

Supplemental Material for “On the Convergence of a Bayesian Algorithm for Joint Dictionary Learning and Sparse Recovery”

Geethu Joseph and Chandra R. Murthy *Senior Member, IEEE*

I. DERIVATION OF DL-SBL ALGORITHM

In this section, we provide the details of the EM-algorithm development, explaining how to obtain (3)-(8), and the γ_k update equations in Algorithm 1 and Algorithm 2. The EM algorithm computes the unknown parameter set Λ by minimizing the negative log likelihood $-\log p(\mathbf{y}^K; \Lambda)$. To compute the likelihood, we first note that the SBL framework imposes a Gaussian prior on the unknown vector $\mathbf{x}_k \sim \mathcal{N}(\mathbf{0}, \Gamma_k)$, where Γ_k is an unknown diagonal matrix. Thus, \mathbf{y}_k also follows a Gaussian distribution: $\mathbf{y}_k \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I} + \mathbf{A} \Gamma_k \mathbf{A}^\top)$ because the noise term $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. Therefore, we have

$$p(\mathbf{y}^K; \Lambda) = \prod_{k=1}^K \frac{1}{\sqrt{(2\pi)^m |\sigma^2 \mathbf{I} + \mathbf{A} \Gamma_k \mathbf{A}^\top|}} \times \exp\left(-\frac{1}{2} \mathbf{y}_k^\top (\sigma^2 \mathbf{I} + \mathbf{A} \Gamma_k \mathbf{A}^\top)^{-1} \mathbf{y}_k\right). \quad (97)$$

Hence, the negative log likelihood is computed as follows:

$$-\log p(\mathbf{y}^K; \Lambda) = \frac{1}{2} \sum_{k=1}^K \left[m \log(2\pi) + \log |\sigma^2 \mathbf{I} + \mathbf{A} \Gamma_k \mathbf{A}^\top| + \frac{1}{2} \mathbf{y}_k^\top (\sigma^2 \mathbf{I} + \mathbf{A} \Gamma_k \mathbf{A}^\top)^{-1} \mathbf{y}_k \right]. \quad (98)$$

Since the $\log(2\pi)$ term is a constant independent of Λ , we omit that term and the scaling factor of $\frac{1}{2}$ to obtain the cost function $T(\Lambda)$ in (3).

The EM algorithm treats the unknowns \mathbf{x}^K as the hidden data and the observations \mathbf{y}^K as the known data. It is an iterative procedure which updates the estimate of the parameters Λ in every iteration using two steps: an expectation step (E-step) and a maximization step (M-step). Let $\Lambda^{(r)}$ be the estimate of Λ at the r^{th} iteration. The E-step computes the marginal log-likelihood of the observed data $Q(\Lambda; \Lambda^{(r-1)})$, and the M-step computes the parameter tuple Λ that maximizes $Q(\Lambda; \Lambda^{(r-1)})$.

E-step: $Q(\Lambda; \Lambda^{(r-1)}) = \mathbb{E}_{\mathbf{x}^K | \mathbf{y}^K; \Lambda^{(r-1)}} \{\log p(\mathbf{y}^K, \mathbf{x}^K; \Lambda)\}$

M-step: $\Lambda^{(r)} = \arg \max_{\Lambda \in \mathbb{O} \times \mathbb{R}_+^{N \times K}} Q(\Lambda; \Lambda^{(r-1)}). \quad (99)$

The authors are with the Dept. of ECE at IISc, Bangalore, India, Emails: {geethu, cmurthy}@iisc.ac.in.

To simplify $Q(\Lambda, \Lambda^{(r-1)})$, we note that

$$p(\mathbf{y}^K, \mathbf{x}^K; \Lambda) = \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{x}_k; \Lambda) p(\mathbf{x}_k; \Lambda). \quad (100)$$

Here, $p(\mathbf{y}_k | \mathbf{x}_k; \Lambda) = \mathcal{N}(\mathbf{A} \mathbf{x}_k, \sigma^2 \mathbf{I})$, and $p(\mathbf{x}_k; \Lambda) = \mathcal{N}(\mathbf{0}, \Gamma_k)$. Thus, we get,

$$\begin{aligned} \log p(\mathbf{y}^K, \mathbf{x}^K; \Lambda) &= \log \left\{ \prod_{k=1}^K \frac{1}{\sqrt{(2\pi\sigma)^{2m}}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y}_k - \mathbf{A} \mathbf{x}_k\|^2\right) \right. \\ &\quad \times \left. \frac{1}{\sqrt{(2\pi)^N |\Gamma_k|}} \exp\left(-\frac{1}{2} \mathbf{x}_k^\top \Gamma_k^{-1} \mathbf{x}_k\right) \right\} \quad (101) \\ &= -\frac{Km}{2} \log((2\pi)^{N+1} \sigma^2) - \frac{1}{2} \sum_{k=1}^K [\log |\Gamma_k| + \text{Tr} \{\Gamma_k^{-1} \mathbf{x}_k \mathbf{x}_k^\top\}] \\ &\quad - \frac{1}{2\sigma^2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k). \quad (102) \end{aligned}$$

Therefore, eliminating the constant terms, we obtain (5) as follows:

$$\begin{aligned} Q(\Lambda; \Lambda^{(r-1)}) &= \\ &= -\frac{1}{2} \sum_{k=1}^K \left[\log |\Gamma_k| + \text{Tr} \left\{ \Gamma_k^{-1} \mathbb{E} \left\{ \mathbf{x}_k \mathbf{x}_k^\top | \mathbf{y}^K; \Lambda^{(r-1)} \right\} \right\} \right] \\ &= -\frac{1}{2\sigma^2} \sum_{k=1}^K \mathbb{E} \left\{ (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k) | \mathbf{y}^K; \Lambda^{(r-1)} \right\}. \quad (103) \end{aligned}$$

We notice that the expectation terms in the above expression depend only on $\Lambda^{(r-1)}$, and are independent of Λ . Thus, the dependence of Γ_k in $Q(\Lambda; \Lambda^{(r-1)})$ is only through the k^{th} term in the first summation, and the dependence on \mathbf{A} is only through the last summation term. Therefore, the optimization in the M-step is separable in its variables Γ_k and \mathbf{A} . Hence, the M-step reduces as follows:

$$\gamma_k^{(r)} = \arg \min_{\gamma \in \mathbb{R}_+^N} \log |\Gamma_k| + \text{Tr} \left\{ \Gamma_k^{-1} \mathbb{E} \left\{ \mathbf{x}_k \mathbf{x}_k^\top | \mathbf{y}^k; \Lambda^{(r-1)} \right\} \right\} \quad (104)$$

$$\mathbf{A}^{(r)} = \arg \min_{\mathbf{A} \in \mathbb{O}} \sum_{k=1}^K \mathbb{E} \left\{ (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k) | \mathbf{y}^k; \Lambda^{(r-1)} \right\}. \quad (105)$$

Here, we note that (105) is same as (7). Further, differentiating the objective function, we get the update equation (6):

$$\gamma_k^{(r)} = \text{Diag} \left\{ \mathbb{E} \left\{ \mathbf{x}_k \mathbf{x}_k^\top | \mathbf{y}^k; \mathbf{\Lambda}^{(r-1)} \right\} \right\} \quad (106)$$

$$= \text{Diag} \left\{ \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top + \boldsymbol{\Sigma}_{(k)} \right\}, \quad (107)$$

where we use the following facts:

$$\boldsymbol{\mu}_k \triangleq \mathbb{E} \left\{ \mathbf{x}_k | \mathbf{y}_k; \mathbf{\Lambda}^{(r-1)} \right\} \quad (108)$$

$$\boldsymbol{\Sigma}_{(k)} \triangleq \mathbb{E} \left\{ (\mathbf{x}_k - \boldsymbol{\mu}_k) (\mathbf{x}_k - \boldsymbol{\mu}_k)^\top | \mathbf{y}_k; \mathbf{\Lambda}^{(r-1)} \right\} \quad (109)$$

$$= \text{cov} \left\{ \mathbf{x}_k | \mathbf{y}^K; \mathbf{\Lambda}^{(r-1)} \right\}. \quad (110)$$

Next, we compute the conditional expectations terms needed to find $\gamma_k^{(r)}$. We start with the following cross-covariance matrix:

$$\begin{aligned} \mathbb{E} \left\{ \mathbf{y}_k \mathbf{x}_k^\top | \gamma_k, \sigma^2 \right\} &= \mathbb{E} \left\{ (\mathbf{A} \mathbf{x}_k + \mathbf{w}_k) \mathbf{x}_k^\top | \gamma_k, \sigma^2 \right\} \\ &= \mathbb{E} \left\{ \mathbf{A} \mathbf{x}_k \mathbf{x}_k^\top | \gamma_k, \sigma^2 \right\} \\ &= \mathbf{A} \boldsymbol{\Gamma}_k. \end{aligned} \quad (111)$$

Thus, the conditional mean and covariance are given as follows:

$$\begin{aligned} &\text{cov} \left\{ \mathbf{x}_k | \mathbf{y}^K; \mathbf{\Lambda} \right\} \\ &= \mathbb{E} \left\{ \mathbf{x}_k \mathbf{x}_k^\top | \gamma_k, \sigma^2 \right\} - \mathbb{E} \left\{ \mathbf{x}_k \mathbf{y}_k^\top | \gamma_k, \sigma^2 \right\} \\ &\quad \times \mathbb{E} \left\{ \mathbf{y}_k \mathbf{y}_k^\top | \gamma_k, \sigma^2 \right\}^{-1} \mathbb{E} \left\{ \mathbf{y}_k \mathbf{x}_k^\top | \gamma_k, \sigma^2 \right\} \\ &= \boldsymbol{\Gamma}_k - \boldsymbol{\Gamma}_k \mathbf{A}^\top \left(\sigma^2 \mathbf{I} + \mathbf{A} \boldsymbol{\Gamma}_k \mathbf{A}^\top \right)^{-1} \mathbf{A} \boldsymbol{\Gamma}_k \\ &\mathbb{E} \left\{ \mathbf{x}_k | \mathbf{y}^K; \mathbf{\Lambda} \right\} \\ &= \mathbb{E} \left\{ \mathbf{x}_k | \gamma_k, \sigma^2 \right\} + \mathbb{E} \left\{ \mathbf{x}_k \mathbf{y}_k^\top | \gamma_k, \sigma^2 \right\} \\ &\quad \times \mathbb{E} \left\{ \mathbf{y}_k \mathbf{y}_k^\top | \gamma_k, \sigma^2 \right\}^{-1} \left(\mathbf{y}_k - \mathbb{E} \left\{ \mathbf{y}_k | \gamma_k, \sigma^2 \right\} \right) \\ &= \boldsymbol{\Gamma}_k \mathbf{A}^\top \left(\sigma^2 \mathbf{I} + \mathbf{A} \boldsymbol{\Gamma}_k \mathbf{A}^\top \right)^{-1} \mathbf{y}_k \\ &= \sigma^{-2} \boldsymbol{\Gamma}_k \mathbf{A}^\top \left(\mathbf{I} - \left(\sigma^2 \mathbf{I} + \mathbf{A} \boldsymbol{\Gamma}_k \mathbf{A}^\top \right)^{-1} \mathbf{A} \boldsymbol{\Gamma}_k \mathbf{A}^\top \right) \mathbf{y}_k \\ &= \sigma^{-2} \text{cov} \left\{ \mathbf{x}_k | \mathbf{y}^K; \mathbf{\Lambda} \right\} \mathbf{A}^\top \mathbf{y}_k. \end{aligned} \quad (112)$$

Therefore, (106), (116) and (121) together gives the update step for γ_k used in Algorithm 1 and Algorithm 2.

Similarly, the optimization problem corresponding the dictionary update (105) reduces as follows:

$$\arg \min_{\mathbf{A} \in \mathbb{O}} \sum_{k=1}^K \mathbb{E} \left\{ (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k) \middle| \mathbf{y}_k; \mathbf{\Lambda}^{(r-1)} \right\} \quad (114)$$

$$= \arg \min_{\mathbf{A} \in \mathbb{O}} \sum_{k=1}^K \mathbb{E} \left\{ -\mathbf{y}_k^\top \mathbf{A} \mathbf{x}_k + \frac{1}{2} \mathbf{x}_k^\top \mathbf{A}^\top \mathbf{A} \mathbf{x}_k \middle| \mathbf{y}_k; \mathbf{\Lambda}^{(r-1)} \right\}$$

$$= \arg \min_{\mathbf{A} \in \mathbb{O}} -\text{Tr} \left\{ \left(\sum_{k=1}^K \boldsymbol{\mu}_k \mathbf{y}_k^\top \right) \mathbf{A} + \frac{1}{2} \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top \right\}$$

$$= \arg \min_{\mathbf{A} \in \mathbb{O}} \text{Tr} \left\{ -\mathbf{M} \mathbf{Y}^\top \mathbf{A} + \frac{1}{2} \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top \right\}. \quad (115)$$

Since $\mathbf{A} \in \mathbb{O}$, we can further simplify the second term here as follows:

$$\begin{aligned} \text{Tr} \left\{ \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top \right\} &= \sum_{i,j=1; i \neq j}^N \boldsymbol{\Sigma}[i, j] \mathbf{A}_i^\top \mathbf{A}_j + \sum_{i=1}^N \boldsymbol{\Sigma}[i, i] \mathbf{A}_i^\top \mathbf{A}_i \\ &= \text{Tr} \left\{ \mathbf{A} (\boldsymbol{\Sigma} - \mathcal{D} \{ \boldsymbol{\Sigma} \}) \mathbf{A}^\top \right\} + \sum_{i=1}^N \boldsymbol{\Sigma}[i, i]. \end{aligned} \quad (116)$$

$$= \text{Tr} \left\{ \mathbf{A} (\boldsymbol{\Sigma} - \mathcal{D} \{ \boldsymbol{\Sigma} \}) \mathbf{A}^\top \right\} + \sum_{i=1}^N \boldsymbol{\Sigma}[i, i]. \quad (117)$$

Here, the second term does not depend on \mathbf{A} , and hence, we remove the term from the objective function to get an equivalent optimization objective function as in (8). Thus, the derivation of algorithm development given by (3)-(8), and the update equations for γ_k in Algorithm 1 and Algorithm 2 are completed.

Learning the noise variance

Following a similar approach as the above, we can learn the noise variance σ^2 along with the dictionary \mathbf{A} and covariance matrices $\boldsymbol{\Gamma}_k$. If σ^2 is unknown, we have to incorporate its update to the M-step by maximizing the Q function defined in (103). Thus, considering the terms that depend on σ^2 , we get

$$\begin{aligned} (\sigma^2)^{(r)} &= \arg \min_{\sigma^2 \in \mathbb{R}_+} K m \log(\sigma^2) \\ &\quad + \frac{1}{\sigma^2} \sum_{k=1}^K \mathbb{E} \left\{ (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k) \middle| \mathbf{y}^K; \mathbf{\Lambda}^{(r-1)} \right\} \\ &= \frac{1}{K m} \sum_{k=1}^K \mathbb{E} \left\{ (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{A} \mathbf{x}_k) \middle| \mathbf{y}^K; \mathbf{\Lambda}^{(r-1)} \right\} \\ &= \frac{1}{K m} \text{Tr} \left\{ \mathbf{Y}^\top \mathbf{Y} - 2 \mathbf{M} \mathbf{Y}^\top \mathbf{A} + \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top \right\}, \end{aligned} \quad (118)$$

where the last step follows because of the same arguments used to derive (125) from (122).

II. PROOF OF KURDYKA-ŁOJASIEWICZ PROPERTY BASED CONVERGENCE RESULT

Theorem 6. A bounded sequence of iterates $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$ generated by the ALS algorithm converges to a stationary point of \tilde{g} if the following four conditions hold:

(i) The objective function $\tilde{g}(\mathbf{A})$ satisfies

$$\inf_{\mathbf{A} \in \mathbb{R}^{m \times N}} \tilde{g}(\mathbf{A}) > -\infty. \quad (119)$$

(ii) There exist constants $\theta \in [0, 1)$, $C, \epsilon > 0$ such that

$$|\tilde{g}(\mathbf{A}) - \tilde{g}(\mathbf{A}^*)|^\theta \leq C \|\mathbf{Z}\| \quad (120)$$

for any stationary point \mathbf{A}^* of \tilde{g} , any \mathbf{A} such that $\|\mathbf{A} - \mathbf{A}^*\| \leq \epsilon$, and any \mathbf{Z} such that $\mathbf{Z} \in \partial g(\mathbf{A})$. The constant θ is called the Łojasiewicz exponent of the Łojasiewicz gradient inequality.

(iii) There exists $C_1 > 0$ such that

$$\tilde{g}(\mathbf{A}^{(r,u-1)}) - \tilde{g}(\mathbf{A}^{(r,u)}) \geq C_1 \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u)} \right\|^2 \quad (121)$$

(iv) There exist $u_0 > 1$, $C_2 > 0$ and $\mathbf{Z} \in \partial g(\mathbf{A}^{(r,u)})$ such that for all $u > u_0$

$$\|\mathbf{Z}\| \leq C_2 \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u)} \right\|. \quad (122)$$

The proof is adapted from the proof of [40, Theorem 2]. At a high level, there are four steps to the proof:

A We first prove that the sequence $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$ converges to a bounded connected set $\mathbb{G} \subseteq \text{crit}(\tilde{g}) \subseteq \mathbb{O}$, where $\text{crit}(\tilde{g})$ is the set of stationary points of \tilde{g} . Moreover, \tilde{g} is constant over the set \mathbb{G} .

B Next, we connect the above result to Condition (ii). To establish the connection, we define a new function $\bar{g} : \mathbb{O} \rightarrow \mathbb{R}_+$ as $\bar{g}(\mathbf{A}) \triangleq \tilde{g}(\mathbf{A}) - \tilde{g}(\mathbf{A}^{(r)})$, where $\mathbf{A}^{(r)}$ is a limit point of the sequence $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$, and \mathbf{A} is any point in the set \mathbb{O} . We note that the definition of \bar{g} is unambiguous because Step A shows that \tilde{g} is constant over the set \mathbb{G} . We then show that there exists a positive integer $U_0 \in \mathbb{N}$ and $\tilde{C} > 0$ such that for all $u \geq U_0$,

$$\left(\bar{g}(\mathbf{A}^{(r,u)}) \right)^\theta \geq \tilde{C} \|\mathbf{Z}\|, \quad (123)$$

for any \mathbf{Z} such that $\mathbf{Z} \in \partial \tilde{g}(\mathbf{A}^{(r,u)})$.

C Finally, using the above relation and other conditions of the theorem, we show that the desired result follows.

Next, we present the details of the above steps:

A. Characterization of \mathbb{G}

From Condition (iii), we get that

$$\begin{aligned} & \sum_{u=1}^{\infty} \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u)} \right\|^2 \\ & \leq \frac{1}{C_1} \left[\lim_{u \rightarrow \infty} \tilde{g}(\mathbf{A}^{(r,u-1)}) - \tilde{g}(\mathbf{A}^{(r,0)}) \right] < \infty, \end{aligned} \quad (124)$$

where the last step follows because $\lim_{u \rightarrow \infty} \tilde{g}(\mathbf{A}^{(r,u-1)}) < \infty$ due to Proposition 1. Further, [45, Theorem 1] states that the set of subsequential limit points of a sequence $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$ in a compact metric space is a connected set if it satisfies the following:

$$\sum_{u=1}^{\infty} \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u)} \right\|^2 < \infty. \quad (125)$$

Consequently, the result applies to any bounded sequence satisfying (137). Since the sequence $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$ generated by the AM procedure belongs to the bounded set \mathbb{O} , it converges to a bounded connected set $\mathbb{G} \subseteq \mathbb{O}$. Also, since the set of subsequential limits is closed, \mathbb{G} is a connected compact set.

Now, for any limit point $\mathbf{A}^{(r)} \in \mathbb{G}$ of the sequence $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$, there exists a sequence $\{u_j\}_{j \in \mathbb{N}}$ of natural numbers such that $\left\{ \left(\mathbf{A}^{(r,u_j)}, \mathbf{Z}^{(r,u_j)}, \tilde{g}(\mathbf{A}^{(r,u_j)}) \right) \right\}_{j \in \mathbb{N}}$ converges to the tuple $(\mathbf{A}^{(r)}, \mathbf{0}, \tilde{g}(\mathbf{A}^{(r)}))$. This is because the subsequence $\left\{ \left(\mathbf{Z}^{(r,u_j)}, \tilde{g}(\mathbf{A}^{(r,u_j)}) \right) \right\}_{j \in \mathbb{N}}$ converges to the same limit point as that of the sequence

$\left\{ \left(\mathbf{Z}^{(r,u)}, \tilde{g}(\mathbf{A}^{(r,u)}) \right) \right\}_{u \in \mathbb{N}}$ which is $(\mathbf{0}, \tilde{g}(\mathbf{A}^{(r)}))$ due to (13) and Proposition 1. Therefore, we conclude that $\mathbb{G} \subseteq \text{crit}(\tilde{g})$ and \tilde{g} is constant over the set \mathbb{G} , completing Step A.

B. Connection to Kurdyka-Łojasiewicz property

The compact set \mathbb{G} can be covered with finite number of closed balls $\mathcal{B}_j = \left\{ \mathbf{A} \in \mathbb{O} : \left\| \mathbf{A} - \mathbf{A}^{*(j)} \right\| \leq \epsilon_j \right\}$ such that Condition (ii) is satisfied by $\mathbf{A}^{(r,j)}$ with constants $C^{(j)}$ and $\epsilon_j > 0$. Therefore, we have the following relation for $\mathbf{A} \in \mathcal{B}_j$:

$$\left| \tilde{g}(\mathbf{A}) - \tilde{g}(\mathbf{A}^{*(j)}) \right|^{\theta_j} \leq C^{(j)} \|\mathbf{Z}\|, \quad (126)$$

for some θ_j and any \mathbf{Z} such that $\mathbf{Z} \in \partial \tilde{g}(\mathbf{A})$. Setting $\epsilon = \min_j \epsilon_j$, $\tilde{C} = \max_j C^{(j)}$, and $\theta = \max_j \theta_j$ we get the following:

$$\left| \tilde{g}(\mathbf{A}) - \tilde{g}(\mathbf{A}^*) \right|^\theta \leq \tilde{C} \|\mathbf{Z}\|, \quad (127)$$

for any $\mathbf{A}^* \in \mathbb{G}$ of \tilde{g} , any \mathbf{A} such that $\|\mathbf{A} - \mathbb{G}\| \leq \epsilon$, and any \mathbf{Z} such that $\mathbf{Z} \in \partial \tilde{g}(\mathbf{A})$. Further, since $\left\{ \mathbf{A}^{(r,u)} \right\}_{u \in \mathbb{N}}$ converges to \mathbb{G} , for any $\epsilon > 0$, there exists a positive integer U_0 such that for all $u \geq U_0$, we have $\left\| \mathbf{A}^{(r,u)} - \mathbb{G} \right\| \leq \epsilon$. Therefore, for all $u \geq U_0$,

$$\left| \bar{g}(\mathbf{A}^{(r,u)}) \right|^\theta = \left| \tilde{g}(\mathbf{A}^{(r,u)}) - \tilde{g}(\mathbf{A}^{(r)}) \right|^\theta \leq \tilde{C} \|\mathbf{Z}\|. \quad (128)$$

Thus, Step B is completed.

C. Convergence to a single point

Since $\left\{ \tilde{g}(\mathbf{A}^{(r,u)}) \right\}_{u \in \mathbb{N}}$ is a non-increasing sequence, we have $\bar{g}(\mathbf{A}^{(r,u)}) \geq 0$, and the following relation holds.

$$\lim_{u \rightarrow \infty} \bar{g}(\mathbf{A}^{(r,u)}) = 0. \quad (129)$$

We first note that the function $h : \mathbb{R}_+ \rightarrow \mathbb{R}$ defined as $h(s) = -s^{1-\theta}$ is convex for all $0 \leq \theta \leq 1$. Thus, for all $u \in \mathbb{N}$ and for θ in Condition (ii), it holds that

$$\begin{aligned} & \left[\bar{g}(\mathbf{A}^{(r,u-1)}) \right]^{1-\theta} - \left[\bar{g}(\mathbf{A}^{(r,u)}) \right]^{1-\theta} \\ & = h\left(\bar{g}(\mathbf{A}^{(r,u-1)})\right) - h\left(\bar{g}(\mathbf{A}^{(r,u)})\right) \end{aligned} \quad (130)$$

$$\begin{aligned} & \geq \frac{dh(s)}{ds} \Big|_{s=\bar{g}(\mathbf{A}^{(r,u-1)})} \left[\bar{g}(\mathbf{A}^{(r,u-1)}) - \bar{g}(\mathbf{A}^{(r,u)}) \right] \\ & = (1-\theta) \left[\bar{g}(\mathbf{A}^{(r,u-1)}) \right]^{-\theta} \left[\bar{g}(\mathbf{A}^{(r,u-1)}) - \bar{g}(\mathbf{A}^{(r,u)}) \right] \end{aligned} \quad (131)$$

$$\begin{aligned} & \geq C_1(1-\theta) \left[\bar{g}(\mathbf{A}^{(r,u)}) \right]^{-\theta} \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u)} \right\|^2, \end{aligned} \quad (132)$$

$$\begin{aligned} & \geq C_1(1-\theta) \left[\bar{g}(\mathbf{A}^{(r,u)}) \right]^{-\theta} \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u)} \right\|^2, \end{aligned} \quad (133)$$

where we use Condition (iii) to obtain the last relation. Further, from Step B, we get that

$$\begin{aligned} & \left[\bar{g}(\mathbf{A}^{(r,u-1)}) \right]^{1-\theta} - \left[\bar{g}(\mathbf{A}^{(r,u)}) \right]^{1-\theta} \\ & \geq \frac{C_1(1-\theta)}{C} \frac{\left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\|^2}{\|\mathbf{Z}\|} \end{aligned} \quad (134)$$

$$\geq \frac{C_1(1-\theta)}{CC_2} \frac{\left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\|^2}{\left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u-2)} \right\|}, \quad (135)$$

where we use Condition (iv).

Next, we fix a constant $0 < \tau < 1$. For some $u \geq U_0$, if $\left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \geq \tau \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u-2)} \right\|$, from (147), we get the following:

$$\begin{aligned} & \frac{CC_2}{rC_1(1-\theta)} \left\{ \left[\bar{g} \left(\mathbf{A}^{(r,u-1)} \right) \right]^{1-\theta} - \left[\bar{g} \left(\mathbf{A}^{(r,u)} \right) \right]^{1-\theta} \right\} \\ & \geq \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\|. \end{aligned} \quad (136)$$

For all other values of $u \geq U_0$, we have the following relation:

$$\left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \leq \tau \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u-2)} \right\|. \quad (137)$$

Combining (148) and (149), for all $u \geq U_0$, we get the upper bound as given below:

$$\begin{aligned} & \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \leq \tau \left\| \mathbf{A}^{(r,u-1)} - \mathbf{A}^{(r,u-2)} \right\| \\ & + \frac{CC_2}{rC_1(1-\theta)} \left\{ \left[\bar{g} \left(\mathbf{A}^{(r,u-1)} \right) \right]^{1-\theta} - \left[\bar{g} \left(\mathbf{A}^{(r,u)} \right) \right]^{1-\theta} \right\}. \end{aligned} \quad (138)$$

Summing both sides, and using (141), we can simplify the expression as follows:

$$\begin{aligned} & \sum_{u=U_0}^{\infty} \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \\ & \leq \frac{\tau}{1-\tau} \left\| \mathbf{A}^{(r,U_0-1)} - \mathbf{A}^{(r,U_0-2)} \right\| \end{aligned}$$

$$+ \frac{CC_2}{rC_1(1-\theta)} \left[\bar{g} \left(\mathbf{A}^{(r,U_0)} \right) \right]^{1-\theta}. \quad (139)$$

Thus, we conclude that the series converges, and there exists a finite constant $\kappa < \infty$ such that the following holds:

$$\sum_{u=1}^{\infty} \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| = \kappa. \quad (140)$$

Hence, for any $\epsilon > 0$, there exists a positive integer U_1 such that for all $U \geq U_1$, we have

$$\kappa - \epsilon/2 \leq \sum_{u=1}^U \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \leq \kappa + \epsilon/2 \quad (141)$$

Thus, for any $U_1 \leq u_1 < u_2$, we have

$$\begin{aligned} & \left\| \left\| \mathbf{A}^{(r,u_2)} \right\| - \left\| \mathbf{A}^{(r,u_1)} \right\| \right\| \\ & \leq \sum_{u=u_1+1}^{u_2} \left\| \left\| \mathbf{A}^{(r,u)} \right\| - \left\| \mathbf{A}^{(r,u-1)} \right\| \right\| \end{aligned} \quad (142)$$

$$\leq \sum_{u=u_1+1}^{u_2} \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \quad (143)$$

$$= \sum_{u=1}^{u_2} \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| - \sum_{u=1}^{u_1} \left\| \mathbf{A}^{(r,u)} - \mathbf{A}^{(r,u-1)} \right\| \quad (144)$$

$$\leq \epsilon. \quad (145)$$

Therefore, the sequence $\left\{ \left\| \mathbf{A}^{(r,u)} \right\| \right\}_{u \in \mathbb{N}}$ is Cauchy, hence it converges.