

Learning Mixtures of Linear Dynamical Systems

Yanxi Chen, H. Vincent Poor



Mixtures of time-series models

Mixture models are powerful in the face of *heterogeneous and complex* time-series data



Mixtures of time-series models

Mixture models are powerful in the face of *heterogeneous and complex* time-series data



- Higher accuracy of fitting the data
- Better interpretability: reveal cluster structures

Numerous applications

(Bulteel et al., 2016) Time-series measurements of certain psychological symptoms for multiple patients → identify *subgroups of patients*, provide tailored treatments.

Numerous applications

(Bulteel et al., 2016) Time-series measurements of certain psychological symptoms for multiple patients → identify *subgroups of patients*, provide tailored treatments.

(Hallac et al., 2017) Sensory data of a car under *a few driving modes* (e.g. “driving straight”, “slowing down”, “turning”, etc.)

Numerous applications

(Bulteel et al., 2016) Time-series measurements of certain psychological symptoms for multiple patients → identify *subgroups of patients*, provide tailored treatments.

(Hallac et al., 2017) Sensory data of a car under *a few driving modes* (e.g. “driving straight”, “slowing down”, “turning”, etc.)

(Brunskill et al., 2009) Sensory data of a robot navigating *a complex environment* (e.g. with areas of grass, sand, carpets, rocks, etc.)

.....

Problem formulation

Recap: linear dynamical system (LDS). A d -dimensional time-series trajectory $\{x_t\}$ generated by an LDS model, i.e. the $d \times d$ state transition matrix \mathbf{A} and noise covariance \mathbf{W} :

Problem formulation

Recap: linear dynamical system (LDS). A d -dimensional time-series trajectory $\{x_t\}$ generated by an LDS model, i.e. the $d \times d$ state transition matrix \mathbf{A} and noise covariance \mathbf{W} :

$$x_{t+1} = \mathbf{A}x_t + w_t, \quad \text{where} \quad \mathbb{E}[w_t] = \mathbf{0}, \quad \text{cov}(w_t) = \mathbf{W} \succ \mathbf{0}.$$

Problem formulation

Recap: linear dynamical system (LDS). A d -dimensional time-series trajectory $\{\mathbf{x}_t\}$ generated by an LDS model, i.e. the $d \times d$ state transition matrix \mathbf{A} and noise covariance \mathbf{W} :

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t, \quad \text{where} \quad \mathbb{E}[\mathbf{w}_t] = \mathbf{0}, \quad \text{cov}(\mathbf{w}_t) = \mathbf{W} \succ \mathbf{0}.$$

Mixed LDSs. K models $\{\mathbf{A}^{(k)}, \mathbf{W}^{(k)}\}_{1 \leq k \leq K}$, M trajectories $\{\mathbf{X}_m\}_{1 \leq m \leq M}$, where $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}$ is generated by the k_m -th model:

$$\mathbf{x}_{m,t+1} = \mathbf{A}^{(k_m)}\mathbf{x}_{m,t} + \mathbf{w}_{m,t}, \quad \text{cov}(\mathbf{w}_{m,t}) = \mathbf{W}^{(k_m)}$$

Problem formulation

Recap: linear dynamical system (LDS). A d -dimensional time-series trajectory $\{\mathbf{x}_t\}$ generated by an LDS model, i.e. the $d \times d$ state transition matrix \mathbf{A} and noise covariance \mathbf{W} :

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t, \quad \text{where} \quad \mathbb{E}[\mathbf{w}_t] = \mathbf{0}, \quad \text{cov}(\mathbf{w}_t) = \mathbf{W} \succ \mathbf{0}.$$

Mixed LDSs. K models $\{\mathbf{A}^{(k)}, \mathbf{W}^{(k)}\}_{1 \leq k \leq K}$, M trajectories $\{\mathbf{X}_m\}_{1 \leq m \leq M}$, where $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}$ is generated by the k_m -th model:

$$\mathbf{x}_{m,t+1} = \mathbf{A}^{(k_m)}\mathbf{x}_{m,t} + \mathbf{w}_{m,t}, \quad \text{cov}(\mathbf{w}_{m,t}) = \mathbf{W}^{(k_m)}$$

Note that the labels $\{k_m\}$ are **unknown!**

Literature: lack of **provable guarantees** for model estimation

Literature: lack of **provable guarantees** for model estimation

Major challenges:

- Latent variables are not observed;

Literature: lack of **provable guarantees** for model estimation

Major challenges:

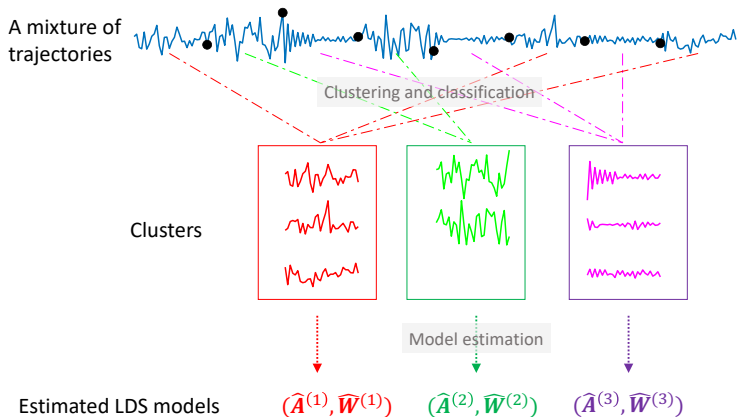
- Latent variables are not observed;
- Short trajectories might have lengths much smaller than the model dimension d ;

Literature: lack of **provable guarantees** for model estimation

Major challenges:

- Latent variables are not observed;
- Short trajectories might have lengths much smaller than the model dimension d ;
- Temporal dependence inherent to time series (in contrast to mixed regression problems).

Outline of our solution



A two-stage approach

Stage 1: coarse estimation

- **Subspace estimation**
- **Clustering** of trajectories (assisted by variance reduction)
- Initial **model estimation** within each cluster

A two-stage approach

Stage 1: coarse estimation

- **Subspace estimation**
- **Clustering** of trajectories (assisted by variance reduction)
- Initial **model estimation** within each cluster

Stage 2: refined estimation

- **Classification** of additional trajectories
- Refined **model estimation** within each cluster

A two-stage approach

Stage 1: coarse estimation

- **Subspace estimation**
- **Clustering** of trajectories (assisted by variance reduction)
- Initial **model estimation** within each cluster

Stage 2: refined estimation

- **Classification** of additional trajectories
- Refined **model estimation** within each cluster

The algorithm outline is largely inspired by the works on meta-learning for mixed linear regression (Kong et al., 2020a;b), but the detailed implementations are substantially different due to temporal dependence in mixed LDSs; see Section 2 of paper for detailed algorithms.

Assumptions for simplification

Initial state: each trajectory starts at $\mathbf{x}_{m,0} = \mathbf{0}$. (Another canonical case is when the short trajectories are segments of a single continuous trajectory; the main results are slightly different, and included in the paper.)

Assumptions for simplification

Initial state: each trajectory starts at $\mathbf{x}_{m,0} = \mathbf{0}$. (Another canonical case is when the short trajectories are segments of a single continuous trajectory; the main results are slightly different, and included in the paper.)

Balanced clusters: each LDS model accounts for (order-wise) $1/K$ proportion of data.

Assumptions for simplification

Initial state: each trajectory starts at $\mathbf{x}_{m,0} = \mathbf{0}$. (Another canonical case is when the short trajectories are segments of a single continuous trajectory; the main results are slightly different, and included in the paper.)

Balanced clusters: each LDS model accounts for (order-wise) $1/K$ proportion of data.

Sample splitting: M sample trajectories,

$$\{1, 2, \dots, M\} = \mathcal{M}_{\text{subspace}} \cup \mathcal{M}_{\text{clustering}} \cup \mathcal{M}_{\text{classification}}.$$

Assume that each trajectory in \mathcal{M}_o has the same length T_o , and denote the total sample size of \mathcal{M}_o as $T_{\text{total},o} = T_o \cdot |\mathcal{M}_o|$.

Essential assumptions

Mixing: for each $\mathbf{A} \in \{\mathbf{A}^{(k)}\}$ and all $t \geq 1$, $\|\mathbf{A}^t\| \leq \kappa_{\mathbf{A}} \cdot \rho^t$ for some $0 \leq \rho < 1$; denote mixing time $t_{\text{mix}} := 1/(1 - \rho)$.

Essential assumptions

Mixing: for each $\mathbf{A} \in \{\mathbf{A}^{(k)}\}$ and all $t \geq 1$, $\|\mathbf{A}^t\| \leq \kappa_{\mathbf{A}} \cdot \rho^t$ for some $0 \leq \rho < 1$; denote mixing time $t_{\text{mix}} := 1/(1 - \rho)$.

Stationary autocovariance matrices $\{\Gamma^{(k)}, \mathbf{Y}^{(k)}\}$, where

$$\Gamma(\mathbf{A}, \mathbf{W}) := \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top | \mathbf{A}, \mathbf{W}],$$

$$\mathbf{Y}(\mathbf{A}, \mathbf{W}) := \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^\top | \mathbf{A}, \mathbf{W}].$$

Essential assumptions

Mixing: for each $\mathbf{A} \in \{\mathbf{A}^{(k)}\}$ and all $t \geq 1$, $\|\mathbf{A}^t\| \leq \kappa_A \cdot \rho^t$ for some $0 \leq \rho < 1$; denote mixing time $t_{\text{mix}} := 1/(1 - \rho)$.

Stationary autocovariance matrices $\{\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)}\}$, where

$$\mathbf{\Gamma}(\mathbf{A}, \mathbf{W}) := \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top | \mathbf{A}, \mathbf{W}],$$

$$\mathbf{Y}(\mathbf{A}, \mathbf{W}) := \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^\top | \mathbf{A}, \mathbf{W}].$$

Model separation: there exist $\delta_{\Gamma, Y}, \delta_{A, W} > 0$ such that

$$\begin{aligned} \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\mathbb{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2 &\geq d \cdot \delta_{\Gamma, Y}^2, \\ \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2 + \|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\mathbb{F}}^2 &\geq d \cdot \delta_{A, W}^2, \end{aligned}$$

for all $1 \leq k < \ell \leq K$

Theorem. With high probability, the proposed two-stage method achieves exact clustering and classification of the sample trajectories, as well as final model estimation errors

$$\|\widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(k)}\| \leq \epsilon, \quad \frac{\|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)}\|}{\|\mathbf{W}^{(k)}\|} \leq \epsilon, \quad 1 \leq k \leq K,$$

Theorem. With high probability, the proposed two-stage method achieves exact clustering and classification of the sample trajectories, as well as final model estimation errors

$$\|\widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(k)}\| \leq \epsilon, \quad \frac{\|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)}\|}{\|\mathbf{W}^{(k)}\|} \leq \epsilon, \quad 1 \leq k \leq K,$$

provided the following sample complexities:

$$\begin{aligned} T_{\text{subspace}} &\gtrsim t_{\text{mix}}, & T_{\text{total,subspace}} &\gtrsim t_{\text{mix}} d \left(\frac{K^4}{\delta_{\Gamma,Y}^4} + 1 \right), \\ T_{\text{clustering}} &\gtrsim t_{\text{mix}} \left(\frac{1}{\delta_{\Gamma,Y}^2} \sqrt{\frac{K}{d}} + 1 \right), & T_{\text{total,clustering}} &\gtrsim K d \left(\frac{1}{\delta_{A,W}^2} + 1 \right), \\ T_{\text{classification}} &\gtrsim \frac{1}{d \delta_{A,W}^2} + 1, & T_{\text{total,clustering}} + T_{\text{total,classification}} &\gtrsim \frac{Kd}{\epsilon^2}. \end{aligned}$$

Theorem. With high probability, the proposed two-stage method achieves exact clustering and classification of the sample trajectories, as well as final model estimation errors

$$\|\widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(k)}\| \leq \epsilon, \quad \frac{\|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)}\|}{\|\mathbf{W}^{(k)}\|} \leq \epsilon, \quad 1 \leq k \leq K,$$

provided the following sample complexities:

$$\begin{aligned} T_{\text{subspace}} &\gtrsim t_{\text{mix}}, & T_{\text{total,subspace}} &\gtrsim t_{\text{mix}} d \left(\frac{K^4}{\delta_{\Gamma,Y}^4} + 1 \right), \\ T_{\text{clustering}} &\gtrsim t_{\text{mix}} \left(\frac{1}{\delta_{\Gamma,Y}^2} \sqrt{\frac{K}{d}} + 1 \right), & T_{\text{total,clustering}} &\gtrsim K d \left(\frac{1}{\delta_{A,W}^2} + 1 \right), \\ T_{\text{classification}} &\gtrsim \frac{1}{d \delta_{A,W}^2} + 1, & T_{\text{total,clustering}} + T_{\text{total,classification}} &\gtrsim \frac{Kd}{\epsilon^2}. \end{aligned}$$

See Section 3 of paper for formal theorems.

Summary

- Problem formulation of mixed LDSs;
- A two-stage approach for solving it;
- Theoretical guarantees with non-asymptotic sample complexities.

Future works

- Strengthening the theoretical analysis and algorithm design.

Future works

- Strengthening the theoretical analysis and algorithm design.
- Learning mixtures of more general time-series models. (Our algorithms essentially require (1) mixing; (2) the existence of stationary autocovariance matrices; (3) well-specified parametric models, and sufficient separation among them.)

Future works

- Strengthening the theoretical analysis and algorithm design.
- Learning mixtures of more general time-series models. (Our algorithms essentially require (1) mixing; (2) the existence of stationary autocovariance matrices; (3) well-specified parametric models, and sufficient separation among them.)
- Applications in real-world problems.

Future works

- Strengthening the theoretical analysis and algorithm design.
- Learning mixtures of more general time-series models. (Our algorithms essentially require (1) mixing; (2) the existence of stationary autocovariance matrices; (3) well-specified parametric models, and sufficient separation among them.)
- Applications in real-world problems.
- Extensions to the cases with controlled inputs, e.g. LQR in control and latent MDP in reinforcement learning.

Future works

- Strengthening the theoretical analysis and algorithm design.
- Learning mixtures of more general time-series models. (Our algorithms essentially require (1) mixing; (2) the existence of stationary autocovariance matrices; (3) well-specified parametric models, and sufficient separation among them.)
- Applications in real-world problems.
- Extensions to the cases with controlled inputs, e.g. LQR in control and latent MDP in reinforcement learning.

Thank you!