

# Nonlinear Systems Identification

Machine Learning in Feedback System #10

Yahya Sattar

Postdoc with Prof. Sarah Dean

# Recap: System Identification

Linear dynamical system with state observations

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{w}_t$$

# Recap: System Identification

Linear dynamical system with state observations

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{w}_t$$

- Identify the dynamics model  $\mathbf{s}_{0:T} \rightarrow \hat{\mathbf{F}}$

# Recap: System Identification

Linear dynamical system with state observations

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{w}_t$$

- Identify the dynamics model  $\mathbf{s}_{0:T} \rightarrow \hat{\mathbf{F}}$

- State, noise  $\mathbf{s}_t, \mathbf{w}_t \in \mathbb{R}^n$

# Recap: System Identification

Linear dynamical system with state observations

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{w}_t$$

- Identify the dynamics model  $\mathbf{s}_{0:T} \rightarrow \hat{\mathbf{F}}$
- State, noise  $\mathbf{s}_t, \mathbf{w}_t \in \mathbb{R}^n$
- System dynamics  $\mathbf{F} \in \mathbb{R}^{n \times n}$

## Recap: System Identification

Least-squares problem with matrix variable

$$\hat{\mathbf{F}} = \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t\|_{\ell_2}^2$$

# Recap: System Identification

Least-squares problem with matrix variable

$$\begin{aligned}\hat{\mathbf{F}} &= \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t\|_{\ell_2}^2 \\ &= \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{w}_t\|_{\ell_2}^2 \text{ s.t. } \underbrace{\begin{bmatrix} \mathbf{s}_1 \\ \vdots \\ \mathbf{s}_T \end{bmatrix}}_{\mathbf{s}_+} = \underbrace{\begin{bmatrix} \mathbf{s}_0 \\ \vdots \\ \mathbf{s}_{T-1} \end{bmatrix}}_{\mathbf{s}_-} \mathbf{F}^\top + \underbrace{\begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{bmatrix}}_{\mathbf{w}}\end{aligned}$$

# Recap: System Identification

Least-squares problem with matrix variable

$$\begin{aligned}\hat{\mathbf{F}} &= \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t\|_{\ell_2}^2 \\&= \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{w}_t\|_{\ell_2}^2 \text{ s.t. } \underbrace{\begin{bmatrix} \mathbf{s}_1 \\ \vdots \\ \mathbf{s}_T \end{bmatrix}}_{\mathbf{S}_+} = \underbrace{\begin{bmatrix} \mathbf{s}_0 \\ \vdots \\ \mathbf{s}_{T-1} \end{bmatrix}}_{\mathbf{S}_-} \mathbf{F}^\top + \underbrace{\begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{bmatrix}}_{\mathbf{W}} \\&= \arg \min_{\mathbf{F}} \|\mathbf{S}_+ - \mathbf{S}_- \mathbf{F}^\top\|_F^2\end{aligned}$$



# Recap: System Identification

Least-squares problem with matrix variable

$$\begin{aligned}\hat{\mathbf{F}} &= \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t\|_{\ell_2}^2 \\&= \arg \min_{\mathbf{F}} \sum_{t=0}^{T-1} \|\mathbf{w}_t\|_{\ell_2}^2 \text{ s.t. } \underbrace{\begin{bmatrix} \mathbf{s}_1 \\ \vdots \\ \mathbf{s}_T \end{bmatrix}}_{\mathbf{S}_+} = \underbrace{\begin{bmatrix} \mathbf{s}_0 \\ \vdots \\ \mathbf{s}_{T-1} \end{bmatrix}}_{\mathbf{S}_-} \mathbf{F}^\top + \underbrace{\begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{bmatrix}}_{\mathbf{W}} \\&= \arg \min_{\mathbf{F}} \|\mathbf{S}_+ - \mathbf{S}_- \mathbf{F}^\top\|_{\mathbf{F}}^2 \\ \hat{\mathbf{F}} &= \mathbf{S}_+^\top \mathbf{S}_- (\mathbf{S}_-^\top \mathbf{S}_-)^{-1}\end{aligned}$$

# Recap: System Identification

Least-squares problem with matrix variable

$$\hat{\mathbf{F}} = \mathbf{S}_+^{\top} \mathbf{S}_- (\mathbf{S}_-^{\top} \mathbf{S}_-)^{-1}$$

- Unlike standard linear regression, the “features”  $(\mathbf{s}_0, \dots, \mathbf{s}_T)$  are not fixed or drawn independently
  - dependence:  $\mathbf{s}_t = \mathbf{F}^t \mathbf{s}_0 + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{w}_{\tau}$

## Recap: System Identification

Least-squares problem with matrix variable

$$\hat{\mathbf{F}} = \mathbf{S}_+^\top \mathbf{S}_- (\mathbf{S}_-^\top \mathbf{S}_-)^{-1}$$

- Unlike standard linear regression, the “features”  $(\mathbf{s}_0, \dots, \mathbf{s}_T)$  are not fixed or drawn independently
  - dependence:  $\mathbf{s}_t = \mathbf{F}^t \mathbf{s}_0 + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{w}_\tau$
- Similar to standard linear regression, the estimation error (when  $\mathbf{S}_-$  is full rank) can be bounded as

$$\|\hat{\mathbf{F}} - \mathbf{F}\| = \|\mathbf{W}^\top \mathbf{S}_- (\mathbf{S}_-^\top \mathbf{S}_-)^{-1}\|$$

# Recap: System Identification

Least-squares problem with matrix variable

$$\hat{\mathbf{F}} = \mathbf{S}_+^T \mathbf{S}_- (\mathbf{S}_-^T \mathbf{S}_-)^{-1}$$

- Unlike standard linear regression, the “features”  $(\mathbf{s}_0, \dots, \mathbf{s}_T)$  are not fixed or drawn independently
  - dependence:  $\mathbf{s}_t = \mathbf{F}^t \mathbf{s}_0 + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{w}_\tau$
- Similar to standard linear regression, the estimation error (when  $\mathbf{S}_-$  is full rank) can be bounded as

$$\begin{aligned} \|\hat{\mathbf{F}} - \mathbf{F}\| &= \|\mathbf{W}^T \mathbf{S}_- (\mathbf{S}_-^T \mathbf{S}_-)^{-1}\| \\ &\leq \frac{\|\mathbf{W}^T \mathbf{S}_-\|}{\lambda_{\min}(\mathbf{S}_-^T \mathbf{S}_-)} \end{aligned}$$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- System dynamics  $\mathbf{F} \in \mathbb{R}^{n \times n}$  and  $\mathbf{G} \in \mathbb{R}^{n \times n}$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- System dynamics  $\mathbf{F} \in \mathbb{R}^{n \times n}$  and  $\mathbf{G} \in \mathbb{R}^{n \times n}$
- Two types of inputs:

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- System dynamics  $\mathbf{F} \in \mathbb{R}^{n \times n}$  and  $\mathbf{G} \in \mathbb{R}^{n \times n}$
- Two types of inputs:
  - the disturbance  $\mathbf{w}_t \in \mathbb{R}^n$



# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- System dynamics  $\mathbf{F} \in \mathbb{R}^{n \times n}$  and  $\mathbf{G} \in \mathbb{R}^{n \times m}$
- Two types of inputs:
  - the disturbance  $\mathbf{w}_t \in \mathbb{R}^n$
  - the control input  $\mathbf{u}_t \in \mathbb{R}^m$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- System dynamics  $\mathbf{F} \in \mathbb{R}^{n \times n}$  and  $\mathbf{G} \in \mathbb{R}^{n \times n}$

- Two types of inputs:

- the disturbance  $\mathbf{w}_t \in \mathbb{R}^n$
- the control input  $\mathbf{u}_t \in \mathbb{R}^m$

- Identify the dynamics model  $\mathbf{s}_{0:T}, \mathbf{u}_{0:T} \rightarrow \hat{\mathbf{F}}, \hat{\mathbf{G}}$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- Suppose  $\rho(\mathbf{F}) > 1$ , can we choose  $\mathbf{u}_t$  to stabilize the system?

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- Suppose  $\rho(\mathbf{F}) > 1$ , can we choose  $\mathbf{u}_t$  to stabilize the system?

- Choose  $\mathbf{u}_t = -\mathbf{K}\mathbf{s}_t$ :

$$\mathbf{s}_{t+1} = (\mathbf{F} - \mathbf{G}\mathbf{K})\mathbf{s}_t + \mathbf{w}_t$$

# Today: Control Inputs

Linear dynamical system with control Inputs

- Observe states  $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T$

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- Suppose  $\rho(\mathbf{F}) > 1$ , can we choose  $\mathbf{u}_t$  to stabilize the system?

- Choose  $\mathbf{u}_t = -\mathbf{K}\mathbf{s}_t$ :

$$\mathbf{s}_{t+1} = (\mathbf{F} - \mathbf{G}\mathbf{K})\mathbf{s}_t + \mathbf{w}_t$$

- closed-loop dynamics  $\tilde{\mathbf{F}} = (\mathbf{F} - \mathbf{G}\mathbf{K}) \in \mathbb{R}^{n \times n}$

# Today: Nonlinear Dynamical Systems

Nonlinear dynamical systems with state observations

- state  $\mathbf{s}_t \in \mathbb{R}^n$
- input  $\mathbf{u}_t \in \mathbb{R}^m$
- noise  $\mathbf{w}_t \in \mathbb{R}^n$
- system dynamics  $\boldsymbol{\theta}_\star \in \mathbb{R}^d$

$$\mathbf{s}_{t+1} = \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta}_\star) + \mathbf{w}_t$$

# Today: Nonlinear Dynamical Systems

Nonlinear dynamical systems with state observations

$$\mathbf{s}_{t+1} = \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta}_\star) + \mathbf{w}_t$$

**Example:**

- Standard linear dynamical system
  - State equation:  $\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$
  - System dynamics:  $\boldsymbol{\theta}_\star = [\mathbf{F} \ \mathbf{G}]$

**Goal:** Learning the dynamics  $\boldsymbol{\theta}_\star$  from data



# Today: Nonlinear Dynamical Systems

Nonlinear dynamical systems with state observations

$$\mathbf{s}_{t+1} = \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta}_\star) + \mathbf{w}_t$$

**Example:**

- RNN type nonlinear dynamical system
  - State equation:  $\mathbf{s}_{t+1} = \phi(\mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t) + \mathbf{w}_t$
  - System dynamics:  $\boldsymbol{\theta}_\star = [\mathbf{F} \ \mathbf{G}]$

**Goal:** Learning the dynamics  $\boldsymbol{\theta}_\star$  from data

Run the system until time  $T$ , collect  $(\mathbf{s}_t, \mathbf{u}_t)_{t=0}^T$

# Learning from Finite Data

Run the system until time  $T$ , collect  $(\mathbf{s}_t, \mathbf{u}_t)_{t=0}^T$

- 1 Set loss function  $\hat{\mathcal{L}}(\boldsymbol{\theta}) = \frac{1}{2(\bar{T}-L)} \sum_{t=L}^{T-1} \|\mathbf{s}_{t+1} - \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta})\|_{\ell_2}^2$ .

# Learning from Finite Data

Run the system until time  $T$ , collect  $(\mathbf{s}_t, \mathbf{u}_t)_{t=0}^T$

- 1 Set loss function  $\hat{\mathcal{L}}(\boldsymbol{\theta}) = \frac{1}{2(\bar{T}-L)} \sum_{t=L}^{T-1} \|\mathbf{s}_{t+1} - \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta})\|_{\ell_2}^2$ .
- 2 Find  $\hat{\boldsymbol{\theta}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \hat{\mathcal{L}}(\boldsymbol{\theta})$  (e.g. via gradient descent)

# Learning from Finite Data

Run the system until time  $T$ , collect  $(\mathbf{s}_t, \mathbf{u}_t)_{t=0}^T$

- 1 Set loss function  $\hat{\mathcal{L}}(\boldsymbol{\theta}) = \frac{1}{2(\bar{T}-L)} \sum_{t=L}^{T-1} \|\mathbf{s}_{t+1} - \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta})\|_{\ell_2}^2$ .
- 2 Find  $\hat{\boldsymbol{\theta}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \hat{\mathcal{L}}(\boldsymbol{\theta})$  (e.g. via gradient descent)
- 3 **Hope that  $\hat{\boldsymbol{\theta}} \approx \boldsymbol{\theta}_\star$**

# Learning from Finite Data

Run the system until time  $T$ , collect  $(\mathbf{s}_t, \mathbf{u}_t)_{t=0}^T$

- 1 Set loss function  $\hat{\mathcal{L}}(\boldsymbol{\theta}) = \frac{1}{2(T-L)} \sum_{t=L}^{T-1} \|\mathbf{s}_{t+1} - \phi(\mathbf{s}_t, \mathbf{u}_t; \boldsymbol{\theta})\|_{\ell_2}^2$ .
- 2 Find  $\hat{\boldsymbol{\theta}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \hat{\mathcal{L}}(\boldsymbol{\theta})$  (e.g. via gradient descent)
- 3 **Hope that  $\hat{\boldsymbol{\theta}} \approx \boldsymbol{\theta}_\star$**

## Challenges:

- temporal dependence
- nonlinearity
- finite samples

# Nonlinear Stability

**Nonlinear dynamical systems:** Use  $\rho$ -stability

## Definition ( $\rho$ -stabilized system)

- (1) Pick inputs  $\mathbf{u}_t = \pi(\mathbf{s}_t) + \mathbf{z}_t$ . Fix  $(\mathbf{z}_\tau)_{\tau=0}^{t-1}$  and  $(\mathbf{w}_\tau)_{\tau=0}^{t-1}$ .
- (2) Denote the state sequence resulting from initial state  $\mathbf{s}_0 = \boldsymbol{\alpha}$  by  $\mathbf{s}_t(\boldsymbol{\alpha})$ .
- (3) There exists  $C_\rho \geq 1$  and  $\rho \in (0, 1)$  such that for all  $\boldsymbol{\alpha}$ ,  $(\mathbf{z}_t)_{t \geq 0}$  and  $(\mathbf{w}_t)_{t \geq 0}$ , we have

$$\|\mathbf{s}_t(\boldsymbol{\alpha}) - \mathbf{s}_t(0)\|_{\ell_2} \leq C_\rho \rho^t \|\boldsymbol{\alpha}\|_{\ell_2},$$

$\rho$  corresponds to *nonlinear spectral radius* (not easy to calculate).

# Nonlinear Stability

**Nonlinear dynamical systems:** Use  $\rho$ -stability

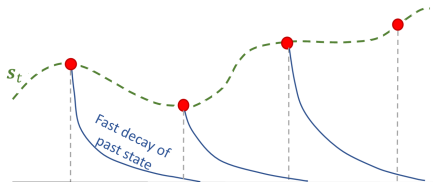
## Definition ( $\rho$ -stabilized system)

- (1) Pick inputs  $\mathbf{u}_t = \pi(\mathbf{s}_t) + \mathbf{z}_t$ . Fix  $(\mathbf{z}_\tau)_{\tau=0}^{t-1}$  and  $(\mathbf{w}_\tau)_{\tau=0}^{t-1}$ .
- (2) Denote the state sequence resulting from initial state  $\mathbf{s}_0 = \alpha$  by  $\mathbf{s}_t(\alpha)$ .
- (3) There exists  $C_\rho \geq 1$  and  $\rho \in (0, 1)$  such that for all  $\alpha$ ,  $(\mathbf{z}_t)_{t \geq 0}$  and  $(\mathbf{w}_t)_{t \geq 0}$ , we have

$$\|\mathbf{s}_t(\alpha) - \mathbf{s}_t(0)\|_{\ell_2} \leq C_\rho \rho^t \|\alpha\|_{\ell_2},$$

$\rho$  corresponds to *nonlinear spectral radius* (not easy to calculate).

**Key observation:** System forgets the past quickly





**Linear dynamical systems:** Gelfand's formula

$$\rho(\mathbf{F}) = \lim_{k \rightarrow \infty} \|\mathbf{F}^k\|^{1/k}$$

**Linear dynamical systems:** Gelfand's formula

$$\rho(\mathbf{F}) = \lim_{k \rightarrow \infty} \|\mathbf{F}^k\|^{1/k}$$

$$\|\mathbf{s}_t(\alpha) - \mathbf{s}_t(0)\|_{\ell_2} = \|\mathbf{F}^t \alpha - \mathbf{F}^t 0\|_{\ell_2}$$

**Linear dynamical systems:** Gelfand's formula

$$\rho(\mathbf{F}) = \lim_{k \rightarrow \infty} \|\mathbf{F}^k\|^{1/k}$$

$$\begin{aligned}\|\mathbf{s}_t(\alpha) - \mathbf{s}_t(0)\|_{\ell_2} &= \|\mathbf{F}^t \alpha - \mathbf{F}^t 0\|_{\ell_2} \\ &\leq \|\mathbf{F}^t\| \|\alpha\|_{\ell_2}\end{aligned}$$

**Linear dynamical systems:** Gelfand's formula

$$\rho(\mathbf{F}) = \lim_{k \rightarrow \infty} \|\mathbf{F}^k\|^{1/k}$$

$$\begin{aligned}\|\mathbf{s}_t(\boldsymbol{\alpha}) - \mathbf{s}_t(0)\|_{\ell_2} &= \|\mathbf{F}^t \boldsymbol{\alpha} - \mathbf{F}^t 0\|_{\ell_2} \\ &\leq \|\mathbf{F}^t\| \|\boldsymbol{\alpha}\|_{\ell_2} \\ &\leq C_\rho \rho(\mathbf{F})^t \|\boldsymbol{\alpha}\|_{\ell_2}\end{aligned}$$

# Assumptions on the System and Inputs

## Assumption (Stability)

*The closed loop system  $\tilde{\phi}$  is  $\rho$ -stable.*

# Assumptions on the System and Inputs

## Assumption (Stability)

*The closed loop system  $\tilde{\phi}$  is  $\rho$ -stable.*

## Assumption (Boundedness)

*There exist scalars  $B, c_w, \sigma > 0$ , such that  $(\mathbf{z}_t)_{t \geq 0} \stackrel{i.i.d.}{\sim} \mathcal{D}_z$  and  $(\mathbf{w}_t)_{t \geq 0} \stackrel{i.i.d.}{\sim} \mathcal{D}_w$  obey  $\|\tilde{\phi}(0, \mathbf{z}_t; \boldsymbol{\theta}_*)\|_{\ell_2} \leq B\sqrt{n}$  and  $\|\mathbf{w}_t\|_{\ell_\infty} \leq c_w\sigma$  for  $0 \leq t \leq T - 1$  with probability at least  $1 - p_0$  over the generation of data.*

# Mixing-time Approach

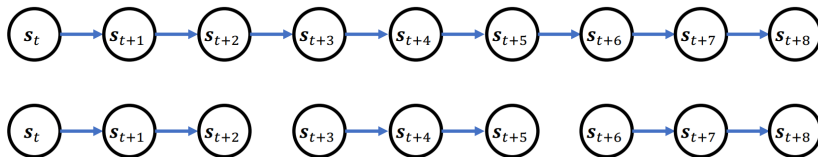
Suppose noise  $(\mathbf{w}_t)_{t \geq 0}$  and excitation  $(\mathbf{z}_t)_{t \geq 0}$  are i.i.d and  $\mathbf{s}_0 = 0$ .

- Leverage mixing-time argument with stability

# Mixing-time Approach

Suppose noise  $(\mathbf{w}_t)_{t \geq 0}$  and excitation  $(\mathbf{z}_t)_{t \geq 0}$  are i.i.d and  $\mathbf{s}_0 = 0$ .

- Leverage mixing-time argument with stability



## Related works

- Mixing-time/stability: Yu (Annals of Prob.'94), Mohri & Rostamizadeh (JMLR'10), Miller & Hardt (ICLR'19), Oymak (COLT'19)
- Martingale-based: Recht, Rakhlin, Tewari & coauthors

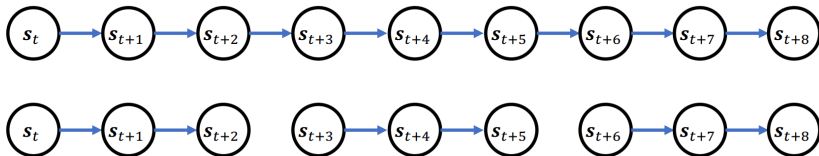


# Mixing-time Approach

## Definition (Multiple trajectory loss)

Let  $\mathbf{s}_{t+1,L}, \mathbf{s}_{t,L-1}$  be  $L$ -truncated and  $L-1$ -truncated states at time  $t+1$  and  $t$  respectively. We define the truncated (empirical) risk as

$$\hat{\mathcal{L}}^{\text{tr}}(\theta) := \frac{1}{2(T-L)} \sum_{t=L}^{T-1} \|\mathbf{s}_{t+1,L} - \tilde{\phi}(\mathbf{s}_{t,L-1}, \mathbf{z}_t; \theta)\|_{\ell_2}^2.$$



## Theorem (Difference between single and multiple trajectories)

*Under  $\rho$ -stability and boundedness assumptions, with probability at least  $1 - p_0$ , for all  $\boldsymbol{\theta} \in \mathcal{B}^d(\boldsymbol{\theta}_\star, r)$ , we have*

$$\|\nabla \hat{\mathcal{L}}(\boldsymbol{\theta}) - \nabla \hat{\mathcal{L}}^{tr}(\boldsymbol{\theta})\|_{\ell_2} \lesssim C_\rho \rho^{L-1} (c_w \sigma + C_{\tilde{\phi}} \|\boldsymbol{\theta} - \boldsymbol{\theta}_\star\|_{\ell_2}).$$

To concretely show how stability helps, we define the following loss function, obtained from i.i.d. samples at time  $L - 1$  and can be used as a proxy for  $\mathbb{E}[\hat{\mathcal{L}}]$ .

# Optimization Landscape

To concretely show how stability helps, we define the following loss function, obtained from i.i.d. samples at time  $L - 1$  and can be used as a proxy for  $\mathbb{E}[\hat{\mathcal{L}}]$ .

## Definition (Population Loss)

Suppose  $\mathbf{s}_0 = 0$ . Let  $(\mathbf{z}_t)_{t \geq 0} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_z$  and  $(\mathbf{w}_t)_{t \geq 0} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_w$ . The auxiliary loss is defined as the expected loss at timestamp  $L - 1$ , that is,

$$\mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}[\|\mathbf{s}_L - \tilde{\phi}(\mathbf{s}_{L-1}, \mathbf{z}_{L-1}; \boldsymbol{\theta})\|_{\ell_2}^2].$$

# Uniform Convergence of Empirical Gradient

Define the empirical and population losses,

$$\hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}_i) \quad \text{and} \quad \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{x})].$$

# Uniform Convergence of Empirical Gradient

Define the empirical and population losses,

$$\hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}_i) \quad \text{and} \quad \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{x})].$$

## Assumption (Lipschitz gradients)

*There exist numbers  $L_{\mathcal{D}}, p_0 > 0$  such that with probability at least  $1 - p_0$  over the generation of data, for all pairs  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathcal{B}^d(\boldsymbol{\theta}_*, r)$ ,*

$$\max(\|\nabla \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) - \nabla \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}')\|_{\ell_2}, \|\nabla \hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}) - \nabla \hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}')\|_{\ell_2}) \leq L_{\mathcal{D}} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_{\ell_2}.$$

# Uniform Convergence of Empirical Gradient

Define the empirical and population losses,

$$\hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}_i) \quad \text{and} \quad \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{x})].$$

## Assumption (Subexponential gradient noise)

*There exist scalars  $K, \sigma_0 > 0$  such that, given  $\mathbf{x} := (\mathbf{h}_L, \mathbf{h}_{L-1}, \mathbf{z}_{L-1}) \sim \mathcal{D}$ , at any point  $\boldsymbol{\theta}$ , the subexponential norm of the gradient is upper bounded as a function of the noise level  $\sigma_0$  and distance to the population minimizer via*

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) - \mathbb{E}[\nabla \mathcal{L}(\boldsymbol{\theta}, \mathbf{x})]\|_{\psi_1} \leq \sigma_0 + K\|\boldsymbol{\theta} - \boldsymbol{\theta}_\star\|_{\ell_2}.$$

# Uniform Convergence of Empirical Gradient

Define the empirical and population losses,

$$\hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}_i) \quad \text{and} \quad \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{x})].$$

## Assumption (Subexponential gradient noise)

*There exist scalars  $K, \sigma_0 > 0$  such that, given  $\mathbf{x} := (\mathbf{h}_L, \mathbf{h}_{L-1}, \mathbf{z}_{L-1}) \sim \mathcal{D}$ , at any point  $\boldsymbol{\theta}$ , the subexponential norm of the gradient is upper bounded as a function of the noise level  $\sigma_0$  and distance to the population minimizer via*

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) - \mathbb{E}[\nabla \mathcal{L}(\boldsymbol{\theta}, \mathbf{x})]\|_{\psi_1} \leq \sigma_0 + K \|\boldsymbol{\theta} - \boldsymbol{\theta}_\star\|_{\ell_2}.$$

For a random variable  $X \in \mathbb{R}$ ,  $\|X\|_{\psi_1} := \sup_{k \geq 1} \frac{(\mathbb{E}[|X|^k])^{1/k}}{k}$

For a random vector  $\mathbf{x} \in \mathbb{R}^n$  is  $\|\mathbf{x}\|_{\psi_1} := \sup_{\mathbf{v} \in \mathcal{S}^{n-1}} \|\mathbf{v}^\top \mathbf{x}\|_{\psi_1}$



# Uniform Convergence of Empirical Gradient

## Theorem (Uniform gradient convergence)

*Suppose the gradients of  $\mathcal{L}_{\mathcal{D}}$  and  $\mathcal{L}_{\mathcal{S}}$  obey the above two Assumptions. Then, for all  $\boldsymbol{\theta} \in \mathcal{B}^d(\boldsymbol{\theta}_*, r)$ , with probability  $1 - p_0 - \log(\frac{Kr}{\sigma_0}) \exp(-100d)$ , we have*

$$\|\nabla \hat{\mathcal{L}}_{\mathcal{S}}(\boldsymbol{\theta}) - \nabla \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta})\|_{\ell_2} \lesssim (\sigma_0 + K\|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_{\ell_2}) C_{\log} \sqrt{\frac{d}{N}}.$$

- Here,  $C_{\log} = \log(3(L_{\mathcal{D}}N/K + 1))$

# Optimization Landscape

A special case of Polyak-Lojasiewicz inequality and provides a generalization of strong convexity to nonconvex functions.

## Assumption (One-point convexity & smoothness)

*There exist scalars  $\beta \geq \alpha > 0$  such that the auxiliary loss  $\mathcal{L}_{\mathcal{D}}(\theta)$  satisfies*

$$\begin{aligned}\langle \theta - \theta_*, \nabla \mathcal{L}_{\mathcal{D}}(\theta) \rangle &\geq \alpha \|\theta - \theta_*\|_{\ell_2}^2, \\ \|\nabla \mathcal{L}_{\mathcal{D}}(\theta)\|_{\ell_2} &\leq \beta \|\theta - \theta_*\|_{\ell_2}.\end{aligned}$$

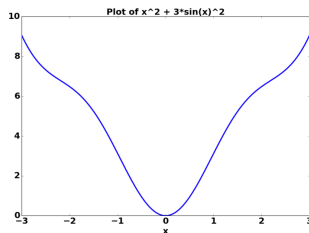


Figure 1: Not a convex function. But it also does not look challenging to optimize. Polyak-Lojasiewicz condition can capture such functions (and beyond).

# Main Result

## Theorem (Finite sample learning)

- *Suppose  $\rho < 1$  and Assumptions above hold*

# Main Result

## Theorem (Finite sample learning)

- *Suppose  $\rho < 1$  and Assumptions above hold*
- *Suppose  $\mathcal{L}_{\mathcal{D}}$  satisfies OPCS with  $\alpha, \beta > 0$*

# Main Result

## Theorem (Finite sample learning)

- Suppose  $\rho < 1$  and Assumptions above hold
- Suppose  $\mathcal{L}_{\mathcal{D}}$  satisfies OPCS with  $\alpha, \beta > 0$
- Suppose trajectory length  $T$  obeys

$$T \gtrsim L(N + 1), \quad \text{where} \quad N \propto d \log(N), \quad L \propto 1 + \log(N) / \log(\rho^{-1})$$

# Main Result

## Theorem (Finite sample learning)

- Suppose  $\rho < 1$  and Assumptions above hold
- Suppose  $\mathcal{L}_{\mathcal{D}}$  satisfies OPCS with  $\alpha, \beta > 0$
- Suppose trajectory length  $T$  obeys

$$T \gtrsim L(N+1), \quad \text{where} \quad N \propto d \log(N), \quad L \propto 1 + \log(N)/\log(\rho^{-1})$$

Pick a proper constant learning rate and starting from  $\theta_0 \in \mathcal{B}^d(\theta_*, r)$ , run gradient iterates  $\theta_{\tau+1} = \theta_{\tau} - \eta \nabla \hat{\mathcal{L}}(\theta_{\tau})$ . All iterates satisfy

$$\|\theta_{\tau} - \theta_{*}\|_{\ell_2} \leq \underbrace{\left(1 - \frac{\alpha^2}{128\beta^2}\right)^{\tau}}_{\text{fast convergence}} \|\theta_0 - \theta_{*}\|_{\ell_2} + \underbrace{\frac{c\sigma_0}{\alpha} \sqrt{\frac{d}{N}}}_{\text{statistical error}}.$$

# Main Result

## Theorem (Finite sample learning)

- Suppose  $\rho < 1$  and Assumptions above hold
- Suppose  $\mathcal{L}_{\mathcal{D}}$  satisfies OPCS with  $\alpha, \beta > 0$
- Suppose trajectory length  $T$  obeys

$$T \gtrsim L(N+1), \quad \text{where} \quad N \propto d \log(N), \quad L \propto 1 + \log(N)/\log(\rho^{-1})$$

Pick a proper constant learning rate and starting from  $\theta_0 \in \mathcal{B}^d(\theta_*, r)$ , run gradient iterates  $\theta_{\tau+1} = \theta_{\tau} - \eta \nabla \hat{\mathcal{L}}(\theta_{\tau})$ . All iterates satisfy

$$\|\theta_{\tau} - \theta_{*}\|_{\ell_2} \leq \underbrace{\left(1 - \frac{\alpha^2}{128\beta^2}\right)^{\tau}}_{\text{fast convergence}} \|\theta_0 - \theta_{*}\|_{\ell_2} + \underbrace{\frac{c\sigma_0}{\alpha} \sqrt{\frac{d}{N}}}_{\text{statistical error}}.$$

**Remark:** Need to calculate  $\alpha, \beta, \sigma_0$  etc. for specific scenario e.g.

$$\mathbf{s}_{t+1} = \phi(\mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t) + \mathbf{w}_t$$

# Case Study

Linear Dynamical System:

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

$$[\mathbf{F} \ \mathbf{G}] = [\boldsymbol{\theta}_1 \ \cdots \ \boldsymbol{\theta}_n]^T \in \mathbb{R}^{n \times (n+m)}$$



# Case Study

Linear Dynamical System:  $\boxed{s_{t+1} = F s_t + G u_t + w_t}$

$$[F \ G] = [\theta_1 \ \cdots \ \theta_n]^T \in \mathbb{R}^{n \times (n+m)}$$

$$s_t = F^t s_0 + \sum_{\tau=0}^{t-1} F^{t-\tau-1} G u_{\tau} + \sum_{\tau=0}^{t-1} F^{t-\tau-1} w_{\tau}$$

# Case Study

Linear Dynamical System:  $\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$

$$[\mathbf{F} \ \mathbf{G}] = [\boldsymbol{\theta}_1 \ \cdots \ \boldsymbol{\theta}_n]^T \in \mathbb{R}^{n \times (n+m)}$$

$$\begin{aligned}\mathbf{s}_t &= \mathbf{F}^t \mathbf{s}_0 + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{G} \mathbf{u}_\tau + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{w}_\tau \\ &= \mathbf{F}^t \mathbf{s}_0 + \boldsymbol{\Gamma}_t^G \begin{bmatrix} \mathbf{u}_0 \\ \vdots \\ \mathbf{u}_{T-1} \end{bmatrix} + \boldsymbol{\Gamma}_t \begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{bmatrix}\end{aligned}$$

# Case Study

Linear Dynamical System:  $\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$

$$[\mathbf{F} \ \mathbf{G}] = [\boldsymbol{\theta}_1 \ \cdots \ \boldsymbol{\theta}_n]^T \in \mathbb{R}^{n \times (n+m)}$$

$$\begin{aligned}\mathbf{s}_t &= \mathbf{F}^t \mathbf{s}_0 + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{G} \mathbf{u}_\tau + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{w}_\tau \\ &= \mathbf{F}^t \mathbf{s}_0 + \boldsymbol{\Gamma}_t^G \begin{bmatrix} \mathbf{u}_0 \\ \vdots \\ \mathbf{u}_{T-1} \end{bmatrix} + \boldsymbol{\Gamma}_t \begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{bmatrix}\end{aligned}$$

where  $\boldsymbol{\Gamma}_t^G := [\mathbf{F}^{t-1} \mathbf{G} \ \mathbf{F}^{t-2} \mathbf{G} \ \cdots \ \mathbf{G}]$  and  $\boldsymbol{\Gamma}_t := [\mathbf{F}^{t-1} \ \mathbf{F}^{t-2} \ \cdots \ \mathbf{I}_n]$

# Case Study

Linear Dynamical System:  $\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$

$$[\mathbf{F} \ \mathbf{G}] = [\boldsymbol{\theta}_1 \ \cdots \ \boldsymbol{\theta}_n]^T \in \mathbb{R}^{n \times (n+m)}$$

$$\begin{aligned}\mathbf{s}_t &= \mathbf{F}^t \mathbf{s}_0 + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{G} \mathbf{u}_\tau + \sum_{\tau=0}^{t-1} \mathbf{F}^{t-\tau-1} \mathbf{w}_\tau \\ &= \mathbf{F}^t \mathbf{s}_0 + \boldsymbol{\Gamma}_t^G \begin{bmatrix} \mathbf{u}_0 \\ \vdots \\ \mathbf{u}_{T-1} \end{bmatrix} + \boldsymbol{\Gamma}_t \begin{bmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{bmatrix}\end{aligned}$$

where  $\boldsymbol{\Gamma}_t^G := [\mathbf{F}^{t-1} \mathbf{G} \ \mathbf{F}^{t-2} \mathbf{G} \ \cdots \ \mathbf{G}]$  and  $\boldsymbol{\Gamma}_t := [\mathbf{F}^{t-1} \ \mathbf{F}^{t-2} \ \cdots \ \mathbf{I}_n]$

$\boldsymbol{\Gamma}_t^G \boldsymbol{\Gamma}_t^{G\top}$  and  $\boldsymbol{\Gamma}_t \boldsymbol{\Gamma}_t^\top$  are the finite time controllability Gramians for the control and noise inputs, respectively.

# Case Study

Linear Dynamical System:

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}u_t + \mathbf{w}_t$$

- Suppose  $\mathbf{s}_0 = 0$

# Case Study

Linear Dynamical System:

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- Suppose  $\mathbf{s}_0 = 0$
- Suppose  $\mathbf{u}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_m)$  and  $\mathbf{w}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_w^2 \mathbf{I}_n)$

# Case Study

Linear Dynamical System:

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

- Suppose  $\mathbf{s}_0 = 0$
- Suppose  $\mathbf{u}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_m)$  and  $\mathbf{w}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_w^2 \mathbf{I}_n)$
- $\mathbf{s}_t \sim \mathcal{N}(0, \mathbf{\Gamma}_t^G \mathbf{\Gamma}_t^{G\top} + \sigma_w^2 \mathbf{\Gamma}_t \mathbf{\Gamma}_t^\top)$

# Case Study

Linear Dynamical System:

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

## Theorem (simplified)

- Suppose  $\rho(\mathbf{F}) < 1$ .
- Suppose  $\mathbf{u}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_p)$  and  $\mathbf{w}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_w^2 \mathbf{I}_n)$ .
- Suppose trajectory length  $T$  obeys

$$T \gtrsim L(N+1), \text{ where } N \propto (n+m)\log(N), \quad L \propto 1 + \log(N)/\log(\rho^{-1})$$



# Case Study

Linear Dynamical System:

$$\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t + \mathbf{G}\mathbf{u}_t + \mathbf{w}_t$$

## Theorem (simplified)

- Suppose  $\rho(\mathbf{F}) < 1$ .
- Suppose  $\mathbf{u}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_p)$  and  $\mathbf{w}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_w^2 \mathbf{I}_n)$ .
- Suppose trajectory length  $T$  obeys

$$T \gtrsim L(N+1), \text{ where } N \propto (n+m) \log(N), \quad L \propto 1 + \log(N)/\log(\rho^{-1})$$

With proper learning rate and the initialization  $[\mathbf{F}^{(0)} \quad \mathbf{G}^{(0)}] = 0$ , all gradient descent iterates satisfy

$$\|\boldsymbol{\theta}_k^{(\tau)} - \boldsymbol{\theta}_k^*\|_{\ell_2} \leq \left(1 - \frac{\gamma_-^2}{128\gamma_+^2}\right)^\tau \|\boldsymbol{\theta}_k^{(0)} - \boldsymbol{\theta}_k^*\|_{\ell_2} + \frac{c\sigma}{\gamma_-} \sqrt{\gamma_+} \log(9N) \sqrt{\frac{n+m}{N}}.$$

# Another Interesting Nonlinear Dynamical System

Bilinear dynamical systems with state observations

- state  $\mathbf{s}_t \in \mathbb{R}^n$
- input  $\mathbf{u}_t \in \mathbb{R}^m$

$$\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m u_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$$

- noise  $\mathbf{w}_t \in \mathbb{R}^n$
- system dynamics  $\{\mathbf{F}_k\}_{k=0}^m \in (\mathbb{R}^{n \times n})^{m+1}$

# Another Interesting Nonlinear Dynamical System

Bilinear dynamical systems with state observations

- state  $\mathbf{s}_t \in \mathbb{R}^n$
- input  $\mathbf{u}_t \in \mathbb{R}^m$

$$\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m u_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$$

- noise  $\mathbf{w}_t \in \mathbb{R}^n$
- system dynamics  $\{\mathbf{F}_k\}_{k=0}^m \in (\mathbb{R}^{n \times n})^{m+1}$

**Goal:** Learning the dynamics  $\{\mathbf{F}_k\}_{k=0}^m$  from data

# Learning Without Mixing

Run the system until time  $T$ , collect a single  $(\mathbf{s}_t, \mathbf{u}_t)_{t=1}^T$

Run the system until time  $T$ , collect a single  $(\mathbf{s}_t, \mathbf{u}_t)_{t=1}^T$

- System dynamics  $\mathbf{F}^* := [\mathbf{F}_0 \quad \sigma_{\mathbf{u}} \mathbf{F}_1 \quad \cdots \quad \sigma_{\mathbf{u}} \mathbf{F}_m] \in \mathbb{R}^{n \times n(m+1)}$

Run the system until time  $T$ , collect a single  $(\mathbf{s}_t, \mathbf{u}_t)_{t=1}^T$

- System dynamics  $\mathbf{F}^* := [\mathbf{F}_0 \quad \sigma_{\mathbf{u}} \mathbf{F}_1 \quad \cdots \quad \sigma_{\mathbf{u}} \mathbf{F}_m] \in \mathbb{R}^{n \times n(m+1)}$
- State  $\tilde{\mathbf{s}}_t := [\mathbf{s}_t^\top \quad \sigma_{\mathbf{u}}^{-1} \mathbf{u}_t[1] \mathbf{s}_t^\top \quad \cdots \quad \sigma_{\mathbf{u}}^{-1} \mathbf{u}_t[m] \mathbf{s}_t^\top]^\top \in \mathbb{R}^{n(m+1)}$

# Learning Without Mixing

Run the system until time  $T$ , collect a single  $(\mathbf{s}_t, \mathbf{u}_t)_{t=1}^T$

- System dynamics  $\mathbf{F}^* := [\mathbf{F}_0 \quad \sigma_{\mathbf{u}} \mathbf{F}_1 \quad \cdots \quad \sigma_{\mathbf{u}} \mathbf{F}_m] \in \mathbb{R}^{n \times n(m+1)}$
- State  $\tilde{\mathbf{s}}_t := [\mathbf{s}_t^\top \quad \sigma_{\mathbf{u}}^{-1} \mathbf{u}_t[1] \mathbf{s}_t^\top \quad \cdots \quad \sigma_{\mathbf{u}}^{-1} \mathbf{u}_t[m] \mathbf{s}_t^\top]^\top \in \mathbb{R}^{n(m+1)}$
- $\mathbf{s}_{t+1} = \mathbf{F}^* \tilde{\mathbf{s}}_t + \mathbf{w}_{t+1}$

# Learning Without Mixing

- Suppose  $\{\mathbf{u}_t\}_{t=0}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{u}}^2 \mathbf{I}_m)$  and  $\{\mathbf{w}_t\}_{t=1}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}_n)$ .



# Learning Without Mixing

- Suppose  $\{\mathbf{u}_t\}_{t=0}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{u}}^2 \mathbf{I}_m)$  and  $\{\mathbf{w}_t\}_{t=1}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}_n)$ .
- Set the loss function  $\hat{\mathcal{L}}(\mathbf{F}) = \frac{1}{2T} \sum_{t=1}^T \|\mathbf{s}_{t+1} - \mathbf{F} \tilde{\mathbf{s}}_t\|_{\ell_2}^2$ .

# Learning Without Mixing

- Suppose  $\{\mathbf{u}_t\}_{t=0}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{u}}^2 \mathbf{I}_m)$  and  $\{\mathbf{w}_t\}_{t=1}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}_n)$ .
- Set the loss function  $\hat{\mathcal{L}}(\mathbf{F}) = \frac{1}{2T} \sum_{t=1}^T \|\mathbf{s}_{t+1} - \mathbf{F} \tilde{\mathbf{s}}_t\|_{\ell_2}^2$ .
- Over-determined problem, solution  $\hat{\mathbf{F}} = \mathbf{S}_+^{\top} \mathbf{S}_- (\mathbf{S}_+^{\top} \mathbf{S}_-)^{-1}$ .

- where  $\mathbf{S}_+ := \begin{bmatrix} \mathbf{s}_2^{\top} \\ \vdots \\ \mathbf{s}_{T+1}^{\top} \end{bmatrix}$ ,  $\mathbf{S}_- := \begin{bmatrix} \tilde{\mathbf{s}}_1^{\top} \\ \vdots \\ \tilde{\mathbf{s}}_T^{\top} \end{bmatrix}$ ,  $\mathbf{W} := \begin{bmatrix} \mathbf{w}_2^{\top} \\ \vdots \\ \mathbf{w}_{T+1}^{\top} \end{bmatrix}$ .

# Learning Without Mixing

- Suppose  $\{\mathbf{u}_t\}_{t=0}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{u}}^2 \mathbf{I}_m)$  and  $\{\mathbf{w}_t\}_{t=1}^{\infty} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}_n)$ .
- Set the loss function  $\hat{\mathcal{L}}(\mathbf{F}) = \frac{1}{2T} \sum_{t=1}^T \|\mathbf{s}_{t+1} - \mathbf{F} \tilde{\mathbf{s}}_t\|_{\ell_2}^2$ .
- Over-determined problem, solution  $\hat{\mathbf{F}} = \mathbf{S}_+^{\top} \mathbf{S}_- (\mathbf{S}_-^{\top} \mathbf{S}_-)^{-1}$ .
- where  $\mathbf{S}_+ := \begin{bmatrix} \mathbf{s}_2^{\top} \\ \vdots \\ \mathbf{s}_{T+1}^{\top} \end{bmatrix}$ ,  $\mathbf{S}_- := \begin{bmatrix} \tilde{\mathbf{s}}_1^{\top} \\ \vdots \\ \tilde{\mathbf{s}}_T^{\top} \end{bmatrix}$ ,  $\mathbf{W} := \begin{bmatrix} \mathbf{w}_2^{\top} \\ \vdots \\ \mathbf{w}_{T+1}^{\top} \end{bmatrix}$ .
- Estimation error  $\|\hat{\mathbf{F}} - \mathbf{F}^*\| \leq \frac{\|\mathbf{W}^{\top} \mathbf{S}_-\|}{\lambda_{\min}(\mathbf{S}_-^{\top} \mathbf{S}_-)}.$

**Challenges:**

## Challenges:

- randomness in dynamical behavior

## Challenges:

- randomness in dynamical behavior
- distribution of states

## Challenges:

- randomness in dynamical behavior
- distribution of states
- temporal dependence
- finite samples

# Stability in Bilinear Dynamical Systems

**State equation:**  $\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m \mathbf{u}_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$

## Definition (Mean-square stability)

The bilinear system is mean-square stable (MSS) if there exists  $\mathbf{s}_\infty \in \mathbb{R}^n$  and  $\mathbf{\Sigma}_\infty \in \mathbb{R}_+^{n \times n}$ , such that for any initial state  $\mathbf{x}_0$ , as  $t \rightarrow \infty$ , we have

$$\|\mathbb{E}[\mathbf{s}_t] - \mathbf{s}_\infty\|_{\ell_2} \rightarrow 0, \quad \|\mathbb{E}[\mathbf{s}_t \mathbf{s}_t^\top] - \mathbf{\Sigma}_\infty\| \rightarrow 0. \quad (1)$$

Here the expectation is over the input sequence  $\{\mathbf{u}_t\}_{t=0}^\infty$ , the noise process  $\{\mathbf{w}_t\}_{t=1}^\infty$  and the initial state  $\mathbf{x}_0$ . In the noise free case ( $\mathbf{w}_t = 0$ ), we have  $\mathbf{x}_\infty = 0$  and  $\mathbf{\Sigma}_\infty = 0$ .

## Related works

- 1 C. Kubrusly and O. Costa, "Mean square stability conditions for discrete stochastic bilinear systems," IEEE Transactions on Automatic Control, vol. 30, no. 11, pp. 1082-1087, 1985.



# BMSB Condition

- Does the random process  $\{\tilde{\mathbf{s}}_t\}_{t \geq 1}$  satisfies the martingale small-ball condition?

# BMSB Condition

- Does the random process  $\{\tilde{\mathbf{S}}_t\}_{t \geq 1}$  satisfies the martingale small-ball condition?

## Definition (Martingale small-ball)

- Let  $\{\mathcal{F}_t\}_{t \geq 1}$  denotes a filtration.

# BMSB Condition

- Does the random process  $\{\tilde{\mathbf{s}}_t\}_{t \geq 1}$  satisfies the martingale small-ball condition?

## Definition (Martingale small-ball)

- Let  $\{\mathcal{F}_t\}_{t \geq 1}$  denotes a filtration.
- Let  $\{Z_t\}_{t \geq 1}$  be  $\{\mathcal{F}_t\}_{t \geq 1}$ -adapted random process taking values in  $\mathbb{R}$ .

# BMSB Condition

- Does the random process  $\{\tilde{\mathbf{s}}_t\}_{t \geq 1}$  satisfies the martingale small-ball condition?

## Definition (Martingale small-ball)

- Let  $\{\mathcal{F}_t\}_{t \geq 1}$  denotes a filtration.
- Let  $\{Z_t\}_{t \geq 1}$  be  $\{\mathcal{F}_t\}_{t \geq 1}$ -adapted random process taking values in  $\mathbb{R}$ .
- $\{Z_t\}_{t \geq 1}$  satisfies the  $(k, \nu, p)$ -block martingale small-ball (BMSB) condition if, for any  $j \geq 0$ , one has  $\frac{1}{k} \sum_{i=1}^k \mathbb{P}(|Z_{j+i}| \geq \nu \mid \mathcal{F}_j) \geq p$  almost surely.

# BMSB Condition

- Does the random process  $\{\tilde{\mathbf{s}}_t\}_{t \geq 1}$  satisfies the martingale small-ball condition?

## Definition (Martingale small-ball)

- Let  $\{\mathcal{F}_t\}_{t \geq 1}$  denotes a filtration.
- Let  $\{Z_t\}_{t \geq 1}$  be  $\{\mathcal{F}_t\}_{t \geq 1}$ -adapted random process taking values in  $\mathbb{R}$ .
- $\{Z_t\}_{t \geq 1}$  satisfies the  $(k, \nu, p)$ -block martingale small-ball (BMSB) condition if, for any  $j \geq 0$ , one has  $\frac{1}{k} \sum_{i=1}^k \mathbb{P}(|Z_{j+i}| \geq \nu \mid \mathcal{F}_j) \geq p$  almost surely.
- A process  $\{\mathbf{x}_t\}_{t \geq 1}$  taking values in  $\mathbb{R}^d$  satisfies the  $(k, \Gamma_{\text{sb}}, p)$ -BMSB condition for  $\Gamma_{\text{sb}} \succ 0$  if, for any fixed  $\mathbf{v} \in \mathcal{S}^{d-1}$ , the process  $\{Z_t = \langle \mathbf{v}, \mathbf{x}_t \rangle\}_{t \geq 1}$  satisfies  $(k, \sqrt{\mathbf{v}^\top \Gamma_{\text{sb}} \mathbf{v}}, p)$ -BMSB.

# Main Result

**State equation:**  $\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m \mathbf{u}_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$

## Theorem (Bilinear system identification)

*Fix  $\delta \in (0, 1)$  and suppose we are given a single trajectory  $(\mathbf{x}_t, \mathbf{u}_t)_{t=1}^T$  of the bilinear dynamical system.*

# Main Result

**State equation:**  $\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m \mathbf{u}_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$

## Theorem (Bilinear system identification)

Fix  $\delta \in (0, 1)$  and suppose we are given a single trajectory  $(\mathbf{x}_t, \mathbf{u}_t)_{t=1}^T$  of the bilinear dynamical system.

- Suppose the inputs and noise are Gaussian.

# Main Result

**State equation:**  $\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m \mathbf{u}_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$

## Theorem (Bilinear system identification)

Fix  $\delta \in (0, 1)$  and suppose we are given a single trajectory  $(\mathbf{x}_t, \mathbf{u}_t)_{t=1}^T$  of the bilinear dynamical system.

- Suppose the inputs and noise are Gaussian.
- Suppose the bilinear system is *marginally mean-square stable*.



# Main Result

**State equation:**  $\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m \mathbf{u}_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$

## Theorem (Bilinear system identification)

Fix  $\delta \in (0, 1)$  and suppose we are given a single trajectory  $(\mathbf{x}_t, \mathbf{u}_t)_{t=1}^T$  of the bilinear dynamical system.

- Suppose the inputs and noise are Gaussian.
- Suppose the bilinear system is *marginally mean-square stable*.
- Suppose  $T \gtrsim n(m+1) + \log(12\Gamma/(\sigma_w^2\delta)) + \log(3/\delta)$ .

# Main Result

**State equation:**  $\mathbf{s}_{t+1} = \mathbf{F}_0 \mathbf{x}_t + \sum_{k=1}^m \mathbf{u}_t[k] \mathbf{F}_k \mathbf{x}_t + \mathbf{w}_{t+1}$

## Theorem (Bilinear system identification)

Fix  $\delta \in (0, 1)$  and suppose we are given a single trajectory  $(\mathbf{x}_t, \mathbf{u}_t)_{t=1}^T$  of the bilinear dynamical system.

- Suppose the inputs and noise are Gaussian.
- Suppose the bilinear system is *marginally mean-square stable*.
- Suppose  $T \gtrsim n(m+1) + \log(12\Gamma/(\sigma_w^2\delta)) + \log(3/\delta)$ .
- With probability at least  $1 - \delta$ , we have

$$\max \left\{ \begin{array}{l} \|\hat{\mathbf{F}}_0 - \mathbf{F}_0\|, \\ \{\sigma_u \|\hat{\mathbf{F}}_k - \mathbf{F}_k\|\}_{k=1}^m \end{array} \right\} \lesssim \sqrt{\frac{n(m+1) + \log(\frac{12\Gamma}{\sigma_w^2\delta}) + \log(\frac{3}{\delta})}{T}}.$$

Thanks

Thank you for your attention!