

Optimization in Machine Learning: Majorization Minimization Method

Ziping Zhao

School of Information Science and Technology
ShanghaiTech University, Shanghai, China

CS182: Introduction to Machine Learning (Spring 2023)
<http://cs182.sist.shanghaitech.edu.cn>

Majorization Minimization

Consider the following problem

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t. } \mathbf{x} \in \mathcal{X}, \quad (1)$$

where \mathcal{X} is a closed convex set; $f(\cdot)$ may be non-convex and/or nonsmooth.

- **Challenge:** For a general $f(\cdot)$, problem (1) can be difficult to solve.
- **Majorization Minimization (MM):** Iteratively generate $\{\mathbf{x}^r\}$ as follows:

$$\mathbf{x}^r \in \min_{\mathbf{x}} u(\mathbf{x}, \mathbf{x}^{r-1}) \quad \text{s.t. } \mathbf{x} \in \mathcal{X}, \quad (2)$$

where $u(\mathbf{x}, \mathbf{x}^{r-1})$ is a surrogate function of $f(\mathbf{x})$, satisfying

1. $u(\mathbf{x}, \mathbf{x}^r) \geq f(\mathbf{x}), \quad \forall \mathbf{x}^r, \mathbf{x} \in \mathcal{X};$
2. $u(\mathbf{x}^r, \mathbf{x}^r) = f(\mathbf{x}^r);$

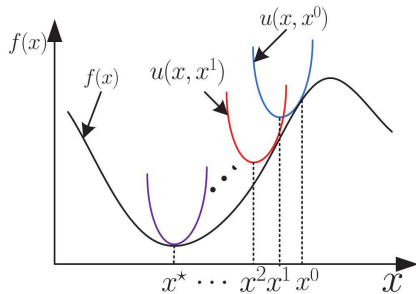


Figure: An pictorial illustration of MM algorithm.

Property 1. $\{f(\mathbf{x}^r)\}$ is nonincreasing, i.e., $f(\mathbf{x}^r) \leq f(\mathbf{x}^{r-1})$, $\forall r = 1, 2, \dots$

Proof.

$$f(\mathbf{x}^r) \leq u(\mathbf{x}^r, \mathbf{x}^{r-1}) \leq u(\mathbf{x}^{r-1}, \mathbf{x}^{r-1}) = f(\mathbf{x}^{r-1}).$$



- The nonincreasing property of $\{f(\mathbf{x}^r)\}$ implies that $f(\mathbf{x}^r) \rightarrow f^\infty$. But how about the convergence of the iterates $\{\mathbf{x}^r\}$?

Technical Preliminaries

- ▶ **Limit point:** \mathbf{x}^∞ is a limit point of $\{\mathbf{x}^r\}$ if there exists a subsequence of $\{\mathbf{x}^r\}$ that converges to \mathbf{x}^∞ .
 - Note that every bounded sequence in \mathbb{R}^n has a limit point (or convergent subsequence).
- ▶ **Directional derivative:** Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be a function where $\mathcal{D} \subseteq \mathbb{R}^m$ is a convex set. The directional derivative of f at point \mathbf{x} in direction \mathbf{d} is defined by

$$f'(\mathbf{x}; \mathbf{d}) \triangleq \liminf_{\lambda \downarrow 0} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda}.$$

- If f is differentiable, then $f'(\mathbf{x}; \mathbf{d}) = \mathbf{d}^\top \nabla f(\mathbf{x})$.
- ▶ **Stationary point:** $\mathbf{x} \in \mathcal{X}$ is a stationary point of $f(\cdot)$ if

$$f'(\mathbf{x}; \mathbf{d}) \geq 0, \quad \forall \mathbf{d} \text{ such that } \mathbf{x} + \mathbf{d} \in \mathcal{D}. \quad (3)$$

- A stationary point may be a local min., a local max. or a saddle point;
 - If $\mathcal{D} = \mathbb{R}^n$ and f is differentiable, then $(3) \iff \nabla f(\mathbf{x}) = \mathbf{0}$.

Convergence of MM

- **Assumption 1** $u(\cdot, \cdot)$ satisfies the following conditions

$$\begin{cases} u(\mathbf{y}, \mathbf{y}) = f(\mathbf{y}), & \forall \mathbf{y} \in \mathcal{X}, \end{cases} \quad (4a)$$

$$\begin{cases} u(\mathbf{x}, \mathbf{y}) \geq f(\mathbf{x}), & \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}, \end{cases} \quad (4b)$$

$$\begin{cases} u'(\mathbf{x}, \mathbf{y}; \mathbf{d})|_{\mathbf{x}=\mathbf{y}} = f'(\mathbf{y}; \mathbf{d}) & \forall \mathbf{d} \text{ with } \mathbf{y} + \mathbf{d} \in \mathcal{X}, \end{cases} \quad (4c)$$

$$\begin{cases} u(\mathbf{x}, \mathbf{y}) \text{ is continuous in } (\mathbf{x}, \mathbf{y}) \end{cases} \quad (4d)$$

- (4c) means the 1st order local behavior of $u(\cdot, \mathbf{x}^{r-1})$ is the same as $f(\cdot)$.

Theorem (Convergence of MM [1])

Assume that Assumption 1 is satisfied. Then every limit point of the iterates generated by MM algorithm is a stationary point of problem (1).

Convergence of MM

Proof of Theorem 1:

From **Property 1**, we know that $f(\mathbf{x}^{r+1}) \leq u(\mathbf{x}^{r+1}, \mathbf{x}^r) \leq u(\mathbf{x}, \mathbf{x}^r)$, $\forall \mathbf{x} \in \mathcal{X}$. Now assume that there exists a subsequence $\{\mathbf{x}^{r_j}\}$ of $\{\mathbf{x}^r\}$ converging to a limit point \mathbf{z} , i.e., $\lim_{j \rightarrow \infty} \mathbf{x}^{r_j} = \mathbf{z}$. Then

$$u(\mathbf{x}^{r_{j+1}}, \mathbf{x}^{r_{j+1}}) = f(\mathbf{x}^{r_{j+1}}) \leq f(\mathbf{x}^{r_{j+1}}) \leq u(\mathbf{x}^{r_{j+1}}, \mathbf{x}^{r_j}) \leq u(\mathbf{x}, \mathbf{x}^{r_j}), \quad \forall \mathbf{x} \in \mathcal{X}.$$

Letting $j \rightarrow \infty$, we obtain $u(\mathbf{z}, \mathbf{z}) \leq u(\mathbf{x}, \mathbf{z})$, $\forall \mathbf{x} \in \mathcal{X}$, which implies that

$$u'(\mathbf{x}, \mathbf{z}; \mathbf{d})|_{\mathbf{x}=\mathbf{z}} \geq \mathbf{0}, \quad \forall \mathbf{z} + \mathbf{d} \in \mathcal{X}.$$

Combining the above inequality with (4c) (i.e., $u'(\mathbf{x}, \mathbf{y}; \mathbf{d})|_{\mathbf{x}=\mathbf{y}} = f'(\mathbf{y}; \mathbf{d})$, $\forall \mathbf{d}$ with $\mathbf{y} + \mathbf{d} \in \mathcal{X}$), we have

$$f'(\mathbf{z}; \mathbf{d}) \geq \mathbf{0}, \quad \forall \mathbf{z} + \mathbf{d} \in \mathcal{X}.$$

Applications — Nonnegative Least Squares

In many engineering applications, we encounter the following problem

$$(\text{NLS}) \quad \min_{\mathbf{x} \geq \mathbf{0}} \|\mathbf{Ax} - \mathbf{b}\|_2^2, \quad (5)$$

where $\mathbf{b} \in \mathbb{R}_+^m$, $\mathbf{b} \neq \mathbf{0}$, and $\mathbf{A} \in \mathbb{R}_{++}^{m \times n}$.

- ▶ It's a least squares (LS) problem with nonnegative constraints, so the conventional LS solution may not be feasible for (5).
- ▶ A simple multiplicative updating algorithm:

$$\mathbf{x}_l^r = c_l^r \mathbf{x}_l^{r-1}, \quad l = 1, \dots, n, \quad (6)$$

where \mathbf{x}_l^r is the l th component of \mathbf{x}^r , and $c_l^r = \frac{[\mathbf{A}^T \mathbf{b}]_l}{[\mathbf{A}^T \mathbf{A} \mathbf{x}^{r-1}]_l}$.

- ▶ Starting with $\mathbf{x}^0 > \mathbf{0}$, then all \mathbf{x}^r generated by (6) are nonnegative.

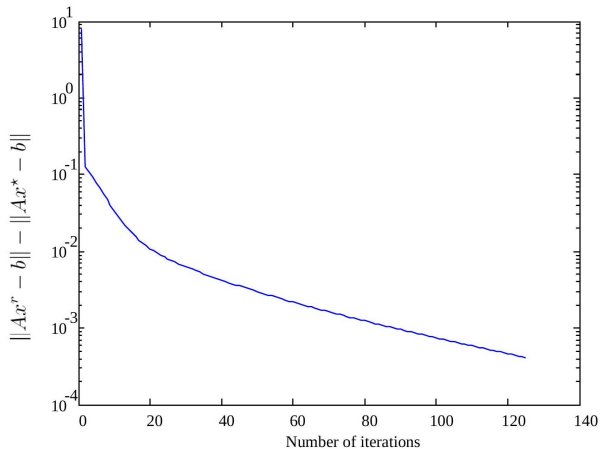


Figure: $\|\mathbf{A}\mathbf{x}^r - \mathbf{b}\|_2$ vs. the number of iterations.

- Usually the multiplicative update converges within a few tens of iterations.

- **MM interpretation:** Let $f(\mathbf{x}) \triangleq \|\mathbf{Ax} - \mathbf{b}\|_2^2$. The multiplicative update essentially solves the following problem

$$\min_{\mathbf{x} \geq \mathbf{0}} u(\mathbf{x}, \mathbf{x}^{r-1}),$$

where

$$u(\mathbf{x}, \mathbf{x}^{r-1}) \triangleq f(\mathbf{x}^{r-1}) + (\mathbf{x} - \mathbf{x}^{r-1})^\top \nabla f(\mathbf{x}^{r-1}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{r-1})^\top \Phi(\mathbf{x}^{r-1}) (\mathbf{x} - \mathbf{x}^{r-1}),$$

$$\Phi(\mathbf{x}^{r-1}) = \left(\frac{[\mathbf{A}^\top \mathbf{Ax}^{r-1}]_1}{x_1^{r-1}}, \dots, \frac{[\mathbf{A}^\top \mathbf{Ax}^{r-1}]_n}{x_n^{r-1}} \right).$$

- Observations:

$$\begin{cases} u(\mathbf{x}, \mathbf{x}^{r-1}) \text{ is quadratic approx. of } f(\mathbf{x}), \\ \Phi(\mathbf{x}^{r-1}) \succeq \mathbf{A}^\top \mathbf{A}, \end{cases} \quad \Rightarrow \quad \begin{cases} u(\mathbf{x}, \mathbf{x}^{r-1}) \geq f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^n, \\ u(\mathbf{x}^{r-1}, \mathbf{x}^{r-1}) = f(\mathbf{x}^{r-1}). \end{cases}$$

- The multiplicative update converges to an optimal solution of NLS (by the MM convergence in **Theorem 1** and convexity of NLS).

Applications — Convex-Concave Procedure / Difference-of-Convex (DC) Programming

- ▶ Suppose that $f(\mathbf{x})$ has the following form $f(\mathbf{x}) = g(\mathbf{x}) - h(\mathbf{x})$, where $g(\mathbf{x})$ and $h(\mathbf{x})$ are convex and differentiable. Thus, $f(\mathbf{x})$ is in general nonconvex.
- ▶ **DC Programming:** Construct $u(\cdot, \cdot)$ as

$$u(\mathbf{x}, \mathbf{x}^r) = g(\mathbf{x}) - \underbrace{\left(h(\mathbf{x}^r) + \nabla_{\mathbf{x}} h(\mathbf{x}^r)^\top (\mathbf{x} - \mathbf{x}^r) \right)}_{\text{linearization of } h \text{ at } \mathbf{x}^r}.$$

- ▶ By the 1st order condition of $h(\mathbf{x})$, it's easy to show that

$$u(\mathbf{x}, \mathbf{x}^r) \geq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{X}, \quad (\mathbf{x}^r, \mathbf{x}^r) = f(\mathbf{x}^r).$$

► **Sparse Signal Recovery by DC Programming**

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0, \quad \text{s.t. } \mathbf{y} = \mathbf{Ax}. \quad (7)$$

- Apart from the popular ℓ_1 approximation, consider the following concave approximation

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^n \log(1 + |x_i|/\epsilon) \quad \text{s.t. } \mathbf{y} = \mathbf{Ax},$$

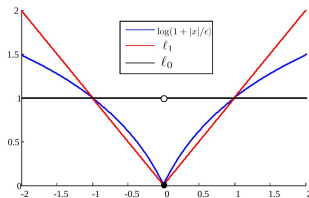


Figure: $\log(1 + |x|/\epsilon)$ promotes more sparsity than ℓ_1

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^n \log(1 + |x_i|/\epsilon) \quad \text{s.t. } \mathbf{y} = \mathbf{Ax},$$

which can be equivalently written as

$$\min_{\mathbf{x}, \mathbf{z} \in \mathbb{R}^n} \sum_{i=1}^n \log(z_i + \epsilon) \quad \text{s.t. } \mathbf{y} = \mathbf{Ax}, |x_i| \leq z_i, i = 1, \dots, n \quad (8)$$

- Problem (8) minimizes a concave objective, so it's a special case of DC programming ($g(\mathbf{x}) = 0$). Linearizing the concave function at $(\mathbf{x}^r, \mathbf{z}^r)$ yields

$$(\mathbf{x}^{r+1}, \mathbf{z}^{r+1}) = \arg \min \sum_{i=1}^n \frac{z_i}{z_i^r + \epsilon} \quad \text{s.t. } \mathbf{y} = \mathbf{Ax}, |x_i| \leq z_i, i = 1, \dots, n$$

- We solve a sequence of reweighted ℓ_1 problems.

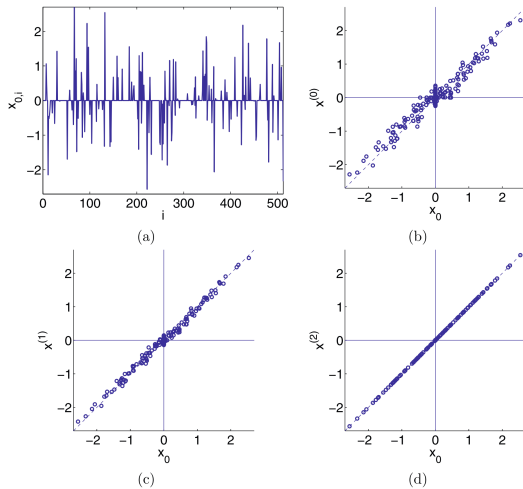


Figure: Sparse signal recovery through reweighted ℓ_1 iterations. **(a)** Original length $n = 512$ signal x_0 with 130 spikes. **(b)** Scatter plot, coefficient-by-coefficient, of x_0 versus its reconstruction $x^{(0)}$ using unweighted ℓ_1 minimization. **(c)** Reconstruction $x^{(1)}$ after the first reweighted iteration. **(d)** Reconstruction $x^{(2)}$ after the second reweighted iteration.

Applications — $\ell_2 - \ell_p$ Optimization

- ▶ Many problems involve solving the following problem (e.g., basis-pursuit denoising)

$$\min_{\mathbf{x}} f(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \mu \|\mathbf{x}\|_p, \quad (9)$$

where $p \geq 1$.

- ▶ If $\mathbf{A} = \mathbf{I}$ or \mathbf{A} is unitary, optimal \mathbf{x}^* is computed in closed-form as

$$\mathbf{x}^* = \mathbf{A}^\top \mathbf{y} - \text{Proj}_C(\mathbf{A}^\top \mathbf{y})$$

where $C \triangleq \{\mathbf{x} : \|\mathbf{x}\|_{p^*} \leq \mu\}$, $\|\cdot\|_{p^*}$ is the dual norm of $\|\cdot\|_p$ and Proj_C denotes the projection operator. In particular, for $p = 1$,

$$x_i^* = \text{soft}(y_i, \mu), \quad i = 1, \dots, n,$$

where $\text{soft}(u, a) \triangleq \text{sign}(u) \max\{|u| - a, 0\}$ denotes a *soft-thresholding* operation.

- ▶ For general \mathbf{A} , there is no simple closed-form solution for (9).

- **MM for $\ell_2 - \ell_p$ Problem:** Consider a modified $\ell_2 - \ell_p$ problem

$$\min_{\mathbf{x}} u(\mathbf{x}, \mathbf{x}^r) \triangleq f(\mathbf{x}) + \text{dist}(\mathbf{x}, \mathbf{x}^r), \quad (10)$$

where $\text{dist}(\mathbf{x}, \mathbf{x}^r) \triangleq \frac{c}{2} \|\mathbf{x} - \mathbf{x}^r\|_2^2 - \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}^r\|_2^2$ and $c > \lambda_{\max}(\mathbf{A}^\top \mathbf{A})$.

- $\text{dist}(\mathbf{x}, \mathbf{x}^r) \geq 0 \ \forall \mathbf{x} \implies u(\mathbf{x}, \mathbf{x}^r)$ majorizes $f(\mathbf{x})$.
- $u(\mathbf{x}, \mathbf{x}^r)$ can be reexpressed as

$$u(\mathbf{x}, \mathbf{x}^r) = \frac{c}{2} \|\mathbf{x} - \bar{\mathbf{x}}^r\|_2^2 + \mu \|\mathbf{x}\|_p + \text{const.},$$

where

$$\bar{\mathbf{x}}^r = \frac{1}{c} \mathbf{A}^\top (\mathbf{y} - \mathbf{A}\mathbf{x}^r) + \mathbf{x}^r.$$

- The modified $\ell_2 - \ell_p$ problem (10) has a simple soft-thresholding solution.
- Repeatedly solving problem (10) leads to an optimal solution of the $\ell_2 - \ell_p$ problem (by the MM convergence in **Theorem 1**).

Applications — Gradient Descent (GD)

► GD:

$$\begin{aligned}\mathbf{x}^{r+1} &= \mathbf{x}^r - \gamma \nabla f(\mathbf{x}^r) \\ &= \arg \min_{\mathbf{x}} \underbrace{\left\{ f(\mathbf{x}^r) + \nabla f(\mathbf{x}^r)^\top (\mathbf{x} - \mathbf{x}^r) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}^r\|^2 \right\}}_{L(\mathbf{x} | \mathbf{x}^r)} \quad [\text{verify it}]\end{aligned}$$

► Choose $\gamma \leq 1/L$: $L(\mathbf{x} | \mathbf{x}^r) \geq f(\mathbf{x})$

Proof of descent

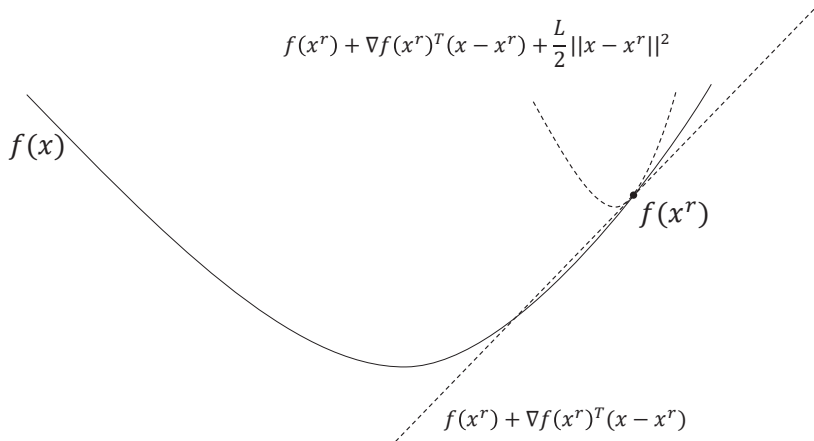
Majorize: by descent lemma and $\gamma \leq 1/L$

$$\begin{aligned} f(\mathbf{x}) &\leq f(\mathbf{x}^r) + \nabla f(\mathbf{x}^r)^\top (\mathbf{x} - \mathbf{x}^r) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}^r\|^2 \\ &\leq f(\mathbf{x}^r) + \nabla f(\mathbf{x}^r)^\top (\mathbf{x} - \mathbf{x}^r) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}^r\|^2 \end{aligned}$$

Minimize: let $\mathbf{x} = \mathbf{x}^r - \gamma \nabla f(\mathbf{x}^r)$

$$\begin{aligned} f(\mathbf{x}^{r+1}) &\leq f(\mathbf{x}^r) - \gamma \|\nabla f(\mathbf{x}^r)\|^2 + \frac{\gamma}{2} \|\nabla f(\mathbf{x}^r)\|^2 \\ &= f(\mathbf{x}^r) - \frac{\gamma}{2} \|\nabla f(\mathbf{x}^r)\|^2. \end{aligned}$$

Quadratic upperbound - a majorization minimization perspective



Applications — Expectation Maximization (EM)

- ▶ Consider an ML estimate of θ , given the random observation w

$$\hat{\theta}_{\text{ML}} = \arg \min_{\theta} -\ln p(w \mid \theta).$$

- ▶ Suppose that there are some missing data or hidden variables z in the model. Then, EM algorithm iteratively compute an ML estimate $\hat{\theta}$ as follows:

- E-step:

$$g(\theta, \theta^r) \triangleq \mathbb{E}_{z|w, \theta^r} \{\ln p(w, z \mid \theta)\}.$$

- M-step:

$$\theta^{r+1} = \arg \max_{\theta} g(\theta, \theta^r).$$


- repeat the above two steps until convergence.

- ▶ EM algorithm generates a nonincreasing sequence of $\{-\ln p(w \mid \theta^r)\}$.
- ▶ EM algorithm can be interpreted by MM.

► MM interpretation of EM algorithm:

$$\begin{aligned} & -\ln p(w \mid \theta) \\ &= -\ln \mathbb{E}_{z \mid \theta} p(w \mid z, \theta) \\ &= -\ln \mathbb{E}_{z \mid \theta} \left[\frac{p(z \mid w, \theta^r) p(w \mid z, \theta)}{p(z \mid w, \theta^r)} \right] \\ &= -\ln \mathbb{E}_{z \mid w, \theta^r} \left[\frac{p(z \mid \theta) p(w \mid z, \theta)}{p(z \mid w, \theta^r)} \right] \quad (\text{interchange the integrations}) \\ &\leq -\mathbb{E}_{z \mid w, \theta^r} \ln \left[\frac{p(z \mid \theta) p(w \mid z, \theta)}{p(z \mid w, \theta^r)} \right] \quad (\text{Jensen's inequality}) \\ &= -\mathbb{E}_{z \mid w, \theta^r} \ln p(w, z \mid \theta) + \mathbb{E}_{z \mid w, \theta^r} \ln p(z \mid w, \theta^r) \\ &\triangleq u(\theta, \theta^r) \end{aligned} \tag{11}$$

- $u(\theta, \theta^r)$ majorizes $-\ln p(w \mid \theta)$, and $-\ln p(w \mid \theta^r) = u(\theta^r, \theta^r)$;
- E-step essentially constructs $u(\theta, \theta^r)$;
- M-step minimizes $u(\theta, \theta^r)$ (note θ appears in the 1st term (11) only).

 Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Luo.
A unified convergence analysis of block successive minimization methods for nonsmooth optimization.
SIAM Journal on Optimization, 23(2):1126–1153, 2013.