# Dynamic Nonlinear Matrix Completion for Time-Varying Data Imputation Supplementary Material

## Jicong Fan

The Chinese University of Hong Kong, Shenzhen, China.
fanjicong@cuhk.edu.cn

## A    Proof for Theorem 1

**Theorem 1.** *Suppose* $\boldsymbol{X}_t = [\boldsymbol{x}_{t-w+1}, \boldsymbol{x}_{t-w+2}, \ldots, \boldsymbol{x}_t]$ *is given by Assumption 1. Let* $\phi : \mathbb{R}^d \mapsto \mathbb{R}^{\binom{d+q}{q}}$ *be a q-order polynomial feature map. Let* $c_t = \max(\|\boldsymbol{z}_{t-w+1}\|, \ldots, \|\boldsymbol{z}_t\|)$. *Then with probability 1, there exists a matrix* $\hat{\boldsymbol{X}}_t$ *with rank at most* $\min\left\{\binom{r+\theta}{\theta}, d, w\right\}$ *such that* $\|\boldsymbol{X}_t - \hat{\boldsymbol{X}}_t\|_F \leq \dfrac{\gamma c_t (w-2)^{1.5}}{3}$ *and* $\mathrm{rank}(\phi(\hat{\boldsymbol{X}}_t)) \leq \min\left\{\binom{r+\theta q}{\theta q}, \binom{d+q}{q}, w\right\}$.

*Proof.* Without loss of generality, we assume that $w$ is an odd number.

$$
\begin{aligned}
&\|g_t(\boldsymbol{z}_t) - g_{t-\frac{w-1}{2}}(\boldsymbol{z}_t)\| \\
\leq &\|g_t(\boldsymbol{z}_t) - g_{t-1}(\boldsymbol{z}_t)\| + \|g_{t-1}(\boldsymbol{z}_t) - g_{t-2}(\boldsymbol{z}_t)\| + \cdots \\
&+ \|g_{t-\frac{w-1}{2}+1}(\boldsymbol{z}_t) - g_{t-\frac{w-1}{2}}(\boldsymbol{z}_t)\| \\
\leq &\frac{w-1}{2}\gamma\|\boldsymbol{z}_t\|.
\end{aligned}
\tag{1}
$$

Similarly, we

$$
\|g_s(\boldsymbol{z}_s) - g_{t-\frac{w-1}{2}}(\boldsymbol{z}_s)\| \leq (\tfrac{w-1}{2} + s - t)\gamma\|\boldsymbol{z}_s\|, \tag{2}
$$

where $s = t - \frac{w-1}{2}, \ldots, t$. We also have

$$
\|g_s(\boldsymbol{z}_s) - g_{t-\frac{w-1}{2}}(\boldsymbol{z}_s)\| \leq (t - s - \tfrac{w-1}{2})\gamma\|\boldsymbol{z}_s\|, \tag{3}
$$

where $s = t - w + 1, \ldots, t - \frac{w-1}{2} - 1$. Putting (2) and (3) together, we get

$$
\begin{aligned}
&\sum_{s=t-w+1}^{t} \left\|g_s(\boldsymbol{z}_s) - g_{t-\frac{w-1}{2}}(\boldsymbol{z}_s)\right\|^2 \\
\leq &2 \sum_{v=1}^{(w-3)/2} v^2\gamma^2 c_t^2 \\
= &\gamma^2 c_t^2 (w-1)(w-2)(w-3)/12 \\
\leq &\gamma^2 c_t^2 (w-2)^3/12,
\end{aligned}
\tag{4}
$$

where $c_t = \max(\|\boldsymbol{z}_{t-w+1}\|, \ldots, \|\boldsymbol{z}_t\|)$. Let

$$
\hat{\boldsymbol{X}}_t = (\hat{\boldsymbol{x}}_{t-w+1}, \hat{\boldsymbol{x}}_{t-w+2}, \ldots, \hat{\boldsymbol{x}}_t),
$$

where $\hat{\boldsymbol{x}}_s = g_{t-\frac{w-1}{2}}(\boldsymbol{z}_s)$, $s = t-w+1, \ldots, t$. According to Lemma 1 of (Fan, Zhang, and Udell 2020), with probability 1, we have

$$
\mathrm{rank}(\hat{\boldsymbol{X}}_t) \leq \min\left\{\binom{r+\theta}{\theta}, d, w\right\}. \tag{5}
$$

On the other hand, according to (4) and the definition of $\hat{\boldsymbol{X}}_t$, we have

$$
\|\boldsymbol{X}_t - \hat{\boldsymbol{X}}_t\|_F \leq \frac{\gamma c_t (w-2)^{1.5}}{3}. \tag{6}
$$

Now combining (5) and (6), we conclude that $\boldsymbol{X}_t$ can be approximated by a matrix $\hat{\boldsymbol{X}}_t$ with rank at most $\min\left\{\binom{r+\theta}{\theta}, d, w\right\}$ and the approximation error is at most $\gamma c_t(w-2)^{1.5}/3$. This finished the proof for the first part of the theorem.

Let $\phi$ be a $q$-order polynomial feature map. According to Lemma 1 of (Fan, Zhang, and Udell 2020), we have

$$
\mathrm{rank}(\phi(\hat{\boldsymbol{X}}_t)) \leq \min\left\{\binom{r+\theta q}{\theta q}, \binom{d+q}{q}, w\right\}. \tag{7}
$$

Then we conclude than $\boldsymbol{X}_t$ can be approximated by a matrix $\hat{\boldsymbol{X}}_t$ satisfying $\mathrm{rank}(\phi(\hat{\boldsymbol{X}}_t)) \leq \min\left\{\binom{r+\theta q}{\theta q}, \binom{d+q}{q}, w\right\}$. Then we finish the proof. $\square$

## B    Gradient related to polynomial kernels

Denote by $\mathcal{L}_t$ the objective function in (5) of the main paper. We have

$$
\frac{\partial \mathcal{L}_t}{\partial \boldsymbol{K}_t} = \frac{p}{2}\boldsymbol{K}_t^{\frac{p}{2}-1} = \frac{p}{2}\boldsymbol{V}_t\boldsymbol{\Lambda}_t^{\frac{p}{2}-1}\boldsymbol{V}_t^\top, \tag{8}
$$

where $\boldsymbol{V}_t$ and $\mathrm{diag}(\boldsymbol{\Lambda}_t)$ are the eigenvectors and eigenvalues of $\boldsymbol{K}_t$ respectively. When $\boldsymbol{K}_t$ is computed by a polynomial kernel $k(\boldsymbol{x}_i, \boldsymbol{x}_j) = (\boldsymbol{x}_i^\top \boldsymbol{x}_j + a)^q$, we have

$$
\begin{aligned}
\frac{\partial \mathcal{L}_t}{\partial [\boldsymbol{x}_t]_{\bar{\omega}}} &= \sum_{i=1}^{w}\sum_{j=1}^{w} \frac{\partial \mathcal{L}_t}{\partial [\boldsymbol{K}_t]_{ij}} \frac{\partial [\boldsymbol{K}_t]_{ij}}{\partial [\boldsymbol{x}_t]_{\bar{\omega}}} \\
&= \left[2q\boldsymbol{X}_t\left(\boldsymbol{\alpha} \odot \left(\boldsymbol{X}_t^\top \boldsymbol{x}_t + a\right)^{\odot(q-1)}\right)\right]_{\bar{\omega}},
\end{aligned}
\tag{9}
$$

where $\boldsymbol{\alpha} = \left[\dfrac{\partial \mathcal{L}_t}{\partial \boldsymbol{K}_t}\right]_{:w}$. Invoking (8) into (9), we arrive at

$$
\frac{\partial \mathcal{L}_t}{\partial [\boldsymbol{x}_t]_{\bar{\omega}}} \left[2q\boldsymbol{X}_t\left(\left(\frac{p}{2}\boldsymbol{V}_t\boldsymbol{\Lambda}_t^{\frac{p}{2}-1}\boldsymbol{v}_t\right) \odot \left(\boldsymbol{X}_t^\top \boldsymbol{x}_t + a\right)^{\odot(q-1)}\right)\right]_{\bar{\omega}}, \tag{10}
$$

where $\boldsymbol{v}_t$ denotes the last columns of $\boldsymbol{V}_t^\top$.

## C  Proof for Theorem 2

**Theorem 2.** *Let $\boldsymbol{K}_t$ be the Gaussian kernel matrix with parameter $\sigma$. There exists a matrix $\tilde{\boldsymbol{K}}_t$ with rank at most $min\left\{\binom{r+\theta q}{\theta q}, \binom{d+q}{q}, w\right\}$ such that*

$$\|\boldsymbol{K}_t - \tilde{\boldsymbol{K}}_t\|_F \leq \frac{C_t \gamma w^2}{2\sigma^2} + \frac{C_t' w^2}{\sigma^{2(q+1)}(q+1)!}, \qquad (11)$$

*where $C_t$ and $C_t'$ are positive values relying on $\theta$, $q$, and $\max(\|\boldsymbol{z}_{t-w+1}\|, \ldots, \|\boldsymbol{z}_t\|)$.*

*Proof.* Let $\tilde{\boldsymbol{K}} = \boldsymbol{\Gamma} \odot \sum_{u=1}^q \sigma^{2u} u! \hat{\boldsymbol{K}}_j$, where $[\boldsymbol{\Gamma}]_{ij} = \exp\left(-\frac{\|\boldsymbol{x}_i\|^2 + \|\boldsymbol{x}_j\|^2 + 2a}{2\sigma^2}\right)$. According to Corollary 1 of (Fan, Zhang, and Udell 2020), we have

$$\left\|\hat{\boldsymbol{K}}_\sigma - \tilde{\boldsymbol{K}}\right\|_F \leq C_1, \qquad (12)$$

where $C_1 = w^2 \exp\left(-\frac{\min_i \|\hat{\boldsymbol{x}}_i\|^2}{\sigma^2}\right) \frac{\max_i \|\hat{\boldsymbol{x}}_i\|^q}{\sigma^{2(q+1)}(q+1)!}$ and $\text{rank}(\tilde{\boldsymbol{K}}) \leq \binom{r+\theta q}{\theta q}$ provided that $w/r$ is large enough. On the other hand, we have

$$\begin{aligned}
&\|\boldsymbol{K}_\sigma - \hat{\boldsymbol{K}}_\sigma\|_F^2 \\
&= \frac{1}{4\sigma^4} \sum_{ij} \left(\|\boldsymbol{x}_i - \boldsymbol{x}_j\|^2 - \|\hat{\boldsymbol{x}}_i - \hat{\boldsymbol{x}}_j\|^2\right)^2 \\
&\leq \frac{1}{4\sigma^4} \sum_{ij} C_{ij} \left(\|\boldsymbol{x}_i - \hat{\boldsymbol{x}}_i\| + \|\boldsymbol{x}_j - \hat{\boldsymbol{x}}_j\|\right)^2 \\
&\leq \frac{\max_{ij} C_{ij}}{4\sigma^4} \sum_{ij} \left(\frac{w-1}{2}\gamma\|\boldsymbol{z}_i\| + \frac{w-1}{2}\gamma\|\boldsymbol{z}_j\|\right)^2 \\
&\leq \frac{\gamma^2 w^2 (w-1)^2 \max_{ij} C_{ij} \max_i \|\boldsymbol{z}_i\|^2}{4\sigma^4},
\end{aligned} \qquad (13)$$

where $C_{ij} = 2\max(\|\boldsymbol{x}_i\|, \|\boldsymbol{x}_j\|, \|\hat{\boldsymbol{x}}_i\|, \|\hat{\boldsymbol{x}}_j\|)$. Combining (12) with (13), we obtain

$$\begin{aligned}
&\|\boldsymbol{K}_\sigma - \tilde{\boldsymbol{K}}\|_F \\
&\leq \|\boldsymbol{K}_\sigma - \hat{\boldsymbol{K}}_\sigma\|_F + \left\|\hat{\boldsymbol{K}}_\sigma - \tilde{\boldsymbol{K}}\right\|_F \\
&\leq \frac{\gamma w^2 C_x C_z}{2\sigma^2} + \frac{w^2 C_x' C_x^q}{\sigma^{2(q+1)}(q+1)!},
\end{aligned} \qquad (14)$$

where $C_x' = \exp\left(-\frac{\min_i \|\hat{\boldsymbol{x}}_i\|^2}{\sigma^2}\right)$, $C_x = \sqrt{2\max(\|\boldsymbol{x}_i\|, \|\boldsymbol{x}_j\|, \|\hat{\boldsymbol{x}}_i\|, \|\hat{\boldsymbol{x}}_j\|)}$, and $C_z = \max_i \|\boldsymbol{z}_i\|$. SInce $g_t$ is polynomial, there exists a constant $C_\theta$ large enough such that $\max_i \|\hat{\boldsymbol{x}}_i\| \leq C_\theta \max_i \|\boldsymbol{z}_i\|$, where $i = t-w+1, \ldots, t$. Letting $C_t = \sqrt{2C_\theta} (\max_i \|\boldsymbol{z}_i\|)^{3/2}$ and $C_t' = \exp(-\frac{C_\theta^2 (\max_i \|\boldsymbol{z}_i\|)^2}{\sigma^2}) (2C_\theta \max_i \|\boldsymbol{z}_i\|)^{q/2}$. It follows from (15) that

$$\|\boldsymbol{K}_\sigma - \tilde{\boldsymbol{K}}\|_F \leq \frac{C_t \gamma w^2}{2\sigma^2} + \frac{C_t' w^2}{\sigma^{2(q+1)}(q+1)!}. \qquad (15)$$

This finished the proof. $\qquad\square$

## D  Rank-one modification for fast EVD

Here we show how to perform rank-one modification (Brand 2006) twice to compute the eigenvalue decomposition of $\boldsymbol{K}_t$. Let $\boldsymbol{e}_w = [0, 0, \ldots, 0, 1]^\top$ and $\tilde{\boldsymbol{k}}' = [\boldsymbol{k}'^\top \; k(\boldsymbol{x}_t, \boldsymbol{x}_t)]^\top$. The method is detailed in Algorithm 1.

---
**Algorithm 1:** Rank-one modification for fast EVD of $\boldsymbol{K}_t$
---
**Input:** $\boldsymbol{V}_{t-1}', \boldsymbol{\Lambda}_{t-1}', \boldsymbol{e}_w, \boldsymbol{k}', \tilde{\boldsymbol{k}}'$
1: $\boldsymbol{U} \leftarrow \boldsymbol{V}_{t-1}', \boldsymbol{V} \leftarrow [\boldsymbol{V}_{t-1}'^\top \; \boldsymbol{0}]^\top, \boldsymbol{a} \leftarrow \tilde{\boldsymbol{k}}', \boldsymbol{b} \leftarrow \boldsymbol{e}_w$
2: $\boldsymbol{m} = \boldsymbol{U}^\top \boldsymbol{a}, \boldsymbol{p} = \boldsymbol{a} - \boldsymbol{U}\boldsymbol{m}, \bar{\boldsymbol{p}} = \boldsymbol{p}/\|\boldsymbol{p}\|$.
3: $\boldsymbol{n} = \boldsymbol{V}^\top \boldsymbol{b}, \boldsymbol{q} = \boldsymbol{b} - \boldsymbol{V}\boldsymbol{n}, \bar{\boldsymbol{q}} = \boldsymbol{q}/\|\boldsymbol{q}\|$.
4: $\boldsymbol{W} := \begin{bmatrix} \boldsymbol{\Lambda}_{t-1}' & \boldsymbol{0} \\ \boldsymbol{0} & 0 \end{bmatrix} + \begin{bmatrix} \boldsymbol{m} \\ \|\boldsymbol{p}\| \end{bmatrix} \begin{bmatrix} \boldsymbol{n} \\ \|\boldsymbol{q}\| \end{bmatrix}^\top$.
5: $\boldsymbol{W} = \boldsymbol{U}'\boldsymbol{\Sigma}'\boldsymbol{V}'^\top$.
6: $\bar{\boldsymbol{U}} \leftarrow \boldsymbol{U} \; \bar{\boldsymbol{p}}]\boldsymbol{U}', \bar{\boldsymbol{V}} \leftarrow \boldsymbol{V} \; \bar{\boldsymbol{q}}]\boldsymbol{V}'$.
7: $\boldsymbol{U} \leftarrow [\bar{\boldsymbol{U}}\top \; \boldsymbol{0}]^\top, \boldsymbol{V} \leftarrow \bar{\boldsymbol{V}}, \boldsymbol{a} \leftarrow \boldsymbol{e}_w, \boldsymbol{b} \leftarrow \tilde{\boldsymbol{k}}'$
8: $\boldsymbol{m} = \boldsymbol{U}^\top \boldsymbol{a}, \boldsymbol{p} = \boldsymbol{a} - \boldsymbol{U}\boldsymbol{m}, \bar{\boldsymbol{p}} = \boldsymbol{p}/\|\boldsymbol{p}\|$.
9: $\boldsymbol{n} = \boldsymbol{V}^\top \boldsymbol{b}, \boldsymbol{q} = \boldsymbol{b} - \boldsymbol{V}\boldsymbol{n}, \bar{\boldsymbol{q}} = \boldsymbol{q}/\|\boldsymbol{q}\|$.
10: $\boldsymbol{W} := \begin{bmatrix} \boldsymbol{\Sigma}' & \boldsymbol{0} \\ \boldsymbol{0} & 0 \end{bmatrix} + \begin{bmatrix} \boldsymbol{m} \\ \|\boldsymbol{p}\| \end{bmatrix} \begin{bmatrix} \boldsymbol{n} \\ \|\boldsymbol{q}\| \end{bmatrix}^\top$.
11: $\boldsymbol{W} = \boldsymbol{U}'\boldsymbol{\Sigma}'\boldsymbol{V}'^\top$.
12: $\boldsymbol{U}_t \leftarrow [\boldsymbol{U} \; \bar{\boldsymbol{p}}]\boldsymbol{U}', \boldsymbol{\Lambda}_t \leftarrow \boldsymbol{\Sigma}', \boldsymbol{V}_t \leftarrow [\boldsymbol{V} \; \bar{\boldsymbol{q}}]\boldsymbol{V}'$.
**Output:** $\boldsymbol{K}_t \approx \boldsymbol{V}_t \boldsymbol{\Lambda}_t \boldsymbol{V}_t^\top$.
---

## E  Proof for Theorem 3

**Theorem 3.** *Suppose $\boldsymbol{X}$ and $\Omega$ are given by Assumption 2. Let $\hat{\boldsymbol{X}}$ be a solution (not necessarily optimal) of (12) with a $q$-order polynomial kernel and let $\hat{\boldsymbol{K}}$ be the corresponding kernel matrix on $\hat{\boldsymbol{X}}$ and $\text{rank}(\hat{\boldsymbol{K}}) = R < \binom{\binom{r+\theta}{\theta}+q}{q}$. Suppose $\|\boldsymbol{X}\|_\infty, \|\hat{\boldsymbol{X}}\|_\infty \leq \beta$ and $\|\hat{\boldsymbol{X}}\|_F \leq \delta$. Then there exists a numerical constant $c$ such that the following inequality holds with probability at least $1 - \frac{2}{dn}$*

$$\begin{aligned}
&\frac{1}{|\bar{\Omega}|} \sum_{(i,j)\in\bar{\Omega}} \left([\boldsymbol{X}]_{ij} - [\hat{\boldsymbol{X}}]_{ij}\right)^2 \\
&\leq \frac{cdn\beta^2}{dn - |\Omega|} \left(\frac{\left(r^\star n + d\binom{r^\star+\theta^\star}{\theta^\star}\right)\log\frac{\delta}{\beta}}{|\Omega|}\right)^{1/2},
\end{aligned} \qquad (16)$$

*where $r^\star = \max\left\{\hat{r} \in \mathbb{Z}^+ : \psi(\theta^\star, \binom{r+\theta}{\theta}) \leq \hat{r} \leq \psi(\theta^\star q, R), \theta^\star \in \mathbb{Z}^+/\{1\}\right\}$.*

Before proving Theorem 3, we give the following lemmas.

**Lemma 1.** *Let $\mathcal{S}$ be the set of matrices $\boldsymbol{X} \in \mathbb{R}^{d \times n}$ whose columns are given by a polynomial function of order at most $\theta$ on a latent variable $\boldsymbol{z} \in \mathbb{R}^r$, where $\|\boldsymbol{X}\|_F \leq \delta$. Then there exists a constant $c$ such that the covering numbers of $\mathcal{S}$ with respect to Frobenius norm satisfy*

$$\mathcal{N}(\mathcal{S}, \|\cdot\|_F, \epsilon) \leq \left(\frac{c\delta}{\epsilon}\right)^{rn + d\binom{r+\theta}{\theta}}.$$

**Lemma 2** (Hoeffding inequality for sampling without replacement (Serfling 1974)). *Let $X_1, X_2, \ldots, X_s$ be a set of samples taken without replacement from a distribution $\{x_1, x_2, \ldots, x_N\}$ of mean $u$ and variance $\sigma^2$. Denote $a = \min_i x_i$ and $b = \max_i x_i$. Then*

$$P\left[\left|\frac{1}{s}\sum_{i=1}^{s} X_i - u\right| \geq t\right] \leq 2\exp\left(-\frac{2st^2}{(1-(s-1)/N)(b-a)^2}\right).$$

Now we start the proof for Theorem 3.

*Proof.* The assumption indicates that there exist $r_j$, $\theta_j$, and $s$ such that

$$\sum_{j=1}^{s}\binom{r_j + \theta_j q}{\theta_j q} \leq R, \qquad (17)$$

and the columns of $X$ can be fitted by $s$ polynomial functions

$$f_j : \mathbb{R}^{r_j} \to \mathbb{R}^d, \quad j = 1, 2, \ldots, s.$$

The difficulty is that we do not know what order and how many polynomials are fitted by the columns of $X$. We consider the following special cases.

*Case 1: highest-order polynomials.* The columns of $X$ lie on polynomials with the possibly highest order., which means $r_1 = \cdots = r_s = 1$. Without loss of generality, let $\theta_1 = \cdots = \theta_s = \theta^+$. We have

$$\theta^+ = \max\left\{\hat{\theta} \in \mathbb{Z}^+ : s\binom{1+\hat{\theta}q}{\hat{\theta}q} \leq R\right\} = \frac{R}{s} - 1.$$

Then the number of parameters (polynomial coefficients and latent variables) required to determine $X$ is

$$\pi_1 = n + s\binom{1+\theta^+}{\theta^+}d = n + Rd.$$

*Case 2: linear functions.* The columns of $X$ lie on lines, which means $\theta_1 = \cdots = \theta_s = 1$. Without loss of generality, let $r_1 = \cdots = r_s = r^+$. We have

$$r^+ = \max\left\{\hat{r} \in \mathbb{Z}^+ : s\binom{\hat{r}+q}{q} \leq R\right\}.$$

Since $\binom{\hat{r}+q}{q} \approx \frac{(\hat{r}+q)^q}{q!}$, we get

$$r^+ \approx \left\lceil\left(\frac{Rq!}{s}\right)^{1/q} - q\right\rceil. \qquad (18)$$

Here the minimum $s$ is 1 and the maximum $s$ is $R/(q+1)$. If $R$ is sufficiently small, we obtain

$$r^+ < \binom{r+\theta}{\theta},$$

which contradicts with the fact $d \geq \text{rank}(X) \geq \binom{r+\theta}{\theta}$. Therefore, Case 2 will not happen if $R$ is sufficiently small, i.e.,

$$R < s\binom{\binom{r+\theta}{\theta}+q}{q}, \qquad (19)$$

or if $R < \binom{\binom{r+\theta}{\theta}+q}{q}$ more strictly.

*Case 3: low-order polynomials.* The columns of $X$ lie on polynomials with order at least 2. Without loss of generality, we assume $\theta_1 = \cdots = \theta_s = \theta^\star \geq 2$ and $r_1 = \cdots = r_s = r^\star$. To ensure that (17) and $s\binom{r^\star+\theta^\star}{\theta^\star} \geq \text{rank}(X)$ hold simultaneously and $r^\star$ is sufficiently large, we get

$$r^\star = \max\left\{\hat{r} \in \mathbb{Z}^+ : s\binom{\hat{r}+\theta^\star q}{\theta^\star q} \leq R, \theta^\star \in \mathbb{Z}^+/\{1\},\right.$$

$$\left. s\binom{\hat{r}+\theta^\star}{\theta^\star} \geq \binom{r+\theta}{\theta}\right\}.$$

$$(20)$$

Let $\psi(v, C)$ be the root of equation $\binom{u+v}{v} = C$ with variable $u$. We have

$$r^\star = \max\left\{\hat{r} \in \mathbb{Z}^+ : \psi(\theta^\star, \binom{r+\theta}{\theta}/s) \leq \hat{r} \leq \psi(\theta^\star q, R/s),\right.$$

$$\left. \theta^\star \in \mathbb{Z}^+/\{1\}\right\}.$$

$$(21)$$

Note that using $\binom{a}{b} \approx \frac{a^b}{b!}$, we have

$$r^\star = \max\left\{\hat{r} \in \mathbb{Z}^+ : r_l \leq \hat{r} \leq r_u\right\}, \qquad (22)$$

where

$$r_l = \left(\frac{\binom{r+\theta}{\theta}\theta^\star!}{s}\right)^{1/\theta^\star} - \theta^\star,$$

$$r_u = \left(\frac{R(\theta^\star q)!}{s}\right)^{1/(\theta^\star q)} - \theta^\star q.$$

But the approximation given by (22) works only when $r$ is much larger than $\theta$ and $r^\star$ is much larger than $\theta^\star q$.

Then the number of parameters required to determine $X$ is

$$\pi_3 = \max_{s \in \mathbb{Z}^+} nr^\star + s\binom{r^\star + \theta^\star}{\theta^\star}d.$$

Since $n \gg d$ and $\binom{r^\star+\theta^\star}{\theta^\star} \ll R$, it suffices to let $s = 1$ and we arrive at

$$\pi_3 = nr^\star + \binom{r^\star + \theta^\star}{\theta^\star}d.$$

It is obvious that

$$\pi_1 < \pi_3.$$

Therefore, we will only consider Case 3 in the remaining context.

Let $\hat{\mathcal{L}}(X) := \frac{1}{|\Omega|}\|\mathcal{P}_\Omega(Y-X)\|_F^2$ and $\mathcal{L}(X) := \frac{1}{N}\|Y-X\|_F^2$. where $N = dn$. Suppose $\max\{\|Y\|_\infty, \|X\|_\infty\} \leq \beta$. According to Lemma 2, we have

$$P\left[|\hat{\mathcal{L}} - \mathcal{L}| \geq t\right] \leq 2\exp\left(-\frac{2|\Omega|t^2}{(1-(|\Omega|-1)/n^d)\eta^2}\right),$$

where $\eta = 4\beta^2$. Using union bound for all $\bar{X} \in \mathcal{S}$ (defined in Lemma 1), we obtain

$$P\left[\sup_{\bar{X}\in\mathcal{S}}|\hat{\mathcal{L}}(\bar{X}) - \mathcal{L}(\bar{X})| \geq t\right]$$

$$\leq 2|\mathcal{S}|\exp\left(-\frac{2|\Omega|t^2}{(1-(|\Omega|-1)/N)\eta^2}\right).$$

Equivalently, with probability at least $1 - 2N^{-1}$, we have

$$\sup_{\bar{\boldsymbol{X}} \in \mathcal{S}} |\hat{\mathcal{L}}(\bar{\boldsymbol{X}}) - \mathcal{L}(\bar{\boldsymbol{X}})| \leq \sqrt{\frac{\eta^2 \log(|\mathcal{S}|N)}{2} \left( \frac{1}{|\Omega|} - \frac{1}{N} + \frac{1}{N|\Omega|} \right)}$$

$$\leq \sqrt{\frac{\eta^2 \log(|\mathcal{S}|N)}{2|\Omega|}} \triangleq \Upsilon.$$

In $\mathcal{S}$ we use the $r^\star$ and $\theta^\star$ given by Case 3 because the corresponding $|\mathcal{S}|$ is largest.

Since $|\sqrt{u} - \sqrt{v}| \leq \sqrt{|u - v|}$ holds for any non-negative $u$ and $v$, we have

$$\sup_{\bar{\boldsymbol{X}} \in \mathcal{S}} \left| \sqrt{\hat{\mathcal{L}}(\bar{\boldsymbol{X}})} - \sqrt{\mathcal{L}(\bar{\boldsymbol{X}})} \right| \leq \sqrt{\Upsilon}.$$

As $\epsilon \geq \|\boldsymbol{X} - \bar{\boldsymbol{X}}\|_F \geq \|\mathcal{P}(\boldsymbol{X} - \bar{\boldsymbol{X}})\|_F$, we have

$$\left| \sqrt{\mathcal{L}(\boldsymbol{X})} - \sqrt{\mathcal{L}(\bar{\boldsymbol{X}})} \right|$$

$$= \frac{1}{\sqrt{N}} \left| \|\boldsymbol{Y} - \boldsymbol{X}\|_F - \|\boldsymbol{Y} - \bar{\boldsymbol{X}}\|_F \right| \leq \frac{\epsilon}{\sqrt{N}}$$

and

$$\left| \sqrt{\hat{\mathcal{L}}(\boldsymbol{X})} - \sqrt{\hat{\mathcal{L}}(\bar{\boldsymbol{X}})} \right|$$

$$= \frac{1}{\sqrt{|\Omega|}} \left| \|\mathcal{P}_\Omega(\boldsymbol{Y} - \boldsymbol{X})\|_F - \|\mathcal{P}_\Omega(\boldsymbol{Y} - \bar{\boldsymbol{X}})\|_F \right| \leq \frac{\epsilon}{\sqrt{|\Omega|}}.$$

It follows that

$$\sup_{\boldsymbol{X} \in \mathcal{S}} \left| \sqrt{\hat{\mathcal{L}}(\boldsymbol{X})} - \sqrt{\mathcal{L}(\boldsymbol{X})} \right|$$

$$\leq \sup_{\boldsymbol{X} \in \mathcal{S}} \left| \sqrt{\hat{\mathcal{L}}(\boldsymbol{X})} - \sqrt{\hat{\mathcal{L}}(\bar{\boldsymbol{X}})} \right| + \left| \sqrt{\hat{\mathcal{L}}(\bar{\boldsymbol{X}})} - \sqrt{\mathcal{L}(l\bar{\boldsymbol{X}})} \right|$$

$$+ \left| \sqrt{\mathcal{L}(\bar{\boldsymbol{X}})} - \sqrt{\mathcal{L}(\boldsymbol{X})} \right|$$

$$\leq \frac{\epsilon}{\sqrt{|\Omega|}} + \sqrt{\Upsilon} + \frac{\epsilon}{\sqrt{N}} \leq \frac{2\epsilon}{\sqrt{|\Omega|}} + \sqrt{\Upsilon}.$$

Now let $\epsilon = \beta$, we arrive at

$$\left| \sqrt{\hat{\mathcal{L}}(\boldsymbol{X})} - \sqrt{\mathcal{L}(\boldsymbol{X})} \right|$$

$$\leq \frac{2\beta}{\sqrt{|\Omega|}} + \beta \left( \frac{8 \log N + 8 \left( r^\star n + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{c\delta}{\beta}}{|\Omega|} \right)^{1/4}$$

$$\leq c'\beta \left( \frac{\left( r^\star n + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{\delta}{\beta}}{|\Omega|} \right)^{1/4},$$

$$\tag{23}$$

where $c'$ is a constant. Equivalently, we have

$$\frac{1}{\sqrt{dn}} \|\boldsymbol{X} - \hat{\boldsymbol{X}}\|_F \leq \frac{1}{\sqrt{|\bar{\Omega}|}} \|\mathcal{P}_{\bar{\Omega}}(\boldsymbol{X} - \hat{\boldsymbol{X}})\|_F$$

$$+ c'\beta \left( \frac{\left( r^\star n + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{\delta}{\beta}}{|\Omega|} \right)^{1/4}$$

$$= c'\beta \left( \frac{\left( r^\star n + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{\delta}{\beta}}{|\Omega|} \right)^{1/4},$$

$$\tag{24}$$

where we have used the fact that $[\hat{\boldsymbol{X}}]_{ij} = [\boldsymbol{X}]_{ij}$ for all $(i,j) \in \Omega$. As $\|\boldsymbol{X} - \hat{\boldsymbol{X}}\|_F^2 = \sum_{(i,j) \in \bar{\Omega}} \left( [\hat{\boldsymbol{X}}]_{ij} - [\boldsymbol{X}]_{ij} \right)^2$, we can write (24) as

$$\frac{\sum_{(i,j) \in \bar{\Omega}} \left( [\hat{\boldsymbol{X}}]_{ij} - [\boldsymbol{X}]_{ij} \right)^2}{|\bar{\Omega}|}$$

$$\leq \frac{c'' dn\beta^2}{dn - |\Omega|} \left( \frac{\left( r^\star n + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{\delta}{\beta}}{|\Omega|} \right)^{1/2}.$$

$$\tag{25}$$

This finished the proof. $\qquad\square$

## F  Proof for Corollary 1

**Corollary 1.** *Suppose $\boldsymbol{X}_t$ is given by Assumption 1. Let $\hat{\boldsymbol{X}}_t$ be the matrix recovered (not necessarily optimal) by (5) with a $q$-order polynomial kernel. Suppose the rank of the kernel matrix of $\hat{\boldsymbol{X}}_t - \boldsymbol{E}_t$ is $R$, where $\|\boldsymbol{E}_t\| \leq \varepsilon_t$. Other assumptions and notations are inherited from Theorem 3. Then there exists a numerical constant $c$ such that the following inequality holds with probability at least $1 - \frac{2}{dw}$*

$$\frac{1}{|\bar{\Omega}_t|} \sum_{(i,j) \in \bar{\Omega}_t} \left( [\boldsymbol{X}_t]_{ij} - [\hat{\boldsymbol{X}}_t]_{ij} \right)^2$$

$$\leq \frac{dw}{dw - |\Omega_t|} \left( \frac{8\varepsilon_t^2}{|\Omega_t|} + c\beta_t^2 \left( \frac{\left( r^\star w + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{\delta_t}{\beta_t}}{|\Omega_t|} \right)^{\frac{1}{2}} \right),$$

*where $r^\star = \max \left\{ \hat{r} \in \mathbb{Z}^+ : \psi(\theta^\star, \binom{r+\theta}{\theta}) \leq \hat{r} \leq \psi(\theta^\star q, R), \theta^\star \in \mathbb{Z}^+/\{1\} \right\}$.*

*Proof.* Let $\bar{\boldsymbol{X}}_t = \hat{\boldsymbol{X}}_t - \boldsymbol{E}_t$. Then we have rank$(\bar{\boldsymbol{X}}) = R$ and

$$\|\boldsymbol{X}_t - \hat{\boldsymbol{X}}_t\| \leq \|\boldsymbol{X}_t - \bar{\boldsymbol{X}}_t\| + \|\boldsymbol{E}_t\|. \tag{26}$$

Denote $\Delta = c\beta \left( \frac{\left( r^\star n + d\binom{r^\star + \theta^\star}{\theta^\star} \right) \log \frac{\delta}{\beta}}{|\Omega|} \right)^{1/4}$. We apply

(23) to $\boldsymbol{X}_t$ and $\hat{\boldsymbol{X}}_t$ and get

$$\frac{1}{\sqrt{dn}}\|\boldsymbol{X}_t - \bar{\boldsymbol{X}}_t\|_F$$
$$\leq \frac{1}{\sqrt{|\Omega|}}\|\mathcal{P}_\Omega(\boldsymbol{X}_t - \bar{\boldsymbol{X}}_t)\|_F + \Delta$$
$$\leq \frac{1}{\sqrt{|\Omega|}}\|\mathcal{P}_\Omega(\boldsymbol{X}_t - \hat{\boldsymbol{X}}_t)\|_F + \frac{1}{\sqrt{|\Omega|}}\|\mathcal{P}_\Omega(\boldsymbol{E}_t)\|_F + \Delta$$
$$\leq \frac{1}{\sqrt{|\Omega|}}\|\mathcal{P}_\Omega(\boldsymbol{E}_t)\|_F + \Delta.$$

(27)

It follows that

$$\frac{1}{\sqrt{dn}}\|\boldsymbol{X}_t - \hat{\boldsymbol{X}}_t\| \leq \frac{1}{\sqrt{dn}}\|\boldsymbol{X}_t - \bar{\boldsymbol{X}}_t\| + \frac{1}{\sqrt{dn}}\|\boldsymbol{E}_t\|$$
$$\leq \frac{1}{\sqrt{|\Omega|}}\|\mathcal{P}_\Omega(\boldsymbol{E}_t)\|_F + \Delta + \frac{1}{\sqrt{dn}}\|\boldsymbol{E}_t\|$$
$$\leq \frac{2\varepsilon_t}{\sqrt{|\Omega|}} + \Delta.$$

(28)

Or equivalently, we have

$$\frac{1}{|\bar{\Omega}|} \sum_{(i,j)\in\bar{\Omega}} \left([\boldsymbol{X}_t]_{ij} - [\hat{\boldsymbol{X}}_t]_{ij}\right)^2$$
$$\leq \frac{dn}{dn - |\Omega|}\left(\frac{2\varepsilon_t}{\sqrt{|\Omega|}} + \Delta\right)^2$$
$$= \frac{dn}{dn - |\Omega|}\left(\frac{2\varepsilon_t}{\sqrt{|\Omega|}} + c\beta\left(\frac{\left(r^\star n + d\binom{r^\star+\theta^\star}{\theta^\star}\right)\log\frac{\delta}{\beta}}{|\Omega|}\right)^{1/4}\right)^2$$
$$\leq \frac{dn}{dn - |\Omega|}\left(\frac{8\varepsilon_t^2}{|\Omega|} + c'\beta^2\left(\frac{\left(r^\star n + d\binom{r^\star+\theta^\star}{\theta^\star}\right)\log\frac{\delta}{\beta}}{|\Omega|}\right)^{1/2}\right),$$

where $c' = 2c$ is a constant. This finished the proof.

$\square$

## G  Proof for Lemma 1

*Proof.* Suppose $\boldsymbol{X}_1, \boldsymbol{X}_2 \in \mathcal{S}$. It means $\boldsymbol{X}_j = g_j(\boldsymbol{Z}_j) = \boldsymbol{P}_j\tilde{\boldsymbol{Z}}_j$, where $\tilde{\boldsymbol{Z}}_j \in \mathbb{R}^{\binom{r+\theta}{\theta}\times n}$ and $\boldsymbol{P}_j \in \mathbb{R}^{d\times\binom{r+\theta}{\theta}}$ denote the binomial terms and the coefficients respectively, $j = 1, 2$. Suppose $g_j$ is $L$-Lipschitz continuous, $\|\boldsymbol{Z}_j\|_F \leq \delta_1$, $\|\tilde{\boldsymbol{Z}}_j\|_F \leq \delta_2$, and $\|\boldsymbol{P}_j\|_F \leq \delta_3$, $j = 1, 2$. We have

$$\|\boldsymbol{X}_1 - \boldsymbol{X}_2\|_F = \|g_1(\boldsymbol{Z}_1) - g_2(\boldsymbol{Z}_2)\|_F$$
$$\leq \|g_1(\boldsymbol{Z}_1) - g_1(\boldsymbol{Z}_2)\|_F + \|g_1(\boldsymbol{Z}_2) - g_2(\boldsymbol{Z}_2)\|_F \quad (29)$$
$$\leq L\|\boldsymbol{Z}_1 - \boldsymbol{Z}_2\| + \|\tilde{\boldsymbol{Z}}_2\|_F\|\boldsymbol{P}_1 - \boldsymbol{P}_2\|_F.$$

Suppose $\|\boldsymbol{Z}_1 - \boldsymbol{Z}_2\| \leq \frac{\epsilon}{2L}$ and $\|\boldsymbol{P}_1 - \boldsymbol{P}_2\|_F \leq \frac{\epsilon}{2\|\tilde{\boldsymbol{Z}}_2\|_F}$. It follows that

$$\|\boldsymbol{X}_1 - \boldsymbol{X}_2\|_F \leq \epsilon. \quad (30)$$

Then we can bound the $\epsilon$-covering number of $\mathcal{S}$ as

$$\mathcal{N}(\mathcal{S}_{ab}, \|\cdot\|_F, \epsilon) \leq \left(\frac{6L\delta_1}{\epsilon}\right)^{rn}\left(\frac{6\delta_2\delta_3}{\epsilon}\right)^{d\binom{r+\theta}{\theta^\star}}$$
$$\leq \left(\frac{6\max(L\delta_1, \delta_2\delta_3)}{\epsilon}\right)^{r^\star n + d\binom{r+\theta}{\theta}}.$$

(31)

Although $L$ and $\{\delta_i\}_{i=1}^3$ are unknown, they are related to $\|\boldsymbol{X}\|_F$. We can bound $6\max(L\delta_1, \delta_2\delta_3)$ by $c\|\boldsymbol{X}\|_F$, where $c$ is a sufficiently large constant. Now we get

$$\mathcal{N}(\mathcal{S}, \|\cdot\|_F, \epsilon) \leq \left(\frac{c\delta}{\epsilon}\right)^{rn + d\binom{r+\theta}{\theta}}.$$

$\square$

## H  More about the experiments

**Data preprocessing** Since the variables in the SML2010 indoor temperature dataset and Air Quality dataset have very different scales, we rescale all variables by their standard deviations.

**Parameter setting of D-NLMC** For the synthetic data, we set $w = 20$, $R = 15$, and $\mu = 1$. For the Chlorine level dataset, we set $w = 100$, $R = 50$, and $\mu = 1$. For the SML2010 indoor temperature dataset, we set $w = 50$, $R = 25$, and $\mu = 1$. For the Air Quality dataset, we set $w = 50$, $R = 25$, and $\mu = 3$. Note that in OL-LRMC, we used the same $w$ as D-NLMC.