

Orthogonally Constrained Spherical Matrix Factorization

Anonymous*

1 Convergence Analysis

Theorem 1 (Subsequence convergence). *Let $\{W_k\}_{k \geq 0} = \{(U_k, V_k)\}_{k \geq 0}$ be the sequence generated by Algorithm 1 with constant step size $\lambda, \mu > L_c$. Then the sequence $\{W_k\}_{k \geq 0}$ is bounded and obeys the following properties:*

(P1) *sufficient decrease:*

$$f(W_{k-1}) - f(W_k) \geq \frac{\min(\lambda, \mu) - L_c}{2} \|W_k - W_{k-1}\|_F^2 \quad (1)$$

implying

$$\lim_{k \rightarrow \infty} \|W^{k-1} - W^k\|_F = 0. \quad (2)$$

(P2) *the sequence $\{f(W_k)\}$ converges to some $\bar{f} \geq 0$.*

(P3) *denote $\mathbb{C}(W_0)$ (depending on W_0) as the set of all limit points of the iterates $\{W_k\}$. Then all the limit points W^* are critical points of f and have the same function value*

$$f(W^*) = \bar{f}. \quad (3)$$

Further, $\mathbb{C}(W_0)$ is a nonempty, compact and connected set and satisfies

$$\lim_{k \rightarrow \infty} \text{dist}(W_k, \mathbb{C}(W_0)) = 0 \quad (4)$$

Proof of Theorem 1. Before proving Theorem 1, we give out some necessary definition.

Definition 1. [Attouch et al., 2013] *Let $f : \mathbb{R}^d \rightarrow (-\infty, \infty]$ be a proper and lower semi-continuous function, whose domain is defined as*

$$\text{dom } f := \{u \in \mathbb{R}^n : f(u) < \infty\}.$$

The (Fréchet) subdifferential ∂f of f at u is defined by

$$\partial f(u) = \left\{ z : \liminf_{v \rightarrow u} \frac{f(v) - f(u) - \langle z, v - u \rangle}{\|u - v\|} \geq 0 \right\}$$

for any $u \in \text{dom } h$ and $\partial f(u) = \emptyset$ if $u \notin \text{dom } f$.

We say u is a limiting critical point, or simply a critical point of f if $0 \in \partial f(u)$.

We now turn to prove Theorem 1.

(P1): First note that for all k , according to our alternating minimization method, we always have $\delta_U(U_k) = \delta_V(V_k) = 0$ and thus $f(W_k) = h(W_k)$.

Since $h(U, V)$ has Lipschitz continuous gradient at $U \in \mathbb{U}, V \in \mathbb{V}$ with Lipschitz gradient L_c and $\lambda > L_c$, we define $h_{L_c}(U, U', V)$ as proximal regularization of $h(U, V)$ linearized at U', V :

$$h(U', V) + \langle \nabla_U h(U', V), U - U' \rangle + \frac{L_c}{2} \|U - U'\|_F^2,$$

By the definition of Lipschitz continuous gradient and Taylor expansion, we have

$$h(U, V) \leq h_{L_c}(U, U', V). \quad (5)$$

Also by the definition of proximal map, we get:

$$U_k = \arg \min_U \delta_U(U) + \frac{\mu}{2} \|U - U_{k-1}\|_F^2 + \langle \nabla_U h(U_{k-1}, V_{k-1}), U - U_{k-1} \rangle \quad (6)$$

and hence we take $U_k = U$, which implies that

$$\delta_U(U_k) + \frac{\mu}{2} \|U_k - U_{k-1}\|_F^2 + \langle \nabla_U h(U_{k-1}, V_{k-1}), U_k - U_{k-1} \rangle \leq \delta_U(U_{k-1}) \quad (7)$$

Combining Eq. (5) to Eq. (7), we have:

$$\begin{aligned} & h(U_k, V_{k-1}) + \delta_U(U_k) \\ & \leq h(U_{k-1}, V_{k-1}) + \langle \nabla_U h(U_{k-1}, V_{k-1}), U_k - U_{k-1} \rangle \\ & \quad + \frac{L_c}{2} \|U_k - U_{k-1}\|_F^2 + \delta_U(U_k) \\ & \leq h(U_{k-1}, V_{k-1}) + \frac{L_c}{2} \|U_k - U_{k-1}\|_F^2 \\ & \quad + \delta_U(U_{k-1}) - \frac{\mu}{2} \|U_k - U_{k-1}\|_F^2 \\ & = h(U_{k-1}, V_{k-1}) + \delta_U(U_{k-1}) - \frac{\mu - L_c}{2} \|U_k - U_{k-1}\|_F^2, \end{aligned} \quad (8)$$

Similarly, we have

$$\begin{aligned} & h(U_k, V_k) - h(U_k, V_{k-1}) + \delta_V(V_k) - \delta_V(V_{k-1}) \\ & \leq -\frac{\lambda - L_c}{2} \|V_k - V_{k-1}\|_F^2 \end{aligned} \quad (9)$$

*This is the supplementary of Paper ID: 5740.

which together with the above equation gives Eq. (1). Now repeating Eq. (1) for all k will give

$$(\min(\lambda, \mu) - L_c) \sum_{k=1}^{\infty} \|W_k - W_{k-1}\|_F^2 \leq 2f(W_0), \quad (10)$$

which gives Eq. (2).

Remark 1. In our proposed algorithm, since in every update, our solution is closed while satisfying the constraints, thus in fact δ_U and δ_V are 0, and ∞ is never achieved.

(P2) It follows from Eq. (16) that $\{f(W_k)\}_{k \geq 0}$ is a decreasing sequence. Due to the fact that f is lower bounded as $f(W_k) \geq 0$ for all k , we conclude that $\{f(W_k)\}_{k \geq 0}$ is convergent to some constant $\bar{f} \geq 0$.

(P3) Extract any convergent subsequence $\{W_{k'} = (U_{k'}, V_{k'})\}$ from $\{W_k\}$ and denote the limit point of this subsequence as W^* . Since $U_{k'} \in \mathbb{U}, V_{k'} \in \mathbb{V}$ for all k' and both of the sets \mathbb{U} and \mathbb{V} are closed, we have $U^* \in \mathbb{U}, V^* \in \mathbb{V}$. Since h is continuous, we have

$$\begin{aligned} \lim_{k' \rightarrow \infty} f(W_{k'}) &= \lim_{k' \rightarrow \infty} h(U_{k'}, V_{k'}) + \delta_U(U_{k'}) + \delta_V(V_{k'}) \\ &= f(W^*), \end{aligned} \quad (11)$$

which together with the fact that $\bar{f} = \lim_{k \rightarrow \infty} f(W_k)$ gives Eq. (3).

To show W^* is a critical point, we first consider Eq. (6) and the optimality condition yields:

$$\nabla_U h(U_{k-1}, V_{k-1}) + \mu(U_k - U_{k-1}) + \partial \delta_U(U_k) = 0. \quad (12)$$

Similarly, we have

$$\nabla_V h(U_k, V_{k-1}) + \lambda(V_k - V_{k-1}) + \partial \delta_V(V_k) = 0. \quad (13)$$

Now, define

$$\underbrace{\nabla_U h(U_k, V_k) + \partial \delta_U(U_k)}_{A_k} \quad \text{and} \quad \underbrace{\nabla_V h(U_k, V_k) + \partial \delta_V(V_k)}_{B_k}.$$

Thus, we have

$$A_k \in \partial_U f(U_k, V_k), B_k \in \partial_V f(U_k, V_k). \quad (14)$$

It follows from the above that

$$\begin{aligned} &\lim_{k \rightarrow \infty} \|A_k\|_F \\ &\leq \lim_{k \rightarrow \infty} \|\nabla_U h(U_k, V_k) - \nabla_U h(U_{k-1}, V_{k-1})\|_F + \mu \|U_k - U_{k-1}\|_F \\ &\leq \lim_{k \rightarrow \infty} (L_c + \mu) \|W_k - W_{k-1}\|_F = 0. \end{aligned} \quad (15)$$

Similarly, we have

$$\lim_{k \rightarrow \infty} \|B_k\|_F \leq \lim_{k \rightarrow \infty} (L_c + \lambda) \|W_k - W_{k-1}\|_F = 0. \quad (16)$$

Then we have:

$$\text{dist}(0, \partial f(W_k)) \leq \mathcal{L}_g \|W_k - W_{k-1}\|_F \quad (17)$$

where $\mathcal{L}_g := (2L_c + \mu + \lambda)$. Owing to the closedness properties of $\partial f(W_{k'})$, we finally obtain $0 \in \partial f(W^*)$. Thus, W^* is a critical point of f .

The remaining proof regarding to the properties of $\mathbb{C}(W_0)$ follows from [Bolte et al., 2014] by using the regularity of $\{W_k\}$ (assertion (P1) of Theorem 1). \square

Theorem 2 (Sequence convergence). *The sequence $\{W_k\}_{k \geq 0}$ generated by Algorithm 1 with a constant step size $\lambda, \mu > L_c$ is global-sequence convergence.*

Remark 2. Theorem 2 is much stronger than Theorem 1, since we are not guaranteed that the iterates $\{W_k\}$ generated by Algorithm 1 would converge to a limit point only by Theorem 1. Theorem 2 fulfills this gap by directly showing the iterates $\{W_k\}$ generated by Algorithm 1 converges to a critical point, and as an consequence, the set of limit point $\mathbb{C}(W_0)$ becomes a singleton.

Before proving Theorem 2, we give out another important definition.

Definition 2 (Kurdyka-Lojasiewicz (KL) property). [Bolte et al., 2007] *We say a proper semi-continuous function $h(\mathbf{u})$ satisfies Kurdyka-Lojasiewicz (KL) property, if $\bar{\mathbf{u}}$ is a critical point of $h(\mathbf{u})$, then there exist $\delta > 0$, $\theta \in [0, 1)$, $C_1 > 0$, s.t.*

$$|h(\mathbf{u}) - h(\bar{\mathbf{u}})|^\theta \leq C_1 \text{dist}(0, \partial h(\mathbf{u})), \quad \forall \mathbf{u} \in B(\bar{\mathbf{u}}, \delta)$$

We mention that the above KL property (also known as KL inequality) states the regularity of $h(\mathbf{u})$ around its critical point \mathbf{u} and the KL inequality trivially holds at non-critical point. There are a very large set of functions satisfying the KL inequality including any semi-algebraic functions [Attouch et al., 2013; Bolte et al., 2014]. Clearly, the objective function f is semi-algebraic as both h , δ_U and δ_V are semi-algebraic [Bolte et al., 2014].

Lemma 1 (Uniform KL property). *There exist $\delta_0 > 0$, $\theta_{KL} \in [0, 1)$, $C_{KL} > 0$ such that for all W s.t. $\text{dist}((W), \mathbb{C}(W_0)) \leq \delta_0$:*

$$|f(W) - \bar{f}|^{\theta_{KL}} \leq C_{KL} \text{dist}(0, \partial f(W)) \quad (18)$$

with \bar{f} denoting the limiting function value defined in P (2) of Theorem 1.

Proof. First we recognize the union $\bigcup_i B(W_i^*, \delta_i)$ forms an open cover of $\mathbb{C}(W_0)$ with W_i^* representing all points in $\mathbb{C}(W_0)$ and δ_i to be chosen so that the the following KL property of f at $W_i^* \in \mathbb{C}(W_0)$ holds:

$$|f(W) - \bar{f}|^{\theta_i} \leq C_i \text{dist}(0, \partial f(W)) \quad \forall (W) \in B(W_i^*, \delta_i)$$

where we have used all $f(W_i^*) = \bar{f}$ by assertion (P3) of Theorem 1. Then due to the compactness of the set $\mathbb{C}(W_0)$, it has a finite subcover $\bigcup_{i=1}^p B(W_{k_i}^*, \delta_{k_i})$ for some positive integer p . Now combining all, we have for all $W \in \bigcup_{i=1}^p B(W_{k_i}^*, \delta_{k_i})$,

$$|f(W) - \bar{f}|^{\theta_{KL}} \leq C_{KL} \text{dist}(0, \partial f(W)) \quad (19)$$

with $\theta_{KL} = \max_{i=1}^p \{\theta_{k_i}\}$ and $C_{KL} = \max_{i=1}^p \{C_{k_i}\}$. Finally, since $\bigcup_{i=1}^p B(W_{k_i}^*, \delta_{k_i})$ is an open cover of $\mathbb{C}(W_0)$, there exists a sufficiently small number δ_0 so that

$$\{(W) : \text{dist}(W, \mathbb{C}(W_0)) \leq \delta_0\} \subset \bigcup_{i=1}^p B(W_{k_i}^*, \delta_{k_i}).$$

Therefore, Eq. (19) holds whenever $\text{dist}(W, \mathbb{C}(W_0)) \leq \delta_0$. \square
Now we are ready to prove Theorem 2.

Proof of Theorem 2. First of all, following from assertion (P3) in Theorem 1, there exists a positive integer k_0 so that $\text{dist}(W_k, \mathbb{C}(W_0)) \leq \delta_0$ for all $k \geq k_0$. Now using Lemma 1, we have

$$[f(W_k) - f(W^*)]^{\theta_{KL}} \leq C_{KL} \text{dist}(0, \partial f(W_k)), \quad \forall k \geq k_0. \quad (20)$$

In the subsequent analysis, we restrict to $k \geq k_0$. Construct a concave function $x^{1-\theta}$ for some $\theta \in [0, 1]$ with domain $x > 0$. Obviously, by the concavity, we have

$$x_2^{1-\theta} - x_1^{1-\theta} \geq (1-\theta)x_2^{-\theta}(x_2 - x_1), \quad \forall x_1 > 0, x_2 > 0$$

Replacing x_1 by $f(W_{k+1}) - f(W^*)$ and x_2 by $f(W_k) - f(W^*)$, choosing $\theta = \theta_{KL}$ and using the sufficient decrease property, we have

$$\begin{aligned} & [f(W_k) - f(W^*)]^{1-\theta_{KL}} - [f(W_{k+1}) - f(W^*)]^{1-\theta_{KL}} \\ & \geq (1-\theta_{KL}) \frac{f(W_k) - f(W_{k+1})}{[f(W_k) - f(W^*)]^{\theta_{KL}}} \\ & \geq \frac{\lambda(1-\theta_{KL})}{2C_{KL}} \frac{\|W_k - W_{k+1}\|_F^2}{\text{dist}(0, \partial f(W_k))} \\ & \geq \frac{\lambda(1-\theta_{KL})}{2C_{KL}\mathcal{L}_g} \frac{\|W_k - W_{k+1}\|_F^2}{\|W_k - W_{k-1}\|_F} \\ & = \kappa \left(\frac{\|W_k - W_{k+1}\|_F^2}{\|W_k - W_{k-1}\|_F} + \|W_k - W_{k-1}\|_F \right) - \kappa \|W_k - W_{k-1}\|_F \\ & \geq \kappa (2\|W_k - W_{k+1}\|_F - \|W_k - W_{k-1}\|_F) \end{aligned}$$

where we have used Eq. (20) in the third line and Eq. (17) in the fourth line. And accordingly, we have:

$$\begin{aligned} & 2\|W_k - W_{k+1}\|_F - \|W_k - W_{k-1}\|_F \\ & \leq \beta ([f(W_k) - f(W^*)]^{1-\theta} - [f(W_{k+1}) - f(W^*)]^{1-\theta_{KL}}) \end{aligned} \quad (21)$$

with $\kappa := \frac{\lambda(1-\theta_{KL})}{2C_{KL}\mathcal{L}_g}$ and $\beta := \left(\frac{\lambda(1-\theta_{KL})}{2C_{KL}\mathcal{L}_g} \right)^{-1}$.

Summing the above inequalities up from some $\tilde{k} > k_0$ to infinity yields

$$\begin{aligned} & \sum_{k=\tilde{k}}^{\infty} \|W_k - W_{k+1}\|_F \\ & \leq \|W_{\tilde{k}} - W_{\tilde{k}-1}\|_F + \beta [f(W_{\tilde{k}}) - f(W^*)]^{1-\theta_{KL}} \end{aligned} \quad (22)$$

implying $\sum_{k=\tilde{k}}^{\infty} \|W_k - W_{k+1}\|_F < \infty$. Following some standard arguments one can see that

$$\limsup_{t \rightarrow \infty, t_1, t_2 \geq t} \|W_{t_1} - W_{t_2}\|_F = 0$$

which implies that the sequence $\{W_k\}$ is Cauchy, and hence convergent. Hence, the limit point set $\mathcal{C}(W_0)$ is singleton W^* . \square

Theorem 3 (Convergence Rate). *The convergence rate is at least sub-linear.*

Towards that end, we first know from the above argument that $\{W_k\}$ converges to some point W^* , i.e.,

$\lim_{k \rightarrow \infty} W^k = W^*$. Then using Eq. (22) and the triangle inequality, we obtain

$$\begin{aligned} \|W_{\tilde{k}} - W^*\|_F & \leq \sum_{k=\tilde{k}}^{\infty} \|W_k - W_{k+1}\|_F \\ & \leq \|W_{\tilde{k}} - W_{\tilde{k}-1}\|_F + \beta [f(W_{\tilde{k}}) - f(W^*)]^{1-\theta_{KL}} \end{aligned} \quad (23)$$

which indicates the convergence rate of $W_{\tilde{k}} \rightarrow W^*$ is at least as fast as the rate that $\|W_{\tilde{k}} - W_{\tilde{k}-1}\|_F + \beta [f(W_{\tilde{k}}) - f(W^*)]^{1-\theta_{KL}}$ converges to 0. In particular, the second term $\beta [f(W_{\tilde{k}}) - f(W^*)]^{1-\theta_{KL}}$ can be controlled:

$$\begin{aligned} \beta [f(W_{\tilde{k}}) - f(W^*)]^{\theta_{KL}} & \leq \beta C_{KL} \text{dist}(0, \partial f(W_{\tilde{k}})) \\ & \leq \underbrace{\beta C_{KL}(2B_0 + \lambda + \|\mathbf{X}\|_F)}_{:=\alpha} \|W_{\tilde{k}} - W_{\tilde{k}-1}\|_F \end{aligned} \quad (24)$$

Plugging Eq. (24) back to Eq. (23), we then have

$$\sum_{k=\tilde{k}}^{\infty} \|W_k - W_{k+1}\|_F \leq \|W_{\tilde{k}} - W_{\tilde{k}-1}\|_F + \alpha \|W_{\tilde{k}} - W_{\tilde{k}-1}\|_F^{\frac{1-\theta_{KL}}{\theta_{KL}}}.$$

We divide the following analysis into two cases based on the value of the KL exponent θ_{KL} .

Case I : $\theta_{KL} \in [0, \frac{1}{2}]$. This case means $\frac{1-\theta_{KL}}{\theta_{KL}} \geq 1$. We define $P_{\tilde{k}} = \sum_{i=\tilde{k}}^{\infty} \|W_{i+1} - W_i\|_F$,

$$P_{\tilde{k}} \leq P_{\tilde{k}-1} - P_{\tilde{k}} + \alpha [P_{\tilde{k}-1} - P_{\tilde{k}}]^{\frac{1-\theta_{KL}}{\theta_{KL}}}. \quad (25)$$

Since $P_{\tilde{k}-1} - P_{\tilde{k}} \rightarrow 0$, there exists a positive integer k_1 such that $P_{\tilde{k}-1} - P_{\tilde{k}} < 1$, $\forall \tilde{k} \geq k_1$. Thus,

$$P_{\tilde{k}} \leq (1 + \alpha) (P_{\tilde{k}-1} - P_{\tilde{k}}), \quad \forall \tilde{k} \geq \max\{k_0, k_1\},$$

which implies that

$$P_{\tilde{k}} \leq \rho \cdot P_{\tilde{k}-1}, \quad \forall \tilde{k} \geq \max\{k_0, k_1\}, \quad (26)$$

where $\rho = \frac{1+\alpha}{2+\alpha} \in (0, 1)$. This together with Eq. (23) gives the linear convergence rate

$$\|W_k - W^*\|_F \leq \mathcal{O}(\rho^{k-\bar{k}}), \quad \forall k \geq \bar{k}. \quad (27)$$

where $\bar{k} = \max\{k_0, k_1\}$.

Case II : $\theta_{KL} \in (1/2, 1)$. This case means $\frac{1-\theta_{KL}}{\theta_{KL}} \leq 1$. Based on the former results, we have

$$P_{\tilde{k}} \leq (1 + \alpha) [P_{\tilde{k}-1} - P_{\tilde{k}}]^{\frac{1-\theta_{KL}}{\theta_{KL}}}, \quad \forall \tilde{k} \geq \max\{k_0, k_1\}$$

which gives

$$P_{\tilde{k}}^{\frac{1-2\theta_{KL}}{1-\theta_{KL}}} - P_{\tilde{k}-1}^{\frac{1-2\theta_{KL}}{1-\theta_{KL}}} \geq \zeta, \quad \forall k \geq \bar{k}$$

for some $\zeta > 0$. Then repeating and summing up the above inequality from $\bar{k} = \max\{k_0, k_1\}$ to any $k > \bar{k}$, we can conclude

$$P_{\tilde{k}} \leq \left[P_{\tilde{k}-1}^{\frac{1-2\theta_{KL}}{1-\theta_{KL}}} + \zeta(\tilde{k} - \bar{k}) \right]^{-\frac{1-\theta_{KL}}{2\theta_{KL}-1}} = \mathcal{O}((\tilde{k} - \bar{k})^{-\frac{1-\theta_{KL}}{2\theta_{KL}-1}})$$

Finally, the following sublinear convergence holds

$$\|W_k - W^*\|_F \leq \mathcal{O}((k - \bar{k})^{-\frac{1-\theta_{KL}}{2\theta_{KL}-1}}), \quad \forall k > \bar{k}.$$

We end this proof by commenting that both linear and sublinear convergence rate are closely related to the KL exponent θ_{KL} at the critical point W^* .

References

- [Attouch and Bolte, 2009] Hedy Attouch and Jérôme Bolte. On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming*, 116(1-2):5–16, 2009.
- [Attouch *et al.*, 2013] Hedy Attouch, Jérôme Bolte, and Benar Fux Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized gauss–seidel methods. *Mathematical Programming*, 137(1-2):91–129, 2013.
- [Bolte *et al.*, 2007] Jérôme Bolte, Aris Daniilidis, and Adrian Lewis. The łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization*, 17(4):1205–1223, 2007.
- [Bolte *et al.*, 2014] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2):459–494, 2014.