

# Non-asymptotic local convergence analysis of alternating minimization

Wenzhi Gao



gwz@stanford.edu ICME, Stanford University

## Abstract

In this blog post, we consider a proof template that allows obtaining non-asymptotic local linear convergence rate of (two-block) alternating minimization (AM). Traditional analysis for alternating minimization often treats AM as a fixed point iteration and resorts to local Jacobian linearization arguments. This type of analysis gives sharp local rates, but it is often asymptotic. We provide an alternative argument that bypasses the Jacobian argument and directly shows convergence of the original function value gap. The resulting rates are non-asymptotic but equally tight.

## 1 Introduction and background

Consider the convex minimization problem

$$\min_{z=(u,v) \in \mathbb{R}^{m+n}} f(u,v),$$

where  $f: \mathbb{R}^{m+n} \rightarrow \mathbb{R}$  is smooth and convex. In particular, I assume that  $f$  is twice continuously differentiable. The optimization variable  $z \in \mathbb{R}^{m+n}$  is partitioned into two blocks

$$z = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{m+n},$$

and when partial minimization is accessible, alternating minimization is probably the simplest algorithm.

---

### Algorithm 1. Two-block alternating minimization

---

```

input  $v^1$ 
for  $k = 1, \dots, K$ 
     $u^{k+1} = \arg \min_u f(u, v^k)$ 
     $v^{k+1} = \arg \min_v f(u^{k+1}, v)$ 
end for
    
```

---

The global convergence analysis of AM has been widely studied in the literature [BT13]. Let  $f$  be  $L$ -smooth. A standard convergence analysis of AM is to reduce the AM update to gradient descent:

Suppose  $z = (u, v)$  satisfies  $\nabla_v f(u, v) = 0$ , then  $\nabla f(u, v) = (\nabla_u f(u, v), \nabla_v f(u, v)) = (\nabla_u f(u, v), 0)$  and

$$\begin{aligned}
 f(u^+, v^+) &\leq f(u^+, v) = \min_u f(u, v) \\
 \text{(By } u \text{ minimizes } f(\cdot, v)) &\leq f\left(u - \frac{1}{L} \nabla_u f(u, v), v\right) \\
 \text{(By } \nabla_v f(u, v) = 0) &= f\left(u - \frac{1}{L} \nabla_u f(u, v), v - \frac{1}{L} \nabla_v f(u, v)\right) \\
 \text{(By descent lemma)} &\leq f\left(z - \frac{1}{L} \nabla f(z)\right) \leq f(z) - \frac{1}{2L} \|\nabla f(z)\|^2,
 \end{aligned}$$

and we can reduce the analysis of AM to that of gradient descent. Alternatively, since AM is invariant with respect to block-diagonal affine transformation  $\begin{pmatrix} A_1 & \\ & A_1 \end{pmatrix}$ , it is possible to improve the dependence on  $L$ . In either case, it is safe to believe that given any target accuracy  $\varepsilon$ , the global convergence analysis ensures that AM can reach this accuracy in

$$T^G(\varepsilon) := \mathcal{O}\left(\text{poly}\left(\frac{1}{\varepsilon}\right)\right)$$

iterations. Given this guarantee, we move on to the local story. We should not expect the global worst-case complexity  $T_\varepsilon^G$  to fully characterize the behavior of AM when  $z \rightarrow z^*$ . Instead, we will resort to local analysis. The principle of local analysis is simple:

1. Analyze the algorithm on quadratics. Given a twice continuously differentiable function, it satisfies

$$f(z) = f(z^*) + \frac{1}{2}\|z - z^*\|_{\nabla^2 f(z^*)}^2 + \mathcal{O}(\|z - z^*\|^3),$$

and analyze the same algorithm on the local approximation

$$\lambda(z) = \lambda(u, v) = f(z^*) + \frac{1}{2}\|z - z^*\|_{\nabla^2 f(z^*)}^2$$

would provide quantitative behavior.

2. Globalize the analysis. Using a global analysis  $T_\varepsilon^G := \mathcal{O}\left(\text{poly}\left(\frac{1}{\varepsilon}\right)\right)$  to ensure that the  $\mathcal{O}(\|z - z^*\|^3)$  error term is dominated by the convergence behavior on  $\lambda(z)$ .

The first step is often clean, while the second step is often messy and involves a lot of algebraic manipulations to control the error. So I'll only cover the first step here and leave the technical second step to **Section 2**. For AM, the local quadratic model has more structures:

$$\begin{aligned} \lambda(u, v) &= f(u^*, v^*) + \frac{1}{2} \left\langle \begin{pmatrix} u - u^* \\ v - v^* \end{pmatrix}, \begin{pmatrix} \nabla_{uu}^2 f(u^*, v^*) & \nabla_{uv}^2 f(u^*, v^*) \\ \nabla_{uv}^2 f(u^*, v^*)^\top & \nabla_{vv}^2 f(u^*, v^*) \end{pmatrix} \begin{pmatrix} u - u^* \\ v - v^* \end{pmatrix} \right\rangle \\ &=: f(u^*, v^*) + \frac{1}{2} \left\langle \begin{pmatrix} u - u^* \\ v - v^* \end{pmatrix}, \begin{pmatrix} P & R \\ R^\top & Q \end{pmatrix} \begin{pmatrix} u - u^* \\ v - v^* \end{pmatrix} \right\rangle, \end{aligned}$$

where  $P := \nabla_{uu}^2 f(u^*, v^*)$ ,  $R := \nabla_{uv}^2 f(u^*, v^*)$ , and  $Q := \nabla_{vv}^2 f(u^*, v^*)$ . Then we can explicitly write the AM update on this quadratic function

$$u^+ = \arg \min_u \lambda(u, v) = u - (u - u^*) - P^{-1}R(v - v^*) = u^* - P^{-1}R(v - v^*)$$

and similarly,  $v^+ = \arg \min_v \lambda(u^+, v) = v^* - Q^{-1}R^\top(u - u^*)$ . Note that this update is exactly preconditioned gradient descent with preconditioner  $S^{-1} := \begin{pmatrix} P & \\ & Q \end{pmatrix}^{-1}$ , and we can immediately obtain

$$\begin{aligned} \|u^+ - u^*\| &= \|P^{-1}R(v - v^*)\| = \|P^{-1}RQ^{-1}R^\top(u - u^*)\| \\ \|v^+ - v^*\| &= \|Q^{-1}R^\top(u - u^*)\| = \|Q^{-1}R^\top P^{-1}R(v - v^*)\|. \end{aligned}$$

Hence asymptotically, the convergence of AM will be governed by the spectral radius:

$$\rho(P^{-1}RQ^{-1}R^\top) = \rho(Q^{-1}R^\top P^{-1}R) = \|P^{-1/2}RQ^{-1}R^\top P^{-1/2}\| =: \sigma_1.$$

**Remark 1.** Note that  $\sigma_1 \in (0, 1)$  since

$$\begin{pmatrix} P & R \\ R^\top & Q \end{pmatrix} \succ 0 \Rightarrow \begin{pmatrix} P^{-1/2} & \\ & Q^{-1/2} \end{pmatrix} \begin{pmatrix} P & R \\ R^\top & Q \end{pmatrix} \begin{pmatrix} P^{-1/2} & \\ & Q^{-1/2} \end{pmatrix} = \begin{pmatrix} I & P^{-1/2}RQ^{-1/2} \\ Q^{-1/2}R^\top P^{-1/2} & I \end{pmatrix} \succ 0$$

implies that the Schur complement  $I - P^{-1/2}RQ^{-1}R^\top P^{-1/2} \succ 0$ .

However, to give a tight non-asymptotic analysis, we need a more explicit control. Since  $\lambda$  is a quadratic, the change in optimality gap can be exactly characterized:

Let  $u, v$  satisfy  $\nabla \lambda_v(u, v) = 0$ . Then  $R^\top(u - u^*) + Q(v - v^*) = 0 \Rightarrow v - v^* = -Q^{-1}R^\top(u - u^*)$  and

$$\begin{aligned}\lambda(u, v) - f(u^*, v^*) &= \frac{1}{2}\|u - u^*\|_P^2 + \langle u - u^*, R(v - v^*) \rangle + \frac{1}{2}\|v - v^*\|_Q^2 \\ &= \frac{1}{2}\|u - u^*\|_{P-RQ^{-1}R^\top}^2.\end{aligned}$$

For  $u^+ = u^* - P^{-1}R(v - v^*)$ , we have

$$\lambda(u^+, v) - f(u^*, v^*) = \frac{1}{2}\|u - u^*\|_{RQ^{-1}R^\top - RQ^{-1}R^\top P^{-1}RQ^{-1}R^\top}^2.$$

Hence it suffices to consider

$$\frac{\lambda(u^+, v) - f(u^*, v^*)}{\lambda(u, v) - f(u^*, v^*)} \leq \frac{\|u - u^*\|_{RQ^{-1}R^\top - RQ^{-1}R^\top P^{-1}RQ^{-1}R^\top}^2}{\|u - u^*\|_{P-RQ^{-1}R^\top}^2} \leq \sigma_1,$$

and the convergence behavior of AM on  $\lambda(u, v)$  is expected.

After going through the messy derivation in **Section 2**, we arrive at the following convergence result.

**Theorem 1.** *Suppose  $f$  is  $L$ -smooth,  $\mu$ -strongly convex and with  $H$ -Lipschitz Hessian. Then it takes*

$$K_\varepsilon = \left\lceil T^G \left( \min \left\{ \sqrt{\frac{1-\sigma_1}{2}} \frac{\mu^3}{H^2}, \sqrt[3]{\frac{(1-\sigma_1)^2}{16} \frac{3\mu^3}{8H^2}}, \frac{\mu^3}{8H^2} \right\} \right) + \frac{2}{1-\sigma_1} \log\left(\frac{1}{\varepsilon}\right) \right\rceil.$$

iterations to find  $(u, v)$  such that  $f(u, v) - f(u^*, v^*) \leq \varepsilon$ . Recall that  $T^G$  is a global complexity bound.

**Remark 2.** Although the theorem is non-asymptotic, a caveat is that it is loose in at least two aspects

- First, as I previously mentioned, AM has certain affine invariant properties that cannot be characterized by standard smoothness assumptions.
- Second, the dependence on  $L$  and  $\mu$  are not ideal.  $L$  and  $\mu$  are global properties of  $f$ , while the iteration trajectory of AM will always satisfy either  $\nabla_u f = 0$  or  $\nabla_v f = 0$ . A tighter analysis should take this fact into account.

## 2 Local convergence analysis of alternating minimization

This section covers the promised AM convergence analysis. The result hinges on the following estimates.

**Lemma 1.** *Suppose  $f$  is  $L$ -smooth,  $\mu$ -strongly convex, and has  $H$ -Lipschitz Hessian, then*

- $|f(u, v) - \lambda(u, v)| \leq \frac{H}{3\mu} [f(u, v) - f(u^*, v^*)]^{3/2}$ ,
- $\|\nabla f(u, v) - \nabla \lambda(u, v)\| \leq \frac{H}{\mu} [f(u, v) - f(u^*, v^*)]$ ,
- $\|\nabla^2 f(u, v) - \nabla^2 \lambda(u, v)\| \leq \frac{\sqrt{2}H}{\sqrt{\mu}} [f(u, v) - f(u^*, v^*)]^{1/2}$ .

**Proof.** By the property of Lipschitz Hessian and strong convexity:  $f(u, v) - f(u^*, v^*) \geq \frac{\mu}{2} \|(u - u^*, v - v^*)\|^2$ ,

$$\begin{aligned}|f(u, v) - \lambda(u, v)| &\leq \frac{H}{6} \|(u - u^*, v - v^*)\|^3 \leq \frac{\sqrt{2}H}{3\mu^{3/2}} [f(u, v) - f(u^*, v^*)]^{3/2} \\ \|\nabla f(u, v) - \nabla \lambda(u, v)\| &\leq \frac{H}{2} \|(u - u^*, v - v^*)\|^2 \leq \frac{H}{\mu} [f(u, v) - f(u^*, v^*)] \\ \|\nabla^2 f(u, v) - \nabla^2 \lambda(u, v)\| &\leq H \|(u - u^*, v - v^*)\| \leq \frac{\sqrt{2}H}{\sqrt{\mu}} [f(u, v) - f(u^*, v^*)]^{1/2},\end{aligned}$$

and this completes the proof.  $\square$

The following lemma characterizes the objective value after a half alternating minimization step.

**Lemma 2.** Suppose  $f(u, v) - f(u^*, v^*) \leq \frac{\mu^3}{8H^4}$ , then

$$f(u^+, v) - f(u^*, v^*) \leq \sigma_1 [f(u, v) - f(u^*, v^*)] + \frac{2\sqrt{2}H}{3\mu^{3/2}} [f(u, v) - f(u^*, v^*)]^{3/2} + \frac{H^2}{2\mu^3} [f(u, v) - f(u^*, v^*)]^2.$$

A direct application of **Lemma 2** will complete the proof.

**Theorem 2. (Theorem 1)** Under the previous assumptions, it takes

$$K_\varepsilon := \left\lceil T^G \left( \min \left\{ \sqrt{\frac{1-\sigma_1}{2} \frac{\mu^3}{H^2}}, \sqrt[3]{\frac{(1-\sigma_1)^2}{16} \frac{3\mu^3}{8H^2}}, \frac{\mu^3}{8H^2} \right\} \right) + \frac{1}{1-\sigma_1} \log \left( \frac{1}{\varepsilon} \right) \right\rceil$$

AM iterations to find an  $\varepsilon$ -optimal solution.

**Proof.** With **Lemma 2**, suppose  $f(u, v) - f(u^*, v^*) \leq \min \left\{ \sqrt{\frac{1-\sigma_1}{2} \frac{\mu^3}{H^2}}, \sqrt[3]{\frac{(1-\sigma_1)^2}{16} \frac{3\mu^3}{8H^2}} \right\}$ , we have

$$\frac{2\sqrt{2}H}{3\mu^{3/2}} [f(u, v) - f(u^*, v^*)]^{3/2} + \frac{H^2}{2\mu^3} [f(u, v) - f(u^*, v^*)]^2 \leq \frac{1-\sigma_1}{2} [f(u, v) - f(u^*, v^*)]$$

and that

$$f(u^+, v) - f(u^*, v^*) \leq \frac{1+\sigma_1}{2} [f(u, v) - f(u^*, v^*)].$$

Putting things together completes the proof.  $\square$

## References

[BT13] Amir Beck and Luba Tetrushvili. On the convergence of block coordinate descent type methods. *SIAM journal on Optimization*, 23(4):2037–2060, 2013.

## A Proof of Lemma 2

**Proof.** For brevity, we denote  $\varepsilon := f(u, v) - f(u^*, v^*)$ .

Define  $\alpha_v(u) := f(u, v)$  and  $\ell_v(u) := \lambda(u, v)$  for fixed  $v$ . Then

$$\begin{aligned} & f(u^+, v) - f(u^*, v^*) \\ &= f(u^+, v) - f(u^*, v^*) \\ &= \min_u \alpha_v(u) - f(u^*, v^*) \\ &\leq \alpha_v(u - P^{-1} \nabla \alpha(u)) - f(u^*, v^*) \\ &= \alpha_v(u - P^{-1} \nabla \alpha(u)) - \ell_v(u - P^{-1} \nabla \alpha(u)) + \ell_v(u - P^{-1} \nabla \alpha(u)) - f(u^*, v^*) \\ &= \underbrace{\alpha_v(u - P^{-1} \nabla \alpha(u)) - \ell_v(u - P^{-1} \nabla \alpha(u))}_{\Sigma_1} + \underbrace{\ell_v(u - P^{-1} \nabla \alpha(u)) - \ell_v(u - P^{-1} \nabla \ell(u))}_{\Sigma_2} \\ &\quad + \underbrace{\ell_v(u - P^{-1} \nabla \ell(u)) - f(u^*, v^*)}_{\Sigma_3}. \end{aligned}$$

and we bound  $\Sigma_1, \Sigma_2, \Sigma_3$  respectively.

For  $\Sigma_1$ , if  $\alpha(u - P^{-1}\nabla\alpha(u)) \leq \alpha(u)$ , then  $f(u - P^{-1}\nabla\alpha(u), v) \leq f(u, v)$ , and

$$\Sigma_1 \leq \frac{\sqrt{2}H}{3\mu^{3/2}}\varepsilon^{3/2}.$$

To show  $\alpha(u - P^{-1}\nabla\alpha(u)) \leq \alpha(u)$ , we first note that  $\nabla^2\alpha(u) = \nabla_{uv}^2 f(u, v)$  and that with

$$f(u, v) - f(u^*, v^*) \leq \frac{\mu^3}{8H^2} \Rightarrow \frac{\sqrt{2}H}{\sqrt{\mu}} \sqrt{f(u, v) - f(u^*, v^*)} \leq \frac{1}{2}\mu$$

we have

$$\begin{aligned} \nabla^2\alpha_v(u) &= \nabla^2\ell_v(u) + \nabla^2\alpha_v(u) - \nabla^2\ell_v(u) \\ &\leq P + \|\nabla^2\alpha_v(u) - \nabla^2\ell_v(u)\| \\ &\leq P + \|\nabla^2 f(u, v) - \nabla^2\lambda(u, v)\| \cdot I \\ &\leq P + \frac{\sqrt{2}H}{\sqrt{\mu}} [f(u, v) - f(u^*, v^*)]^{1/2} \cdot I \\ &\leq P + \frac{\mu}{2} I \leq \frac{3}{2}P. \end{aligned}$$

Consider scalar function  $\gamma(\theta) := \alpha_v(u - \theta P^{-1}\nabla\alpha(u))$  and let  $\theta_{\max} = \sup_{\theta > 0} \{\theta > 0: \gamma(\theta) \leq \alpha_v(u)\}$ .  $\theta_{\max}$  is well-defined since  $P^{-1} > 0$  and  $P^{-1}\nabla\alpha_v(u)$  is a descent direction of  $\alpha_v$  at  $u$ .

By convexity, we have  $\alpha_v(u - \theta P^{-1}\nabla\alpha(u)) \leq \alpha_v(u)$  for all  $\theta \leq \theta_{\max}$ . Then we deduce that

$$\begin{aligned} \alpha_v(u) &= \gamma(\theta_{\max}) \\ &= \alpha_v(u) - \theta_{\max} \langle \nabla\alpha_v(u), P^{-1}\nabla\alpha_v(u) \rangle \\ &\quad + \theta_{\max}^2 \int_0^1 \langle \nabla\alpha_v(u), P^{-1}\nabla^2\alpha_v(u - t\theta_{\max}P^{-1}\nabla\alpha_v(u)) P^{-1}\nabla\alpha_v(u) \rangle (1-t) dt \\ &\leq \alpha_v(u) - \theta_{\max} \langle \nabla\alpha_v(u), P^{-1}\nabla\alpha_v(u) \rangle + \theta_{\max}^2 \int_0^1 \langle \nabla\alpha_v(u), P^{-1}(\frac{3}{2}P) P^{-1}\nabla\alpha_v(u) \rangle (1-t) dt, \\ &= \alpha_v(u) - \theta_{\max} \|\nabla\alpha_v(u)\|_{P^{-1}}^2 + \frac{3}{4}\theta_{\max}^2 \|\nabla\alpha_v(u)\|_{P^{-1}}^2, \end{aligned}$$

where the inequality holds since  $\alpha_v(u - t\theta_{\max}P^{-1}\nabla\alpha_v(u)) \leq \alpha_v(u)$ . Given

$$\frac{3}{4}\theta_{\max}^2 \|\nabla\alpha_v(u)\|_{P^{-1}}^2 \geq \theta_{\max} \|\nabla\alpha_v(u)\|_{P^{-1}}^2,$$

we have  $\theta_{\max} \geq \frac{4}{3}$ . Therefore,  $\alpha_v(u - P^{-1}\nabla\alpha(u)) = \gamma(1) \leq \gamma(\theta_{\max}) = \alpha_v(u)$ .

For  $\Sigma_2$ , since  $\ell_v(u)$  is a quadratic function, we have

$$\begin{aligned} \Sigma_2 &= \ell_v(u - P^{-1}\nabla\alpha(u)) - \ell_v(u - P^{-1}\nabla\ell(u)) \\ &= \frac{1}{2} \|\nabla\alpha_v(u) - \nabla\ell_v(u)\|_{P^{-1}}^2 \\ &\leq \frac{H^2}{2\mu^3} [f(u, v) - f(u^*, v^*)]^2 = \frac{H^2}{2\mu^3} \varepsilon^2 \end{aligned}$$

For  $\Sigma_3$ , we have

$$\begin{aligned} \ell_v(u - P^{-1}\nabla\ell(u)) - f(u^*, v^*) &\leq \sigma_1 [\ell_v(u) - f(u^*, v^*)] \\ &= \sigma_1 [\ell_v(u) - \alpha_v(u) + \alpha_v(u) - f(u^*, v^*)] \\ &\leq \sigma_1 [\alpha_v(u) - f(u^*, v^*)] + |\ell_v(u) - \alpha_v(u)| \\ &\leq \sigma_1 [\alpha_v(u) - f(u^*, v^*)] + \frac{H}{3\mu} \varepsilon^{3/2} \\ &= \sigma_1 \varepsilon + \frac{\sqrt{2}H}{3\mu^{3/2}} \varepsilon^{3/2} \end{aligned}$$

Putting things together completes the proof.  $\square$