{GD} & & \pmb{x}_{k+1} = \pmb{x}_k - t_k \nabla f(\pmb{x}_k) \\ \mathsf{HBM} & & \pmb{x}_{k+1} = \pmb{x}_k - t_k \nabla f(\pmb{x}_k) + \beta_k(\pmb{x}_k - \pmb{x}_{k-1}) \\ \mathsf{Nesterov} & & \pmb{x}_{k+1} = \pmb{x}_k - t_k \nabla f(\pmb{x}_k + \beta_k(\pmb{x}_k - \pmb{x}_{k-1})) + \beta_k(\pmb{x}_k - \pmb{x}_{k-1}) \\ \mathsf{Nesterov-2} & & \pmb{y}_{k+1} = \pmb{x}_k - t_k \nabla f(\pmb{x}_k) \\ & & \pmb{x}_{k+1} = \pmb{y}_{k+1} + \beta_k(\pmb{y}_{k+1} - \pmb{y}_k) \end{array}$

- We saw HBM with fixed β .
- Nesterov gave the update scheme with *close-form formula* for β_k (in 1983)

$$\alpha_1 \in [0,1], \ \alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2} - \alpha_k^2}{2}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$$

Note: β_k is not fix, it is a function of α_1 . We need to guess α_1 .

▶ How to get Nesterov-2 from Nesterov : set $x_{-1} = x_0$, $y_{-1} = y_0$

Algorithm 5: Nesterov's accelerated gradient for (\mathcal{P})

Result: A solution x that approximately solves (\mathcal{P}) 1 Initialization Set $x_0 \in \mathbb{R}^n$, $p = A^\top b$, $Q = A^\top A$, $\alpha_1 \in (0 \ 1)$ 2 while stopping condition is not met do3Compute $\nabla f(x_k)$ and step size t4Compute $\alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}$, $\beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$ 5 $y_{k+1} = x_k - t_k \nabla f(x_k)$ 6 $x_{k+1} = y_{k+1} + \beta_k(y_{k+1} - y_k)$ 7end
Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Constrained problem and projected gradient descent

$$(\mathcal{P})$$
: given $A \in \mathbb{R}^{m imes n}$, $b \in \mathbb{R}^m$, find $x \in \mathcal{C} \subset \mathbb{R}^n$ by solving
 $x \coloneqq \operatorname*{argmin}_{x \in \mathcal{C}} f(x) = rac{1}{2} \|Ax - b\|_2^2.$

The basic GD becomes *Projected GD*.

Algorithm 6: PGD for (\mathcal{P}) , with step size t

Result: A solution $x \in C$ that approximately solves (\mathcal{P})

- 1 Initialization Set $x_0 \in C$
- 2 while stopping condition is not met do
- 3 Compute $\nabla f(\boldsymbol{x})$ and step size t_k

4 Compute gradient update
$$oldsymbol{y} = oldsymbol{x} - t_k
abla f(oldsymbol{x})$$

5 Compute projection $\boldsymbol{x} = P_{\mathcal{C}}(\boldsymbol{y})$

6 end

Gradient and projection in 1 line:
$$oldsymbol{x}=P_{\mathcal{C}}ig(oldsymbol{x}-t_k
abla f(oldsymbol{x})ig).$$

Note: now the distance traveled between consecutive points is not $t \|\nabla f(x)\|$ but $\|x - P_{\mathcal{C}}(x - t_k \nabla f(x))\|$.

Nonnegative least squares

 $(\mathcal{P}): \text{ given } \boldsymbol{A} \in \mathbb{R}^{m \times n}, \, \boldsymbol{b} \in \mathbb{R}^{m}, \, \text{find } \boldsymbol{x} \in \mathbb{R}^{n}_{+} \text{ by solving } \boldsymbol{x} \coloneqq \operatorname*{argmin}_{\boldsymbol{x} \in \mathbb{R}^{n}_{+}} f(\boldsymbol{x}) = \frac{1}{2} \| \boldsymbol{A} \boldsymbol{x} - \boldsymbol{b} \|_{2}^{2}.$

Algorithm 7: PGD for (\mathcal{P}) , with step size t

Result: A solution $\boldsymbol{x} \in \mathbb{R}^n_+$ that approximately solves (\mathcal{P})

- 1 Initialization Set $oldsymbol{x}_0\in\mathbb{R}^n_+$
- 2 while stopping condition is not met do
- 3 Compute $abla f(oldsymbol{x})$ and step size t_k
- 4 Compute projected gradient update $m{x} = [m{x} t_k
 abla f(m{x})]_+ = \max(m{x}, 0).$
- 5 end

Algorithm 8: Nesterov's accelerated projected gradient for (\mathcal{P})

 $\begin{array}{l|l} \textbf{Result: A solution } \boldsymbol{x} \in \mathbb{R}^n_+ \text{ that approximately solves } (\mathcal{P}) \\ \textbf{1 Initialization Set } \boldsymbol{x}_0 \in \mathbb{R}^n_+, \ \boldsymbol{p} = \boldsymbol{A}^\top \boldsymbol{b}, \ \boldsymbol{Q} = \boldsymbol{A}^\top \boldsymbol{A}, \ \alpha_1 \in (0 \ 1) \\ \textbf{2 while stopping condition is not met } \textbf{do} \\ \textbf{3 } & \text{Compute } \nabla f(\boldsymbol{x}_k) \text{ and step size } t \\ \textbf{4 } & \text{Compute } \alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}} \\ \textbf{5 } & \boldsymbol{y}_{k+1} = [\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)]_+ \\ \textbf{6 } & \boldsymbol{x}_{k+1} = \boldsymbol{y}_{k+1} + \beta_k(\boldsymbol{y}_{k+1} - \boldsymbol{y}_k) \\ \textbf{7 end} \end{array}$

PGD vs Accelerated PG on NNLS



Figure: m=n=50, $oldsymbol{A}\in\mathbb{R}^{m imes n}$, $oldsymbol{x}^*=0$, $oldsymbol{x}_0\in\mathbb{R}^n$, $lpha_1=0.9$

- ► APGD is much faster.
- ▶ PGD is *monotonic*, APGD is not.
- Choice of step size makes small difference.
- ► Important: APGD and PGD have a very different convergence rate
 - the y-axis is in log-scale.
 - the curve of PGD and that of APGD have different slopes.

Other β_k schemes

► Nesterov's parameters looks so complicated

$$\alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$$

Another Nesterov's parameters

$$\alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2 + \kappa^{-1}\alpha_{k+1}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$$

► Yet another Nesterov's parameters

$$\alpha_{k+1} = \frac{1 + \sqrt{1 + 4\alpha_k^2}}{2}, \ \beta_k = \frac{1 - \alpha_k}{\alpha_{k+1}}.$$

- ▶ Paul Tseng parameter $\beta_k = \frac{k-1}{k+2}, \text{ or } \frac{1}{k+2},$
- Using conditional number

$$\beta_k = \beta = \frac{1 - \sqrt{\kappa'}}{1 + \sqrt{\kappa'}}, \quad \kappa' = \frac{1}{\kappa}, \quad \kappa = \frac{\sigma_{\max}(\boldsymbol{Q})}{\sigma_{\min}(\boldsymbol{Q})} = \frac{\lambda_{\max}(\boldsymbol{Q})}{\lambda_{\min}(\boldsymbol{Q})}$$

Different β on NNLS



Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE



- Function values of GD always go down.
- ► Function values of APGD sometimes go up.
- ► APGD exhibits some periodic behavior.

Acceleration with restarts

- When you run GD $x_{k+1} = x_k t \nabla f(x_k)$, you always get $f(x_{k+1}) \leq f(x_k)$. Why :
 - GD is the global minimizer of a local approximator F of f at x_k. More technical, see theorem 2 in https://angms.science/doc/CVX/CVX_GD_Convergence.pdf
 - ► GD is a local decision maker, a *greedy algorithm*.
- ► We "twist" the Accelerated GD algorithm:



Go back to normal GD when function value \uparrow

Acceleration with restarts



But go back to GD = go back to a slow algorithm! \rightarrow Re-run acceleration right after back to GD.

Effect of restarts on APGD



47 / 92

Algorithm 9: APGD (using Nesterov's parameter) for (\mathcal{P}) , no restart

 $\begin{array}{c|c} \textbf{Result: A solution } \boldsymbol{x} \text{ that approximately solves } (\mathcal{P}) \\ \textbf{1} \quad \textbf{Initialization Set } \boldsymbol{x}_0 \in \mathbb{R}^n_+, \boldsymbol{p} = \boldsymbol{A}^\top \boldsymbol{b}, \boldsymbol{Q} = \boldsymbol{A}^\top \boldsymbol{A}, \, \alpha_1 \in (0 \ 1) \\ \textbf{2} \quad \textbf{while stopping condition is not met do} \\ \textbf{3} \quad \hline \textbf{Compute } \nabla f(\boldsymbol{x}_k), \, \textbf{step size } t_k \\ \textbf{4} \quad \hline \textbf{Compute } \alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}, \, \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}} \\ \textbf{5} \quad \boldsymbol{y}_{k+1} = [\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)]_+ \\ \textbf{6} \quad \boldsymbol{x}_{k+1} = \boldsymbol{y}_{k+1} + \beta_k(\boldsymbol{y}_{k+1} - \boldsymbol{y}_k) \\ \textbf{7} \quad \textbf{end} \end{array}$

Algorithm 10: APGD (using Nesterov's parameter) for (\mathcal{P}) with restart

Result: A solution x that approximately solves (\mathcal{P}) Initialization Set $\boldsymbol{x}_0 \in \mathbb{R}^n_+$, $\boldsymbol{p} = \boldsymbol{A}^\top \boldsymbol{b}$, $\boldsymbol{Q} = \boldsymbol{A}^\top \boldsymbol{A}$, $\alpha_1 \in (0 \ 1)$ while stopping condition is not met do 2 Compute $\nabla f(\boldsymbol{x}_k)$, step size t_k Compute $\alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$ 4 $\boldsymbol{y}_{k+1} = [\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)]_+$ 5 $\boldsymbol{x}_{k+1} = \boldsymbol{y}_{k+1} + \beta_k (\boldsymbol{y}_{k+1} - \boldsymbol{y}_k)$ If error increase do $\boldsymbol{x}_{k+1} = [\boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)]_+$ (go back to gradient descent) $y_{k+1} = x_{k+1}, \ \alpha_k = \alpha_1$ (reset all parameters) a end if 10 11 end

Algorithm 15: APGD (using Paul Tseng's parameter) for (\mathcal{P}) with restart

Result: A solution x that approximately solves (\mathcal{P}) 1Initialization Set $x_0 \in \mathbb{R}_+^n$, $p = A^\top b$, $Q = A^\top A$, $\hat{k} = 0$ 2while stopping condition is not met do3Compute $\nabla f(x_k)$, step size t_k 4Compute $\beta_k = \frac{(k - \hat{k}) - 1}{(k - \hat{k}) + 2}$ 5 $y_{k+1} = [x_k - t_k \nabla f(x_k)]_+$ 6 $x_{k+1} = y_{k+1} + \beta_k (y_{k+1} - y_k)$ 7If error increase do8 $x_{k+1} = [x_k - t_k \nabla f(x_k)]_+$ (go back to gradient descent)9 $y_{k+1} = x_{k+1}$, $\hat{k} = \hat{k} + k$ (reset all parameters)10end if

Algorithm 16: APGD (using Conditional number) for (\mathcal{P}) with restart

Result: A solution x that approximately solves (\mathcal{P}) 1 Initialization Set $x_0 \in \mathbb{R}^n_+$, $p = \mathbf{A}^\top b$, $Q = \mathbf{A}^\top \mathbf{A}$, $\beta = \frac{1 - \sqrt{\kappa'}}{1 + \sqrt{\kappa'}}$ 2 while stopping condition is not met do 3 Compute $\nabla f(\mathbf{x}_k)$, step size t_k 4 $y_{k+1} = [\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)]_+$ 5 $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \beta(\mathbf{y}_{k+1} - \mathbf{y}_k)$ 6 If error increase do 7 $\mathbf{x}_{k+1} = [\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)]_+$ (go back to gradient descent) 8 $y_{k+1} = \mathbf{x}_{k+1}$ (reset all parameters) 10 end

Comparing all schemes

NNLS(50, 50)



Warning: restart is not "always good"

NNLS(100,100)



50 / 92

Summary so far

$$(\mathcal{P})$$
: given $A \in \mathbb{R}^{m imes n}$, $b \in \mathbb{R}^m$, find $x \in \mathcal{C}$ by solving
 $x \coloneqq \operatorname*{argmin}_{x \in \mathcal{C}} f(x) = rac{1}{2} \|Ax - b\|_2^2 = rac{1}{2} x^\top Q x - p^\top x$,
 $Q = A^\top A$, $p = A^\top b$, $L = \|Q\|_2$.

Algorithm: Gradient Descent $\boldsymbol{x} = \boldsymbol{x} - t \nabla f(\boldsymbol{x})$

• Minimizes a local model of f, slow if level sets of f is elliptic.

Acceleration : add momentum $eta_k(oldsymbol{x}_k-oldsymbol{x}_{k-1})$

$$\blacktriangleright \text{ HBM } \boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t \nabla f(\boldsymbol{x}_k) + \beta_k (\boldsymbol{x}_k - \boldsymbol{x}_{k-1})$$

$$\blacktriangleright \text{ Nesterov } \boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t \nabla f(\boldsymbol{x}_k + \beta_k(\boldsymbol{x}_k - \boldsymbol{x}_{k-1})) + \beta_k(\boldsymbol{x}_k - \boldsymbol{x}_{k-1})$$

► Nesterov's parameter
$$\alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$$

• Other parameters
$$\beta = \frac{k-1}{k+2}$$
, $\beta = \frac{1-\sqrt{\kappa'}}{1+\sqrt{\kappa'}}$

► Adaptive restarts

Why the acceleration works?



Pourquoi l'acceleration fonctionne?

Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on $\left(\mathcal{P}\right)$

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Convergence rate

Recall the observation : different slope \implies different convergence rate



We are going to study the convergence properties of GD.

- \blacktriangleright For simplicity, consider the model (\mathcal{P}) with no constraint.
- ▶ The theory applies for general any smooth convex *f*.

Convergence condition of GD on (\mathcal{P})

$$(\mathcal{P})$$
: given $A \in \mathbb{R}^{m imes n}$, $b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ by solving
 $x \coloneqq \operatorname*{argmin}_{x} f(x) = rac{1}{2} \|Ax - b\|_2^2 = rac{1}{2} x^\top Q x - p^\top x$,
where $Q = A^\top A$, $p = A^\top b$, assumes A is full rank.

Theorem:

If stepsize $t_k > 0$ fulfills: $\max \left\{ |1 - t_k \lambda_{\max}(\boldsymbol{Q})|, |1 - t_k \lambda_{\min}(\boldsymbol{Q})| \right\} < 1 \text{ for all } k \in \mathbb{N},$ the sequence $\{\boldsymbol{x}_k\}_{k \in \mathbb{N}}$ produced by GD iteration $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k)$ converges to the minimizer \boldsymbol{x}^* of (\mathcal{P}) .

Why: reading exercise in the next three pages.

Convergence condition of gradient descent \dots 1/2

$$f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^{\top} \boldsymbol{Q} \boldsymbol{x} - \boldsymbol{p}^{\top} \boldsymbol{x}, \ \nabla f(\boldsymbol{x}) = \boldsymbol{Q} \boldsymbol{x} - \boldsymbol{p} \text{ and } \boldsymbol{x}_{k+1} = \boldsymbol{x}_k - t_k \nabla f(\boldsymbol{x}_k) \\ = \boldsymbol{x}_k - t_k (\boldsymbol{Q} \boldsymbol{x}_k - \boldsymbol{p}) = \boldsymbol{x}_k - t_k \boldsymbol{Q} \boldsymbol{x}_k + t_k \boldsymbol{p} \\ = (\boldsymbol{I}_n - t_k \boldsymbol{Q}) \boldsymbol{x}_k + t_k \boldsymbol{p}$$

Recall 1st-order optimality condition $\nabla f(\boldsymbol{x}^*) = 0$ so $\boldsymbol{Q}\boldsymbol{x}^* - \boldsymbol{p} = 0$.

$$\begin{array}{rcl} \boldsymbol{x}_{k+1} &=& (\boldsymbol{I}_n - t_k \boldsymbol{Q}) \boldsymbol{x}_k + t_k \boldsymbol{Q} \boldsymbol{x}^* \\ \boldsymbol{x}_{k+1} - \boldsymbol{x}^* &=& (\boldsymbol{I}_n - t_k \boldsymbol{Q}) \boldsymbol{x}_k + t_k \boldsymbol{Q} \boldsymbol{x}^* - \boldsymbol{x}^* \\ &=& (\boldsymbol{I}_n - t_k \boldsymbol{Q}) \boldsymbol{x}_k + (t_k \boldsymbol{Q} - \boldsymbol{I}_n) \boldsymbol{x}^* \\ &=& (\boldsymbol{I}_n - t_k \boldsymbol{Q}) (\boldsymbol{x}_k - \boldsymbol{x}^*) \\ \| \boldsymbol{x}_{k+1} - \boldsymbol{x}^* \|_2 &=& \| (\boldsymbol{I}_n - t_k \boldsymbol{Q}) (\boldsymbol{x}_k - \boldsymbol{x}^*) \|_2 \\ & \stackrel{\mathsf{oni}}{\leq} & \| \boldsymbol{I}_n - t_k \boldsymbol{Q} \|_2 \| \boldsymbol{x}_k - \boldsymbol{x}^* \|_2. \quad (\mathsf{oni}: \mathsf{operator norm inequality}) \end{array}$$

Let $\mu I_n \preceq \lambda_{\min}(\boldsymbol{Q}) I_n \preceq \boldsymbol{Q} \preceq \lambda_{\max}(\boldsymbol{Q}) I_n \preceq L I_n$, then $\|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \leq \max\left\{|1 - t_k L|, |1 - t_k \mu|\right\} \|\boldsymbol{x}_k - \boldsymbol{x}^*\|_2$

Derivation

56 / 92

Convergence of gradient descent $\dots 2/2$

$$egin{aligned} \|m{x}_{k+1} - m{x}^*\|_2 &\leq & \maxig\{|1 - t_k L|, |1 - t_k \mu|ig\}\|m{x}_k - m{x}^*\|_2 \ &\leq & \maxig\{|1 - t_k L|, |1 - t_k \mu|ig\}\Big(\maxig\{|1 - t_{k-1} L|, |1 - t_{k-1} \mu|ig\}\|m{x}_{k-1} - m{x}^*\|_2\Big) \ &dots \ & dots \ & dots \ & igg(\prod\limits_k \maxig\{|1 - t_k L|, |1 - t_k \mu|ig\}\Big)\|m{x}_0 - m{x}^*\|_2. \end{aligned}$$

Hence GD converge for (\mathcal{P}) if

$$\max\{|1 - t_k L|, |1 - t_k \mu|\} < 1 \quad \forall k. \quad \Box$$

Convergence of GD on (\mathcal{P}) with fix stepsize

$$\begin{split} \|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 &\leq \left(\prod_k \max\left\{|1 - t_k L|, |1 - t_k \mu|\right\}\right) \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2 \\ &= \left(\max\left\{|1 - tL|, |1 - t\mu|\right\}\right)^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2. \\ \\ \text{Case A} & \|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \leq |1 - tL|^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2. \\ \\ \text{Case B} & \|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \leq |1 - t\mu|^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2. \end{split}$$

No matter which case, we need $|\cdot| < 1$. Case A $|1 - tL| < 1 \iff 1 - tL < 1$ or 1 - tL > -1

Similarly we get $0 < t < \frac{2}{\mu}$ for case B.

Convergence of GD on (\mathcal{P}) for fix t

$$(\mathcal{P})$$
: given full rank $A \in \mathbb{R}^{m imes n}$, $b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ by solving
 $x \coloneqq \operatorname*{argmin}_{x} f(x) = rac{1}{2} \|Ax - b\|_2^2 = rac{1}{2} x^\top Q x - p^\top x$,
where $Q = A^\top A$, $p = A^\top b$, $L = \|Q\|_2$.

We have

$$\|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \le \left(\max\left\{|1 - tL|, |1 - t\mu|\right\}\right)^{\kappa} \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2.$$

. L

Suppose case A occurs

$$\|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \le |1 - tL|^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2.$$

We have

Theorem (Convergence of GD with fixed step size, case A) For fixed step size $0 < t < \frac{2}{L},$

the sequence
$$\{x_k\}_{k\in\mathbb{N}}$$
 produced by GD converges to x^* of (\mathcal{P}) .

Convergence of GD on (\mathcal{P}) with $t = \frac{1}{L}$

$$\begin{split} \| \boldsymbol{x}_{k+1} - \boldsymbol{x}^* \|_2 &\leq \left(\max\left\{ |1 - tL|, |1 - t\mu| \right\} \right)^k \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|_2. \\ \text{If } t &= \frac{1}{L} \\ \| \boldsymbol{x}_{k+1} - \boldsymbol{x}^* \|_2 &\leq \left(\max\left\{ \left| 1 - \frac{L}{L} \right|, \left| 1 - \frac{\mu}{L} \right| \right\} \right)^k \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|_2 \\ &= \left(1 - \frac{1}{\kappa} \right)^k \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|_2, \text{ where } \kappa = \frac{L}{\mu} > 1 \end{split}$$

• If $\kappa = 1$, $1 - \frac{1}{\kappa} = 0$, and GD converge in 1 step. • If $\kappa = 9999$, $1 - \frac{1}{\kappa} = 0.999$..., and GD converge very slowly.

Convergence of GD on (\mathcal{P}) with the "best" \boldsymbol{t}

$$\begin{array}{l} \text{The "best" } t = \frac{2}{L+\mu} \\ \| \pmb{x}_{k+1} - \pmb{x}^* \|_2 \leq \left(1 - \frac{2}{\kappa+1}\right)^k \| \pmb{x}_0 - \pmb{x}^* \|_2. \\ \end{array}$$

$$\begin{array}{l} \text{Why: put } t = \frac{2}{L+\mu} \text{ into max } \{|1 - tL|, |1 - t\mu|\} : \\ \max\left\{ \left|1 - \frac{2L}{L+\mu}\right|, \left|1 - \frac{2\mu}{L+\mu}\right|\right\} = \max\left\{ \left|\frac{-L+\mu}{L+\mu}\right|, \left|\frac{L-\mu}{L+\mu}\right|\right\} = \frac{L-\mu}{L+\mu} = \frac{\kappa-1}{\kappa+1} = 1 - \frac{2}{\kappa+1}. \end{array}$$

$$\text{r If } \kappa = 1, \ 1 - \frac{2}{\kappa+1} = 0, \text{ and GD converge in 1 step} \\ \text{r If } \kappa = 9999, \ 1 - \frac{2}{\kappa+1} = 0.999..., \text{ and GD converge very slow} \\ \text{Compared with } t = \frac{1}{L} \text{ with } \left(1 - \frac{1}{\kappa}\right)^k : \text{ let } \kappa = 2, \text{ we have } \frac{1}{2^k} \text{ vs } \frac{1}{3^k}. \end{array}$$

Summary on the convergence rate of GD on (\mathcal{P})

Theorem (Convergence condition of GD) If step size $t_k > 0$ satisfies the following condition

$$\max\left\{|1-t_k\lambda_{\max}(\boldsymbol{Q})|, |1-t_k\lambda_{\min}(\boldsymbol{Q})|\right\} < 1 \; \forall k \in \mathbb{N},$$

the sequence $\{x_k\}_{k\in\mathbb{N}}$ produced by GD converges to the minimizer x^* .

And the theorem with different stepsize t:

If $t = \frac{1}{L}$ then $\|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \leq \left(1 - \frac{1}{\kappa}\right)^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2.$ If $t = \frac{2}{L + \mu}$ then $\|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|_2 \leq \left(1 - \frac{2}{\kappa + 1}\right)^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2.$

Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Convergence rate of GD on general smooth \boldsymbol{f}

▶ Previously we studied the convergence rate of GD on minimizing

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2$$

That is, we restrict f to be the least squares function.

- Now we consider general f that is smooth and convex.
- ► We now study the convergence rate of GD on

 $\mathop{\rm argmin}_{{\boldsymbol{x}}} \, f({\boldsymbol{x}})$

where

• f is convex

• f is L-smooth: the gradient of f is L-Lipschitz

$$\|\nabla f(\boldsymbol{a}) - \nabla f(\boldsymbol{b})\| \leq L \|\boldsymbol{a} - \boldsymbol{b}\|.$$

A very important inequality for L-smooth f

• If f is L-smooth: the gradient of f is L-Lipschitz

$$\|\nabla f(\boldsymbol{a}) - \nabla f(\boldsymbol{b})\| \le L \|\boldsymbol{a} - \boldsymbol{b}\|.$$

Then

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{L}{2} \| \boldsymbol{y} - \boldsymbol{x} \|^2.$$

▶ Why and how: see next two pages (reading exercise).

Proof (1/2)

We show for L > 0, $\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\| \le L \|\boldsymbol{x} - \boldsymbol{y}\|$ implies $\left|f(\boldsymbol{y}) - f(\boldsymbol{x}) - \left\langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x}\right\rangle\right| \le \frac{L}{2} \|\boldsymbol{y} - \boldsymbol{x}\|_2^2.$

Recall from calculus $G(b) - G(a) = \int_a^b g(\theta) d\theta$. Next, a smart step, let $g(\theta)$ as $g(\tau) = \langle \nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})), \boldsymbol{y} - \boldsymbol{x} \rangle$ be a function in τ and $d\theta = d\tau$. Consider the definite integral of $g(\tau)$ from 0 to 1, let $G(b) = f(\boldsymbol{y})$ and $G(a) = f(\boldsymbol{x})$, hence

$$\begin{split} f(\boldsymbol{y}) - f(\boldsymbol{x}) &= \int_0^1 \left\langle \nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})), \boldsymbol{y} - \boldsymbol{x} \right\rangle d\tau \\ &= \int_0^1 \left\langle \nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x}) + \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \right\rangle d\tau. \end{split}$$

As $\nabla f(\boldsymbol{x})$ is independent of τ , can take out from the integral

$$f(\boldsymbol{y}) - f(\boldsymbol{x}) = \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \int_0^1 \left\langle \nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \right\rangle d\tau.$$

The idea is to create the term $\langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle$ so that we can move it to the left and get $|f(\boldsymbol{y}) - f(\boldsymbol{x}) - \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle|$.

Proof (2/2)

$$\begin{array}{ll} |f(\boldsymbol{y}) - f(\boldsymbol{x}) - \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle| &= & |\int_0^1 \langle \nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle \, d\tau| \\ &\leq & \int_0^1 |\langle \nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle \, |d\tau| \\ &\leq & \int_0^1 ||\nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x})|| \cdot ||\boldsymbol{y} - \boldsymbol{x}|| d\tau. \end{array}$$

c.s. means Cauchy - Schwarz inequality.

Now look at $\|\nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x})\|$, this is exactly where we can apply the Lipschitz gradient inequality

$$\|\nabla f(\boldsymbol{x} + \tau(\boldsymbol{y} - \boldsymbol{x})) - \nabla f(\boldsymbol{x})\| \le L \|\tau(\boldsymbol{y} - \boldsymbol{x})\| \le L |\tau| \|\boldsymbol{y} - \boldsymbol{x}\| = L \tau \|\boldsymbol{y} - \boldsymbol{x}\|$$

where $\|\tau(\boldsymbol{y} - \boldsymbol{x})\| = |\tau| \|\boldsymbol{y} - \boldsymbol{x}\|$ as norm is non-negative. Note that the integral range is from 0 to 1 so the absolute sign in τ can be removed.

Lastly

$$\left|f(\boldsymbol{y}) - f(\boldsymbol{x}) - \left\langle
abla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x}
ight
angle
ight| \leq \int_{0}^{1} L au d au \cdot \| \boldsymbol{y} - \boldsymbol{x} \|^{2} = rac{L}{2} \| \boldsymbol{y} - \boldsymbol{x} \|^{2}.$$

Remove the absolute value sign gives

$$f(oldsymbol{y}) \leq f(oldsymbol{x}) + \langle
abla f(oldsymbol{x}), oldsymbol{y} - oldsymbol{x}
angle + rac{L}{2} \|oldsymbol{y} - oldsymbol{x}\|^2.$$

Meaning of L-smoothness: quadratic upper bound

A function f is L-smooth if for any two points $x, y \in \text{dom} f$, $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||y - x||_2^2$



Interpretation: f is globally bounded above by a quadratic function. i.e. f cannot be "grow too fast" than the quadratic upper bound. Polyak-Lojasiewicz inequality

A function f : ℝⁿ → ℝ that is differentiable (i.e. ∇f(x) exists for x ∈ domf) satisfies Polyak-Lojasiewicz (PL) inequality if there exists a positive scalar µ > 0 such that

$$\frac{1}{2} \|\nabla f(\boldsymbol{x})\|^2 \ge \mu \big(f(\boldsymbol{x}) - f^* \big)$$

for all $\boldsymbol{x} \in \mathrm{dom} f$, where $f^* = f(\boldsymbol{x}^*)$ and \boldsymbol{x}^* is a minimizer of f.

- It means the gradient grows faster than a quadratic function (scaled by a scalar $\mu > 0$) as we move x away from x^* .
- ▶ Why an inequality with crazy name suddenly jumps out from nowhere: because it is USEFUL!
- ► Good news: this inequality is true for many many many functions you see in daily life. So you can assume the *f* you work with is PL.

Poof of linear convergence of GD

$$f$$
 has L -Lips. grad $f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{L}{2} \| \boldsymbol{y} - \boldsymbol{x} \|^2$
GD update $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \frac{1}{L} \nabla f(\boldsymbol{x}_k)$

• Put $y = x_{k+1}$, $x = x_k$ in the first inequality, then plug in the second equation gives the *descent lemma*

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) - rac{1}{2L} \|
abla f(\boldsymbol{x}_k) \|^2.$$

• PL inequality: $-\frac{1}{2} \|\nabla f(\boldsymbol{x}_k)\|^2 \leq -\mu (f(\boldsymbol{x}_k) - f^*)$, so

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) - \frac{\mu}{L} (f(\boldsymbol{x}_k) - f^*).$$

► Subtract both side by f^*

$$f(\boldsymbol{x}_{k+1}) - f^* \leq f(\boldsymbol{x}_k) - \frac{\mu}{L} (f(\boldsymbol{x}_k) - f^*) - f^* = (1 - \frac{\mu}{L}) (f(\boldsymbol{x}_k) - f^*).$$

Recursion:

$$f(\boldsymbol{x}_{k+1}) - f^* \leq \left(1 - \frac{\mu}{L}\right)^k \left(f(\boldsymbol{x}_0) - f^*\right).$$

Remarks

► The proof also applies to optimal stepsize, since

$$f(\boldsymbol{x}_{k+1}) = \min_{\alpha} f\left(\boldsymbol{x}_{k} - \alpha \nabla f(\boldsymbol{x}_{k})\right) \leq f\left(\boldsymbol{x}_{k} - \frac{1}{L} \nabla f(\boldsymbol{x}_{k})\right),$$

where the \leq is by definition of the optimal stepsize.

► Another approach to show the convergence of GD for smooth-*f* is to impose *f* is *µ*-strongly convex. The proof is longer and more complicated.
Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on $\left(\mathcal{P}\right)$

Convergence of accelerated gradient method explained using ODE

The convergence analysis of GD is so hard-core!



The convergence analysis of AGD is even more hard-core.

Convergence rates

► (One of the) Definition:

$$o \coloneqq \lim_{k \to \infty} rac{f(oldsymbol{x}_{k+1}) - f^*}{f(oldsymbol{x}_k) - f^*}.$$

What is it: limit of ratio of successive errors.

- Convergence rate
 - $\rho = 1$: sublinear convergence rate
 - ▶ $\rho \in]0,1[:$ linear convergence rate
 - $\rho = 0$: super-linear convergence rate
- ► Quadratic convergence rate

$$\lim_{k \to \infty} \frac{f(\boldsymbol{x}_{k+1}) - f^*}{\left(f(\boldsymbol{x}_k) - f^*\right)^2} < M,$$

for some constant M > 0.

• Other definition: $\|\boldsymbol{x}_k - \boldsymbol{x}^*\|$, $\|\nabla f(\boldsymbol{x}_k)\|$.

Convergence rate Recall the observation : different slope \implies different convergence rate



Convergence rate of accelerated gradient is a order higher

We have the following rate of GD on solving (\mathcal{P}) for the best $t=\frac{2}{L+\mu}$:

$$\|m{x}_{k+1} - m{x}^*\|_2 \le \left(1 - rac{2}{\kappa + 1}
ight)^k \|m{x}_0 - m{x}^*\|_2$$

It can be shown, for accelerated gradient (e.g. the heavy ball method), we have

$$\|m{x}_{k+1} - m{x}^*\|_2 \le \left(1 - rac{2}{\sqrt{\kappa} + 1}
ight)^k \|m{x}_0 - m{x}^*\|_2.$$

Example: $\kappa = 2$, we have $\frac{1}{3^k} = 0.3333^k$ vs 0.1716^k . For 4 iterations : GD : 0.33, 0.11, 0.03, 0.01 AGD : 0.17, 0.02, 0.0049, 0.0008

Let's make it even easier to understand

Original problem :

$$(\mathcal{P})$$
: given full rank $A \in \mathbb{R}^{m imes n}$, $b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ by solving
 $x \coloneqq \operatorname*{argmin}_x f(x) = rac{1}{2} \|Ax - b\|_2^2 = rac{1}{2} x^\top Q x - p^\top x$,
where $Q = A^\top A$, $p = A^\top b$, $L = \|Q\|_2$.

To make it easier to understand, set $\boldsymbol{b} = 0$

 (\mathcal{P}') : given symmetric positive definite $oldsymbol{Q} \in \mathbb{R}^{n imes n}$, find $oldsymbol{x} \in \mathbb{R}^n$ by solving $oldsymbol{x} \coloneqq \operatorname*{argmin}_{oldsymbol{x}} f(oldsymbol{x}) = rac{1}{2} oldsymbol{x}^\top oldsymbol{Q} oldsymbol{x}.$

Again, let $\lambda_{\max}(\boldsymbol{Q}) = L$ and $\lambda_{\min}(\boldsymbol{Q}) = \mu$.

Nesterov's accelerated gradient on (\mathcal{P}')

 $(\mathcal{P}'): ext{ given sym. p.d.} oldsymbol{Q} \in \mathbb{R}^{n imes n}, ext{ find } oldsymbol{x} \in \mathbb{R}^n ext{ by } oldsymbol{x} \coloneqq rgmin_{oldsymbol{x}} rac{1}{2} oldsymbol{x}^ op oldsymbol{Q} x.$

The update steps in Nesterov's AGD with step size $t = \frac{1}{L}$ and fixed β

$$y_{k+1} = x_k - \frac{1}{L}Qx_k$$
 $x_{k+1} = y_{k+1} + \beta(y_{k+1} - y_k)$

In one line :

$$\begin{split} \boldsymbol{x}_{k+1} &= \boldsymbol{x}_k - \frac{1}{L} \boldsymbol{Q} \Big(\boldsymbol{x}_k + \beta (\boldsymbol{x}_k - \boldsymbol{x}_{k-1}) \Big) + \beta (\boldsymbol{x}_k - \boldsymbol{x}_{k-1}) \\ &= \left((1+\beta) \boldsymbol{x}_k - \beta \boldsymbol{x}_{k-1} \right) - \frac{1}{L} \boldsymbol{Q} \Big((1+\beta) \boldsymbol{x}_k - \beta \boldsymbol{x}_{k-1} \Big) \\ &= \left(\boldsymbol{I}_n - \frac{1}{L} \boldsymbol{Q} \right) \Big((1+\beta) \boldsymbol{x}_k - \beta \boldsymbol{x}_{k-1} \Big) \\ &\stackrel{\text{e.d.}}{=} \left(\boldsymbol{V} \boldsymbol{V}^\top - \frac{1}{L} \boldsymbol{V} \boldsymbol{\Lambda} \boldsymbol{V}^\top \right) \Big((1+\beta) \boldsymbol{x}_k - \beta \boldsymbol{x}_{k-1} \Big) \\ \boldsymbol{V}^\top \boldsymbol{x}_{k+1} &= \left(\boldsymbol{V}^\top - \frac{1}{L} \boldsymbol{\Lambda} \boldsymbol{V}^\top \right) \Big((1+\beta) \boldsymbol{x}_k - \beta \boldsymbol{x}_{k-1} \Big) \\ &= \left(\boldsymbol{I}_n - \frac{1}{L} \boldsymbol{\Lambda} \right) \Big((1+\beta) \boldsymbol{V}^\top \boldsymbol{x}_k - \beta \boldsymbol{V}^\top \boldsymbol{x}_{k-1} \Big) \\ \boldsymbol{w}_{k+1} &= \left(\boldsymbol{I}_n - \frac{1}{L} \boldsymbol{\Lambda} \right) \Big((1+\beta) \boldsymbol{w}_k - \beta \boldsymbol{w}_{k-1} \Big). \end{split}$$

Nesterov's accelerated gradient method on (\mathcal{P}')

$$\boldsymbol{w}_{k+1} = \left(\boldsymbol{I}_n - \frac{1}{L}\Lambda\right) \left((1+\beta)\boldsymbol{w}_k - \beta\boldsymbol{w}_{k-1}\right),$$
$$\begin{bmatrix} \boldsymbol{w}_{k+1}^{[1]} \\ \vdots \\ \boldsymbol{w}_{k+1}^{[n]} \end{bmatrix} = \left(\begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} - \frac{1}{L} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \right) \left((1+\beta) \begin{bmatrix} \boldsymbol{w}_k^{[1]} \\ \vdots \\ \boldsymbol{w}_k^{[n]} \end{bmatrix} - \beta \begin{bmatrix} \boldsymbol{w}_{k-1}^{[1]} \\ \vdots \\ \boldsymbol{w}_{k-1}^{[n]} \end{bmatrix} \right)$$

That is, we have the decoupled element-wise expression

$$\begin{aligned} w_{k+1}^{[i]} &= \left(1 - \frac{\lambda_i}{L}\right) \left((1+\beta) w_k^{[i]} - \beta w_{k-1}^{[i]}\right) \\ &= \left(1 - \frac{\lambda_i}{L}\right) (1+\beta) w_k^{[i]} - \left(1 - \frac{\lambda_i}{L}\right) \beta w_{k-1}^{[i]}, \quad i = 1, 2, ..., n \end{aligned}$$

(* Recall AGD reduces to GD if $\beta = 0$, this expression also applies to GD)

Second-order dynamical system

$$w_{k+1}^{[i]} = \left(1 - \frac{\lambda_i}{L}\right)(1+\beta)w_k^{[i]} - \left(1 - \frac{\lambda_i}{L}\right)\beta w_{k-1}^{[i]}, \quad i = 1, 2, ..., n$$

The characteristic equation is

$$r^{2} = \left(1 - \frac{\lambda_{i}}{L}\right)(1 + \beta)r - \left(1 - \frac{\lambda_{i}}{L}\right)\beta.$$

The value $\beta=\beta^*$ that the equation has repeated roots are

$$\beta_{i=1,2}^* = \frac{1 - \sqrt{\lambda_i/L}}{1 + \sqrt{\lambda_i/L}}, \ r(\beta^*) = 1 - \sqrt{\lambda_i/L}.$$

If $\beta \leq \beta^*$ the equation has two distinctive real roots. If $\beta > \beta^*$, the equation has two complex roots. The characteristic equation $r^2 = (1 - \frac{\lambda_i}{L})(1 + \beta)r - (1 - \frac{\lambda_i}{L})\beta$. At the value β^* , the equation has repeated roots



 $\blacktriangleright \ \beta = \beta^*$

- Best β
- Best amount of momentum
- System is critically damped
- $\beta \leq \beta^*$
 - β too small
 - Not enough momentum
 - System is over-damped
- $\blacktriangleright \ \beta > \beta^*$
 - β too high
 - Too much momentum
 - System is under-damped
 - System is oscillatory

The acceleration comes from damping!!

The periodic ripples is due to $\beta > \beta^*$ When $\beta > \beta^*$, the system is under-damped \implies the periodic ripples



In fact we have

$$w_k^{[i]} = c_1^{[i]} \left(\beta \left(1 - \frac{\lambda_i}{L} \right) \right)^{k/2} \cos(k\psi^{[i]} - c_2^{[i]})$$

where $c_{1,2}$ are some unimportant constants and

$$\psi^{[i]} = \cos^{-1}\left(\left(1 - \frac{\lambda_i}{L}\right)\frac{1 + \beta}{2}\sqrt{\beta\left(1 - \frac{\lambda_i}{L}\right)}\right) \approx \sqrt{\frac{\lambda_i}{L}}, \quad \psi^{[i]} \approx \sqrt{\frac{\lambda_{\min}}{L}} = \sqrt{\kappa^{-1}}$$

(Detail derivations : take home exercise)

Table of Contents

Introduction & preliminaries

Gradient Descent as an iterative algorithm

Gradient Descent is slow when level set is elliptic

Accelerated gradient method

Nonnegative least squaress

Adaptive restarts

Convergence rate of Gradient Descent on (\mathcal{P})

Convergence rate of GD on smooth f via Polyak-Lojasiewicz inequality

Convergence rate of accelerated gradient method on (\mathcal{P})

Convergence of accelerated gradient method explained using ODE

Nesterov's accelerated gradient

▶ One of the form of Nesterov's accelerated gradient using Paul Tseng's parameter¹

$$egin{array}{rll} m{x}_{k+1} &=& m{y}_k - t_k
abla f(m{y}_k) \ m{y}_{k+1} &=& m{x}_{k+1} + rac{k}{k+3} (m{x}_{k+1} - m{x}_k) \end{array}$$

• Theorem: if f is convex and L-smooth, picking stepsize $t_k = \frac{1}{L}$, the Nesterov's accelerated gradient has the convergence rate as

$$f(\boldsymbol{x}_k) - f^* \le rac{c \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|}{(k+1)^2},$$

where c is some (unimportant) constant.

The accelerated gradient actually associated to an ODE

$$\ddot{\boldsymbol{X}} + rac{3}{ au}\dot{\boldsymbol{X}} +
abla f(\boldsymbol{X}) = 0$$

which we will show it in the coming slides.

¹This is NOT the one proposed by Nesterov in 1983.

Derive the ODE ... (1/4)

Accelerated gradient (constant stepsize)

$$x_{k+1} = y_k - t \nabla f(y_k), \quad y_{k+1} = x_{k+1} + \beta_k (x_{k+1} - x_k)$$

so $oldsymbol{y}_k = oldsymbol{x}_k + eta_{k-1}(oldsymbol{x}_k - oldsymbol{x}_{k-1})$, put it to the gradient step gives

$$egin{array}{rcl} oldsymbol{x}_{k+1}&=&oldsymbol{x}_k+eta_{k-1}(oldsymbol{x}_k-oldsymbol{x}_{k-1})-t
abla f(oldsymbol{y}_k)\ &\Leftrightarrow&oldsymbol{x}_{k+1}-oldsymbol{x}_k&=η_{k-1}(oldsymbol{x}_k-oldsymbol{x}_{k-1})-t
abla f(oldsymbol{y}_k). \end{array}$$

Why do this: 1) to cancel the y_k and 2) get difference terms such as $x_{k+1} - x_{k-1}$.

We now have

$$\frac{\boldsymbol{x}_{k+1} - \boldsymbol{x}_k}{\sqrt{t}} = \beta_{k-1} \frac{\boldsymbol{x}_k - \boldsymbol{x}_{k-1}}{\sqrt{t}} - \sqrt{t} \nabla f(\boldsymbol{y}_k).$$
(odel)

Our goal: derive an ODE from (ode1). How: link discrete iteration with continuous time.

• Let
$$k = \frac{\tau}{\sqrt{t}}$$
 : discrete iteration = $\frac{\text{continuous time}}{\text{stepsize}}$ so

$$\boldsymbol{X}(\tau) = \boldsymbol{X}(k\sqrt{t}) \approx \boldsymbol{x}_k, \quad \boldsymbol{X}(\tau+1\cdot\sqrt{t}) = \boldsymbol{X}\big((k+1)\sqrt{t}\big) \approx \boldsymbol{x}_{k+1},$$

and the " \approx " becomes "=" if we take $\lim_{t\to 0}.$

We also have

$$\frac{\boldsymbol{x}_{k+1} - \boldsymbol{x}_k}{\sqrt{t}} \approx \frac{\boldsymbol{X}(\tau + \sqrt{t}) - \boldsymbol{X}(\tau)}{\sqrt{t}}, \qquad \frac{\boldsymbol{x}_k - \boldsymbol{x}_{k-1}}{\sqrt{t}} \approx \frac{\boldsymbol{X}(\tau) - \boldsymbol{X}(\tau - \sqrt{t})}{\sqrt{t}}$$

Derive the ODE ... (2/4)

Recall Taylor's series

$$u(x_0 + \Delta x) = u(x_0) + \Delta x \frac{\partial u}{\partial x}\Big|_{x=x_0} + \frac{(\Delta x)^2}{2!} \frac{\partial^2 u}{\partial^2 x}\Big|_{x=x_0} + o(\Delta x)$$

Let u = X, $x_0 = \tau$, $\Delta x = \sqrt{t}$ and $\frac{\partial u}{\partial x} = \frac{\partial X}{\partial \tau} = \dot{X}$ then

$$\boldsymbol{X}(x_0 + \Delta x) = \boldsymbol{X}(\tau) + \sqrt{t} \dot{\boldsymbol{X}}(\tau) + \frac{t}{2} \ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t}).$$

$$\blacktriangleright \text{ So } \frac{\boldsymbol{x}_{k+1} - \boldsymbol{x}_k}{\sqrt{t}} \approx \frac{\boldsymbol{X}(\tau + \sqrt{t}) - \boldsymbol{X}(\tau)}{\sqrt{t}} \text{ becomes } \frac{\boldsymbol{x}_{k+1} - \boldsymbol{x}_k}{\sqrt{t}} = \dot{\boldsymbol{X}}(\tau) + \frac{\sqrt{t}}{2} \ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t}).$$

• Similarly, using Taylor's expansion on $u(x_0 - \Delta x)$ with $\Delta x = -\sqrt{t}$, we get

$$\boldsymbol{X}(\boldsymbol{x}_0 - \Delta) = \boldsymbol{X}(\tau) - \sqrt{t} \dot{\boldsymbol{X}}(\tau) + \frac{t}{2} \ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t}),$$

so $oldsymbol{x}_k - oldsymbol{x}_{k-1} pprox oldsymbol{X}(au) - oldsymbol{X}(au - \sqrt{t})$ and

$$\frac{\boldsymbol{x}_k - \boldsymbol{x}_{k-1}}{\sqrt{t}} = \dot{\boldsymbol{X}}(\tau) - \frac{\sqrt{t}}{2}\ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t}).$$

Put them into (ode1)

$$\dot{\boldsymbol{X}}(\tau) + \frac{\sqrt{t}}{2}\ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t}) = \beta_{k-1}\left(\dot{\boldsymbol{X}}(\tau) - \frac{\sqrt{t}}{2}\ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t})\right) - \sqrt{t}\nabla f(\boldsymbol{y}_k).$$

Derive the ODE ... (3/4)

We have

$$\dot{\boldsymbol{X}}(\tau) + \frac{\sqrt{t}}{2}\ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t}) = \beta_{k-1}\left(\dot{\boldsymbol{X}}(\tau) - \frac{\sqrt{t}}{2}\ddot{\boldsymbol{X}}(\tau) + o(\sqrt{t})\right) - \sqrt{t}\nabla f(\boldsymbol{y}_k).$$

Rearrange

$$\frac{\sqrt{t}}{2} \left(1 + \beta_{k-1}\right) \ddot{\boldsymbol{X}}(\tau) + \left(1 - \beta_{k-1}\right) \dot{\boldsymbol{X}}(\tau) + \sqrt{t} \nabla f(\boldsymbol{y}_k) + o(\sqrt{t}) = 0.$$

• For y_k , as $x_k = y_k$ in the long run, we can take $y_k = X(\tau)$.

$$\frac{\sqrt{t}}{2} \left(1 + \beta_{k-1}\right) \ddot{\boldsymbol{X}}(\tau) + \left(1 - \beta_{k-1}\right) \dot{\boldsymbol{X}}(\tau) + \sqrt{t} \nabla f \left(\boldsymbol{X}(\tau)\right) + o(\sqrt{t}) = 0.$$

$$\blacktriangleright \text{ Hide } \tau \text{ in } \mathbf{X}$$

$$\frac{\sqrt{t}}{2} \left(1 + \beta_{k-1}\right) \ddot{\mathbf{X}} + \left(1 - \beta_{k-1}\right) \dot{\mathbf{X}} + \sqrt{t} \nabla f(\mathbf{X}) + o(\sqrt{t}) = 0.$$

• What next: time to plug in β_k .

Derive the ODE ... (4/4)

• If
$$\beta_k = \frac{k}{k+3}$$
 then $\beta_{k-1} = \frac{k-1}{k+2}$ and
 $\frac{k-1}{k+2} = 1 - \frac{-3}{k+2} \overset{k \gg 2}{\approx} 1 - \frac{3}{k} \overset{k = \frac{\tau}{\sqrt{t}}}{=} 1 - \frac{3\sqrt{t}}{\tau}.$
Then $1 + \beta_{k-1} = 2 - \frac{3\sqrt{t}}{\tau}, \ 1 - \beta_{k-1} = \frac{3\sqrt{t}}{\tau}$
 $\frac{\sqrt{t}}{2} \left(2 - \frac{3\sqrt{t}}{\tau}\right) \ddot{\mathbf{X}} + \frac{3\sqrt{t}}{\tau} \dot{\mathbf{X}} + \sqrt{t} \nabla f(\mathbf{X}) + o(\sqrt{t}) = 0.$

• Divide the whole equation by \sqrt{t}

$$\left(1 - \frac{3\sqrt{t}}{2\tau}\right)\ddot{\boldsymbol{X}} + \frac{3}{\tau}\dot{\boldsymbol{X}} + \nabla f(\boldsymbol{X}) + o(\sqrt{t}) = 0.$$

Review of the derivation

Step-1. Nesterov's accelerated gradient gives

$$\frac{\boldsymbol{x}_{k+1} - \boldsymbol{x}_k}{\sqrt{t}} = \beta_{k-1} \frac{\boldsymbol{x}_k - \boldsymbol{x}_{k-1}}{\sqrt{t}} - \sqrt{t} \nabla f(\boldsymbol{y}_k).$$

► Step-2. Apply Taylor's approximation

$$\frac{\sqrt{t}}{2} \left(1 + \beta_{k-1}\right) \ddot{\boldsymbol{X}} + \left(1 - \beta_{k-1}\right) \dot{\boldsymbol{X}} + \sqrt{t} \nabla f(\boldsymbol{X}) + o(\sqrt{t}) = 0.$$

or

$$\frac{1+\beta_{k-1}}{2}\ddot{\boldsymbol{X}} + \frac{1-\beta_{k-1}}{\sqrt{t}}\dot{\boldsymbol{X}} + \nabla f(\boldsymbol{X}) + o(\sqrt{t}) = 0.$$

• Step-3. Plugin β_{k-1} and take limit $t \to 0$

$$\ddot{\boldsymbol{X}} + \frac{3}{\tau}\dot{\boldsymbol{X}} + \nabla f(\boldsymbol{X}) = 0.$$

About the
$$\beta_k = \frac{3}{k+3}$$

• It is instrumental such that taking limit $t \to 0$ will not cancel \ddot{X} or blow up the ODE.

• Selection of
$$\beta_k = \frac{3}{k+3}$$
 is legit because it satisfies

$$\frac{1-\beta_{k+1}}{\beta_{k+1}^2} \le \frac{1}{\beta_k^2}.$$

Convergence rate using ODE

$$\ddot{\boldsymbol{X}}(\tau) + \frac{3}{\tau} \dot{\boldsymbol{X}}(\tau) + \nabla f \Big(\boldsymbol{X}(\tau) \Big) = 0.$$

Standard ODE theory gives

$$f\left(\boldsymbol{X}(\tau)\right) - f^* \le \mathcal{O}\left(\frac{1}{\tau^2}\right)$$

• As ODE \iff Nesterov's accelerated gradient, so this (partially) explains

$$f(\boldsymbol{x}_k) - f^* \le rac{c \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|}{(k+1)^2}.$$

After stories / What is not discussed

The proof of convergence of AGD for general problem https://angms.science/doc/CVX/CVX_NAGD.pdf https://angms.science/doc/CVX/CVX_NAGDalpha.pdf https://angms.science/doc/CVX/fista_convergence.pdf

► Is AGD really always good?

AGD is not suitable when gradient has noise. In other words, AGD scarifies robustness for speed.

 Extension to other gradient scheme Accelerated Stochastic Gradient Descent Accelerated Primal-Dual Gradient

Extension to other model

What about $\min_{x} ||Ax - b||_1$: non-smooth optimization. What about f non-convex: current research topics.

What actually leads to the acceleration?

ODE, Variational perspective, principle of least action, damped Lagrangian ... in the most recent research within last 7 years!

C'est bon bon. End of lecture