Introduction
Matrix completion : the independent case
Time series completion

# Tight Risk Bound for High Dimensional Time Series Completion

## Pierre Alquier

RIKEN

AIP
Center for
Advanced Intelligence Project

EcoDep 2021 Conference
September 16, 2021

Introduction
Matrix completion : the independent case
Time series completion

# Co-authors

Alquier, P., Marie, N. and Rosier, A. (2021). Tight Risk Bound for High Dimensional Time Series Completion. *Preprint arXiv :2102.08178*.
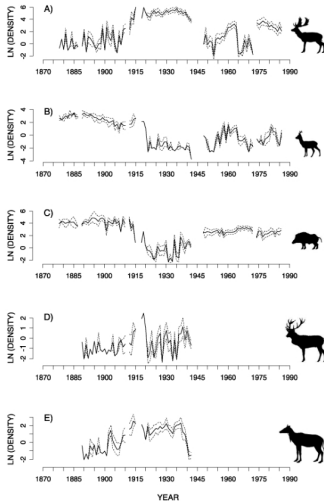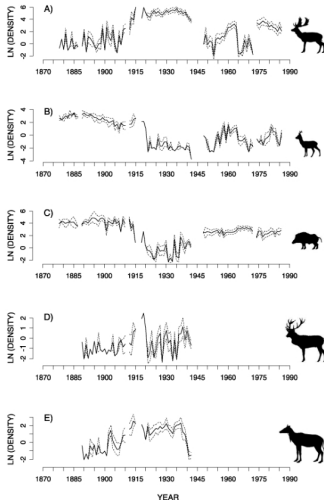


## Nicolas Marie

**Université Paris Nanterre**



## Amélie Rosier

**ESME Sudria and Université Paris Nanterre**

**Introduction**
Matrix completion : the independent case
Time series completion

# Multivariate time series

**Introduction**
Matrix completion : the independent case
Time series completion

# Multivariate time series



S. Imperio *et al*. (2010). Investigating population dynamics in ungulates : Do hunting statistics make up a good index of population abundance ? *Wildlife Biology*.

- multivariate series
- correlations
- noisy observations
- missing entries

**Introduction**
Matrix completion : the independent case
Time series completion

# Partially observed multivariate time series

| $i$ | $\ldots$ | $t-3$ | $t-2$ | $t-1$ | $t$ | $t+1$ | $t+2$ | $t+3$ | $\ldots$ |
|---|---|---|---|---|---|---|---|---|---|
| 1 | $\ldots$ |  | 12.5 |  |  | 17 |  |  | $\ldots$ |
| 2 | $\ldots$ | 1.2 |  |  | 3.8 |  |  | 2.9 | $\ldots$ |
| 3 | $\ldots$ |  |  | 0 |  | 7.2 |  |  | $\ldots$ |
| 4 | $\ldots$ |  |  |  | 4.2 | 3.1 | 2.4 | 2.3 | $\ldots$ |
| 5 | $\ldots$ | 23.1 |  |  | 45.1 | 39.9 |  |  | $\ldots$ |
| 6 | $\ldots$ |  | 4.1 | 4.1 |  | 6.3 |  | 2.9 | $\ldots$ |
| 7 | $\ldots$ | 0.1 |  | 0.9 | 0 |  |  |  | $\ldots$ |
| 8 | $\ldots$ |  |  |  | 34.7 |  |  |  | $\ldots$ |
| $\vdots$ | $\vdots$ |  |  |  | $\vdots$ |  |  |  | $\ddots$ |

**Introduction**
Matrix completion : the independent case
Time series completion

## Examples

- econometrics : panel data with missing entries,
- industry : data from sensors at multiple locations,
- ecology : spatial data with observations from a few sites only at each date,
- . . .
- more generally, any situation where we have multivariate time series and each measurement is expensive.

**Introduction**
Matrix completion : the independent case
Time series completion

# Matrix completion methods

- matrix completion algorithms exist, and were successful in many applications.
- many of them are based on a low-rank assumption and on matrix factorization.
- however, the theory was developed only in the independent case.

**Introduction**
Matrix completion : the independent case
Time series completion

# Contents

Introduction
Matrix completion : the independent case
Time series completion

Matrix completion model
Minimax rate of estimation

# Contents

Introduction
Matrix completion : the independent case
Time series completion

Matrix completion model
Minimax rate of estimation

# Classical example : collaborative filtering

Introduction
Matrix completion : the independent case
Time series completion

Matrix completion model
Minimax rate of estimation

# A statistical model

There is a $d \times T$ matrix $M$ and $n$ i.i.d observations $Y_1, \ldots, Y_n$ drawn as :

- $(i_\ell, j_\ell)$ drawn uniformly on $\{1, \ldots, d\} \times \{1, \ldots, T\}$,
- $Y_\ell = M_{i_\ell j_\ell} + \varepsilon_\ell$

where $\varepsilon_\ell$ is some noise ($= 0$ in the first papers on the topic, subgaussian with variance $\sigma^2$ later).

Introduction
Matrix completion : the independent case
Time series completion

Matrix completion model
Minimax rate of estimation

# A statistical model

There is a $d \times T$ matrix $M$ and $n$ i.i.d observations $Y_1, \ldots, Y_n$ drawn as :

- $(i_\ell, j_\ell)$ drawn uniformly on $\{1, \ldots, d\} \times \{1, \ldots, T\}$,
- $Y_\ell = M_{i_\ell j_\ell} + \varepsilon_\ell$

where $\varepsilon_\ell$ is some noise ($= 0$ in the first papers on the topic, subgaussian with variance $\sigma^2$ later).

Key assumption : $k := \operatorname{rank}(M) \ll \min(d, T) = K$.

Introduction
Matrix completion : the independent case
Time series completion

Matrix completion model
Minimax rate of estimation

# SVD & matrix factorization

$$
M = \underbrace{\left( U_1 \middle| \ldots \middle| U_k \middle| \ldots \right)}_{=U\ (d \times K)}
\underbrace{\begin{pmatrix} \sigma_1 & 0 & \ldots & \\ 0 & \ddots & 0 & \ldots \\ \vdots & & \sigma_k & \\ & & & 0 & \\ & & & & \ddots \end{pmatrix}}_{=\Sigma\ (K \times K)}
\underbrace{\left( \begin{array}{c} V_1^T \\ \hline \vdots \\ \hline V_k^T \\ \hline \vdots \end{array} \right)}_{=V^T\ (K \times T)}
$$

Introduction
Matrix completion : the independent case
Time series completion

Matrix completion model
Minimax rate of estimation

# SVD & matrix factorization

$$M = \underbrace{\left( U_1 \Big| \ldots \Big| U_k \Big| \ldots \right)}_{=U \ (d \times K)} \underbrace{\begin{pmatrix} \sigma_1 & 0 & \ldots & \\ 0 & \ddots & 0 & \ldots \\ \vdots & & \sigma_k & \\ & & & 0 & \\ & & & & \ddots \end{pmatrix}}_{=\Sigma \ (K \times K)} \underbrace{\left( \frac{\overline{V_1^T}}{\frac{\vdots}{\overline{V_k^T}}} \right)}_{=V^T \ (K \times T)}$$

$$M = \underbrace{\left( U_1 \Big| \ldots \Big| U_k \right) \underbrace{\begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_k \end{pmatrix}}_{=A \ (d \times k)}}_{} \underbrace{\left( \frac{\overline{V_1^T}}{\frac{\vdots}{\overline{V_k^T}}} \right)}_{=B \ (k \times T)}$$

Introduction
**Matrix completion : the independent case**
Time series completion

Matrix completion model
**Minimax rate of estimation**

## Estimation

$$
\hat{M}^\lambda = \arg\min_X \left\{ \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell j_\ell})^2 + \lambda \sum_{h=1}^{\min(d,T)} \sigma_h(X) \right\}.
$$

Introduction
**Matrix completion : the independent case**
Time series completion

Matrix completion model
**Minimax rate of estimation**

# Estimation

$$
\hat{M}^\lambda = \underset{X}{\arg\min} \left\{ \sum_{\ell=1}^n (Y_\ell - X_{i_\ell j_\ell})^2 + \lambda \sum_{h=1}^{\min(d,T)} \sigma_h(X) \right\}.
$$

### Theorem

For a well chosen $\lambda$ that does not depend on $k$, and under minimal assumptions on $M$, with large probability

$$
\frac{1}{dT} \sum_{i,j} \left( \hat{M}_{i,j}^\lambda - M_{i,j} \right)^2 \leq \mathrm{Cst} \frac{\sigma k(d+T) \log(d+T)}{n}
$$

📄 Koltchinskii, V., Lounici, K. and Tsybakov, A. (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics*.

Introduction
Matrix completion : the independent case
**Time series completion**

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Contents

Introduction
Matrix completion : the independent case
**Time series completion**

**Generalization of the results for i.i.d data to time series**
Using the time series structure : faster rates

# Time series completion : the model

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Time series completion : the model



- low-rank trend :

$$M = \underbrace{A}_{d \times k} \underbrace{B}_{k \times T}$$

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Time series completion : the model



- low-rank trend :

$$M = \underbrace{A}_{d \times k} \underbrace{B}_{k \times T}$$

- temporal correlated noise $\varepsilon$ :

  $\varepsilon_{i,t}$ indep. $\varepsilon_{j,t'}$ ($i \neq j$)

  $\varepsilon_{i,t}$ not indep. $\varepsilon_{i,t'}$

Introduction
Matrix completion : the independent case
**Time series completion**

**Generalization of the results for i.i.d data to time series**
Using the time series structure : faster rates

# Time series completion : the model



- low-rank trend :

$$M = \underbrace{A}_{d \times k} \underbrace{B}_{k \times T}$$

- temporal correlated noise $\varepsilon$ :

  $\varepsilon_{i,t}$ indep. $\varepsilon_{j,t'}$ $(i \neq j)$

  $\varepsilon_{i,t}$ not indep. $\varepsilon_{i,t'}$

- $(i_\ell, t_\ell)$ i.i.d uniform, $\xi_\ell$ observation noise :

  $Y_\ell = M_{i_\ell, t_\ell} + \varepsilon_{i_\ell, t_\ell} + \xi_\ell.$

Introduction
Matrix completion : the independent case
**Time series completion**

**Generalization of the results for i.i.d data to time series**
Using the time series structure : faster rates

# Assumptions

### Reminder : the model

$$Y_\ell = M_{i_\ell, t_\ell} + \varepsilon_{i_\ell, t_\ell} + \xi_\ell.$$

Introduction
Matrix completion : the independent case
**Time series completion**

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Assumptions

### Reminder : the model

$$Y_\ell = M_{i_\ell, t_\ell} + \varepsilon_{i_\ell, t_\ell} + \xi_\ell.$$

- $M = \underbrace{A}_{d \times k} \underbrace{B}_{k \times T}$ and $|A_{i,h}|, |B_{h,t}| \leq c_{A,B}/\sqrt{k}$.
- $(i_\ell, t_\ell)$ i.i.d uniform on $\{1, \ldots, d\} \times \{1 \ldots, T\}$ ;
- $(\varepsilon_{i,t})_{t=1,\ldots,T}$ is a bounded, $\phi$-mixing time series :

$$|\varepsilon_{i,t}| \leq m_\varepsilon \text{ and } \sum_{t=1}^{\infty} \phi_{\varepsilon_{i,.}}(t) \leq \Phi_\varepsilon.$$

- $(\xi_\ell)$ are i.i.d, sub-exponential variables : for $k \geq 2$,

$$\mathbb{E}(|\xi_\ell|^q) \leq \frac{v_\xi c_\xi^{q-2} q!}{2}.$$

Introduction
Matrix completion : the independent case
**Time series completion**

**Generalization of the results for i.i.d data to time series**
Using the time series structure : faster rates

# Estimator and risk bound

$$\hat{M}^{(k)} = \underset{\underbrace{X}_{d \times T} = \underbrace{A}_{d \times k} \underbrace{B}_{k \times T}}{\arg\min} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2.$$

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Estimator and risk bound

$$\hat{M}^{(k)} = \underset{\underbrace{X}_{d \times T} = \underbrace{A}_{d \times k} \underbrace{B}_{k \times T}}{\arg\min} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2.$$

## Theorem

With probability at least $1 - \eta$,

$$\frac{1}{dT} \sum_{i,j} \left( \hat{M}^{(k)}_{i,j} - M_{i,j} \right)^2 \leq C \frac{k(d + T)\log(n) + \log\left(\frac{1}{\eta}\right)}{n}$$

where $C = C(c_{A,B}, m_\varepsilon, \Phi_\varepsilon, v_\xi, c_\xi)$ is known.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Remarks on the proof

1. decompose the difference between *empirical risk* and *expected risk* $\frac{1}{n}\sum_{\ell=1}^{n}(Y_\ell - X_{i_\ell,j_\ell})^2 - \frac{1}{dT}\sum_{i,j}(M_{i,j} - X_{i,j})^2$ in elementary terms.

2. some of these terms are sums of i.i.d variables. Bound them via Bernstein inequality. Some are sums of $\phi$-mixing variables, use :

📄 Samson, P.-M. (2000). Concentration of measure inequalities for Markov chains and $\Phi$-mixing processes. *The Annals of Probability*.

3. union bound.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
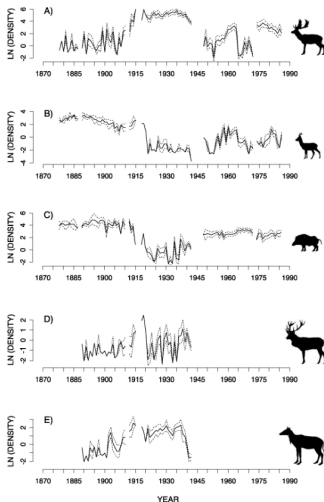Using the time series structure : faster rates

# Remarks on the proof

1. decompose the difference between *empirical risk* and *expected risk* $\frac{1}{n} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2 - \frac{1}{dT} \sum_{i,j} (M_{i,j} - X_{i,j})^2$ in elementary terms.

2. some of these terms are sums of i.i.d variables. Bound them via Bernstein inequality. Some are sums of $\phi$-mixing variables, use :

Samson, P.-M. (2000). Concentration of measure inequalities for Markov chains and Φ-mixing processes. *The Annals of Probability*.

3. union bound.

REMARK : if the $\varepsilon_{i,\cdot}$ satisfy another notion of mixing or weak-dependence, we can use alternative versions of Bernstein inequality but this lead to slower rates of convergence, in $1/\sqrt{n}$.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Rank selection

$$\hat{k} = \underset{1 \leq k \leq K}{\arg\min} \left\{ \frac{1}{n} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2 + C' \frac{k(d+T)\log(n)}{n} \right\}$$

where $C' = C'(c_{A,B}, m_\varepsilon, \Phi_\varepsilon, v_\xi, c_\xi)$ is **known** but too large.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Rank selection

$$\hat{k} = \arg\min_{1 \le k \le K} \left\{ \frac{1}{n} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell,j_\ell})^2 + C' \frac{k(d+T)\log(n)}{n} \right\}$$

where $C' = C'(c_{A,B}, m_\varepsilon, \Phi_\varepsilon, v_\xi, c_\xi)$ is **known** but too large.

In practice : we use the slope heuristic to calibrate a better $C'$.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Rank selection

$$\hat{k} = \underset{1 \leq k \leq K}{\arg\min} \left\{ \frac{1}{n} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2 + C' \frac{k(d+T)\log(n)}{n} \right\}$$

where $C' = C'(c_{A,B}, m_\varepsilon, \Phi_\varepsilon, v_\xi, c_\xi)$ is **known** but too large.

In practice : we use the slope heuristic to calibrate a better $C'$.

### Theorem

With probability at least $1 - \eta$,

$$\frac{1}{dT} \sum_{i,j} \left( \hat{M}_{i,j}^{(\hat{k})} - M_{i,j} \right)^2 \leq C'' \frac{k(d+T)\log(n) + \log\left(\frac{1}{\eta}\right)}{n}.$$

Introduction
Matrix completion : the independent case
**Time series completion**

Generalization of the results for i.i.d data to time series
**Using the time series structure : faster rates**

# Time series with a structure

**Example :** assume that the trends in $M$ are $p$-periodic. This means that

$$\underbrace{M}_{d \times T} = \underbrace{C}_{d \times p} \underbrace{(I_p | \dots | I_p)}_{= \Lambda \ (p \times T)}.$$

Introduction
Matrix completion : the independent case
**Time series completion**

Generalization of the results for i.i.d data to time series
**Using the time series structure : faster rates**

# Time series with a structure

**Example :** assume that the trends in $M$ are $p$-periodic. This means that

$$\underbrace{M}_{d \times T} = \underbrace{C}_{d \times p} \underbrace{(I_p | \ldots | I_p)}_{= \Lambda \ (p \times T)}.$$

More generally, we can assume that there is a known structure in $M$ :

$$\underbrace{M}_{d \times T} = \underbrace{C}_{d \times p} \underbrace{\Lambda}_{p \times T}$$

and still add the initial "low-rank decomposition" to ensure correlations in the rows :

$$\underbrace{M}_{d \times T} = \underbrace{A}_{d \times k} \underbrace{B}_{k \times p} \underbrace{\Lambda}_{p \times T}.$$

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Faster rates

$$\hat{M}^{(k)} = \underset{\underbrace{X}_{d \times T} = \underbrace{A}_{d \times k} \underbrace{B}_{k \times p} \underbrace{\Lambda}_{p \times T}}{\arg\min} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2.$$

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Faster rates

$$\hat{M}^{(k)} = \underbrace{\arg\min}_{\underbrace{X}_{d \times T} = \underbrace{A}_{d \times k} \underbrace{B}_{k \times p} \underbrace{\Lambda}_{p \times T}} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2.$$

### Theorem

With probability at least $1 - \eta$,

$$\frac{1}{dT} \sum_{i,j} \left( \hat{M}^{(k)}_{i,j} - M_{i,j} \right)^2 \leq C \frac{k(d+p)\log(n) + \log\left(\frac{1}{\eta}\right)}{n}.$$

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# Faster rates

$$\hat{M}^{(k)} = \underset{\underbrace{X}_{d \times T} = \underbrace{A}_{d \times k} \underbrace{B}_{k \times p} \underbrace{\Lambda}_{p \times T}}{\arg\min} \sum_{\ell=1}^{n} (Y_\ell - X_{i_\ell, j_\ell})^2.$$

### Theorem

With probability at least $1 - \eta$,

$$\frac{1}{dT} \sum_{i,j} \left( \hat{M}_{i,j}^{(k)} - M_{i,j} \right)^2 \leq C \frac{k(d+p)\log(n) + \log\left(\frac{1}{\eta}\right)}{n}.$$

We also have a similar rank-selection procedure.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates

# RIKEN AIP : position in the ABI team



Approximate Bayesian
Inference team (ABI), lead
by Emtiyaz Khan

## Please visit the team website

https ://team-approx-bayes.github.io/

Open Position : Research Scientist (1 position, Female only)

- research only (= chargé de recherches),
- indefinite-term,
- located in Tokyo center.

Introduction
Matrix completion : the independent case
Time series completion

Generalization of the results for i.i.d data to time series
Using the time series structure : faster rates