



POLYTECHNIQUE
MONTRÉAL

UNIVERSITÉ
D'INGÉNIERIE



Spatiotemporal Traffic Data Imputation and Forecasting with Tensor Learning

Ph.D. Research Project

Xinyu Chen (Ph.D. candidate)

Polytechnique Montréal

May 24, 2022



Supervisor
Prof. Nicolas Saunier
Polytechnique Montréal



Co-supervisor
Prof. Lijun Sun
McGill University

Committee members:

- Prof. Francesco Ciari (Polytechnique Montréal)
- Prof. Nicolas Saunier (Polytechnique Montréal)
- Prof. Lijun Sun (McGill University)
- Prof. Guillaume Rabusseau (University of Montréal)

Sources:

- ➊ GitHub repository: <https://github.com/xinychen/transdim>
- ➋ Data: <https://transdim.github.io/data>

Outline

- **Motivation**

- Multivariate traffic time series
- Multidimensional traffic time series
- Multiple data behaviors

- **Literature Review**

- Spatiotemporal traffic data imputation
- Spatiotemporal traffic forecasting
- Low-rank tensor learning

- **Objectives**

- **Methodology**

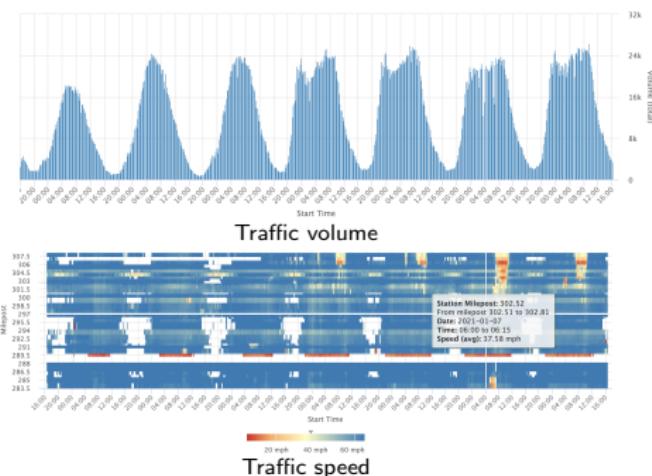
- Spatiotemporal traffic data imputation
- High-dimensional traffic forecasting
- Multidimensional traffic forecasting
- Traffic forecasting on sparse data

- **Conclusion**

Multivariate Traffic Time Series

Many spatiotemporal traffic time series data are in the form of **matrix**.

- Example: Portland highway traffic data¹.



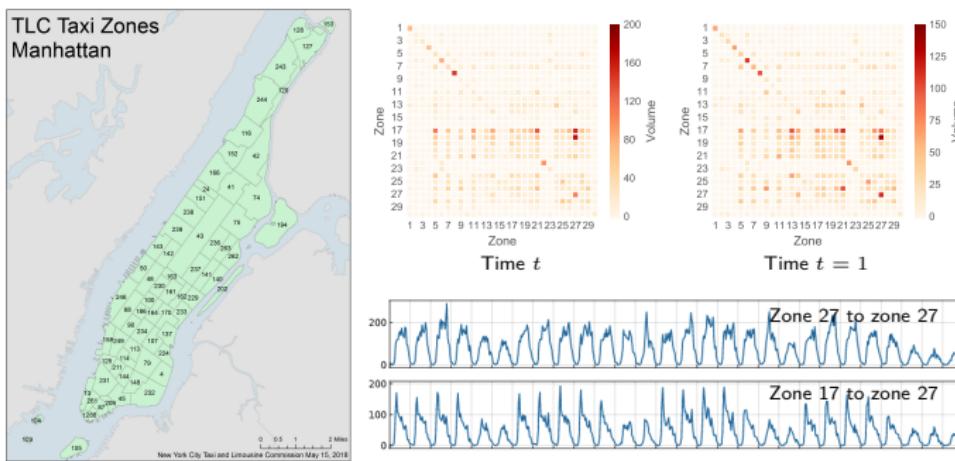
- $X \in \mathbb{R}^{N \times T}$ with N spatial locations $\times T$ time steps

¹ <https://portal.its.pdx.edu/home>

Multidimensional Traffic Time Series

Many spatiotemporal traffic time series data are in the form of **tensor**.

- Example: NYC (hourly) taxi flow data².



- $\mathcal{X} \in \mathbb{R}^{M \times N \times T}$ with M zones $\times N$ zones $\times T$ time steps

²<https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

Multiple Data Behaviors

Spatiotemporal traffic data are time series, but they involve multiple data behaviors.

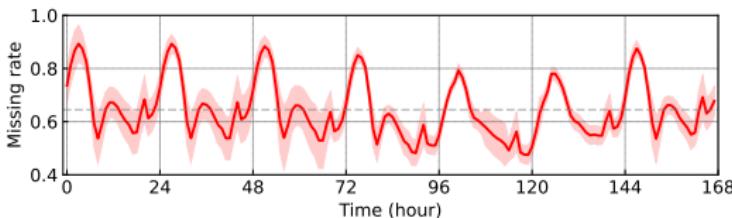
- Incompleteness & sparsity
 - High-dimensionality
 - Multidimensionality
 - Noises & outliers
 - Nonstationarity
 - ...

In addition, spatiotemporal correlations are also very important.

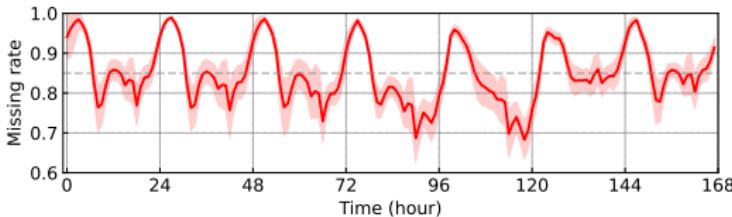
Multiple Data Behaviors

Sparsity & high dimensionality

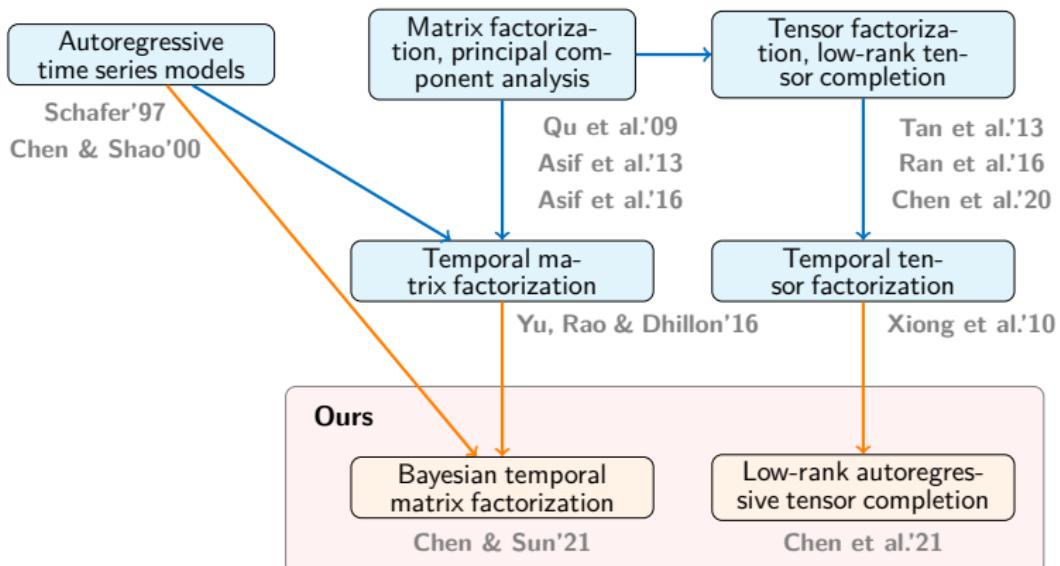
- **NYC** movement speed data (2019)
 - **98,210** road segments & 8,760 time steps (hours)
 - Overall missing rate: **64.43%**



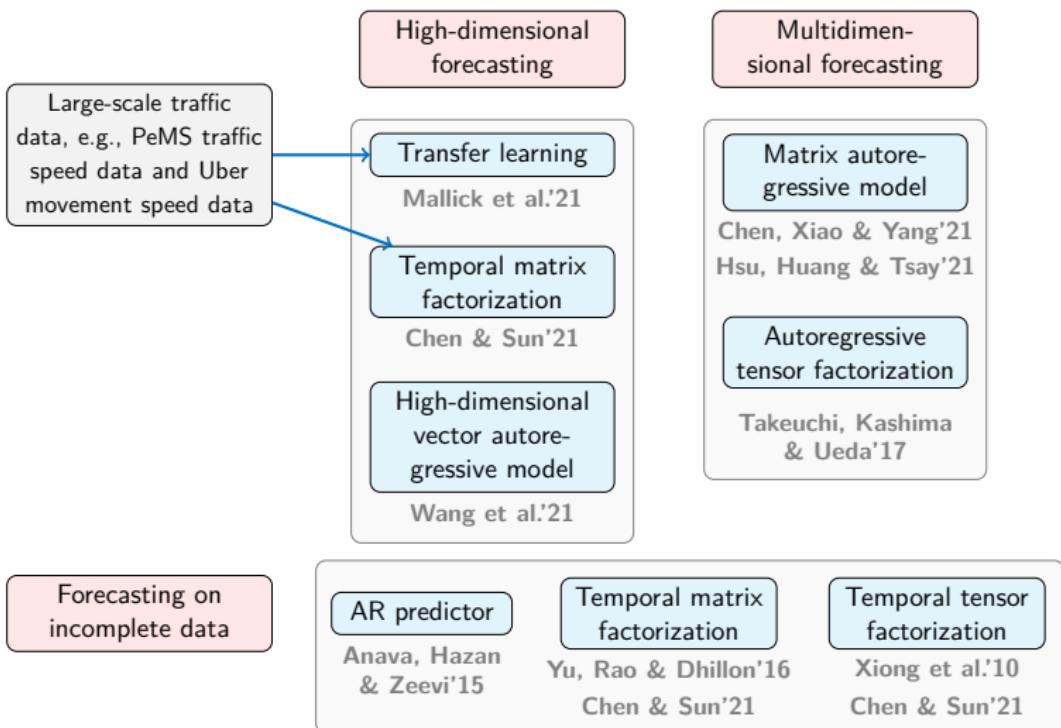
- Seattle movement speed data (2019)
 - 63,490 road segments & 8,760 time steps (hours)
 - Overall missing rate: 84.95%



Spatiotemporal Traffic Data Imputation



Spatiotemporal Traffic Forecasting

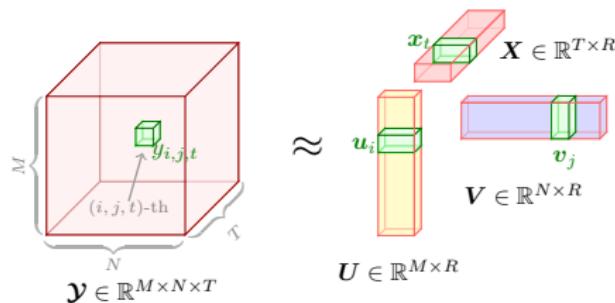


Low-Rank Tensor Learning

From data to model

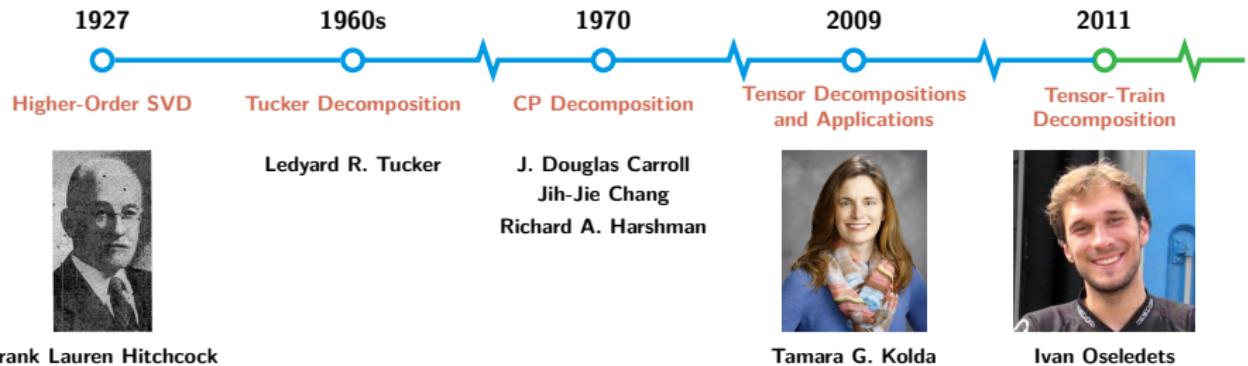
- #### ■ Matrix factorization:

$$\min_{\mathbf{W}, \mathbf{X}} \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Y}) - \mathcal{P}_\Omega(\mathbf{W}\mathbf{X})\|_F^2 + \frac{\lambda}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{X}\|_F^2)$$



- Tensor factorization:

$$\min_{\boldsymbol{U}, \boldsymbol{V}, \boldsymbol{X}} \frac{1}{2} \|\mathcal{P}_\Omega(\boldsymbol{\mathcal{Y}}) - \mathcal{P}_\Omega\left(\sum_{r=1}^R \boldsymbol{u}_r \circ \boldsymbol{v}_r \circ \boldsymbol{x}_r\right)\|_F^2 + \frac{\lambda}{2} (\|\boldsymbol{U}\|_F^2 + \|\boldsymbol{V}\|_F^2 + \|\boldsymbol{X}\|_F^2)$$



Low-Rank Tensor Learning

- Low-rank matrix/tensor completion

Candès & Recht'09: Convex nuclear norm minimization for matrix completion.

$$\begin{aligned} \min_{\mathbf{X}} \quad & \|\mathbf{X}\|_* \\ \text{s.t. } & \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{Y}) \end{aligned}$$

Cai, Candès & Shen'10: Singular value thresholding algorithm.

$$\begin{cases} \mathbf{X}^\ell = \mathcal{D}_\tau(\mathbf{Z}^{\ell-1}) \\ \mathbf{Z}^\ell = \mathbf{Z}^{\ell-1} + \delta_\ell \mathcal{P}_\Omega(\mathbf{Y} - \mathbf{X}^\ell) \end{cases}$$

Zhang et al.'12: Nonconvex truncated nuclear norm minimization.

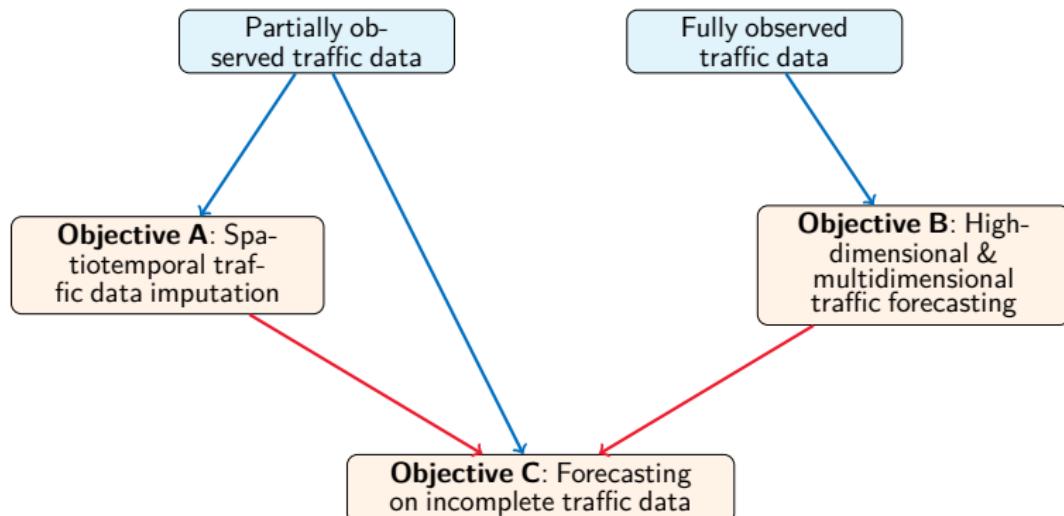
Liu et al.'13: Convex nuclear norm minimization for tensor completion.

$$\begin{aligned} \min_{\mathcal{X}} \quad & \| \mathcal{X} \|_* \\ \text{s.t. } & \mathcal{P}_\Omega(\mathcal{X}) = \mathcal{P}_\Omega(\mathcal{Y}) \end{aligned}$$

Lu, Peng & Wei'19: Tensor nuclear norm induced by linear transform.

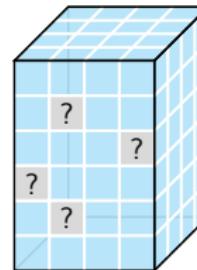
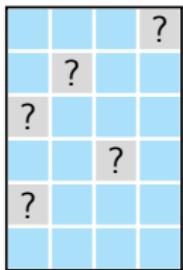
Objectives: A Whole Structure

We are working on **spatiotemporal traffic data modeling**.



Spatiotemporal Traffic Data Imputation

- **Objective A:** Given a multivariate time series data like $\mathbf{Y} \in \mathbb{R}^{N \times T}$ or a multidimensional time series data like $\mathcal{Y} \in \mathbb{R}^{M \times N \times T}$, impute the missing values of the data.



[Q]

- How to learn and reconstruct missing values from observed data?
 - How to make use of spatiotemporal correlations?
 - How to make use of traffic time series dynamics?

Spatiotemporal Traffic Data Imputation

Low-rank matrix completion (Candès & Recht'09)

For any partially observed data matrix $\mathbf{Y} \in \mathbb{R}^{N \times T}$ with observed index set Ω , then low-rank matrix completion takes the form of

$$\begin{aligned} & \min_{\mathbf{X}} \|\mathbf{X}\|_* \\ & \text{s.t. } \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{Y}). \end{aligned} \tag{1}$$

Low-rank tensor completion (Liu et al.'13)

For any partially observed data matrix $\mathcal{Y} \in \mathbb{R}^{M \times N \times T}$ with observed index set Ω , then low-rank matrix completion takes the form of

$$\begin{aligned} \min_{\boldsymbol{\chi}} \quad & \| \boldsymbol{\chi} \|_* \\ \text{s.t. } & \mathcal{P}_\Omega(\boldsymbol{\chi}) = \mathcal{P}_\Omega(\boldsymbol{y}). \end{aligned} \tag{2}$$

- **Limitation:** Only cover the global consistency.
 - **Comment:** For modeling spatiotemporal traffic data, local consistency (e.g., temporal correlations) is also important.

Spatiotemporal Traffic Data Imputation

Low-rank autoregressive matrix completion

Given the multivariate traffic time series data $\mathbf{Y} \in \mathbb{R}^{N \times T}$ with the observed index set Ω , the low-rank matrix completion combines both nuclear norm and univariate autoregressive process:

$$\begin{aligned} & \min_{\mathbf{X}} \|\mathbf{X}\|_* + \frac{\lambda}{2} \sum_{n=1}^N \sum_{t=d+1}^T (z_{n,t} - \sum_{k=1}^d a_{n,k} z_{n,t-k})^2 \\ & \text{s.t. } \begin{cases} \mathbf{X} = \mathbf{Z}, \\ \mathcal{P}_\Omega(\mathbf{Z}) = \mathcal{P}_\Omega(\mathbf{Y}). \end{cases} \end{aligned} \tag{3}$$

Low-rank autoregressive tensor completion

If T time steps can be separated into I time steps per day and J days, i.e., $T = IJ$, then we can define an operator $\mathcal{Q}(\cdot)$ to convert multivariate data to third-order tensor:

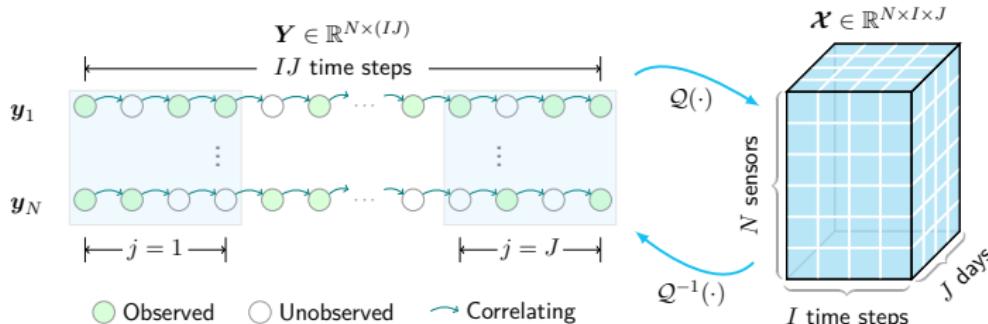
$$\begin{aligned} \min_{\boldsymbol{\mathcal{X}}} \quad & \| \boldsymbol{\mathcal{X}} \|_* + \frac{\lambda}{2} \sum_{n=1}^N \sum_{t=d+1}^T (z_{n,t} - \sum_{k=1}^d a_{n,k} z_{n,t-k})^2 \\ \text{s.t.} \quad & \begin{cases} \boldsymbol{\mathcal{X}} = \mathcal{Q}(\boldsymbol{Z}), \\ \mathcal{P}_{\Omega}(\boldsymbol{Z}) = \mathcal{P}_{\Omega}(\boldsymbol{Y}). \end{cases} \end{aligned} \tag{4}$$

Spatiotemporal Traffic Data Imputation

Low-rank autoregressive tensor completion

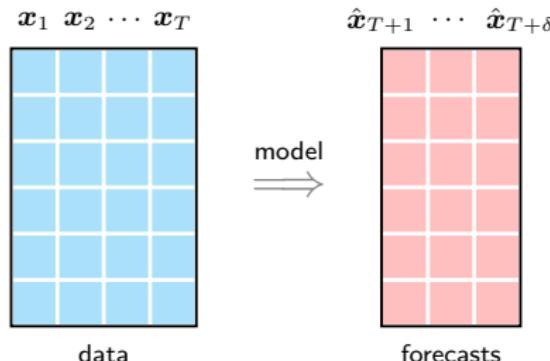
$$\begin{aligned} \min_{\boldsymbol{\chi}} \quad & \| \boldsymbol{\chi} \|_* + \frac{\lambda}{2} \sum_{n=1}^N \sum_{t=d+1}^T (z_{n,t} - \sum_{k=1}^d a_{n,k} z_{n,t-k})^2 \\ \text{s.t.} \quad & \begin{cases} \boldsymbol{\chi} = \mathcal{Q}(\mathbf{Z}), \\ \mathcal{P}_{\Omega}(\mathbf{Z}) = \mathcal{P}_{\Omega}(\mathbf{Y}). \end{cases} \end{aligned} \quad (5)$$

- **Advantage:** Global consistency + local consistency.



High-Dimensional Traffic Forecasting

- **Objective B-1:** Given a multivariate traffic time series $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^N$ with $N \gg T$ ("tall-skinny"), forecast data points $\mathbf{x}_{T+\delta}, \delta \in \mathbb{N}^+$.

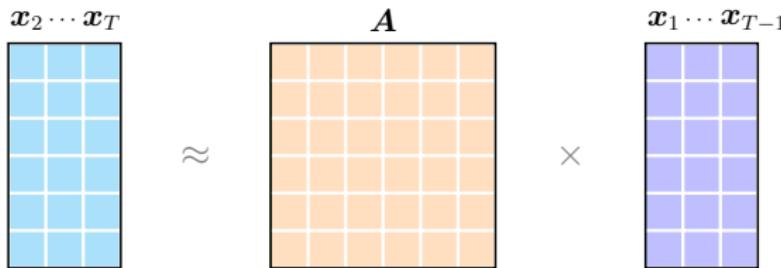


- **Solution:** For time series $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^N$, the d th-order vector autoregressive (VAR(d)) model: $\mathbf{x}_t = \sum_{k=1}^d \mathbf{A}_k \mathbf{x}_{t-k} + \boldsymbol{\epsilon}_t$.
 - **Advantage:** Co-evolution patterns
 - **Limitation:** Over-parameterization

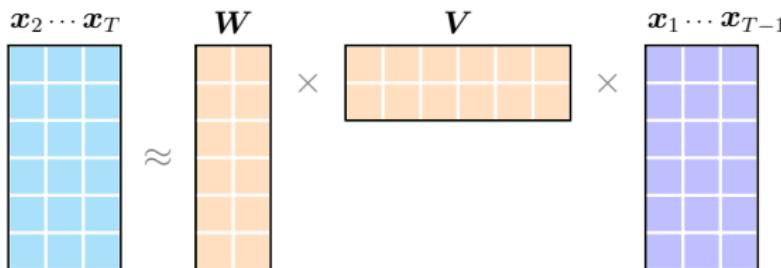
High-Dimensional Traffic Forecasting

VAR(1) model

- Over-parameterization in the case of $N \gg T$.



- Reduced-rank autoregression: $A = WV$ with $W \in \mathbb{R}^{N \times R}$, $V \in \mathbb{R}^{R \times N}$.



High-Dimensional Traffic Forecasting

Multivariate reduced-rank regression

- Data: $Z \in \mathbb{R}^{N \times T}$ (input) & $Y \in \mathbb{R}^{M \times T}$ (output)
 - Component matrices $W \in \mathbb{R}^{M \times R}, V \in \mathbb{R}^{R \times N}$ such that $A = WV$
 - Optimization problem:

$$\mathbf{W}^*, \mathbf{V}^* \triangleq \arg \min_{\mathbf{W}, \mathbf{V}} \frac{1}{2} \|\mathbf{Y} - \mathbf{WVZ}\|_F^2 \quad (6)$$

Theorem 1 (Izenman'75)

Suppose the optimization problem of multivariate reduced-rank regression trained on the centered data matrices Z and Y , then for any positive definite matrix Γ , the solutions are

$$W^* = \Gamma^{-1/2} \xi, V^* = \xi^\top \Gamma^{1/2} \Sigma_{yz} \Sigma_{zz}^{-1}, \quad (7)$$

where $\xi \in \mathbb{R}^{N \times R}$ consists of the R eigenvectors corresponding to the R largest eigenvalues of the matrix $\Gamma^{1/2} \Sigma_{yz} \Sigma_{zz}^{-1} \Sigma_{zy} \Gamma^{1/2}$.

High-Dimensional Traffic Forecasting

VAR(d) model

- Recall that $\mathbf{x}_t = \sum_{k=1}^d \mathbf{A}_k \mathbf{x}_{t-k} + \boldsymbol{\epsilon}_t$.
 - Coefficients $\mathbf{A}_k \in \mathbb{R}^{N \times N}$, $k = 1, \dots, d$ are tensor, e.g., $\mathbf{A} \in \mathbb{R}^{N \times N \times d}$.

VAR(d) with Tucker decomposition (Wang et al.'21)

For VAR(d) on the multivariate time series $\mathbf{x}_t \in \mathbb{R}^N, t = 1, \dots, T$, the reduced-rank VAR via Tucker decomposition is given by

$$\min_{\mathcal{G}, U_1, U_2, U_3} \frac{1}{2} \sum_{t=d+1}^T \| \mathbf{x}_t - (\mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3)_{(1)} \mathbf{z}_t \|_2^2 \quad (8)$$

where $\mathbf{z}_t = (\mathbf{x}_{t-1}^\top, \dots, \mathbf{x}_{t-d}^\top)^\top \in \mathbb{R}^{dN}$. The multilinear rank is (R_1, R_2, R_3) .

$\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$ is the core tensor, while $\mathbf{U}_1 \in \mathbb{R}^{N \times R_1}$, $\mathbf{U}_2 \in \mathbb{R}^{N \times R_2}$, and $\mathbf{U}_3 \in \mathbb{R}^{d \times R_3}$ are the component matrices.

Advantage: High compression rate.

Limitations: Nonconvex optimization; not scalable to large problems.

High-Dimensional Traffic Forecasting

High-dimensional VAR(d)

Given the (high-dimensional) multivariate time series $\mathbf{x}_t \in \mathbb{R}^N, t = 1, \dots, T$, for the data $\mathbf{x}_t \in \mathbb{R}^N, \mathbf{z}_t \in \mathbb{R}^{dN}, t = d+1, \dots, T$, we have a lower-dimensional problem:

$$\min_{\{\mathbf{U}, \mathbf{W}, \mathbf{V} | \mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_R\}} \frac{1}{2} \sum_{t=d+1}^T \|\mathbf{x}_t - \mathbf{U} \mathbf{W} \mathbf{V}^\top \mathbf{z}_t\|_2^2 \quad (9)$$

or equivalently, we have a canonical form:

$$\min_{\{\mathbf{U}, \mathbf{W}, \mathbf{V} | \mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_R\}} \frac{1}{2} \sum_{t=d+1}^T \|\mathbf{U}^\top \mathbf{x}_t - \mathbf{W} \mathbf{V}^\top \mathbf{z}_t\|_2^2 \quad (10)$$

where $\mathbf{W} \in \mathbb{R}^{R \times R}$ is the coefficient matrix.

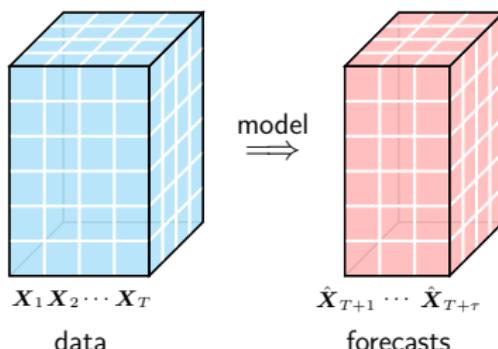
VARMA(d, q)

$$\mathbf{x}_t = \sum_{k=1}^d \mathbf{A}_k \mathbf{x}_{t-k} + \boldsymbol{\epsilon}_t + \sum_{p=1}^q \mathbf{B}_p \boldsymbol{\epsilon}_{t-p} \quad (11)$$

Moving average (MA) processes smooth out the data noises.

Multidimensional Traffic Forecasting

- **Objective B-2:** Given a multidimensional traffic time series $X_1, \dots, X_T \in \mathbb{R}^{M \times N}$, forecast data points $\hat{X}_{T+\tau}, \tau \in \mathbb{N}_+$.



[Q]

- How to perform forecasting on this kind of data?
- How to preserve the intrinsic tensor representation of data?

Multidimensional Traffic Forecasting

Matrix autoregressive model (Chen, Xiao & Yang'21)

Given matrix-variate time series $\mathbf{X}_t \in \mathbb{R}^{M \times N}$, $t = 1, \dots, T$, then the d th-order matrix autoregressive (MAR(d)) model takes the form of

$$\mathbf{X}_t = \sum_{k=1}^d \mathbf{A}_k \mathbf{X}_{t-k} \mathbf{B}_k^\top + \mathbf{E}_t \quad (12)$$

where $\mathbf{A}_k \in \mathbb{R}^{M \times M}$, $\mathbf{B}_k \in \mathbb{R}^{N \times N}$, $k = 1, \dots, d$ are the coefficient matrices.

Advantages:

- Preserve the intrinsic tensor representation.
- Reduce parameters in autoregressive models (if $n = \max\{M, N\}$), e.g.,

$$\mathcal{O}(n^4) \text{ in VAR(1)} \quad \text{vs.} \quad \mathcal{O}(n^2) \text{ in MAR(1)}$$

Limitation: Not scalable to large problems.

Multidimensional Traffic Forecasting

Borrowing the idea of Sylvester equation.

Sylvester-type MAR(d)

Given matrix-variate time series $\mathbf{X}_t \in \mathbb{R}^{M \times N}, t = 1, \dots, T$, then we define the MAR(d) as follows,

$$\mathbf{X}_t = \sum_{k=1}^d \left(\mathbf{A}_k \mathbf{X}_{t-k} + \mathbf{X}_{t-k} \mathbf{B}_k^\top \right) + \mathbf{E}_t \quad (13)$$

where $\mathbf{A}_k \in \mathbb{R}^{M \times M}, \mathbf{B}_k \in \mathbb{R}^{N \times N}, k = 1, \dots, d$ are the coefficient matrices.

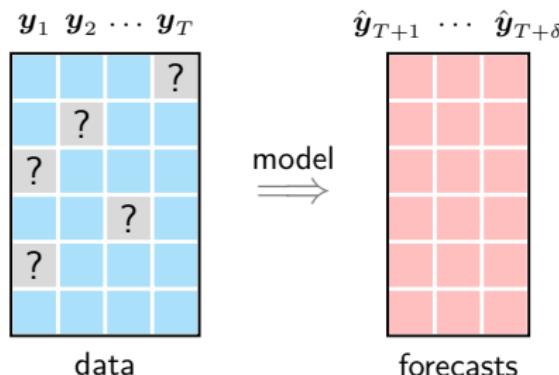
Advantage:

- Instead of the bilinear form of MAR(d), solving Sylvester-type MAR(d) is computationally economic.

Traffic Forecasting on Sparse Data

Multivariate traffic time series

- **Objective C-1:** Given a partially observed data $\mathbf{Y} \in \mathbb{R}^{N \times T}$ consisting of time series $y_1, \dots, y_T \in \mathbb{R}^N$, forecast data points $\hat{y}_{T+\delta}, \delta \in \mathbb{N}^+$.



[Q]

- How to learn from *high-dimensional* and *sparse* data?
- How to model *nonstationarity* in time series?
- How to perform forecasting on these time series?

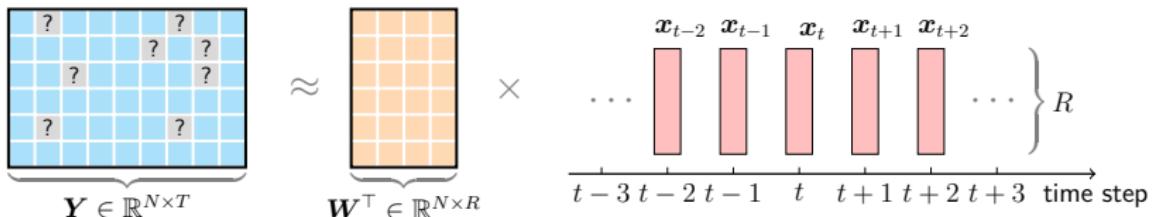
Traffic Forecasting on Sparse Data

Multivariate traffic time series

Temporal matrix factorization (Yu et al.'16; Chen & Sun'21)

Given any partially observed time series data $\mathbf{Y} \in \mathbb{R}^{N \times T}$ with observed index set Ω , then temporal matrix factorization assumes a d th-order vector autoregressive (VAR) process on the temporal factor matrix:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{X}, \{\mathbf{A}_k\}_{k=1}^d} \quad & \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{W}^\top \mathbf{X})\|_F^2 + \frac{\rho}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{X}\|_F^2) \\ & + \frac{\lambda}{2} \sum_{t=d+1}^T \|\mathbf{x}_t - \sum_{k=1}^d \mathbf{A}_k \mathbf{x}_{t-k}\|_2^2 \end{aligned} \quad (14)$$



VAR is usually built on stationary time series (temporal factors).

Traffic Forecasting on Sparse Data

Multivariate traffic time series

Nonstationary temporal matrix factorization (NoTMF)

Given any partially observed time series data $\mathbf{Y} \in \mathbb{R}^{N \times T}$ with observed index set Ω , then we assume a season- m differencing on the latent temporal factors:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{X}, \{\mathbf{A}_k\}_{k=1}^d} & \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{W}^\top \mathbf{X})\|_F^2 + \frac{\rho}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{X}\|_F^2) \\ & + \frac{\lambda}{2} \sum_{t=d+m+1}^T \|(\mathbf{x}_t - \mathbf{x}_{t-m}) - \sum_{k=1}^d \mathbf{A}_k (\mathbf{x}_{t-k} - \mathbf{x}_{t-m-k})\|_2^2 \end{aligned} \quad (15)$$

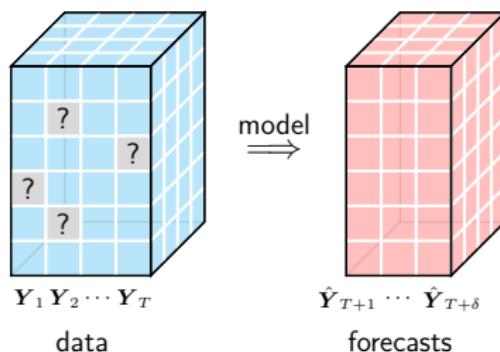
- First-order differencing $\mathbf{x}'_t = \mathbf{x}_t - \mathbf{x}_{t-1}$.
 - Second-order differencing $\mathbf{x}''_t = (\mathbf{x}_t - \mathbf{x}_{t-1}) - (\mathbf{x}_{t-1} - \mathbf{x}_{t-2})$.
 - Twice-differenced series $\mathbf{x}'''_t = (\mathbf{x}_t - \mathbf{x}_{t-m}) - (\mathbf{x}_{t-1} - \mathbf{x}_{t-m-1})$.
- 😊 Stationarizing a time series with differencing can improve the prediction.⁴

⁴ Stationarity and differencing: <https://otexts.com/fpp2/stationarity.html>

Traffic Forecasting on Sparse Data

Multidimensional traffic time series

- **Objective C-2:** Given a partially observed data $\mathcal{Y} \in \mathbb{R}^{M \times N \times T}$ consisting of time series $\mathbf{Y}_1, \dots, \mathbf{Y}_T \in \mathbb{R}^{M \times N}$, forecast data points $\hat{\mathbf{Y}}_{T+\delta}, \delta \in \mathbb{N}^+$.



- **Solution:** Temporal tensor factorization, e.g., CP factorization + VAR (on latent temporal factors).

Conclusion

Contributions

- *Objective A: Spatiotemporal traffic data imputation.* Develop a low-rank temporal modeling framework and improve the imputation accuracy, efficiency, and scalability.
 - *Objective B: High-dimensional and multidimensional forecasting.* Fast and accurate forecasting approach for high-dimensional and large-scale data; tensor representation based autoregressive model for multidimensional data.
 - *Objective C: Forecasting on sparse data.* Low-rank temporal modeling framework for traffic time series forecasting in the presence of missing values.

Research Work during Ph.D. Research

- Publications
 - [J1] X. Chen, M. Lei, N. Saunier, L. Sun (2021). Low-rank autoregressive tensor completion for spatiotemporal traffic data imputation. *IEEE Transactions on Intelligent Transportation Systems*. (Early access)
 - [J2] X. Chen, Y. Chen, N. Saunier, L. Sun (2021). Scalable low-rank tensor learning for spatiotemporal traffic data imputation. *Transportation Research Part C: Emerging Technologies*, 129: 103226.
 - [C1] X. Chen, M. Lei, N. Saunier, L. Sun (2021). Low-rank autoregressive tensor completion for spatiotemporal traffic data imputation. *The 7th SIGKDD Workshop on Mining and Learning from Time Series (MiLeTS)* at KDD 2021.
- Preprint (under review)
 - [P1] X. Chen, C. Zhang, X.L. Zhao, N. Saunier, L. Sun (2022). Nonstationary temporal matrix factorization for multivariate time series forecasting. *IEEE Transactions on Knowledge and Data Engineering*.
- Open-source projects
 - **transdim**: Machine learning for spatiotemporal traffic data imputation and forecasting. (780 stars & 240 forks on GitHub)
<https://github.com/xinchen/transdim>
 - **tracebase**: Multivariate time series forecasting on high-dimensional and sparse Uber movement speed data. (19 stars & 5 forks on GitHub)
<https://github.com/xinchen/tracebase>



Thanks for your attention!

Any Questions?

About me:

-  Homepage: <https://xinychen.github.io>
-  GitHub: <https://github.com/xinychen> (2.4K+ stars)
-  Blog: <https://medium.com/@xinyu.chen> (30K+ views)
-  How to reach me: chenxy346@gmail.com