

An Inertial Non-smooth Non-convex Bregman Minimization Framework



Peter Ochs
Mathematical Optimization Group

— 22.11.2021 —



[Nesterov, 2004]:

Optimization is unsolvable!!!

[Nesterov, 2004]:

Optimization is unsolvable!!!

- ◆ This statement is meant “**in general**”.
- ◆ Nature is just too complicated!
- ◆ This is a **mathematical fact** rather than intuition.

↪ study classes of optimization problems.

[Nesterov, 2004]:

Optimization is unsolvable!!!

- ◆ This statement is meant “**in general**”.
- ◆ Nature is just too complicated!
- ◆ This is a **mathematical fact** rather than intuition.

↪ study classes of optimization problems.

Ultimate Goal:

**meta algorithm that automatically detects
and adapts to the problem structure**

Goal of this Talk

Unified framework for the design and analysis of many available algorithms and thereby to explore synergies and to understand the essentials.

Unified framework for the design and analysis of many available algorithms and thereby to explore synergies and to understand the essentials.

Setup: finite dimensional, non-smooth, non-convex, first-order optimization.

Part 1: Model Framework and Global Convergence

- ◆ unified framework for algorithm design
- ◆ convergence analysis
- ◆ examples

Part 2: Inertial Variant based on Nesterov Extrapolation

- ◆ introduce in standard Euclidean setting
- ◆ and generalize to model framework

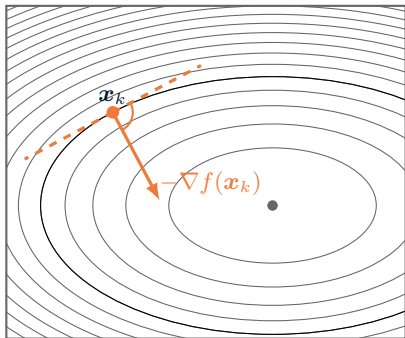
A Simple Illustrative Setting

Smooth optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (\nabla f \text{ is } L\text{-Lipschitz continuous})$$

◆ Gradient Descent update step with step size $\tau > 0$:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)$$



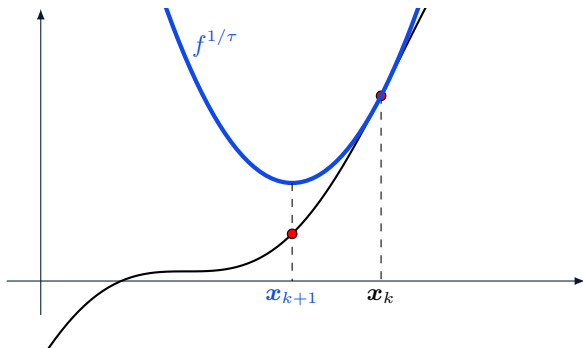
Equivalent Form of Gradient Descent

Equivalent to sequential minimization of quadratic functions:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)$$

$$\Leftrightarrow \mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{=: f^{1/\tau}}$$

(majorizer if $\tau \leq \frac{1}{L}$)



Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (\nabla f \text{ is } L\text{-Lipschitz})$$

Update step:

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle}_{\text{linearization of } f} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (\nabla f \text{ is } L\text{-Lipschitz})$$

Update step:

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle}_{\text{linearization of } f} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Classic setting: (Descent Lemma)

$$|f(\mathbf{x}) - f(\mathbf{x}_k) - \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x})$$

Update step:

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle}_{\text{linearization of } f} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Classic setting: (Descent Lemma)

$$|f(\mathbf{x}) - f(\mathbf{x}_k) - \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x})$$

Update step:

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \underbrace{f_{\mathbf{x}_k}(\mathbf{x})}_{\text{model function}} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Classic setting: (Descent Lemma)

$$|f(\mathbf{x}) - f(\mathbf{x}_k) - \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x})$$

Update step:

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \underbrace{f_{\mathbf{x}_k}(\mathbf{x})}_{\text{model function}} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Model Approximation Property: [Drusvyatskiy et al. 21], [O. et al. 18], [Mukkamala et al. 21], ...

$$|f(\mathbf{x}) - f_{\mathbf{x}_k}(\mathbf{x})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (f \text{ non-smooth})$$

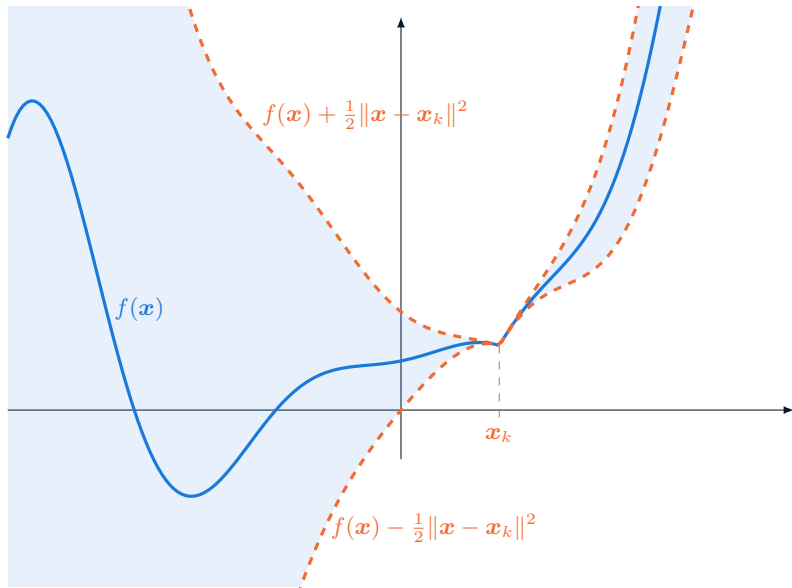
Update step:

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \underbrace{f_{\mathbf{x}_k}(\mathbf{x})}_{\text{model function}} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Model Approximation Property: [Drusvyatskiy et al. 21], [O. et al. 18], [Mukkamala et al. 21], ...

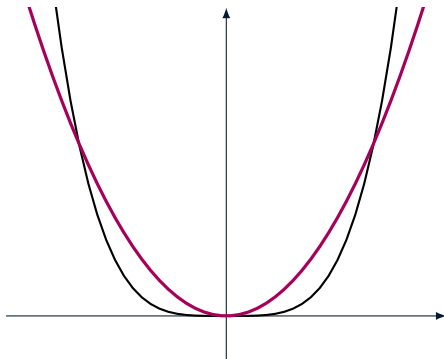
$$|f(\mathbf{x}) - f_{\mathbf{x}_k}(\mathbf{x})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Model Approximation Property



Failure of Quadratic Model Error: \rightsquigarrow relative smoothness

◆ Simple functions like x^4 **do not** allow for quadratic bounds:



◆ Same situation in Matrix Factorization:

$$Q(\mathbf{X}, \mathbf{X}^\top) = \frac{1}{2} \|\mathbf{A} - \mathbf{X}\mathbf{X}^\top\|_F^2 \quad \text{polynomial of degree 4 in } \mathbf{X}$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (f \text{ non-smooth})$$

Update step:

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \underbrace{f_{\mathbf{x}_k}(\mathbf{x})}_{\text{model function}} + \underbrace{\frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{\text{distance to } \mathbf{x}_k}$$

Model Approximation Property: [Drusvyatskiy et al. 21], [O. et al. 18], [Mukkamala et al. 21], ...

$$|f(\mathbf{x}) - f_{\mathbf{x}_k}(\mathbf{x})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (f \text{ non-smooth})$$

Update step:

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \underbrace{f_{\mathbf{x}_k}(\mathbf{x})}_{\text{model function}} + \underbrace{\frac{1}{\tau} D_h(\mathbf{x}, \mathbf{x}_k)}_{\text{Bregman distance to } \mathbf{x}_k}$$

Model Approximation Property: [Drusvyatskiy et al. 21], [O. et al. 18], [Mukkamala et al. 21], ...

$$|f(\mathbf{x}) - f_{\mathbf{x}_k}(\mathbf{x})| \leq L D_h(\mathbf{x}, \mathbf{x}_k)$$

Optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (f \text{ non-smooth})$$

Update step:

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \underbrace{f_{\mathbf{x}_k}(\mathbf{x})}_{\text{model function}} + \underbrace{\frac{1}{\tau} D_h(\mathbf{x}, \mathbf{x}_k)}_{\text{Bregman distance to } \mathbf{x}_k}$$

Model Approximation Property: [Drusvyatskiy et al. 21], [O. et al. 18], [Mukkamala et al. 21], ...

$$|f(\mathbf{x}) - f_{\mathbf{x}_k}(\mathbf{x})| \leq L D_h(\mathbf{x}, \mathbf{x}_k)$$

Generalizes L-smad property (smooth adaptable functions) / relative smoothness

[Birnbaum et al., 2011], [Bauschke et al., 2017], [Bolte et al., 2018], [Lu et al. 2018].

Global Convergence to a Stationary Point

Algorithm: Model BPG ($0 < \inf_k \tau_k \leq \sup_k \tau_k < 1/L$)

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} f_{\mathbf{x}_k}(\mathbf{x}) + \frac{1}{\tau_k} D_h(\mathbf{x}, \mathbf{x}_k)$$

Global Convergence to a Stationary Point

Algorithm: Model BPG ($0 < \inf_k \tau_k \leq \sup_k \tau_k < 1/L$)

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} f_{\mathbf{x}_k}(\mathbf{x}) + \frac{1}{\tau_k} D_h(\mathbf{x}, \mathbf{x}_k)$$

Theorem:

- $(\mathbf{x}_k)_{k \in \mathbb{N}}$ has finite length and converges a critical point of f .
- Standard convergence rates depending on KL-exponent of a Lyapunov function.

Global Convergence to a Stationary Point

Algorithm: Model BPG ($0 < \inf_k \tau_k \leq \sup_k \tau_k < 1/L$)

$$\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} f_{\mathbf{x}_k}(\mathbf{x}) + \frac{1}{\tau_k} D_h(\mathbf{x}, \mathbf{x}_k)$$

Theorem:

- $(\mathbf{x}_k)_{k \in \mathbb{N}}$ has finite length and converges to a critical point of f .
- Standard convergence rates depending on KL-exponent of a Lyapunov function.

Assumptions:

- ◆ D_h generated by Legendre function h that is C^2 on $\operatorname{int} \operatorname{dom} h$. Moreover h is σ_B -strongly convex and $\nabla^2 h$ is bounded on compact convex subsets $B \subset \operatorname{int} \operatorname{dom} h$.
- ◆ $(\mathbf{x}_k)_{k \in \mathbb{N}}$ is bounded (e.g. by coercivity) and set of limit points $\omega(\mathbf{x}_0) \subset \operatorname{int} \operatorname{dom} h$.
- ◆ model consistency: $\|\partial f_{\mathbf{y}}(\mathbf{x})\|_- \leq c \|\mathbf{x} - \mathbf{y}\|$.
- ◆ Functions are definable. (to apply KL-framework; details omitted)

Analysis based on **Lyapunov Function**: $F_L^h(\mathbf{x}, \bar{\mathbf{x}}) := f_{\bar{\mathbf{x}}}(\mathbf{x}) + LD_h(\mathbf{x}, \bar{\mathbf{x}})$

- ◆ Optimization problem:

$$\min_{\mathbf{x}} \underbrace{f_0(\mathbf{x})}_{\text{simple}} + \underbrace{f_1(\mathbf{x})}_{\substack{\text{smooth} \\ \text{non-convex}}}$$

- ◆ Model function:

$$f_{\bar{\mathbf{x}}}(\mathbf{x}) = f_0(\mathbf{x}) + f_1(\bar{\mathbf{x}}) + \langle \mathbf{x} - \bar{\mathbf{x}}, \nabla f_1(\bar{\mathbf{x}}) \rangle$$

- ◆ Update step: (**Proximal Gradient Step**)

$$\begin{aligned} \mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x}} f_0(\mathbf{x}) + f_1(\mathbf{x}_k) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f_1(\mathbf{x}_k) \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2 \\ &= \operatorname{prox}_{\tau f_0}(\mathbf{x}_k - \tau \nabla f_1(\mathbf{x}_k)) \end{aligned}$$

◆ Optimization problem:

$$\min_{\mathbf{x}} \underbrace{f_0(\mathbf{x})}_{\text{simple}} + \underbrace{f_1(\mathbf{x})}_{\substack{\text{smooth} \\ \text{non-convex}}}$$

◆ Model function:

$$f_{\bar{\mathbf{x}}}(\mathbf{x}) = f_0(\mathbf{x}) + f_1(\bar{\mathbf{x}}) + \langle \mathbf{x} - \bar{\mathbf{x}}, \nabla f_1(\bar{\mathbf{x}}) \rangle$$

◆ Update step: (**Bregman Proximal Gradient Method**)

$$\begin{aligned} \mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x}} f_0(\mathbf{x}) + f_1(\mathbf{x}_k) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f_1(\mathbf{x}_k) \rangle + D_h(\mathbf{x}, \mathbf{x}_k) \\ &= \operatorname{prox}_{\tau f_0}^h(\mathbf{x}_k - \tau \nabla f_1(\mathbf{x}_k)) \end{aligned}$$

Bregman distance generated by a Legendre function h :

$$D_h(\mathbf{x}, \bar{\mathbf{x}}) = h(\mathbf{x}) - h(\bar{\mathbf{x}}) - \langle \nabla h(\bar{\mathbf{x}}), \mathbf{x} - \bar{\mathbf{x}} \rangle$$

Examples: (in terms of the Legendre function / kernel function)

◆ Energy: $h(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2$

◆ Burg's entropy: $h(\mathbf{x}) = - \sum_i \log(\mathbf{x}_i)$

◆ Boltzmann–Shannon entropy: $h(\mathbf{x}) = \sum_i \mathbf{x}_i \log(\mathbf{x}_i)$

◆ Polynomials of type: $h(\mathbf{x}) = \sum_{j=1}^m a_j \|\mathbf{x}\|_2^{2j}$

For example in Phase retrieval: $h(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\|_2^4 + \frac{1}{2} \|\mathbf{x}\|_2^2$; (e.g. [\[Bolte et al. 2018\]](#)).

◆ Fermi-Dirac: $h(x) = x \log(x) + (1 - x) \log(1 - x)$

◆ Hellinger function: $h(x) = -\sqrt{1 - x^2}$

Matrix Factorization (non-smooth, non-convex)

$$\min_{\mathbf{X}, \mathbf{Y}} \frac{1}{2} \|\mathbf{A} - \mathbf{XY}\|^2 + R_1(\mathbf{X}) + R_2(\mathbf{Y})$$

- ◆ $\frac{1}{2} \|\mathbf{A} - \mathbf{XY}\|^2$ does not have Lipschitz gradient.
- ◆ Line search can be too expensive (requires matrix multiplication).
- ◆ *Alternating Minimization*: limited convergence guarantees, sometimes biased, not applicable for $\mathbf{X} = \mathbf{Y}$.

Matrix Factorization (non-smooth, non-convex)

$$\min_{\mathbf{X}, \mathbf{Y}} \frac{1}{2} \|\mathbf{A} - \mathbf{XY}\|^2 + R_1(\mathbf{X}) + R_2(\mathbf{Y})$$

- ◆ $\frac{1}{2} \|\mathbf{A} - \mathbf{XY}\|^2$ does not have Lipschitz gradient.
- ◆ Line search can be too expensive (requires matrix multiplication).
- ◆ *Alternating Minimization*: limited convergence guarantees, sometimes biased, not applicable for $\mathbf{X} = \mathbf{Y}$.

Apply Model BPG: (= Bregman Proximal Gradient Method) ([Bolte et al. 2018])

Matrix Factorization (non-smooth, non-convex)

$$\min_{\mathbf{X}, \mathbf{Y}} \frac{1}{2} \|\mathbf{A} - \mathbf{XY}\|^2 + R_1(\mathbf{X}) + R_2(\mathbf{Y})$$

- ◆ $\frac{1}{2} \|\mathbf{A} - \mathbf{XY}\|^2$ does not have Lipschitz gradient.
- ◆ Line search can be too expensive (requires matrix multiplication).
- ◆ *Alternating Minimization*: limited convergence guarantees, sometimes biased, not applicable for $\mathbf{X} = \mathbf{Y}$.

Apply Model BPG: (= Bregman Proximal Gradient Method) ([Bolte et al. 2018])

- ◆ **New Bregman Distance:** [Mukkamala, O. 19] (see also [Teboulle, Vaisbourd 20])

$$h(\mathbf{X}, \mathbf{Y}) = 3 \left(\frac{\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2}{2} \right)^2 + \|\mathbf{A}\|_F \left(\frac{\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2}{2} \right).$$

- ◆ Symmetric MF setting developed in [Dragomir, Bolte, d'Aspremont 19].

Deep Linear Neural Networks or Deep Matrix Factorization:

$$\min_{W_1, \dots, W_N} \frac{1}{2} \|\mathbf{Y} - W_1 W_2 \cdots W_N \mathbf{X}\|_F^2 + \sum_{i=1}^N \mathcal{R}_i(W_i)$$

Deep Linear Neural Networks or Deep Matrix Factorization:

$$\min_{W_1, \dots, W_N} \frac{1}{2} \|\mathbf{Y} - W_1 W_2 \cdots W_N \mathbf{X}\|_F^2 + \sum_{i=1}^N \mathcal{R}_i(W_i)$$

◆ Write: $\mathbf{W} := (W_1, \dots, W_N)$ and $\|\mathbf{W}\|_F^2 := \sum_{i=1}^N \|W_i\|_F^2$.

◆ [Mukkamala et al. 2021] shows relative smoothness w.r.t.

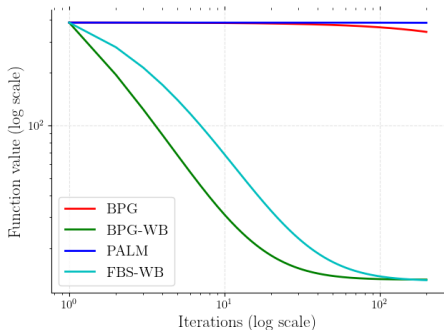
$$h(\mathbf{W}) = \begin{cases} \frac{\|\mathbf{X}\|_F^2}{N^{N-2}} \|\mathbf{W}\|_F^{2N} + \frac{\|\mathbf{Y}\|_F \|\mathbf{X}\|_F}{(N-2)^{\frac{N-2}{2}}} \|\mathbf{W}\|_F^N, & \text{if } N \text{ is even} \\ \frac{\|\mathbf{X}\|_F^2}{N^{N-2}} \|\mathbf{W}\|_F^{2N} + \frac{\|\mathbf{Y}\|_F \|\mathbf{X}\|_F}{(N-1)^{\frac{N-1}{2}}} \left(\|\mathbf{W}\|_F^2 + 1 \right)^{\frac{N+1}{2}}, & \text{if } N \text{ is odd.} \end{cases}$$

optimization / training with **constant step size rule** using BPG.

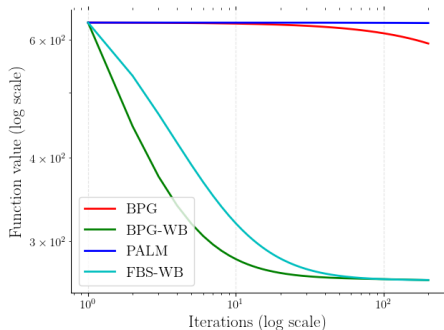
◆ **Non-linear Neural Networks:** see [Mukkamala, PhD Thesis, 2021].

Deep Matrix Completion on MovieLens Datasets

Deep Matrix Completion on MovieLens Datasets:



L2-Regularization ($N = 4$)
MovieLens-100K



L1-Regularization ($N = 4$)
MovieLens-100K

Clipped Proximal Gradient Method

◆ Optimization problem:

$$\min_{\mathbf{x}} \underbrace{f_0(\mathbf{x})}_{\text{simple}} + \underbrace{f_1(\mathbf{x})}_{\substack{\text{smooth} \\ \text{non-convex} \\ \text{non-negative}}}$$

◆ Model function:

$$f_{\bar{\mathbf{x}}}(\mathbf{x}) = f_0(\mathbf{x}) + \max\left(0, f_1(\bar{\mathbf{x}}) + \langle \mathbf{x} - \bar{\mathbf{x}}, \nabla f_1(\bar{\mathbf{x}}) \rangle\right)$$

Proximal Newton Method

◆ Optimization problem:

$$\min_{\mathbf{x}} \underbrace{f_0(\mathbf{x})}_{\text{simple}} + \underbrace{f_1(\mathbf{x})}_{\substack{C^2\text{-smooth} \\ \text{strongly convex}}}$$

◆ Model function:

$$f_{\bar{\mathbf{x}}}(\mathbf{x}) = f_0(\mathbf{x}) + f_1(\bar{\mathbf{x}}) + \langle \mathbf{x} - \bar{\mathbf{x}}, \nabla f_1(\bar{\mathbf{x}}) \rangle + \langle \mathbf{x} - \bar{\mathbf{x}}, \nabla^2 f_1(\bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}}) \rangle$$

ProxDescent [Lewis and Wright, 2016], [Drusvyatskiy and Lewis, 2016]

◆ **Optimization problem:**

$$\min_{\mathbf{x}} \underbrace{f_0(\mathbf{x})}_{\substack{\text{non-smooth} \\ \text{convex}}} + \underbrace{g\left(\underbrace{F(\mathbf{x})}_{\text{smooth}}\right)}_{\substack{\text{non-smooth} \\ \text{convex} \\ \text{finite-valued}}}$$

◆ **Model function:** ($DF(\bar{\mathbf{x}})$ is the Jacobian matrix of F at $\bar{\mathbf{x}}$)

$$f_{\bar{\mathbf{x}}}(\mathbf{x}) = f_0(\mathbf{x}) + g\left(F(\bar{\mathbf{x}}) + DF(\bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})\right)$$

Example: Levenberg–Marquardt Algorithm for Non-linear least-squares problem

More Examples:

- ◆ Outer-linearization of the composite problem:

$$\min_{\mathbf{x}} \underbrace{f_0(\mathbf{x})}_{\substack{\text{non-smooth} \\ \text{convex}}} + \underbrace{g}_{\substack{\text{smooth} \\ \text{non-convex} \\ \text{non-decreasing}}}\left(\underbrace{F(\mathbf{x})}_{\substack{\text{non-smooth} \\ \text{coordinate-wise} \\ \text{convex}}}\right)$$

- ◆ Combine previous concepts of model functions, e.g. for block problems.
- ◆ Higher order approximations.

Be creative! Design good model functions for **your** problem.

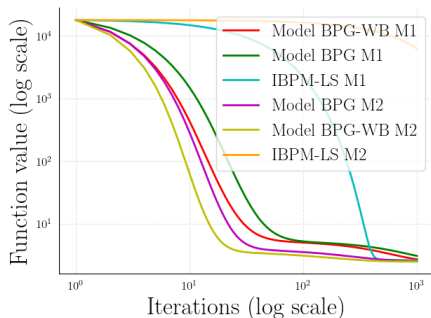
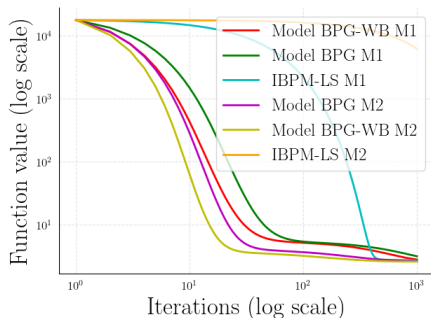
Comparison of Model Functions for a Phase Retrieval problem

$$\min_{\mathbf{x}} \frac{1}{M} \sum_{i=1}^M (\mathbf{x}^\top A_i \mathbf{x} - \mathbf{b}_i)^2 + \mathcal{R}(\mathbf{x})$$

◆ **Model 1:** Bregman Proximal Gradient Model.

◆ **Model 2:** Non-linear composite (\sim ProxDescent): Linearize inner part of $|(\mathbf{x}^\top A_i \mathbf{x} - \mathbf{b}_i)^2|$.

For both the model functions we use $h(\mathbf{x}) = \frac{1}{4}\|\mathbf{x}\|^4 + \frac{1}{2}\|\mathbf{x}\|^2$.



Model BPG with Model 2 is faster in both the settings.

Conclusion of Part 1:

- ◆ unified flexible algorithm with several interesting special cases
- ◆ convergence in non-smooth non-convex setting to stationary point
- ◆ do not require line search

However:

If problem size permits, line search / backtracking can yield a significant speed up.

Motivation for Part 2: Improve practical performance

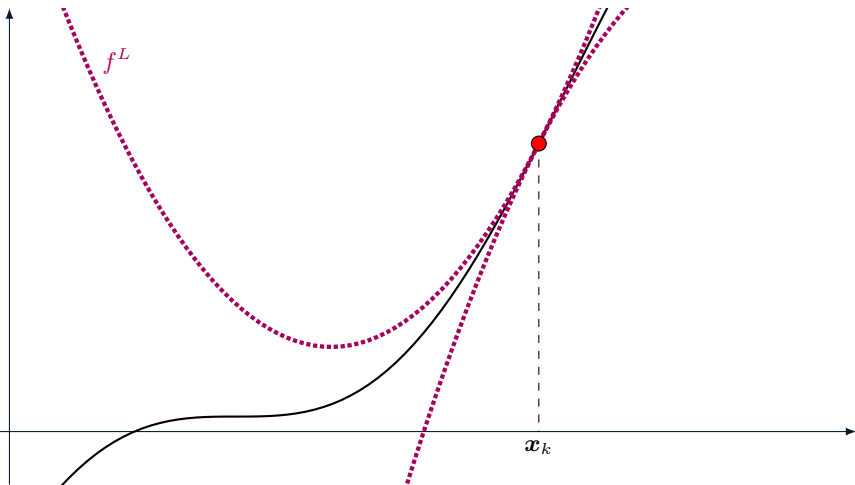
- ◆ Adaptive Bregman distances (similar to quasi-Newton methods).
- ◆ Acceleration / Momentum (**this talk**) based on Nesterov extrapolation:

$$\begin{aligned}\mathbf{y}_k &= \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{x}_{k+1} &= \mathbf{y}_k - \tau \nabla f(\mathbf{y}_k)\end{aligned}$$

with novel CONvex ConcAve INertial backtracking procedure.

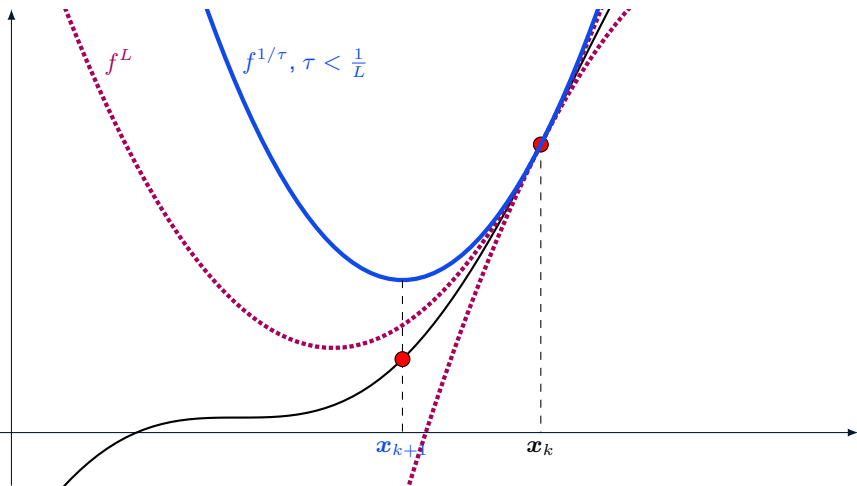
Gradient Descent: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{x}_k)}.$$



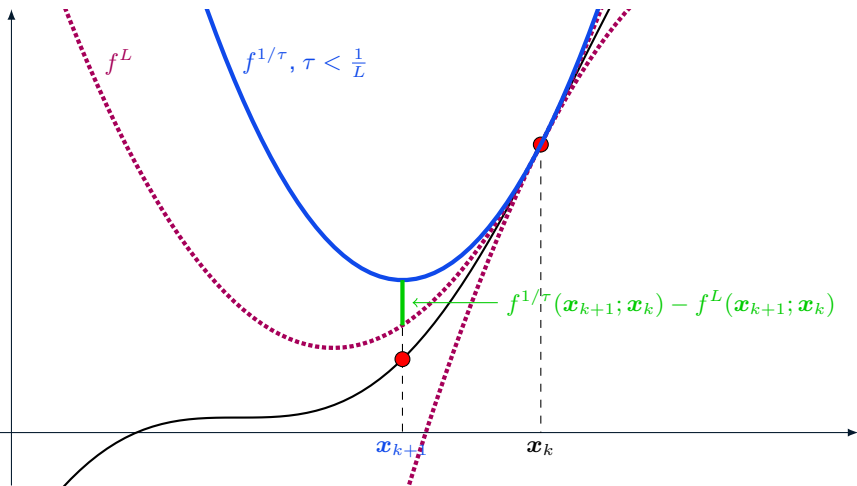
Gradient Descent: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{x}_k)}.$$



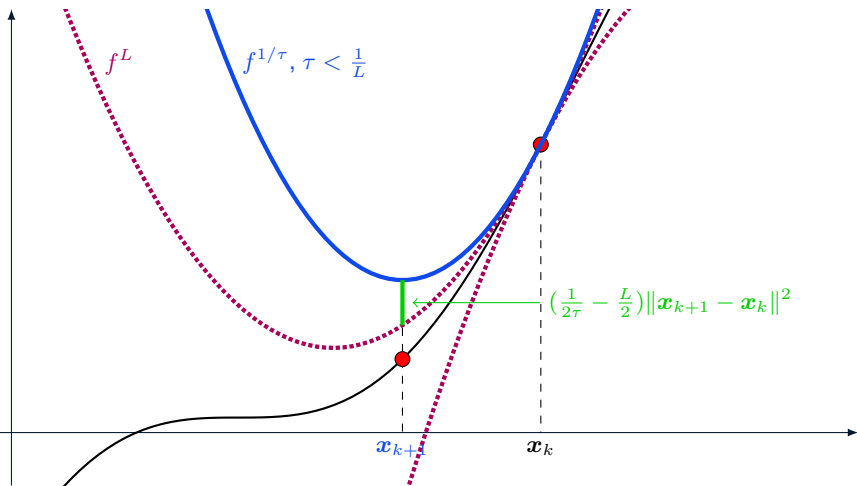
Gradient Descent: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{x}_k)}.$$



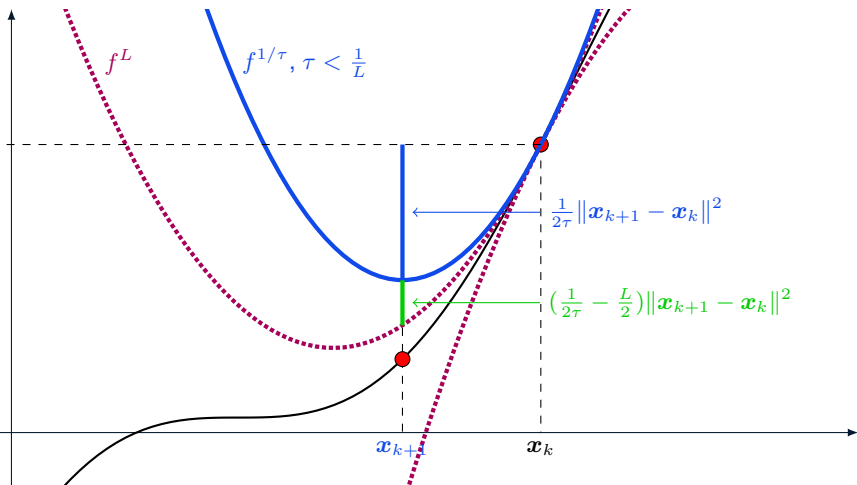
Gradient Descent: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{x}_k)}.$$

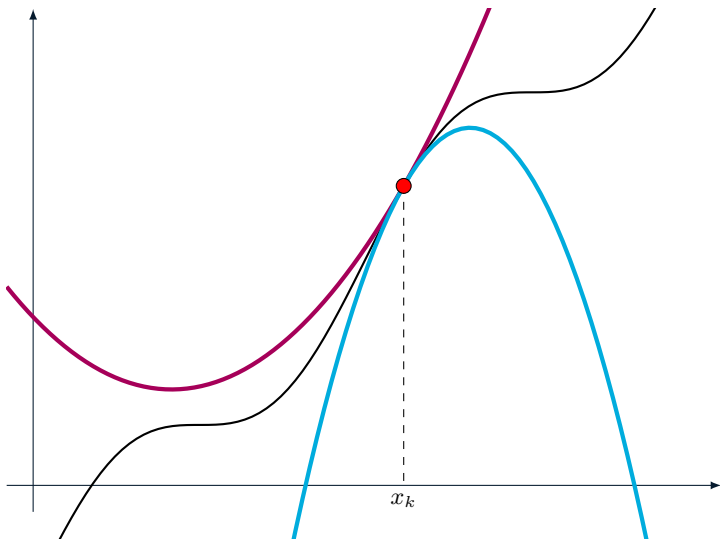


Gradient Descent: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{x}_k)}.$$



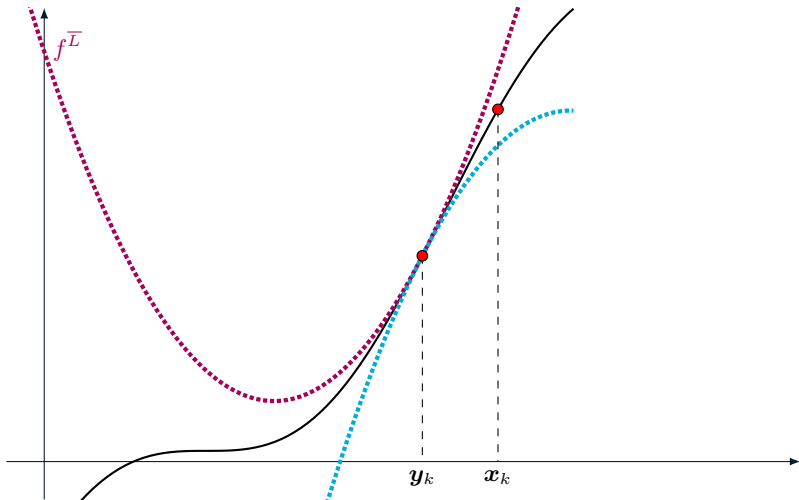
Different Constants for Tight Upper and Lower Quadratic Bounds



$$-\frac{\underline{L}}{2}\|\mathbf{x} - \mathbf{x}_k\|^2 \leq f(\mathbf{x}) - f(\mathbf{x}_k) - \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle \leq \frac{\overline{L}}{2}\|\mathbf{x} - \mathbf{x}_k\|^2$$

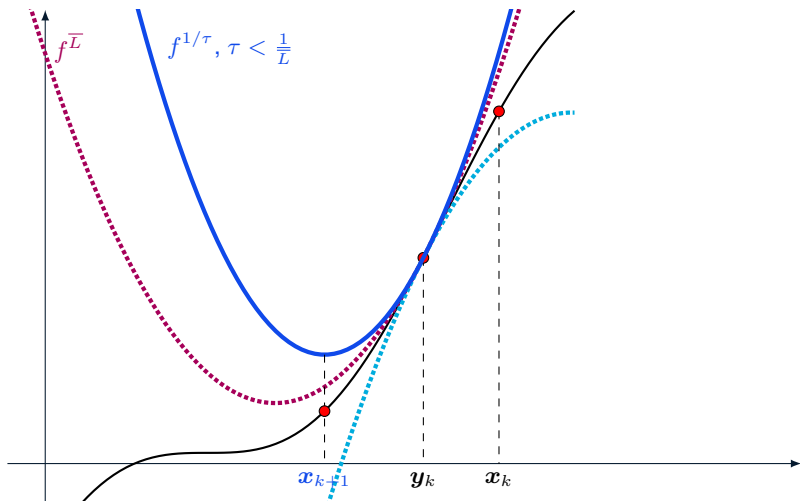
Nesterov's Extrapolation: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{y}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{y}_k)}, \quad \mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$$



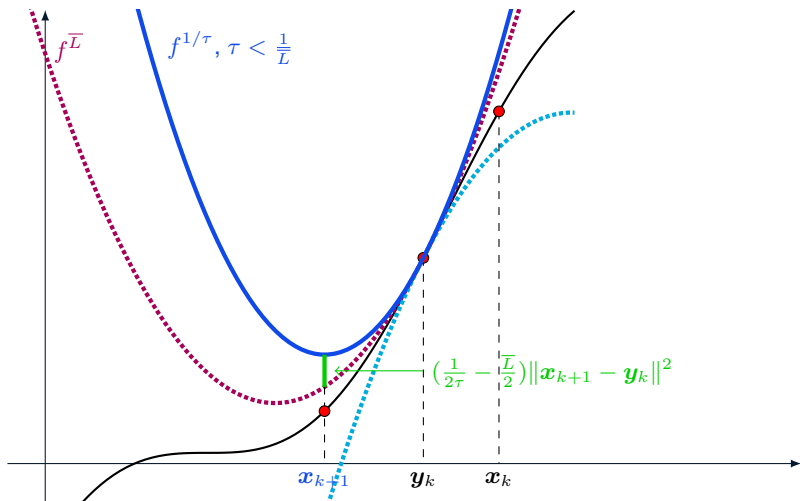
Nesterov's Extrapolation: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{y}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{y}_k)}, \quad \mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$$



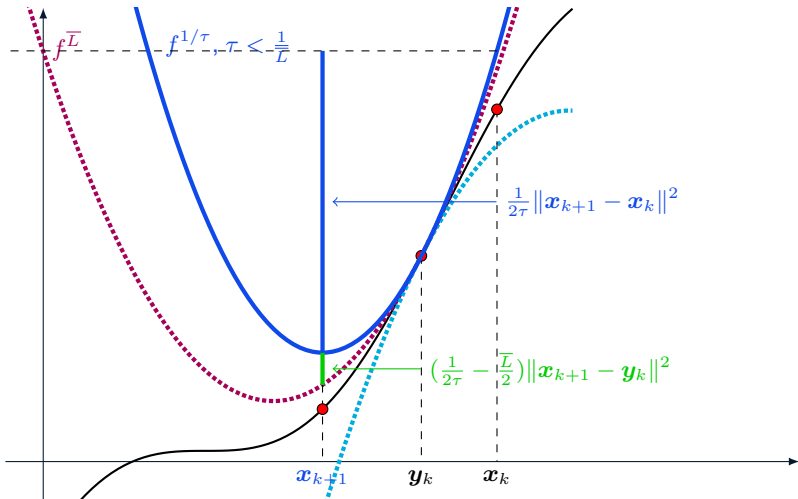
Nesterov's Extrapolation: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{y}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{y}_k)}, \quad \mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$$



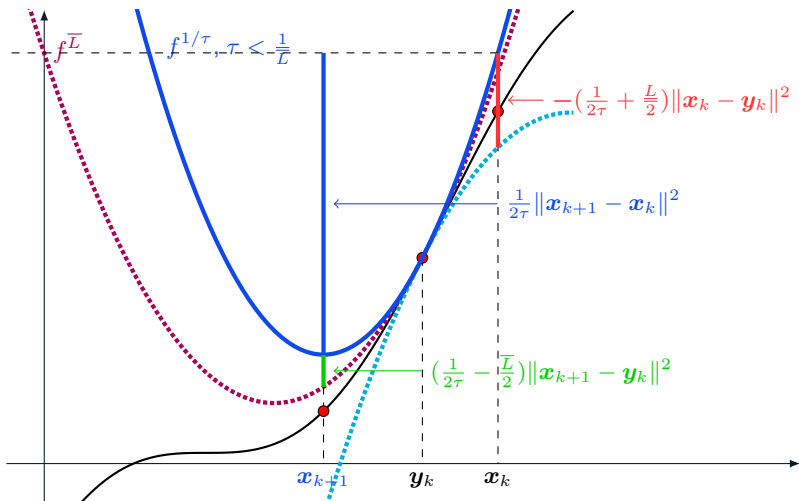
Nesterov's Extrapolation: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{y}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{y}_k)}, \quad \mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$$



Nesterov's Extrapolation: Sufficient decrease

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{y}_k\|^2}_{=: f^{1/\tau}(\mathbf{x}; \mathbf{y}_k)}, \quad \mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$$



◆ From geometric derivation:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2.$$

- ◆ From geometric derivation:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2.$$

- ◆ Define $F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) := f(\mathbf{x}_{k+1}) + \frac{1}{2\tau} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$, then this is equivalent to

$$F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) \leq F_\tau(\mathbf{x}_k, \mathbf{x}_{k-1}) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2$$

Proof of Convergence for Nesterov's Extrapolation

◆ From geometric derivation:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2.$$

◆ Define $F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) := f(\mathbf{x}_{k+1}) + \frac{1}{2\tau} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$, then this is equivalent to

$$F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) \leq F_\tau(\mathbf{x}_k, \mathbf{x}_{k-1}) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2$$

Decrease Condition:	$\left(\frac{L}{2} + \frac{1}{2\tau}\right) \ \mathbf{x}_k - \mathbf{y}_k\ ^2 - \frac{1}{2\tau} \ \mathbf{x}_{k-1} - \mathbf{x}_k\ ^2 \leq 0$
------------------------	---

Proof of Convergence for Nesterov's Extrapolation

- From geometric derivation:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2.$$

- Define $F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) := f(\mathbf{x}_{k+1}) + \frac{1}{2\tau} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$, then this is equivalent to

$$F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) \leq F_\tau(\mathbf{x}_k, \mathbf{x}_{k-1}) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2$$

Decrease Condition: $\left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2 \leq 0$

- For $\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$:

$$\sup_k \gamma_k^2 \left(\frac{L}{2} + \frac{1}{2\tau}\right) < \frac{1}{2\tau} \iff \sup_k \gamma_k^2 < \frac{\tau^{-1}}{L + \tau^{-1}} \stackrel{\tau^{-1} = \bar{L}}{=} \frac{\bar{L}}{L + \bar{L}} \stackrel{L = \bar{L}}{=} \frac{1}{2}.$$

[Wen, Chen, Pong 2017]

- For f convex, i.e., $L = 0$, we obtain convergence for $\sup_{k \in \mathbb{N}} \gamma_k < 1$.

Proof of Convergence for Nesterov's Extrapolation

- From geometric derivation:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|^2.$$

- Define $F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) := f(\mathbf{x}_{k+1}) + \frac{1}{2\tau} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$, then this is equivalent to

$$F_\tau(\mathbf{x}_{k+1}, \mathbf{x}_k) \leq F_\tau(\mathbf{x}_k, \mathbf{x}_{k-1}) + \left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2$$

Decrease Condition: $\left(\frac{L}{2} + \frac{1}{2\tau}\right) \|\mathbf{x}_k - \mathbf{y}_k\|^2 - \frac{1}{2\tau} \|\mathbf{x}_{k-1} - \mathbf{x}_k\|^2 \leq 0$

- For $\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1})$:

$$\sup_k \gamma_k^2 \left(\frac{L}{2} + \frac{1}{2\tau}\right) < \frac{1}{2\tau} \iff \sup_k \gamma_k^2 < \frac{\tau^{-1}}{L + \tau^{-1}} \stackrel{\tau^{-1} = \bar{L}}{=} \frac{\bar{L}}{L + \bar{L}} \stackrel{L = \bar{L}}{=} \frac{1}{2}.$$

[Wen, Chen, Pong 2017]

- For f convex, i.e., $L = 0$, we obtain convergence for $\sup_{k \in \mathbb{N}} \gamma_k < 1$.

\rightsquigarrow tight lower bounds L allow for large extrapolation.

Convex–Concave Inertial Backtracking (CoCaln)

Convex–Concave Inertial Backtracking: (simple version) [Mukkamala, O., Sabach, Pock, 2021]

Convex–Concave Inertial Backtracking: (simple version) [Mukkamala, O., Sabach, Pock, 2021]

◆ Find \underline{L}_k (via backtracking) such that

$$\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \text{with } \gamma_k = \sqrt{\frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k}} \quad \begin{array}{l} \text{[extrapolation} \\ \text{step]} \end{array}$$

satisfies

$$f(\mathbf{x}_k) \geq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle - \frac{\underline{L}_k}{2} \|\mathbf{x}_k - \mathbf{y}_k\|^2. \quad \text{[lower bound]}$$

Convex–Concave Inertial Backtracking: (simple version) [Mukkamala, O., Sabach, Pock, 2021]

◆ Find \underline{L}_k (via backtracking) such that

$$\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \text{with } \gamma_k = \sqrt{\frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k}} \quad \text{[extrapolation step]}$$

satisfies

$$f(\mathbf{x}_k) \geq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle - \frac{\underline{L}_k}{2} \|\mathbf{x}_k - \mathbf{y}_k\|^2. \quad \text{[lower bound]}$$

◆ Find $\bar{L}_k \geq \bar{L}_{k-1}$ (via backtracking) such that

$$\mathbf{x}_{k+1} = \underset{\mathbf{x}}{\operatorname{argmin}} \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\bar{L}_k}{2} \|\mathbf{x} - \mathbf{y}_k\|^2 \quad \text{[update]}$$

satisfies

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_{k+1} - \mathbf{y}_k \rangle + \frac{\bar{L}_k}{2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2. \quad \text{[upper bound]}$$

Convex–Concave Inertial Backtracking: $\min_{\mathbf{x}} f(\mathbf{x})$

◆ Find \underline{L}_k and γ_k (via backtracking) such that

$$\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \text{[extrapolation step]}$$

satisfies

$$\left\{ \begin{array}{l} f(\mathbf{x}_k) \geq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_k - \mathbf{y}_k \rangle - \frac{\underline{L}_k}{2} \|\mathbf{x}_k - \mathbf{y}_k\|^2 \quad \text{[lower bound]} \\ \|\mathbf{x}_k - \mathbf{y}_k\|^2 \leq \frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2. \quad \text{[extrapolation bound]} \end{array} \right.$$

◆ Find $\bar{L}_k \geq \bar{L}_{k-1}$ (via backtracking) such that

$$\mathbf{x}_{k+1} = \underset{\mathbf{x}}{\operatorname{argmin}} \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \frac{\bar{L}_k}{2} \|\mathbf{x} - \mathbf{y}_k\|^2 \quad \text{[update]}$$

satisfies

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_{k+1} - \mathbf{y}_k \rangle + \frac{\bar{L}_k}{2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2. \quad \text{[upper bound]}$$

Convex–Concave Inertial Backtracking: $\min_{\mathbf{x}} f(\mathbf{x})$

◆ Find \underline{L}_k and γ_k (via backtracking) such that

$$\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \text{[extrapolation step]}$$

satisfies

$$\left\{ \begin{array}{l} f(\mathbf{x}_k) \geq f_{\mathbf{y}_k}(\mathbf{x}_k) - \frac{\underline{L}_k}{2} \|\mathbf{x}_k - \mathbf{y}_k\|^2 \quad \text{[lower bound]} \\ \|\mathbf{x}_k - \mathbf{y}_k\|^2 \leq \frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2. \quad \text{[extrapolation bound]} \end{array} \right.$$

◆ Find $\bar{L}_k \geq \bar{L}_{k-1}$ (via backtracking) such that

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} f_{\mathbf{y}_k}(\mathbf{x}) + \frac{\bar{L}_k}{2} \|\mathbf{x} - \mathbf{y}_k\|^2 \quad \text{[update]}$$

satisfies

$$f(\mathbf{x}_{k+1}) \leq f_{\mathbf{y}_k}(\mathbf{x}_{k+1}) + \frac{\bar{L}_k}{2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|^2. \quad \text{[upper bound]}$$

Convex–Concave Inertial Backtracking: $\min_{\mathbf{x}} f(\mathbf{x})$

◆ Find \underline{L}_k and γ_k (via backtracking) such that

$$\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \text{[extrapolation step]}$$

satisfies

$$\left\{ \begin{array}{l} f(\mathbf{x}_k) \geq f_{\mathbf{y}_k}(\mathbf{x}_k) - \underline{L}_k D_h(\mathbf{x}_k, \mathbf{y}_k) \quad \text{[lower bound]} \\ D_h(\mathbf{x}_k, \mathbf{y}_k) \leq \frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k} D_h(\mathbf{x}_{k-1}, \mathbf{x}_k). \quad \text{[extrapolation bound]} \end{array} \right.$$

◆ Find $\bar{L}_k \geq \bar{L}_{k-1}$ (via backtracking) such that

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} f_{\mathbf{y}_k}(\mathbf{x}) + \bar{L}_k D_h(\mathbf{x}, \mathbf{y}_k) \quad \text{[update]}$$

satisfies

$$f(\mathbf{x}_{k+1}) \leq f_{\mathbf{y}_k}(\mathbf{x}_{k+1}) + \bar{L}_k D_h(\mathbf{x}_{k+1}, \mathbf{y}_k). \quad \text{[upper bound]}$$

Here, parameters for convex model functions. Can be extended to weakly convex model function.

Convex–Concave Inertial Backtracking: $\min_{\mathbf{x}} f(\mathbf{x})$

- Find \underline{L}_k and γ_k (via backtracking) such that

$$\mathbf{y}_k = \mathbf{x}_k + \gamma_k(\mathbf{x}_k - \mathbf{x}_{k-1}) \quad \text{[extrapolation step]}$$

satisfies

$$\left\{ \begin{array}{l} f(\mathbf{x}_k) \geq f_{\mathbf{y}_k}(\mathbf{x}_k) - \underline{L}_k D_h(\mathbf{x}_k, \mathbf{y}_k) \quad \text{[lower bound]} \\ D_h(\mathbf{x}_k, \mathbf{y}_k) \leq \frac{\bar{L}_{k-1}}{\bar{L}_{k-1} + \underline{L}_k} D_h(\mathbf{x}_{k-1}, \mathbf{x}_k). \quad \text{[extrapolation bound]} \end{array} \right.$$

- Find $\bar{L}_k \geq \bar{L}_{k-1}$ (via backtracking) such that

$$\mathbf{x}_{k+1} = \underset{\mathbf{x}}{\operatorname{argmin}} f_{\mathbf{y}_k}(\mathbf{x}) + \bar{L}_k D_h(\mathbf{x}, \mathbf{y}_k) \quad \text{[update]}$$

satisfies

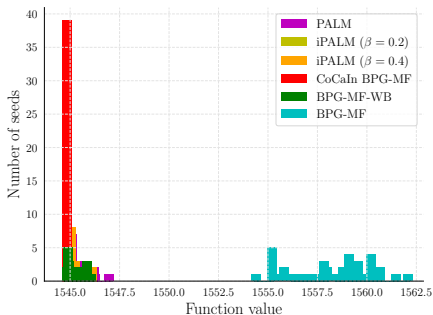
$$f(\mathbf{x}_{k+1}) \leq f_{\mathbf{y}_k}(\mathbf{x}_{k+1}) + \bar{L}_k D_h(\mathbf{x}_{k+1}, \mathbf{y}_k). \quad \text{[upper bound]}$$

Global convergence to a stationary point in KL-framework.

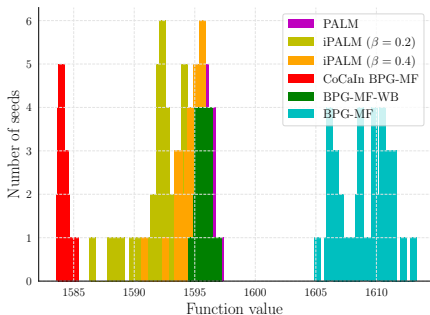
Matrix Factorization

$$\min_{\mathbf{X}, \mathbf{Y}} \frac{1}{2} \|\mathbf{A} - \mathbf{X}\mathbf{Y}\|^2 + R_1(\mathbf{X}) + R_2(\mathbf{Y})$$

Statistical evaluation on simple matrix factorization ($N = 2$).



L2-regularization

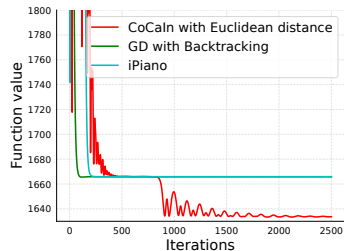
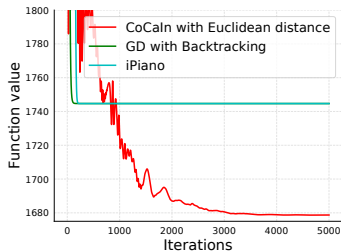


L1-regularization

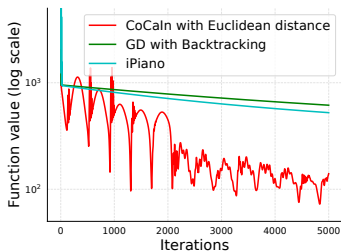
Inertial strategy often results in a smaller objective value.

Matrix Factorization

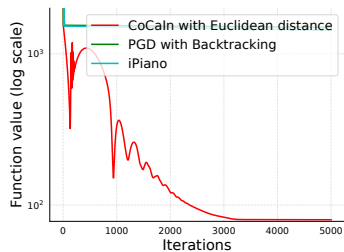
toy experiment



Medulloblastoma dataset
[Brunet, Tamayo,
Golub, Mesirov 2004]



L2 regularization



L1 regularization

Conclusion:

- ◆ Model BPG and CoCaIn Model BPG:

unified framework for design and analysis of algorithms.

- ◆ Global Convergence to a stationary point
- ◆ Good performance in experiments.
- ◆ Requires expert knowledge of the problem and models.
- ◆ CoCaIn Backtracking seeks to locally adapt to the convexity of the objective.