

Non-asymptotic and Accurate Learning of Nonlinear Dynamical Systems

Yahya Sattar Samet Oymak

Department of ECE, University of California Riverside

- **Dynamical systems** appear in
 - physical systems,
 - reinforcement learning & control
 - natural language processing (i.e. RNN, LSTM)



- **Goal:** Efficient learning guarantees for nonlinear systems
- **Challenge:** spatio-temporal dependencies, nonlinear state equation, single trajectory ...

State Equation

Nonlinear systems with state observations

- state $\mathbf{h}_t \in \mathbb{R}^n$
- input $\mathbf{u}_t \in \mathbb{R}^p$
- noise $\mathbf{w}_t \in \mathbb{R}^n$
- system dynamics $\boldsymbol{\theta}_* \in \mathbb{R}^d$

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

Example: A nonlinear linear dynamical system

State equation: $\mathbf{h}_{t+1} = \phi(\mathbf{A}\mathbf{h}_t + \mathbf{B}\mathbf{u}_t) + \mathbf{w}_t$.

Learning from Finite Data

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

Learning from Finite Data

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

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

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

Challenges:

- temporal dependence
- nonlinearity
- finite samples

Nonlinear systems: Use ρ -stability

Definition (ρ -stabilized system)

- (1) Pick inputs $\mathbf{u}_t = \boldsymbol{\pi}(\mathbf{h}_t) + \mathbf{z}_t$. 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 $\mathbf{h}_0 = \boldsymbol{\alpha}$ by $\mathbf{h}_t(\boldsymbol{\alpha})$.
- (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

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

ρ corresponds to *nonlinear spectral radius* (not easy to calculate).

Nonlinear systems: Use ρ -stability

Definition (ρ -stabilized system)

- (1) Pick inputs $\mathbf{u}_t = \boldsymbol{\pi}(\mathbf{h}_t) + \mathbf{z}_t$. 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 $\mathbf{h}_0 = \boldsymbol{\alpha}$ by $\mathbf{h}_t(\boldsymbol{\alpha})$.
- (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

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

ρ corresponds to *nonlinear spectral radius* (not easy to calculate).

Key observation: System forgets the past quickly

Assumptions on the System and Inputs

Assumption (Stability)

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

Assumption (Stability)

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

- a linear dynamical system is ρ -stable if spectral radius $\rho(\mathbf{A}_*) < 1$.

Assumption (Stability)

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

- a linear dynamical system is ρ -stable if spectral radius $\rho(\mathbf{A}_*) < 1$.
- a ρ -stable nonlinear system is a contractive system.

Assumptions on the System and Inputs

Assumption (Stability)

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

- a linear dynamical system is ρ -stable if spectral radius $\rho(\mathbf{A}_\star) < 1$.
- a ρ -stable nonlinear system is a contractive system.

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 $\|\tilde{\phi}(0, \mathbf{z}_t; \boldsymbol{\theta}_\star)\|_{l_2} \leq B\sqrt{n}$ and $\|\mathbf{w}_t\|_{l_\infty} \leq c_w\sigma$ for $0 \leq t \leq T - 1$ with probability at least $1 - p_0$ over the generation of data.

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 (Auxiliary Loss)

Suppose $\mathbf{h}_0 = 0$. 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,

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

Assumption (One-point convexity & smoothness (OPCS))

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

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

Assumption (One-point convexity & smoothness (OPCS))

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

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

- aka restricted secant inequality, and implies Polyak-Lojasiewicz condition.

Assumption (One-point convexity & smoothness (OPCS))

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

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

- aka restricted secant inequality, and implies Polyak-Lojasiewicz condition.
- use OPC with one-point smoothness (rather than global smoothness).

Assumption (One-point convexity & smoothness (OPCS))

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

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

- aka restricted secant inequality, and implies Polyak-Lojasiewicz condition.
- use OPC with one-point smoothness (rather than global smoothness).
- example, nonlinear state equation $\mathbf{h}_{t+1} = \phi(\mathbf{A}_* \mathbf{h}_t + \mathbf{B}_* \mathbf{z}_t) + \mathbf{w}_t$, with γ -increasing activation (i.e. $\phi'(x) \geq \gamma > 0$ for all $x \in \mathbb{R}$).

Main Result

Theorem (Main result – informal)

Suppose we run gradient descent algorithm, $\theta_{\tau+1} = \theta_{\tau} - \eta \nabla \hat{\mathcal{L}}(\theta_{\tau})$ to solve the ERM problem. Suppose $T \gtrsim \frac{d}{\alpha^2(1-\rho)}$ and $r \gtrsim \frac{\sigma}{\alpha} \sqrt{\frac{d}{T(1-\rho)}}$. Under certain assumptions, the following statements hold with high probability over the trajectory.

Main Result

Theorem (Main result – informal)

Suppose we run gradient descent algorithm, $\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_{\tau} - \eta \nabla \hat{\mathcal{L}}(\boldsymbol{\theta}_{\tau})$ to solve the ERM problem. Suppose $T \gtrsim \frac{d}{\alpha^2(1-\rho)}$ and $r \gtrsim \frac{\sigma}{\alpha} \sqrt{\frac{d}{T(1-\rho)}}$. Under certain assumptions, the following statements hold with high probability over the trajectory.

- **Uniform convergence of gradient:** For all $\boldsymbol{\theta} \in \mathcal{B}^d(\boldsymbol{\theta}_{\star}, r)$, $\nabla \hat{\mathcal{L}}(\boldsymbol{\theta})$ satisfies

$$\|\nabla \hat{\mathcal{L}}(\boldsymbol{\theta}) - \nabla \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta})\|_{\ell_2} \lesssim (\sigma + \|\boldsymbol{\theta} - \boldsymbol{\theta}_{\star}\|_{\ell_2}) \sqrt{\frac{d}{T(1-\rho)}}$$

Main Result

Theorem (Main result – informal)

Suppose we run gradient descent algorithm, $\theta_{\tau+1} = \theta_{\tau} - \eta \nabla \hat{\mathcal{L}}(\theta_{\tau})$ to solve the ERM problem. Suppose $T \gtrsim \frac{d}{\alpha^2(1-\rho)}$ and $r \gtrsim \frac{\sigma}{\alpha} \sqrt{\frac{d}{T(1-\rho)}}$. Under certain assumptions, the following statements hold with high probability over the trajectory.

- **Uniform convergence of gradient:** For all $\theta \in \mathcal{B}^d(\theta_*, r)$, $\nabla \hat{\mathcal{L}}(\theta)$ satisfies

$$\|\nabla \hat{\mathcal{L}}(\theta) - \nabla \mathcal{L}_{\mathcal{D}}(\theta)\|_{\ell_2} \lesssim (\sigma + \|\theta - \theta_*\|_{\ell_2}) \sqrt{\frac{d}{T(1-\rho)}}$$

- **Convergence of gradient descent:** Set the learning rate $\eta = \alpha/(16\beta^2)$ and fix $\theta_0 \in \mathcal{B}^d(\theta_*, r)$. All gradient descent iterates θ_{τ} on $\hat{\mathcal{L}}(\theta)$ satisfy

$$\|\theta_{\tau} - \theta_*\|_{\ell_2} \lesssim \left(1 - \frac{\alpha^2}{\beta^2}\right)^{\tau} \|\theta_0 - \theta_*\|_{\ell_2} + \frac{\sigma}{\alpha} \sqrt{\frac{d}{T(1-\rho)}}.$$

Case Study

Entrywise nonlinearity: $\mathbf{h}_{t+1} = \phi(\mathbf{A}_* \mathbf{h}_t) + \mathbf{z}_t + \mathbf{w}_t$

- $\mathbf{A}_* = [\mathbf{a}_1^* \cdots \mathbf{a}_n^*]^T \in \mathbb{R}^{n \times n}$.
- Assume $\phi' \geq \gamma > 0$, $|\phi'|, |\phi''| \leq 1$ and $\phi(0) = 0$.

Theorem (simplified)

- Suppose the system satisfies ρ -stability.
- Let $\mathbf{w}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and $\mathbf{z}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_p)$.
- Suppose trajectory length T obeys $T \gtrsim n \log(T)/(1 - \rho)$

With proper learning rate and the initialization $\mathbf{A}^{(0)} = 0$, all gradient descent iterates satisfy

$$\|\mathbf{a}_k^{(\tau)} - \mathbf{a}_k^*\|_{\ell_2} \lesssim \left(1 - \frac{\gamma^4(1 - \rho)^4}{C_\rho^4 n^2}\right)^\tau \|\mathbf{a}_k^{(0)} - \mathbf{a}_k^*\|_{\ell_2} + \frac{\sigma}{\gamma^2} \sqrt{\frac{n}{T(1 - \rho)}}.$$

Case Study

Linear Dynamical System: $\mathbf{h}_{t+1} = \mathbf{A}_* \mathbf{h}_t + \mathbf{B}_* \mathbf{z}_t + \mathbf{w}_t$

- $[\mathbf{A}_* \ \mathbf{B}_*] = [\boldsymbol{\theta}_1^* \ \dots \ \boldsymbol{\theta}_n^*]^T \in \mathbb{R}^{n \times (n+p)}$.
- $\gamma_- := 1 \wedge \lambda_{\min}(\boldsymbol{\Gamma}_L^{\mathbf{B}_*} + \sigma^2 \boldsymbol{\Gamma}_L)$, and $\gamma_+ := 1 \vee \lambda_{\max}(\boldsymbol{\Gamma}_L^{\mathbf{B}_*} + \sigma^2 \boldsymbol{\Gamma}_L)$.

Theorem (simplified)

- Suppose $\rho(\mathbf{A}_*) < 1$.
- Let $\mathbf{w}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and $\mathbf{z}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2 \mathbf{I}_p)$.
- Suppose trajectory length T obeys $T \gtrsim (n+p) \log(T)/(1-\rho)$

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

$$\|\boldsymbol{\theta}_k^{(\tau)} - \boldsymbol{\theta}_k^*\|_{\ell_2} \lesssim \left(1 - \frac{\gamma_-^2}{\gamma_+^2}\right)^\tau \|\boldsymbol{\theta}_k^{(0)} - \boldsymbol{\theta}_k^*\|_{\ell_2} + \frac{\sigma \sqrt{\gamma_+}}{\gamma_-} \sqrt{\frac{n+p}{T(1-\rho)}}.$$

Possible Extensions:

- Partial state observations $\mathbf{y}_t = \mathbf{C}\mathbf{h}_t$
- NARMAX: $\mathbf{y}_{t+1} = \phi(\mathbf{y}_t, \dots, \mathbf{y}_{t-T}, \mathbf{u}_t, \dots, \mathbf{u}_{t-T}; \boldsymbol{\theta}_\star)$.
- Better dependence on spectral radius ρ (e.g. by using martingales)