

High-Dimensional Vector Autoregressive Time Series Modeling via Tensor Decomposition

Ilia Lomasov

UIUC

11/28/2023

Outline

- 1 Introduction
- 2 Tensor
- 3 MLR VAR
- 4 Low-dim. TS
- 5 High-dim. TS
- 6 Simulation
- 7 Conclusion

Table of Contents

- 1 Introduction
- 2 Tensor
- 3 MLR VAR
- 4 Low-dim. TS
- 5 High-dim. TS
- 6 Simulation
- 7 Conclusion

VAR model (Lütkepohl 2005; Tsay 2010)

Consider the VAR model of the form:

$$y_t = A_1 y_{t-1} + \dots + A_P y_{t-P} + \epsilon_t \quad (1)$$

$$\{y_t\}_{t=1}^T \in \mathbb{R}^N, \quad \{A_j\}_{j=1}^P \in \mathbb{R}^{N \times N}$$

$$\epsilon_t \in \mathbb{R}^N \stackrel{iid}{\sim} (0, \Sigma_\epsilon), \Sigma_\epsilon \succ 0, \Sigma_\epsilon < \inf$$

- Difficult to perform estimation as number of parameters is $N^2 P$ - very large.
- Even in low-dimensional setting: $N = 5, P = 2 \Rightarrow npar = 50$

Alternatives to estimating N^2P parameters:

- ① PCA & Factor Models – in the STAT556 course
- ②
 - **Approach:** Assume sparsity of A_j and apply regularization (e.g., ℓ_1 , LASSO or Dantzig selector)
 - **Papers:** Basu & Michailidis 2015; Han, Lu & Liu 2015 etc.
 - **Drawbacks:**
 - Sacrifices temporal and cross-sectional dependencies
 - Average magnitude of parameters is bounded by $O(N^{-1/2})$, which limits sparsity-inducing regularization

- ③
 - **Approach:** Reduced-rank regression ($A^{(C)}$ – low-rank):

$$y_t = (A_1 \dots A_P) (y'_{t-1} \dots y'_{t-P})' + \epsilon_t =: A^{(C)} x_t + \epsilon_t \quad (2)$$

- **Papers:** Velu & Reinsel 2013; Carriero, Kapetanios, & Marcellino 2011 (Bayesian extension) etc.
- **Limitation:** Allows for only 1 type of low-rank structure

Table of Contents

- 1 Introduction
- 2 Tensor**
- 3 MLR VAR
- 4 Low-dim. TS
- 5 High-dim. TS
- 6 Simulation
- 7 Conclusion

3-Dimensional Tensor $\mathcal{X} \in \mathbb{R}^{l_1 \times l_2 \times l_3}$

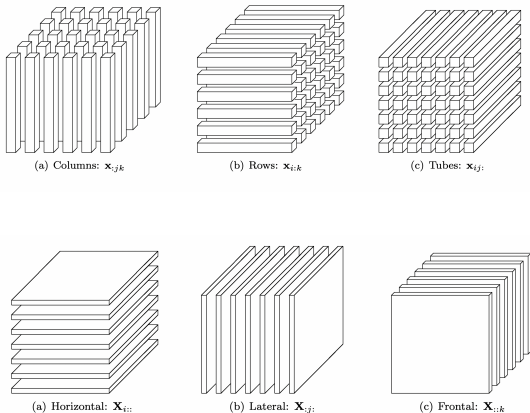


Figure 1: Fibers (top) and Slices (bottom) of \mathcal{X} (Kolda, 2006)

Visualization 1

Let's consider a simple tensor $\mathcal{X} \in \mathbb{R}^{2 \times 3 \times 2}$ defined by $x_{ijk} = 100i + 10j + k$

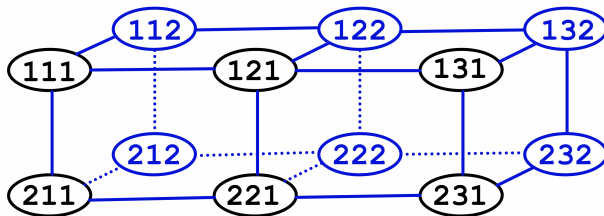


Figure 2: Tensor \mathcal{X}

Matricization

$\{\mathcal{X}_{(n)}\}_{n=1}^3$ – mode- n matricization of \mathcal{X} – rearrangement of mode- n fibers into a matrix. $\mathcal{X}_{(n)} \in \mathbb{R}^{I_n \times (I_1 I_2 I_3 / I_n)}$.

- $\mathcal{X}_{(1)} = (X_{::1} \ X_{::2}) = \begin{pmatrix} 111 & 121 & 131 & 112 & 122 & 132 \\ 211 & 221 & 231 & 212 & 222 & 232 \end{pmatrix}$
- $\mathcal{X}_{(2)} = (X'_{::1} \ X'_{::2}) = \begin{pmatrix} 111 & 211 & 112 & 212 \\ 121 & 221 & 122 & 222 \\ 131 & 231 & 132 & 232 \end{pmatrix}$
- $\mathcal{X}_{(3)} = (\text{vec}(X_{::1}) \ \text{vec}(X_{::2}))' = \begin{pmatrix} 111 & 211 & 121 & 221 & 131 & 231 \\ 112 & 212 & 122 & 222 & 132 & 232 \end{pmatrix}$
- $r_n = \text{rank}_n(\mathcal{X}) = \text{rank}(\mathcal{X}_{(n)}) \leq I_n$

Mode-n multiplication

- Mode-n multiplication:

$$\mathcal{X}_{I_1 \times I_2 \times I_3} \times_n U_{J \times I_n} = \sum_{i_n=1}^{I_n} x_{i_1 i_2 i_3} y_{j i_n}, \quad n \in \{1, 2, 3\}$$

- Let's take $U = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}$, \mathcal{X} defined by $x_{ijk} = 100i + 10j + k$

$$\mathcal{Y}_{2 \times 3 \times 3} = \mathcal{X}_{2 \times 3 \times 2} \times_3 U_{3 \times 2} \Rightarrow y_{ijl} = \sum_{k=1}^3 x_{ijk} u_{lk}$$

- $l \in \{1, 2\} \Rightarrow u_{lk} = \mathbb{1}(l = k) \Rightarrow y_{ijl} = x_{ijk}$
- $l = 3 \Rightarrow u_{lk} = 1 \Rightarrow y_{ij3} = x_{ij1} + x_{ij2}$

Visualization 2

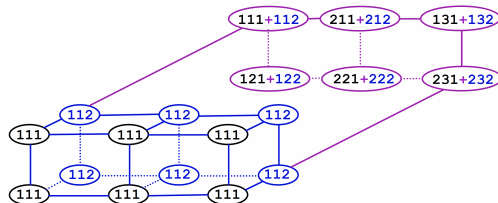


Figure 3: Tensor \mathcal{Y}

$$\begin{aligned}
 \mathcal{Y}_{(3)} &= (\text{vec}(Y_{::1}) \text{vec}(Y_{::2}) \text{vec}(Y_{::3}))' \\
 &= (\text{vec}(X_{::1}) \text{vec}(X_{::2}) \text{vec}(X_{::1}) + \text{vec}(X_{::2}))' \\
 &= U (\text{vec}(X_{::1}) \text{vec}(X_{::2}))' = U \mathcal{X}_{(3)}
 \end{aligned}$$

In general, $\mathcal{Y} = \mathcal{X} \times_n U \Rightarrow \mathcal{Y}_{(n)} = U \mathcal{X}_{(n)}$

Some facts

1 For $m \neq n$, $\mathcal{X} \times_m A \times_n B = \mathcal{X} \times_n B \times_m A$

$$\begin{aligned} \text{Proof: } [\mathcal{X} \times_m A \times_m B]_{i_{-n}, -mjk} &= \sum_{i_n=1}^{I_n} \left[\sum_{i_m=1}^{I_m} x_{i_1 i_2 i_3} a_{ji_m} \right] b_{ki_n} = \\ &= \sum_{i_m=1}^{I_m} \left[\sum_{i_n=1}^{I_n} x_{i_1 i_2 i_3} b_{ki_n} \right] a_{ji_m} = [\mathcal{X} \times_n B \times_m A]_{i_{-n}, -mjk} \end{aligned}$$

2 $\mathcal{X} \times_n A \times_n B = \mathcal{X} \times_n (BA)$

$$\begin{aligned} \text{Proof: } [\mathcal{X} \times_n A \times_n B]_{i_{-n}, k} &= \sum_{j=1}^k \left[\sum_{i_n=1}^{I_n} x_{i_1 i_2 i_3} a_{ji_n} \right] b_{kj} = \\ &= \sum_{i_n=1}^{I_n} x_{i_1 i_2 i_3} \sum_{j=1}^k b_{kj} a_{ji_n} = \sum_{i_n=1}^{I_n} x_{i_1 i_2 i_3} (BA)_{ki_n} = [\mathcal{X} \times_n (BA)]_{i_{-n}, k} \end{aligned}$$

Some facts 2

- 3 If $\text{rank}_n(\mathcal{Y}) < I_n$ for some $n \in \{1, 2, 3\}$, then there exists a Tucker decomposition $\mathcal{Y} = \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3$. (Notation: $\llbracket \mathcal{G}; U_1, U_2, U_3 \rrbracket$)
Here, one option is $\mathcal{G} = \mathcal{X}$, $U_1 = \mathbb{I}_2$, $U_2 = \mathbb{I}_3$, $U_3 = U$
- 4 (Following from 1 & 2): Tucker decomposition is not unique: for any non-singular matrices $\{O_n \in \mathbb{R}^{I_n \times I_n}\}_{n=1}^3$

$$\begin{aligned} \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3 &= \\ &= (\mathcal{G} \times_1 O_1 \times_2 O_2 \times_3 O_3) \times_1 U_1 O_1^{-1} \times_2 U_2 O_2^{-1} \times_3 U_3 O_3^{-1} \end{aligned}$$

- 5 (Kolda, 2006)

$$(\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} = U_1 \mathcal{G}_{(1)} (U_3 \otimes U_2)'$$

Table of Contents

- 1 Introduction
- 2 Tensor
- 3 MLR VAR**
- 4 Low-dim. TS
- 5 High-dim. TS
- 6 Simulation
- 7 Conclusion

Setup

- We can rearrange transition matrices A_1, \dots, A_P into a tensor $\mathcal{A} \in \mathbb{R}^{N \times N \times P}$, assuming those are its frontal slices (as we can see in Figure 4):

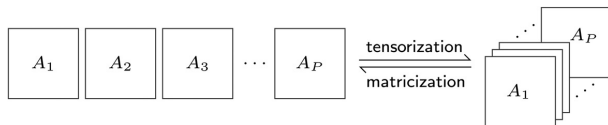


Figure 4: Tensorization from A_1, \dots, A_P to \mathcal{A}

- Recall representation (2): $y_t = (A_1, \dots, A_P) x_t + \epsilon_t$. It's equivalent to $y_t = \mathcal{A}_{(1)} x_t + \epsilon_t$
- Key idea:** assume $\mathcal{A} = \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3$

Some equivalent representations (under certain assumptions)

$$y_t = (\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} x_t + \epsilon_t \quad (3)$$

$$y_t = U_1 \mathcal{G}_{(1)} (U_3 \otimes U_2)' x_t + \epsilon_t = U_1 \mathcal{G}_{(1)} \text{vec}(U_2' X_t U_3) + \epsilon_t \quad (4)$$

Where $X_t = (y_{t-1} \dots y_{t-P})$

Furthermore, consider a special Tucker decomposition – high-order SVD (HOSVD) (De Lathauwer, DeMoor & Vandewalle 2000):

Given $n \in \{1, 2, 3\}$ and $r_n = \text{rank}_n(\mathcal{A}) \leq I_n$, construct $U_n \in \mathbb{R}^{I_n \times r_n}$ as a matrix of top- r_n left singular vectors of $\mathcal{A}_{(n)}$ (i.e., eigenvectors of $\mathcal{A}_{(n)} \mathcal{A}_{(n)}'$). $\mathcal{A}_{(n)} \mathcal{A}_{(n)}'$ – symmetric $\Rightarrow U_n' U_n = \mathbb{I}_{I_n}$, $U_n U_n' = \mathbb{I}_{r_n}$

Number of parameters comparison

Model	Number of parameters
VAR	N^2P
RRR	$(NP + N - r_1)r_1$ $(r_1 \text{ independent rows with } NP \text{ elements each plus } N - r_1 \text{ dependent rows, with } r_1 \text{ dependency coefficients each})$
MLR	$r_1 r_2 r_3 + (N - r_1)r_1 + (N - r_2)r_2 + (P - r_3)r_3$

* – Since \mathcal{A} is to be estimated, r_n has to be assumed?

Connection with Factor Models

- Recall unknown factor model with r_1 common factors: $Y = F\Lambda' + E$
Where $Y = (y_1, \dots, y_T)'$, $X = (x_1, \dots, x_T)'$, $E = (e_1, \dots, e_T)'$,
 $F'F/T = \mathbb{I}_{r_1}$, $\Lambda'\Lambda \in \mathbb{R}^{r_1 \times r_1}$ – full-rank and diagonal.
- Rewrite $y_t = U_1 \mathcal{G}_{(1)} (U_3 \otimes U_2)' x_t + \epsilon_t$ in matrix form:

$$Y = X (U_3 \otimes U_2) \mathcal{G}_{(1)}' U_1' + E \quad (5)$$

- Consider SVD: $(X (U_3 \otimes U_2) \mathcal{G}_{(1)}')$ = $U_x D_x V_x'$, where $D_x \in \mathbb{R}^{r_1 \times r_1}$ is diagonal, and U_x and V_x are orthonormal.
- Define $F = \sqrt{T} U_x$ and $\Lambda = U_1 V_x D_x / \sqrt{T}$. Note that $F'F/T = \mathbb{I}_{r_1}$ and that $\Lambda'\Lambda$ is diagonal.
- Thus, $Y = X (U_3 \otimes U_2) \mathcal{G}_{(1)}' U_1' + E = F\Lambda' + E$

Key difference: MLR can be used directly for predictions

Table of Contents

- 1 Introduction
- 2 Tensor
- 3 MLR VAR
- 4 Low-dim. TS**
- 5 High-dim. TS
- 6 Simulation
- 7 Conclusion

$\widehat{\mathcal{A}}_{\text{MLR}}, \widehat{A}_{\text{RRR}}, \widehat{A}_{\text{OLS}}$

$$\widehat{\mathcal{A}}_{\text{MLR}} = \widehat{\mathcal{G}} \times_1 \widehat{U}_1 \times_2 \widehat{U}_2 \times_3 \widehat{U}_3 = \arg \min L(\mathcal{G}, U_1, U_2, U_3), \text{ where}$$

$$L(\mathcal{G}, U_1, U_2, U_3) = \frac{1}{T} \sum_{t=1}^T \left\| y_t - (\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} x_t \right\|_2^2$$

$$\left(\widehat{\mathcal{A}}_{\text{OLS}} \right)_{(1)} = \widehat{A}_{\text{OLS}} = \arg \min_{B \in \mathbb{R}^{N \times NP}} \sum_{t=1}^T \|y_t - Bx_t\|_2^2$$

$$\left(\widehat{\mathcal{A}}_{\text{RRR}} \right)_{(1)} = \widehat{A}_{\text{RRR}} = \arg \min_{B \in \mathbb{R}^{N \times NP}, \text{rank}(B) \leq r_1} \sum_{t=1}^T \|y_t - Bx_t\|_2^2$$

Asymptotic Properties of $\widehat{\mathcal{A}}_{\text{MLR}}, \widehat{\mathcal{A}}_{\text{RRR}}, \widehat{\mathcal{A}}_{\text{OLS}}$

- Assume true (r_1, r_2, r_3) are known, N, P – fixed (**low-dim. setup**)
- Also assume $\mathbb{E} \|\epsilon_t\|_2^4 < \infty$, and that all roots of the matrix polynomial $A(z) = \mathbb{I}_N - A_1 z - \dots - A_P z^P, z \in \mathbb{C}$ lie outside unit circle.

Then for method $\in \{\text{"MLR"}, \text{"RRR"}, \text{"OLS"}\}$

$$\sqrt{T} \left\{ \text{vec} \left(\left(\widehat{\mathcal{A}}_{\text{method}} \right)_{(1)} \right) - \text{vec} \left(\mathcal{A}_{(1)} \right) \right\} \xrightarrow[T \rightarrow \infty]{D} N(0, \Sigma_{\text{method}})$$

Where Σ_{MLR} is a function of $\mathcal{G}, U_1, U_2, U_3, \Sigma_{\text{OLS}}, \Sigma_{\text{RRR}}$ - of A_1, \dots, A_P

Moreover, $\Sigma_{\text{MLR}} \preceq \Sigma_{\text{RRR}} \preceq \Sigma_{\text{OLS}}$

Alternating Least Squares Estimation

- $L(\mathcal{G}, U_1, U_2, U_3) = \frac{1}{T} \sum_{t=1}^T \left\| y_t - (\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} x_t \right\|_2^2$
- L – convex w.r.t any of \mathcal{G} , U_1 , U_2 , and U_3 when the other three are fixed
- Hence, an ALS algorithm can be implemented. **Idea:**
 - 1 Initialize $\mathcal{A}^{(0)}$
 - 2 Perform HOSVD to obtain $U_1^{(0)}, U_2^{(0)}, U_3^{(0)}, \mathcal{G}^{(0)}$
 - 3 Update individually $U_1^{(k+1)}, U_2^{(k+1)}, U_3^{(k+1)}, \mathcal{G}^{(k+1)}$ (in that order), other 3 fixed
 - 4 When convergence reached, obtain $\widehat{\mathcal{A}}$
- Authors recommend to initialize $\mathcal{A}^{(0)} = \widehat{\mathcal{A}}_{\text{prelim}} + T^{-1/2} \mathcal{T}$, where $\widehat{\mathcal{A}}_{\text{prelim}}$ is $\widehat{\mathcal{A}}_{\text{OLS}}$ for large T , $\widehat{\mathcal{A}}_{\text{RRR}}$ for small T , and $\text{vec}(\mathcal{T}) \sim N(0, \mathbb{I}_{NMP})$. Global minimum is not guaranteed

ALS update equations

$$U_1^{(k+1)} \leftarrow \arg \min_{U_1} \sum_{t=1}^T \|y_t - \left((X_t' (U_3^{(k)} \otimes U_2^{(k)}) \mathcal{G}_{(1)}^{(k)'} \otimes I_N) \text{vec}(U_1) \right) \|_2^2$$

$$U_2^{(k+1)} \leftarrow \arg \min_{U_2} \sum_{t=1}^T \|y_t - U_1^{(k+1)} \mathcal{G}_{(1)}^{(k)} \left((X_t U_3^{(k)})' \otimes I_{r_2} \right) \text{vec}(U_2') \|_2^2$$

$$U_3^{(k+1)} \leftarrow \arg \min_{U_3} \sum_{t=1}^T \|y_t - U_1^{(k+1)} \mathcal{G}_{(1)}^{(k)} \left(I_{r_3} \otimes (U_2^{(k+1)'} X_t) \right) \text{vec}(U_3) \|_2^2$$

$$\mathcal{G}^{(k+1)} \leftarrow \arg \min_{\mathcal{G}} \sum_{t=1}^T \|y_t - \left(\left((U_3^{(k+1)} \otimes U_2^{(k+1)})' x_t \right)' \otimes U_1^{(k+1)} \right) \text{vec}(\mathcal{G}_{(1)}) \|_2^2$$

- **Remark:** Let $h(U_1, U_2, U_3, \mathcal{G}) = \text{vec} \left((\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} x_t \right) = \text{vec} \left(U_1 \mathcal{G}_{(1)} (U_3 \otimes U_2)' x_t \right)$. Consider $\partial h / \partial \text{vec}(U_i)$, $\partial h / \partial \text{vec}(\mathcal{G})$.

Table of Contents

- 1 Introduction
- 2 Tensor
- 3 MLR VAR
- 4 Low-dim. TS
- 5 High-dim. TS**
- 6 Simulation
- 7 Conclusion

Sparse Higher-Order Reduced-Rank VAR (SHORR)

- Same $L = \frac{1}{T} \sum_{t=1}^T \left\| y_t - (\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} \mathbf{x}_t \right\|_2^2$
- Introduce regularization and all-orthogonality constraint:

$$\begin{aligned} \widehat{\mathcal{A}}_{\text{SHORR}} \equiv [\widehat{\mathcal{G}}; \widehat{U}_1, \widehat{U}_2, \widehat{U}_3] = & \arg \min_{\mathcal{G}, U_1, U_2, U_3} \{L(\mathcal{G}, U_1, U_2, U_3) \\ & + \lambda \|U_3 \otimes U_2 \otimes U_1\|_1\} \quad \text{subject to } U_i' U_i = \mathbb{I}_{r_i} \text{ and} \\ & \mathcal{G} \in \left\{ \mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times r_3} : (\mathcal{G}_{(i)})_{i=1}^3 \text{ -- row-orthogonal} \right\} \end{aligned}$$

- Sparsity assumption: each column of U_i has at most s_i nonzero entries
- Under these and certain extra assumptions, non-asymptotic UB's for $\left\| \widehat{\mathcal{A}}_{\text{SHORR}} - \mathcal{A} \right\|_{\text{F}}$ and $T^{-1} \sum_{t=1}^T \left\| \left(\widehat{\mathcal{A}}_{\text{SHORR}} - \mathcal{A} \right)_{(1)} \mathbf{x}_t \right\|_2^2$ were derived by the authors
- Difference with LASSO: \mathcal{A} – not necessarily sparse

Developing an algorithm

- **Issue:** ℓ_1 regularization – non-smooth, orthogonality constraint – non-convex
- **Solution:** alternating direction method of multipliers (ADMM) algorithm (Boyd et al. 2011)
- **Idea:** assume a decomposition $\mathcal{G}_{(i)} = D_i V_i'$ exists, where $D_i \in \mathbb{R}^{r_i \times r_i}$, $V_i \in \mathbb{R}^{(r_1 r_2 r_3 / r_i) \times r_i}$, $V_i' V_i = \mathbb{I}_{r_i}$
- Augmented Lagrangian:

$$\begin{aligned} \mathcal{L}_\varrho(\mathcal{G}, \{U_i\}, \{D_i\}, \{V_i\}; \{\mathcal{C}_i\}) &= L(\mathcal{G}, U_1, U_2, U_3) + \lambda \|U_3 \otimes U_2 \otimes U_1\|_1 \\ &+ 2 \sum_{i=1}^3 \varrho_i \left\langle (\mathcal{C}_i)_{(i)}, \mathcal{G}_{(i)} - D_i V_i' \right\rangle + \sum_{i=1}^3 \varrho_i \|\mathcal{G}_{(i)} - D_i V_i'\|_F^2 \end{aligned}$$

Where ϱ_i – regularization constants, $\mathcal{C}_i \in \mathbb{R}^{r_1 \times r_2 \times r_3}$ – dual variables.

ADMM algorithm for SHORR

```

1: Initialize:  $\mathcal{A}^{(0)}$ 
2: HOSVD:  $\mathcal{A}^{(0)} \approx \mathcal{G}^{(0)} \times_1 \mathbf{U}_1^{(0)} \times_2 \mathbf{U}_2^{(0)} \times_3 \mathbf{U}_3^{(0)}$  with multilinear ranks  $(r_1, r_2, r_3)$ .
3: repeat
4:    $\mathbf{U}_1^{(k+1)} \leftarrow \arg \min_{\mathbf{U}_1' \mathbf{U}_1 = \mathbf{I}_{r_1}} \left\{ L(\mathcal{G}^{(k)}, \mathbf{U}_1, \mathbf{U}_2^{(k)}, \mathbf{U}_3^{(k)}) + \lambda \|\mathbf{U}_1\|_1 \|\mathbf{U}_2^{(k)}\|_1 \|\mathbf{U}_3^{(k)}\|_1 \right\}$ 
5:    $\mathbf{U}_2^{(k+1)} \leftarrow \arg \min_{\mathbf{U}_2' \mathbf{U}_2 = \mathbf{I}_{r_2}} \left\{ L(\mathcal{G}^{(k)}, \mathbf{U}_1^{(k+1)}, \mathbf{U}_2, \mathbf{U}_3^{(k)}) + \lambda \|\mathbf{U}_1^{(k+1)}\|_1 \|\mathbf{U}_2\|_1 \|\mathbf{U}_3^{(k)}\|_1 \right\}$ 
6:    $\mathbf{U}_3^{(k+1)} \leftarrow \arg \min_{\mathbf{U}_3' \mathbf{U}_3 = \mathbf{I}_{r_3}} \left\{ L(\mathcal{G}^{(k)}, \mathbf{U}_1^{(k+1)}, \mathbf{U}_2^{(k+1)}, \mathbf{U}_3) + \lambda \|\mathbf{U}_1^{(k+1)}\|_1 \|\mathbf{U}_2^{(k+1)}\|_1 \|\mathbf{U}_3\|_1 \right\}$ 
7:    $\mathcal{G}^{(k+1)} \leftarrow \arg \min \left\{ L(\mathcal{G}, \mathbf{U}_1^{(k+1)}, \mathbf{U}_2^{(k+1)}, \mathbf{U}_3^{(k+1)}) + \sum_{i=1}^3 \varrho_i \|\mathcal{G}_{(i)} - \mathbf{D}_i^{(k)} \mathbf{V}_i^{(k)'} + (\mathcal{C}_i^{(k)})_{(i)}\|_F^2 \right\}$ 
8:   for  $i \in \{1, 2, 3\}$  do
9:      $\mathbf{D}_i^{(k+1)} \leftarrow \arg \min_{\mathbf{D}_i = \text{diag}(\mathbf{d}_i)} \|\mathcal{G}_{(i)}^{(k+1)} - \mathbf{D}_i \mathbf{V}_i^{(k)'} + (\mathcal{C}_i^{(k)})_{(i)}\|_F^2$ 
10:     $\mathbf{V}_i^{(k+1)} \leftarrow \arg \min_{\mathbf{V}_i' \mathbf{V}_i = \mathbf{I}_{r_i}} \|\mathcal{G}_{(i)}^{(k+1)} - \mathbf{D}_i^{(k+1)} \mathbf{V}_i' + (\mathcal{C}_i^{(k)})_{(i)}\|_F^2$ 
11:     $(\mathcal{C}_i^{(k+1)})_{(i)} \leftarrow (\mathcal{C}_i^{(k)})_{(i)} + \mathcal{G}_{(i)}^{(k+1)} - \mathbf{D}_i^{(k+1)} \mathbf{V}_i^{(k+1)'}$ 
12:   end for
13:    $\mathcal{A}^{(k+1)} \leftarrow \mathcal{G}^{(k+1)} \times_1 \mathbf{U}_1^{(k+1)} \times_2 \mathbf{U}_2^{(k+1)} \times_3 \mathbf{U}_3^{(k+1)}$ 
14: until convergence

```

Updating \mathcal{G} , \mathbf{D}_i , \mathbf{V}_i – LS problem, updating \mathbf{U}_i – very complicated

Updating U_i

- Original problem:

$$U_i = \arg \min_B \{n^{-1} \|y - X \text{vec}(B)\|_2^2 + \lambda \|B\|_1\} \quad \text{s.t.} \quad B'B = \mathbb{I}$$

- Idea: separate orthogonality and regularization

$$\min_B \{n^{-1} \|y - X \text{vec}(B)\|_2^2 + \lambda \|W\|_1\} \quad \text{s.t.} \quad B'B = \mathbb{I}, B = W$$

- Augmented Lagrangian (M – dual variable)

$$n^{-1} \|y - X \text{vec}(B)\|_2^2 + \lambda \|W\|_1 + 2\kappa \langle M, B - W \rangle + \kappa \|B - W\|_F^2$$

- Apply ADMM to find $B = W = U_i$:

-
- 1: Initialize: $B^{(0)} = W^{(0)}, M^{(0)} = \mathbf{0}$
 - 2: **repeat**
 - 3: $B^{(k+1)} \leftarrow \arg \min_{B'B=I} \{n^{-1} \|y - X \text{vec}(B)\|_2^2 + \kappa \|B - W^{(k)} + M^{(k)}\|_F^2\}$
 - 4: $W^{(k+1)} \leftarrow \arg \min_W \{\kappa \|B^{(k+1)} - W + M^{(k)}\|_F^2 + \lambda \|W\|_1\}$
 - 5: $M^{(k+1)} \leftarrow M^{(k)} + B^{(k+1)} - W^{(k+1)}$
 - 6: **until convergence**
-

Convergence and initialization

- Under certain conditions on \mathcal{L}_θ , algorithm converges to local minimum of our objective function:

$$L(\mathcal{G}, U_1, U_2, U_3) + \lambda \|U_3 \otimes U_2 \otimes U_1\|_1$$

- Authors recommend to choose $\mathcal{A}^{(0)} = \widehat{\mathcal{A}}_{\text{NN}} + (NP/T)^{1/2}\mathcal{T}$, where:
 - $\text{vec}(\mathcal{T}) \sim N(0, \mathbb{I}_{NPN})$, $\|\mathcal{T}\|_F = O_p(1)$
 - $\widehat{\mathcal{A}}_{\text{NN}} = \arg \min \frac{1}{T} \sum_{t=1}^T \|y_t - \mathcal{A}_{(1)}x_t\|_2^2 + \lambda \|\mathcal{A}_{(1)}\|_*$
 - $\|\mathcal{A}_{(1)}\|_*$ – nuclear norm, or sum of all singular values of $\mathcal{A}_{(1)}$

Rank selection

- Let $\widehat{\mathcal{A}}$ be a consistent initial estimator of \mathcal{A} (e.g., $\widehat{\mathcal{A}}_{\text{NN}}$)
- Ridge-type ratio estimator (Xia, Xu, and Zhu 2015):

$$\widehat{r}_i = \arg \min_{1 \leq j \leq p_i - 1} \frac{\sigma_{j+1}(\widehat{\mathcal{A}}_{(i)}) + c}{\sigma_j(\widehat{\mathcal{A}}_{(i)}) + c} \quad \text{where } p_1 = p_2 = N, p_3 = P$$

- Denote $\zeta_i = \frac{1}{\sigma_{r_i}(\mathcal{A}_{(i)})} \cdot \max_{1 \leq j < r_i} \frac{\sigma_j(\mathcal{A}_{(i)})}{\sigma_{j+1}(\mathcal{A}_{(i)})}$
- $c > 0$ is chosen such that:
 - 1 $\|\widehat{\mathcal{A}} - \mathcal{A}\|_F = o_p(c)$
 - 2 $\max_{1 \leq i \leq 3} \zeta_i = o(1/c)$
- Authors recommend $c = \sqrt{NP \ln(T)/(10T)}$

Table of Contents

- 1 Introduction
- 2 Tensor
- 3 MLR VAR
- 4 Low-dim. TS
- 5 High-dim. TS
- 6 Simulation**
- 7 Conclusion

Rank selection consistency – simulation setup

- $(N, P) = (10, 5)$, $(r_1, r_2, r_3) = (3, 3, 3)$, and $\epsilon_t \stackrel{\text{iid}}{\sim} N(0, \mathbb{I}_N)$
- \mathcal{G} – a diagonal cube with $(\mathcal{G}_{111}, \mathcal{G}_{222}, \mathcal{G}_{333}) = (2, 2, 2)$ (case a), $(4, 3, 2)$ (case b), $(1, 1, 1)$ (case c), or $(2, 1, 0.5)$ (case d).
- Then nonzero singular values of $\mathcal{A}_{(i)}$ are \mathcal{G}_{111} , \mathcal{G}_{222} , and \mathcal{G}_{333}
- Generate U_i 's as the first r_i left singular vectors of Gaussian random matrices while ensuring the stationarity.
- $c = \sqrt{NP \ln(T)/(10T)}$ was used
- 1000 replications for each $T \in \{50, 100, \dots, 400\}$

Rank selection consistency – simulation results

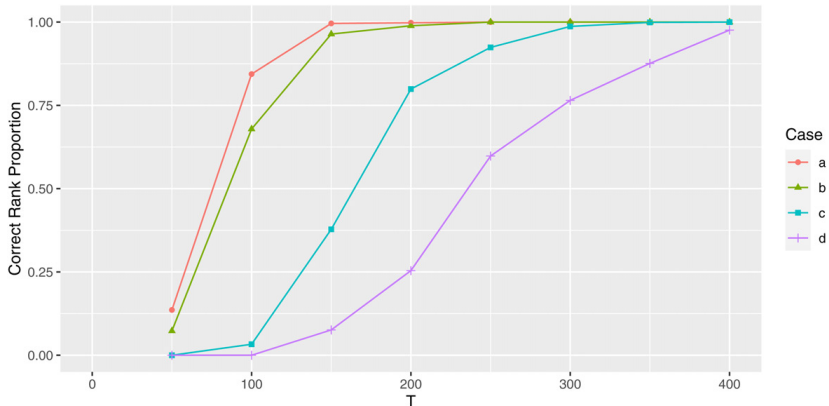


Figure 5: Proportion of correct rank selection when the nonzero singular values of each $\mathcal{A}_{(i)}$ are (2, 2, 2) (case a), (4, 3, 2) (case b), (1, 1, 1) (case c), or (2, 1, 0.5) (case d)

OLS vs. RRR vs. MLR – setup

- N, P, U_i – same. Number of replications – same
- $r_1 = r_2 = 3$, and $r_3 \in \{2, 3, 4\}$
- Generate \mathcal{G} by scaling a random iid Gaussian tensor s.t.
$$\min_{1 \leq i \leq 3} \sigma_{r_i}(\mathcal{G}_{(i)}) = 1$$

OLS vs. RRR vs. MLR – results

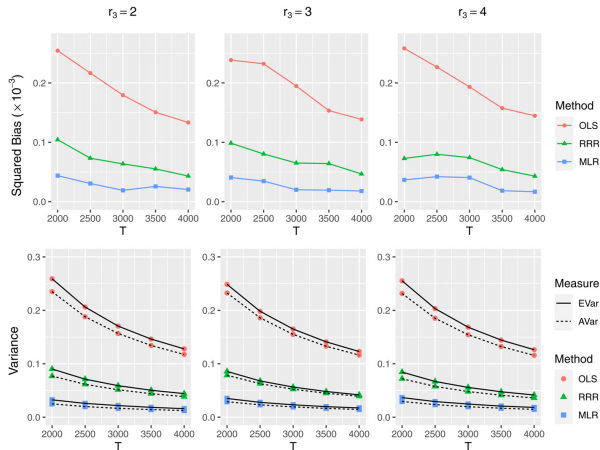


Figure 6: Squared bias, empirical variance (EVar) and asymptotic variance (AVar) for $\widehat{\mathcal{A}}_{\text{OLS}}$, $\widehat{\mathcal{A}}_{\text{RRR}}$, and $\widehat{\mathcal{A}}_{\text{MLR}}$ under various multilinear ranks.

Comparison with existing methods – setup

- $(N, P) = (10, 5)$ (case a), $(15, 8)$ (case b)
- $(r_1, r_2, r_3) = (3, 3, 3)$, $(s_1, s_2, s_3) = (3, 3, 2)$
- For case a, \mathcal{G} and U_i 's are generated by the same methods as in RRR vs. MLR vs. OLS
- For case b, zeros rows are added below the U_i 's in case a
- In both cases, $\|\mathcal{A}\|_0 = 500$. Hence, \mathcal{A} is not sparse in case a, but is sparse in case b

Comparison with existing methods – results

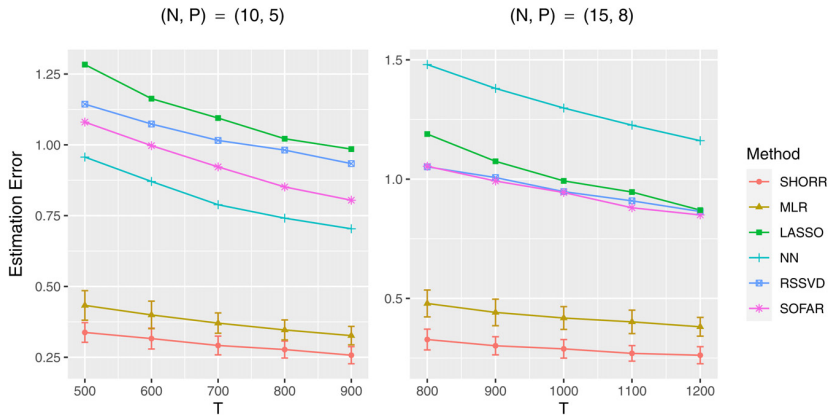


Figure 7: Plots of the estimation error $\|\widehat{\mathcal{A}} - \mathcal{A}\|_F$ against T for six estimation methods under two settings of (N, P) .

- **Issue:** NN performs the worst in the sparse case

Modeling real data

- **Data:** $N = 40$ quarterly macroeconomic sequences of the United States from 1959 to 2007 (from Koop, 2013)
- Lag $P = 4$ for the VAR model is suggested by Koop (2013).
 $N \gg P \Rightarrow$ penalty on U_3 is not needed. New penalty $-\|U_2 \otimes U_1\|_1$
- $(r_1, r_2, r_3) = (4, 3, 2)$ are selected by the ridge-type ratio estimator
- Tuning parameter λ is selected by BIC

Criterion	Unregularized methods				Regularized methods				
	OLS	RRR	DFM	MLR	SHORR	LASSO	NN	RSSVD	SOFAR
ℓ_2 norm	20.16	13.31	6.36	5.81	5.35	6.72	8.16	6.33	6.28
ℓ_∞ norm	8.32	4.55	2.85	2.56	2.44	3.06	3.36	3.02	3.02

Figure 8: Forecasting errors for different methods

- Again, NN performs the worst

Table of Contents

- 1 Introduction
- 2 Tensor
- 3 MLR VAR
- 4 Low-dim. TS
- 5 High-dim. TS
- 6 Simulation
- 7 Conclusion**

Conclusions, issues and improvements

- The novelty of the approach is in its ability to jointly enforce three different reduced-rank structures at the same time
- **order P of VAR is not estimated.** Possible solution: IC-based selection
- **Selecting r_i is dependent on initialization $\widehat{\mathcal{A}}_0$, derived from other methods and which can even be consistent but biased/inefficient and hence make low-T estimation incorrect.**
IC-based selection or hypothesis testing – problematic (3 parameters, too many combinations)
- **NN estimator perform the worst in both simulations and real data**
Possible solution: use other estimators at initialization $\widehat{\mathcal{A}}_0$ (e.g., SOFAR or RSSVD)
- Despite those limitations, MLR and SHORR perform the best on real macroeconomic data