

Polynomial Matrix Completion for Missing Data Imputation and Transductive Learning (Supplementary Material)

Jicong Fan, Yuqian Zhang, Madeleine Udell

Cornell University, Ithaca, NY 14853, USA
{jf577, yz2557, udell}@cornell.edu

Proof for Lemmas and Theorems

Proof for Theorem 1

Theorem 1. Suppose \mathbf{X} is given by Assumption 1 with $\|\mathbf{Z}\|_\infty \leq c_z$. Then for any q' obeying $\binom{d+q'}{q'} \leq \min\{l, n\}$, there exists a matrix Φ with rank at most $\binom{d+q'}{q'}$, such that

$$\|\phi(\mathbf{X}) - \Phi\|_\infty \leq c_{q'}, \quad (1)$$

where $c_{q'} = \frac{c_z^{q'+1}}{(q'+1)!} \max_{\|\mathbf{z}\|_\infty \leq c_z} \|(\phi \circ f)^{(q'+1)}(\mathbf{z})\|_\infty$.

Proof. Denoting $\phi \circ f$ by h , we have $\phi(\mathbf{x}) = h(\mathbf{z}) = [h_1(\mathbf{z}), h_2(\mathbf{z}), \dots, h_l(\mathbf{z})]^T$. h is analytic and hence has convergent Taylor expansion. The q' order Taylor expansion of $h_i(\mathbf{z})$ at \mathbf{a} is given by

$$\begin{aligned} h_i(\mathbf{z}) &= h_i(\mathbf{a}) + \sum_{j=1}^d \frac{\partial h_i^{(1)}(\mathbf{a})}{\partial z_j} (z_j - a_j) \\ &+ \frac{1}{2!} \sum_{j_1=1}^d \sum_{j_2=1}^d \frac{\partial h_i^{(2)}(\mathbf{a})}{\partial z_{j_1} \partial z_{j_2}} (z_{j_1} - a_{j_1})(z_{j_2} - a_{j_2}) + \dots \\ &+ \frac{1}{q'!} \sum_{j_1=1}^d \dots \sum_{j_{q'}=1}^d \frac{\partial h_i^{(q')}(\mathbf{a})}{\partial z_{j_1} \dots \partial z_{j_{q'}}} (z_{j_1} - a_{j_1}) \dots (z_{j_{q'}} - a_{j_{q'}}) \\ &+ c_{i,q'}, \end{aligned} \quad (2)$$

where $c_{i,q'}$ denotes the residual of Taylor expansion. In addition,

$$|c_{i,q'}| \leq \frac{\|\mathbf{z}\|_\infty^{q'+1}}{(q'+1)!} \max_{\|\mathbf{z}\|_\infty \leq c_z, \|\mathbf{a}\|_\infty} |(\phi \circ f)_{ij}^{(q'+1)}(\mathbf{z})|.$$

For convenience, we set $\mathbf{a} = 0$. Let $\tilde{\mathbf{z}}$ be the vector consisting of $\{z_1^{\mu_1} z_2^{\mu_2} \dots z_d^{\mu_d}\}_{|\mu| \leq q'}$. Then $\tilde{\mathbf{z}} \in \mathbb{R}^{\tilde{d}'}$, where $\tilde{d}' = 1 + \sum_{i=1}^{q'} \binom{d+i-1}{i} = \binom{d+q'}{q'}$. Now (2) can be rewritten as

$$h_i(\mathbf{z}) = \boldsymbol{\theta}_i^T \tilde{\mathbf{z}} + c_{i,q'}, \quad (3)$$

where $\boldsymbol{\theta}_i \in \mathbb{R}^{\tilde{d}'}$ consists of the corresponding coefficients. Let $\Theta = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_l]^T \in \mathbb{R}^{l \times \tilde{d}'}$ and $\tilde{\mathbf{Z}} = [\tilde{\mathbf{z}}_1, \tilde{\mathbf{z}}_2, \dots, \tilde{\mathbf{z}}_n] \in$

$\mathbb{R}^{\tilde{d}' \times n}$. Then for all $1 \leq i \leq l$ and $1 \leq j \leq n$,

$$|[h(\mathbf{Z}) - \Theta \tilde{\mathbf{Z}}]_{ij}| \leq c_{q'}, \quad (4)$$

where $c_{q'} = \frac{c_z^{q'+1}}{(q'+1)!} \max_{\|\mathbf{z}\|_\infty \leq c_z} \|(\phi \circ f)^{(q'+1)}(\mathbf{z})\|_\infty$.

Denote $\Phi = \Theta \tilde{\mathbf{Z}}$ and $\text{rank}(\Phi) = \tilde{d}$. We have

$$\|\phi(\mathbf{X}) - \Phi\|_\infty \leq c_{q'}. \quad (5)$$

Since \mathbf{Z} is of full-rank, it is easy to show that $\tilde{\mathbf{Z}}$ is of full-rank. For a specific $\{f_i\}_{i=1}^m$, the rank of Θ may be lower than \tilde{d}' especially when some of the partial derivatives in (2) vanish or/and repeated. Therefore, $\tilde{d} \leq \tilde{d}' = \binom{d+q'}{q'}$. This finished the proof. \square

Proof for Lemma 1

Lemma 1. Suppose \mathbf{X} is given by Assumption 2. Then

$$\text{rank}(\mathbf{X}) \leq \min\left\{\binom{d+\alpha}{\alpha}, m, n\right\}$$

and

$$\text{rank}(\phi(\mathbf{X})) \leq \min\left\{\binom{d+\alpha q}{\alpha q}, \binom{m+q}{q}, n\right\}.$$

Proof. It is easy to have $\text{rank}(\mathbf{X}) \leq \min\left\{\binom{d+\alpha}{\alpha}, m, n\right\}$. In addition, $\phi(\mathbf{X}) \in \mathbb{R}^{l \times n}$, where $l = \binom{m+q}{\alpha}$. As $\phi(\mathbf{x}) = \phi(f(\mathbf{z}))$ is a αq order polynomial mapping of \mathbf{z} , the number of polynomial features of \mathbf{z} is at most $\binom{d+\alpha q}{\alpha q}$. Therefore $\text{rank}(\phi(\mathbf{X})) \leq \min\left\{\binom{d+\alpha q}{\alpha q}, \binom{m+q}{q}, n\right\}$. \square

Proof for Lemma 2

Lemma 2. For any \mathbf{X} given by Assumption 1 and any positive integer α , there exists an $\tilde{\mathbf{X}}$ obeying Assumption 2, such that

$$\|\mathbf{X} - \tilde{\mathbf{X}}\|_\infty \leq c_\alpha, \quad (6)$$

where $c_\alpha = \frac{c_z^{\alpha+1}}{(\alpha+1)!} \max_{\|\mathbf{z}\|_\infty \leq c_z} \|f^{(\alpha+1)}(\mathbf{z})\|_\infty$.

Proof. Lemma 2 is a special case of Theorem 1. The proof can be obtained via letting $\phi(\mathbf{x}) = \mathbf{x}$ in the proof of Theorem 1. \square

Proof for Theorem 2

Theorem 2. Let $\text{mnp}_\phi(\mathbf{X})$ be the minimum number of parameters required to determine \mathbf{X} uniquely among all matrices in the set $\{\mathcal{X} \in \mathbb{R}^{m \times n} : \text{rank}(\phi(\mathcal{X})) = \tilde{d}\}$. Define $\tilde{r} := \min\{o : \binom{o+q}{q} \geq \tilde{d}\}$. Then

$$\text{mnp}_\phi(\mathbf{X}) = (m - \tilde{r}) \times \tilde{d} + n\tilde{r}. \quad (7)$$

Proof. As the first-order feature in $\phi(\mathbf{X})$ is \mathbf{X} itself and the rank of $\phi(\mathbf{X})$ is \tilde{d} , then there exists a $\mathbf{W} \in \mathbb{R}^{m \times \tilde{d}}$ such that

$$\mathbf{X} = \mathbf{W}\phi_{\tilde{d}}(\mathbf{X}), \quad (8)$$

where $\phi_{\tilde{d}}(\mathbf{X})$ denotes the corresponding \tilde{d} rows of $\phi(\mathbf{X})$. Suppose \mathbf{X}_o consists of o rows of \mathbf{X} and \mathbf{X}_{m-o} consists of the remainders. Let o be sufficiently large such that the number of rows of $\phi(\mathbf{X}_o)$ is no less than \tilde{d} :

$$\binom{o+q}{q} \geq \tilde{d}. \quad (9)$$

Then it follows from (8) that

$$\mathbf{X}_{m-o} = \mathbf{W}_{m-o}\phi(\mathbf{X}_o), \quad (10)$$

for some $\mathbf{W}_{m-o} \in \mathbb{R}^{(m-o) \times \tilde{d}}$. We see that \mathbf{X}_{m-o} can be reconstructed from \mathbf{X}_o . In addition, since $\mathbf{X}_o = f_o(\mathbf{Z})$ (according to Assumption 2), we have

$$\mathbf{X}_{m-o} = \mathbf{W}_{m-o}(\phi(f_o(\mathbf{Z})) := \hat{f}_{m-o}(\mathbf{Z}). \quad (11)$$

Therefore, minimizing the rank of $\phi(\mathbf{X})$ implicitly approximates the unknown mapping f .

According to (10), we can count the minimum number of parameters required to determine \mathbf{X} uniquely among all matrices in the set $\{\mathcal{X} \in \mathbb{R}^{m \times n} : \text{rank}(\phi(\mathcal{X})) = \tilde{d}\}$. First, we denote the minimum o satisfying (9) as

$$\tilde{r} := \min\{o : \binom{o+q}{q} \geq \tilde{d}\}, \quad (12)$$

and have

$$(q!\tilde{d})^{1/q} - q \leq \tilde{r} \leq (q!\tilde{d})^{1/q}, \quad (13)$$

which means we can compute \tilde{r} via at most $q+1$ trials. Then $\mathbf{X}_{\tilde{r}}$ counts $n\tilde{r}$ parameters. Second, there are $(m - \tilde{r}) \times \tilde{d}$ parameters for $\mathbf{W}_{m-\tilde{r}}$ to reconstruct $\mathbf{X}_{m-\tilde{r}}$ from $\mathbf{X}_{\tilde{r}}$ and ϕ has been given. Therefore, the minimum number of parameters we required is

$$\text{mnp}_\phi(\mathbf{X}) := (m - \tilde{r}) \times \tilde{d} + n\tilde{r}. \quad \square$$

Proof for Lemma 3

Lemma 3. For any $\mathbf{X} \in \mathbb{R}^{m \times n}$,

$$\begin{aligned} \|\phi(\mathbf{X})\|_{S_p|s}^p &= \text{Tr}(\mathcal{K}(\mathbf{X})^{p/2}) \\ &- \max_{\mathbf{P}^T \mathbf{P} = \mathbf{I}_s, \mathbf{P} \in \mathbb{R}^{n \times s}} \text{Tr}\left((\mathbf{P}^T \mathcal{K}(\mathbf{X}) \mathbf{P})^{p/2}\right). \end{aligned} \quad (14)$$

Proof. Denote $\tilde{\mathcal{K}}(\mathbf{X}) = \mathbf{P}^T \mathcal{K}(\mathbf{X}) \mathbf{P}$. Let λ_i be the i -th eigenvalue of $\mathcal{K}(\mathbf{X})$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. Let $\tilde{\lambda}_i$ be the i -th eigenvalue of $\tilde{\mathcal{K}}(\mathbf{X})$ and $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_n$. Then

$$\text{Tr}\left((\mathbf{P}^T \mathcal{K}(\mathbf{X}) \mathbf{P})^{p/2}\right) = \sum_{i=1}^s \tilde{\lambda}_i^{p/2}.$$

Using the Lemma 3.31(a) of (Horn and Johnson 1991), we have

$$\tilde{\lambda}_i \leq \lambda_i,$$

where $i = 1, 2, \dots, s$. It follows that

$$\sum_{i=1}^s \tilde{\lambda}_i^{p/2} \leq \sum_{i=1}^s \lambda_i^{p/2}.$$

Therefore,

$$\max_{\mathbf{P}^T \mathbf{P} = \mathbf{I}_s} \text{Tr}\left((\mathbf{P}^T \mathcal{K}(\mathbf{X}) \mathbf{P})^{p/2}\right) = \sum_{i=1}^s \lambda_i^{p/2},$$

and

$$\begin{aligned} \|\phi(\mathbf{X})\|_{S_p|s}^p &= \text{Tr}(\mathcal{K}(\mathbf{X})^{p/2}) - \sum_{i=1}^s \lambda_i^{p/2} \\ &= \sum_{i=s+1}^n \lambda_i^{p/2} = \sum_{i=s+1}^n \sigma_i^p. \end{aligned}$$

This finished the proof. \square

Proof for Lemma 4

Lemma 4. For any $\mathbf{X} \in \mathbb{R}^{m \times n}$,

$$\|\phi(\mathbf{X})\|_{S_p|w}^p = \min_{\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I}_n} \text{Tr}\left((\mathbf{W}^{1/p} \mathbf{Q}^T \mathcal{K}(\mathbf{X}) \mathbf{Q} \mathbf{W}^{1/p})^{p/2}\right).$$

Proof. Denote $\bar{\mathcal{K}}(\mathbf{X}) = \mathbf{Q}^T \mathcal{K}(\mathbf{X}) \mathbf{Q}$. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ and $\bar{\lambda}_1 \geq \bar{\lambda}_2 \geq \dots \geq \bar{\lambda}_n$ be the eigenvalues of $\mathcal{K}(\mathbf{X})$ and $\bar{\mathcal{K}}(\mathbf{X})$ respectively. Since \mathbf{Q} is an orthogonal matrix, we have $\bar{\lambda}_i = \lambda_i$ for all $i = 1, 2, \dots, n$. It follows that

$$\begin{aligned} &\text{Tr}\left((\mathbf{W}^{1/p} \mathbf{Q}^T \mathcal{K}(\mathbf{X}) \mathbf{Q} \mathbf{W}^{1/p})^{p/2}\right) \\ &\geq \text{Tr}(\mathbf{W}^{1/2} \bar{\mathcal{K}}(\mathbf{X})^{p/2} \mathbf{W}^{1/2}) \\ &= \text{Tr}(\mathbf{W} \bar{\mathcal{K}}(\mathbf{X})^{p/2}) \\ &\geq \sum_{i=1}^n w_i \bar{\lambda}_i^{p/2} = \sum_{i=1}^n w_i \lambda_i^{p/2} \\ &= \sum_{i=1}^n w_i \sigma_i^p = \|\phi(\mathbf{X})\|_{S_p|w}^p, \end{aligned} \quad (15)$$

where the first inequality holds due to the Araki–Lieb–Thirring inequality (Araki 1990) and the second inequality holds due to the product trace inequality in (Marshall, Olkin, and Arnold 2011) (pages 340-341). This finished the proof. \square

Table 1: Datasets for classification

Dataset	# features	# samples	# classes	missing value
Mice protein	82	1080	8	1.68%
Shuttle	9	2175	5	10%
Dermatology	34	366	6	1%
Satimage	36	6435	7	10%

Table 2: Classification errors

Data set	θ	SVM	LRMC	LRMC+SVM	SRMC	VMC-2	VMC-3	NLMC	PMC-S	PMC-W
Mice protein	10%	8.96	6.33	0.8	2.96	0.54	0.46	0.44	0.41	0.39
	30%	21.41	10.44	2.59	5.65	0.61	0.48	0.44	0.44	0.33
	50%	32.5	18.78	5.26	13.19	1.24	0.87	0.81	0.71	0.63
	70%	46.7	34.02	19.2	29.8	8.28	5.15	4.06	2.43	2.63
Shuttle	10%	12.18	24.7	2.48	21.06	9.7	7.72	4.8	2.66	3.86
	30%	14.74	25.9	6.4	22.8	11.04	8.7	6.12	4.32	5.56
	50%	17.82	28.7	10.6	27.3	13.4	11.1	9.58	8.02	9.16
	70%	19.7	38.12	19.76	40.28	18.08	16.16	16.24	14.8	15.82
Dermatology	10%	4.48	4.54	3.28	4.21	3.17	3.12	3.08	2.84	2.84
	30%	6.83	6.56	5.96	6.01	5.41	5.3	5.14	5.03	4.92
	50%	13.07	9.95	8.83	9.08	8.64	8.74	8.31	8.16	7.98
	70%	26.67	21.09	20.87	21.31	19.67	20.08	19.13	18.5	18.87
Satimage	10%	39.6	23.34	14.38	20.62	17.32	15.5	14.7	13.06	14.24
	30%	41.46	23.58	16.22	22.52	17.86	16.12	15.2	14.56	14.94
	50%	44.2	24.24	16.96	24.1	18.44	16.18	15.84	14.82	15.18
	70%	53.58	27.24	22.62	26.5	21.92	18.86	17.86	16.64	17.3

More details about the experiments

Parameter setting

In our PMC-S and PMC-W, we use RBF kernel because it is often more effective than polynomial kernel. In VMC-2, VMC-3 (Ongie et al. 2017), NLMC (Fan and Chow 2018), PMC-S, and PMC-W, for the Schatten- p norm, we set $p = 1$ or 0.5 and report the better result. In VMC-2 and VMC-3, the parameter α of polynomial kernel is chosen from $\{1, 10\}$. In NLMC, PMC-S, and PMC-W, the hyper-parameter σ of RBF kernel is set as the average distance of pair-wise data points multiplied by 1 or 3. The hyper-parameters in LRMC, SRMC (Fan and Chow 2017) and LADMC (Ongie et al. 2018) (the iterative algorithm) are carefully determined to provide the best performance as possible.

Datasets with missing values for classification

As shown in Table 1, we consider four real datasets with missing values¹. For each dataset, we randomly remove a fraction (0.1, 0.3, 0.5, or 0.7) of the observed entries of the feature matrix and perform classification. In this study, we use a subset of Shuttle/Satimage dataset, which consists of 1000 randomly chosen samples. In addition, the proportion of training (labeled) data is 50%. Since the features of each dataset often have different scale, we normalize the features to have zero mean and unit variance. We use matrix completion to recover the missing entries and classify the data simultaneously. As a well-known baseline of classification

technique, SVM with RBF kernel is compared, where the missing entries are replaced by zeros. Moreover, we also perform SVM on the data recovered by LRMC.

The classification results are shown in Table 2. Our PMC-S and PMC-W outperform other methods in almost all cases.

References

- Araki, H. 1990. On an inequality of Lieb and Thirring. *Letters in Mathematical Physics* 19(2):167–170.
- Fan, J., and Chow, T. W. 2017. Matrix completion by least-square, low-rank, and sparse self-representations. *Pattern Recognition* 71:290 – 305.
- Fan, J., and Chow, T. W. 2018. Non-linear matrix completion. *Pattern Recognition* 77:378 – 394.
- Horn, R. A., and Johnson, C. R. 1991. *Topics in Matrix Analysis*. Cambridge University Press.
- Marshall, A. W.; Olkin, I.; and Arnold, B. C. 2011. *Inequalities: Theory of Majorization and Its Applications*. Springer, New York.
- Ongie, G.; Willett, R.; Nowak, R. D.; and Balzano, L. 2017. Algebraic variety models for high-rank matrix completion. In *Proceedings of the 34th International Conference on Machine Learning*, 2691–2700. Sydney, Australia: PMLR.
- Ongie, G.; Balzano, L.; Pimentel-Alarcón, D.; Willett, R.; and Nowak, R. D. 2018. Tensor Methods for Nonlinear Matrix Completion. *ArXiv e-prints*.

¹<https://archive.ics.uci.edu/ml/datasets/Mice+Protein+Expression>
<https://sci2s.ugr.es/keel/missing.php>