

Summary on basic time series studies

tensor data analysis with different data types

Haofan Zheng

Contents

1	High-dimentional α-PCA method	3
1.1	Overall Summary	3
1.2	Main model	3
1.3	Main Statistics	3
1.4	Transformation	4
1.5	Interpretation	4
1.6	Relative estimators	5
1.7	Theoretical Properties	6
1.8	Simulation	6
1.9	Application	6
2	High-Dimensional GLM with Binary Outcomes	6
2.1	Overall Summary	6
3	Ultra-High Dimensional GFM	6
3.1	Overall Summary	6
4	Matrix-variate Logistic Regression with Measurement Error	6
5	A Likelihood-Based Approach for Multivariate Categorical Response Regression in High Dimensions	6
6	A likelihood-Based Approach for Semiparametric Regression with Panel Count Data	6
7	Time Series Latent Gaussian Count	6
8	Time Series Factor Models(tensor)	6

1 High-dimentional α -PCA method

1.1 Overall Summary

This article considers the estimation and inference of the **low rank** components in high-dimentional matrix-variate models(tensor), and we propose an estimation method called α -PCA and it has some benefits with the high dimensions data favorably compared with other methods(traditional PCA, etc) based on the performance in the simulation.

1.2 Main model

The model is shown as the following:

$$\mathbf{Y}_t = \underbrace{\mathbf{R}\mathbf{F}_t\mathbf{C}^T}_{\text{signal part}} + \underbrace{\mathbf{E}_t}_{\text{noise part}} \quad (1)$$

$\mathbf{Y}_t : \mathbf{Y}_t \in \mathbb{R}^{p \times q}$, $1 \leq t \leq T$, observations,

$\mathbf{F}_t : \mathbf{F}_t \in \mathbb{R}^{k \times r}$, where $k \ll p$ and $r \ll q$ (**low rank**), latent matrix,

$\mathbf{E}_t : \mathbf{E}_t \in \mathbb{R}^{p \times q}$, noise matrix.

1.3 Main Statistics

An estimation procedure, namely α -PCA, aggregates the information in both first and second moments. Specifically, the two statistics are defined:

$$\widehat{\mathbf{M}}_R \triangleq \frac{1}{pq} \left((1 + \alpha) \cdot \overline{\mathbf{Y}}\overline{\mathbf{Y}}^T + \frac{1}{T} \sum_{t=1}^T (\mathbf{Y}_t - \overline{\mathbf{Y}})(\mathbf{Y}_t - \overline{\mathbf{Y}})^T \right) \quad (2)$$

$$\widehat{\mathbf{M}}_C \triangleq \frac{1}{pq} \left((1 + \alpha) \cdot \overline{\mathbf{Y}}^T \overline{\mathbf{Y}} + \frac{1}{T} \sum_{t=1}^T (\mathbf{Y}_t - \overline{\mathbf{Y}})^T (\mathbf{Y}_t - \overline{\mathbf{Y}}) \right) \quad (3)$$

$\alpha : \alpha \in [-1, +\infty)$, a hyperparameter,

$\overline{\mathbf{Y}} = \frac{1}{T} \sum_{i=1}^T \mathbf{Y}_t$, the sample mean.

Based on these two statistics, estimation of \mathbf{R} and \mathbf{C} can be obtained as \sqrt{p} times the top k eigenvectors of $\widehat{\mathbf{M}}_R$ and \sqrt{q} times the top q eigenvectors of $\widehat{\mathbf{M}}_C$ respectively, in descending order by corresponding eigenvalues.

1.4 Transformation

To simplify the estimator, we can transform the parameters, let the $\tilde{\alpha} = \sqrt{\alpha + 1} - 1$, $\tilde{\mathbf{Y}}_t \triangleq \mathbf{Y}_t + \tilde{\alpha}\bar{\mathbf{Y}}$, $\tilde{\mathbf{F}}_t \triangleq \mathbf{F}_t + \tilde{\alpha}\bar{\mathbf{F}}$, and $\tilde{\mathbf{E}}_t \triangleq \mathbf{E}_t + \tilde{\alpha}\bar{\mathbf{E}}$, Then we have

$$\tilde{\mathbf{Y}}_t = \mathbf{R}\tilde{\mathbf{F}}_t\mathbf{C}^T + \tilde{\mathbf{E}}_t \quad (4)$$

The equation 2 and 3 can be rewritten as:

$$\widehat{\mathbf{M}}_R = \frac{1}{pqT} \sum_{t=1}^T \tilde{\mathbf{Y}}_t \tilde{\mathbf{Y}}_t^T, \text{ and } \widehat{\mathbf{M}}_C = \frac{1}{pqT} \sum_{t=1}^T \tilde{\mathbf{Y}}_t^T \tilde{\mathbf{Y}}_t \quad (5)$$

Same as in section 1.3, $\widehat{\mathbf{R}}$ and $\widehat{\mathbf{C}}$ can be obtained as \sqrt{p} times the top k eigenvectors of $\widehat{\mathbf{M}}_R$ and \sqrt{q} times the top q eigenvectors of $\widehat{\mathbf{M}}_C$ respectively, in descending order by corresponding eigenvalues.

1.5 Interpretation

The estimator in Section 1.2 approximately minimized jointly the unexplained variation and bias

$$\begin{aligned} & \underset{\mathbf{R}, \mathbf{C}, \{\mathbf{F}_t\}_{t=1}^T}{\text{minimize}} \quad (1 + \alpha) \underbrace{\frac{1}{pq} \|\bar{\mathbf{Y}} - \mathbf{R}\bar{\mathbf{F}}\mathbf{C}^T\|_F^2}_{\text{sample bias}} + \underbrace{\frac{1}{pqT} \sum_{t=1}^T \|\mathbf{Y}_t - \mathbf{R}\mathbf{F}_t\mathbf{C}^T\|_F^2}_{\text{sample variance}} \\ & \text{subject to} \quad \frac{1}{p} \mathbf{R}^T \mathbf{R} = \mathbf{I}, \frac{1}{q} \mathbf{C}^T \mathbf{C} = \mathbf{I} \end{aligned} \quad (6)$$

The special case for $\alpha = -1$ corresponds to the least-square estimator. (*not convex*)

Projecting on \mathbf{R} :

$$\begin{aligned} & \underset{\mathbf{R}}{\text{maximize}} \quad Tr \left(\mathbb{E} \left[(1 + \alpha)(\mathbf{R}^T \bar{\mathbf{Y}})(\mathbf{R}^T \bar{\mathbf{Y}})^T + (\mathbf{R}^T \mathbf{Y}_t - \mathbb{E}[\mathbf{R}^T \mathbf{Y}_t])(\mathbf{R}^T \mathbf{Y}_t - \mathbb{E}[\mathbf{R}^T \mathbf{Y}_t])^T \right] \right) \\ & \text{subject to} \quad \frac{1}{p} \mathbf{R}^T \mathbf{R} = \mathbf{I}, \frac{1}{q} \mathbf{C}^T \mathbf{C} = \mathbf{I} \end{aligned} \quad (7)$$

Where $\mathbf{M}_R \triangleq (1 + \alpha)\mathbf{M}_R^{(1)} + \mathbf{M}_R^{(2)}$, $\mathbf{M}_R^{(1)} \triangleq \frac{1}{pq} \mathbb{E}[\bar{\mathbf{Y}}\bar{\mathbf{Y}}^T]$, and $\mathbf{M}_R^{(2)} \triangleq \frac{1}{pq} \mathbb{E}[(\mathbf{Y}_t - [\bar{\mathbf{Y}}])(\mathbf{Y}_t - [\bar{\mathbf{Y}}])^T]$

Then a solution by maximizing row and column variances respectively after projection is considered, projecting on \mathbf{C} is similar. (*convex*)

1.6 Relative estimators

Based on the section 1.4,

$$\widehat{\mathbf{F}}_t = \frac{1}{pq} \widehat{\mathbf{R}}^T \widehat{\mathbf{Y}}_t \widehat{\mathbf{C}}, \text{ and the signal part } \widehat{\mathbf{S}}_t = \frac{1}{pq} \widehat{\mathbf{R}} \widehat{\mathbf{R}}^T \widehat{\mathbf{Y}}_t \widehat{\mathbf{C}} \widehat{\mathbf{C}}^T$$

Dimensions k and r are need to be determined:

1. the eigenvalue ratio-based estimator, proposed by Ahn and Horestein(2013)
2. the Scree plot which is standard in principal component analysis.

Let $\hat{\lambda}_1 \geq \hat{\lambda} \geq \dots \geq \hat{\lambda}_k \geq 0$ be the ordered eigenvalues of $\widehat{\mathbf{M}}_R$. The ratio-based estimator for k is defined as follows:

$$\widehat{k} = \arg \max_{1 \leq j \leq k_{max}} \frac{\widehat{\lambda}_j}{\widehat{\lambda}_{j+1}}$$

where k_{max} is the upper bound, usually taken as $\left\lceil \frac{p}{2} \right\rceil$ or $\left\lceil \frac{p}{3} \right\rceil$, according to Ahn and Horestein(2013), similarly for \widehat{r} with respect to $\widehat{\mathbf{M}}_C$.

1.7 Theoretical Properties

1.7.1 Assumptions

1.8 Simulation

1.9 Application

2 High-Dimensional GLM with Binary Outcomes

2.1 Overall Summary

3 Ultra-High Dimensional GFM¹

3.1 Overall Summary

4 Matrix-variate Logistic Regression with Measurement Error

5 A Likelihood-Based Approach for Multivariate Categorical Response Regression in High Dimensions

6 A likelihood-Based Approach for Semiparametric Regression with Panel Count Data

7 Time Series Latent Gaussian Count

8 Time Series Factor Models(tensor)

¹Generalized Factor Model