

# Exponential Ergodicity in a Sobolev Space

Applications to Reinforcement Learning, and ...

A. M. Devraj, I. Kontoyiannis and S. Meyn



Laboratory for Cognition & Control in Complex Systems  
University of Florida, Gainesville

PDE and Probability Methods for Interactions

March 30-31, 2017

Thanks to NSF

# Outline

- 1 Differential Exponential Ergodicity
- 2 Value Function Approximation
- 3 Conclusions

## Goals

Markov process  $X$  on state space  $X = \mathbb{R}^\ell$

Transition semigroup: for  $t \geq 0$ ,  $x \in X$ ,  $A \in \mathcal{B}$ ,

$$P^t(x, A) := P_x\{X(t) \in A\} := \Pr\{X(t) \in A \mid X(0) = x\}$$

Operator notation: for  $f : X \rightarrow \mathbb{R}$ , signed measure  $\nu$  on  $(X, \mathcal{B})$ :

$$P^t f(x) = \int f(y) P^t(x, dy)$$

$$\nu P^t(A) = \int \nu(dx) P^t(x, A)$$

Generator  $\mathcal{D}$  ( $= P - I$  in discrete time)

# Goals

Markov process on  $X = \mathbb{R}^\ell$  (continuous or discrete time)

Generator  $\mathcal{D}$  ( $= P - I$  in discrete time)

# Goals

Markov process on  $X = \mathbb{R}^\ell$  (continuous or discrete time)

Generator  $\mathcal{D}$  ( $= P - I$  in discrete time)

Compute or bound solution  $h: X \rightarrow \mathbb{C}$

- 1 Eigenfunction:  $\mathcal{D}h = \lambda h$

# Goals

Markov process on  $X = \mathbb{R}^\ell$  (continuous or discrete time)

Generator  $\mathcal{D}$  ( $= P - I$  in discrete time)

Compute or bound solution  $h: X \rightarrow \mathbb{C}$

- 1 Eigenfunction:  $\mathcal{D}h = \lambda h$
- 2 Poisson eqn:  $\mathcal{D}h = -\tilde{c}$  [appl. to simulation and control]
- 3 Dirichlet+:  $\mathcal{D}h = \mathcal{G}(h, \nabla h)$  [appl. to simulation and control]

# Goals

Markov process on  $X = \mathbb{R}^\ell$  (continuous or discrete time)

Generator  $\mathcal{D}$  ( $= P - I$  in discrete time)

Compute or bound solution  $h: X \rightarrow \mathbb{C}$

- 1 Eigenfunction:  $\mathcal{D}h = \lambda h$
- 2 Poisson eqn:  $\mathcal{D}h = -\tilde{c}$  [appl. to simulation and control]
- 3 Dirichlet+:  $\mathcal{D}h = \mathcal{G}(h, \nabla h)$  [appl. to simulation and control]
- 4 Its gradient  $K = \nabla h$  [appl. to nonlinear filtering – Mehta & M.]

# Goals

Markov process on  $X = \mathbb{R}^\ell$  (continuous or discrete time)

Generator  $\mathcal{D}$  ( $= P - I$  in discrete time)

Compute or bound solution  $h: X \rightarrow \mathbb{C}$

- 1 Eigenfunction:  $\mathcal{D}h = \lambda h$
- 2 Poisson eqn:  $\mathcal{D}h = -\tilde{c}$  [appl. to simulation and control]
- 3 Dirichlet+:  $\mathcal{D}h = \mathcal{G}(h, \nabla h)$  [appl. to simulation and control]
- 4 Its gradient  $K = \nabla h$  [appl. to nonlinear filtering – Mehta & M.]

**Approach:** New operator norm for spectral theory.

*Stick to discrete time here*



## Notation and Assumptions

Operator notation for geometric ergodicity (M&T and K&M):

$v: X \rightarrow [1, \infty)$  continuous “weighting function”.

For  $f: X \rightarrow \mathbb{R}$ ,

$$\|f\|_v := \sup_x \frac{|f(x)|}{v(x)}.$$

Corresponding Banach space:

$$L_\infty^v := \{f: X \rightarrow \mathbb{R} : \|f\|_v < \infty\}$$

## Notation and Assumptions

Operator notation for geometric ergodicity (M&T and K&M):

$v: X \rightarrow [1, \infty)$  continuous “weighting function”.

For  $f: X \rightarrow \mathbb{R}$ ,

$$\|f\|_v := \sup_x \frac{|f(x)|}{v(x)}.$$

Corresponding Banach space:

$$L_\infty^v := \{f: X \rightarrow \mathbb{R} : \|f\|_v < \infty\}$$

**Geometric ergodicity:** There is  $b_0 < \infty$ ,  $\rho_0 < 1$  such that

$$\|\tilde{P}^t\|_v \leq b_0 \rho_0^t, \quad t \geq 0; \quad \tilde{P}^t = P^t - \mathbf{1} \otimes \pi.$$

$\iff$  spectral gap in  $L_\infty^v$

# Notation and Assumptions

$$\|f\|_v := \sup_x |f(x)|/v(x)$$

$$\text{New norm: } \|f\|_{v,k} = \max_{|\alpha| \leq k} \|\partial^\alpha f\|_v, \quad k \geq 1$$

# Notation and Assumptions

$$\|f\|_v := \sup_x |f(x)|/v(x)$$

$$\text{New norm: } \|f\|_{v,k} = \max_{|\alpha| \leq k} \|\partial^\alpha f\|_v, \quad k \geq 1$$

New Banach spaces:

$$L_\infty^{v,0} = \{f \in L_\infty^v : f \text{ is continuous}\}$$

$$L_\infty^{v,k} = \{g: \mathbf{X} \rightarrow \mathbb{R} : \partial^\alpha f \in L_\infty^{v,0} \text{ for all } |\alpha| \leq k\}$$

Induced operator norm, for any kernel  $\hat{P}$ ,

$$\|\hat{P}\|_{v,k} := \sup \left\{ \frac{\|\hat{P}h\|_{v,k}}{\|h\|_{v,k}} : h \in L_\infty^{v,k}, \|h\|_{v,k} \neq 0 \right\}.$$

## Notation and Assumptions

$$\|f\|_v := \sup_x |f(x)|/v(x)$$

$$\text{New norm: } \|f\|_{v,k} = \max_{|\alpha| \leq k} \|\partial^\alpha f\|_v, \quad k \geq 1$$

New Banach spaces:

$$L_\infty^{v,0} = \{f \in L_\infty^v : f \text{ is continuous}\}$$

$$L_\infty^{v,k} = \{g: \mathbf{X} \rightarrow \mathbb{R} : \partial^\alpha f \in L_\infty^{v,0} \text{ for all } |\alpha| \leq k\}$$

Induced operator norm, for any kernel  $\widehat{P}$ ,

$$\|\widehat{P}\|_{v,k} := \sup \left\{ \frac{\|\widehat{P}h\|_{v,k}}{\|h\|_{v,k}} : h \in L_\infty^{v,k}, \|h\|_{v,k} \neq 0 \right\}.$$

We will stick to  $k = 1$ :

$$\|f\|_{v,1} = \max \left\{ \|f\|_v, \|\partial^1 f\|_v, \dots, \|\partial^\ell f\|_v \right\}$$

# Notation and Assumptions

## Markovian system dynamics

$$X(t+1) = a(X(t), N(t+1)), \quad t \in \mathbb{Z}_+, \mathbf{N} \text{ i.i.d.}$$

$$P(x, A) = \mathbf{P}(a(x, N(1)) \in A)$$

# Notation and Assumptions

## Markovian system dynamics

$$X(t+1) = a(X(t), N(t+1)), \quad t \in \mathbb{Z}_+, \mathbf{N} \text{ i.i.d.}$$

$$P(x, A) = \mathbf{P}(a(x, N(1)) \in A)$$

A1 **Smooth dynamics:**  $a : \mathbb{R}^{d \times m} \rightarrow \mathbb{R}^d$  is  $C^1$  and Lipschitz in  $x$

A2 **Densities:** For some  $t_0 \geq 1$  and  $C^1$  function  $p_{t_0}$ ,

$$P^{t_0}(x, A) = \int_A p_{t_0}(x, y) dy, \quad x \in \mathbf{X}, A \in \mathcal{B}$$

A3  **$\psi$ -irreducibility:** for some  $x_0 \in \mathbf{X}$ ,

$$P^t(x, O) > 0 \quad \text{all } t = t_x \geq 0 \text{ sufficiently large, each nbd } O \text{ of } x_0$$

# Notation and Assumptions

## Markovian system dynamics

$$X(t+1) = a(X(t), N(t+1)), \quad t \in \mathbb{Z}_+, \mathbf{N} \text{ i.i.d.}$$

$$P(x, A) = \mathbb{P}(a(x, N(1)) \in A)$$

A1 **Smooth dynamics:**  $a : \mathbb{R}^{d \times m} \rightarrow \mathbb{R}^d$  is  $C^1$  and Lipschitz in  $x$

A2 **Densities:** For some  $t_0 \geq 1$  and  $C^1$  function  $p_{t_0}$ ,

$$P^{t_0}(x, A) = \int_A p_{t_0}(x, y) dy, \quad x \in X, A \in \mathcal{B}$$

A3  **$\psi$ -irreducibility:** for some  $x_0 \in X$ ,

$$P^t(x, O) > 0 \quad \text{all } t = t_x \geq 0 \text{ sufficiently large, each nbd } O \text{ of } x_0$$

A4 **Donsker-Varadhan drift condition:**

$$\mathcal{H}(V) := \log(Pe^V) - V \leq -\delta W + b\mathbb{I}_C, \quad (\text{DV3})$$

$$V = \log(v) \text{ and } W \in L_\infty^{v,0} \text{ coercive, } C \text{ compact.} \quad [\text{K\&M, L. Wu}]$$



Main Result: Separability in  $L_\infty^{v,1}$ 

- A1 Smooth dynamics
- A2 Densities
- A3  $\psi$ -irreducibility
- A4 Donsker-Varadhan drift condition:

$$\mathcal{H}(V) \leq -\delta W + b\mathbb{I}_C. \quad (\text{DV3})$$

Theorem: *Separability in  $L_\infty^{v,1}$*

The kernel  $P^t$  is *separable* in  $L_\infty^{v,1}$  for some  $t_1$  and all  $t \geq t_1$ :

$$P^t \approx \sum_{k=1}^n s_k \otimes \nu_k, \quad \text{approximation in } L_\infty^{v,1}$$

Main Result: Separability in  $L_\infty^{v,1}$ 

Theorem: *Separability in  $L_\infty^{v,1}$*

The kernel  $P^t$  is *separable* in  $L_\infty^{v,1}$  for some  $t_1$  and all  $t \geq t_1$ :

$$P^t \approx \sum_{k=1}^n s_k \otimes \nu_k, \quad \text{approximation in } L_\infty^{v,1}$$

Corollaries: Discrete spectrum and ... there is  $b_0 < \infty$  and  $\varrho_0 < 1$  s.t.

$$\|\tilde{P}^t\|_{v,1} \leq b_0 \varrho_0^t, \quad t \geq t_1.$$

Main Result: Separability in  $L_\infty^{v,1}$ Theorem: *Separability in  $L_\infty^{v,1}$* The kernel  $P^t$  is *separable* in  $L_\infty^{v,1}$  for some  $t_1$  and all  $t \geq t_1$ :

$$P^t \approx \sum_{k=1}^n s_k \otimes \nu_k, \quad \text{approximation in } L_\infty^{v,1}$$

Corollaries: Discrete spectrum and ... there is  $b_0 < \infty$  and  $\varrho_0 < 1$  s.t.

$$\|\tilde{P}^t\|_{v,1} \leq b_0 \varrho_0^t, \quad t \geq t_1.$$

Interpretation: for  $f \in L_\infty^{v,1}$ ,  $P^t f(x) \rightarrow \pi(f)$ ,  $\partial^i P^t f(x) \rightarrow 0$ ,  
*uniform geometric convergence rate.*

**Proof:** Truncation of  $P^{t_1}$  to compact domain, as in [K&M 200X];  
 Spectrum $_{L_\infty^{v,1}}(P^t) \subseteq$  Spectrum $_{L_\infty}(P^t)$

# Connection with Lyapunov Exponents

Remains a mystery

Sensitivity process:  $\mathcal{S}^T(t) = \frac{\partial}{\partial X(0)} X(t)$

$$\mathcal{S}(t+1) = \mathcal{A}(t+1)\mathcal{S}(t), \quad \mathcal{S}(0) = I$$

where  $\mathcal{A}^T(t) := \nabla_x a(X(t-1), N(t))$ .

# Connection with Lyapunov Exponents

Remains a mystery

Sensitivity process:  $\mathcal{S}^T(t) = \frac{\partial}{\partial X(0)} X(t)$

$$\mathcal{S}(t+1) = \mathcal{A}(t+1)\mathcal{S}(t), \quad \mathcal{S}(0) = I$$

where  $\mathcal{A}^T(t) := \nabla_x a(X(t-1), N(t))$ .

Lyapunov exponents:  $\Lambda = \lim_{t \rightarrow \infty} \frac{1}{t} \log \|\mathcal{S}(t)\| \quad a.s.$

$$\Lambda_p = \lim_{t \rightarrow \infty} \frac{1}{t} \log \mathbf{E}[\|\mathcal{S}(t)\|^p]$$

Gradient representation:

$$\nabla P^t = Q^t \nabla$$

$$Q^t g(x) := \mathbf{E}_x [\mathcal{S}^T(t) g(X(t))], \quad g = \nabla c.$$

# Discounted cost value function

Cost function:  $c : \mathbb{R}^d \rightarrow \mathbb{R}$

Discount factor:  $\alpha < 1$

Value function: 
$$h_\alpha(x) = \sum_{t=0}^{\infty} \alpha^t \mathbf{E}[c(X(t)) \mid X(0) = x]$$

## Discounted cost value function

Cost function:  $c : \mathbb{R}^d \rightarrow \mathbb{R}$

Discount factor:  $\alpha < 1$

Value function: 
$$h_\alpha(x) = \sum_{t=0}^{\infty} \alpha^t \mathbf{E}[c(X(t)) \mid X(0) = x]$$

**Goal of TD learning:** Approximate  $h_\alpha$  in parameterized class  $\{h_\alpha^\theta : \theta \in \mathbb{R}^\ell\}$

## Discounted cost value function

Cost function:  $c : \mathbb{R}^d \rightarrow \mathbb{R}$

Discount factor:  $\alpha < 1$

Value function: 
$$h_\alpha(x) = \sum_{t=0}^{\infty} \alpha^t \mathbb{E}[c(X(t)) \mid X(0) = x]$$

**Goal of TD learning:** Approximate  $h_\alpha$  in parameterized class  $\{h_\alpha^\theta : \theta \in \mathbb{R}^\ell\}$

**Goal of  $\nabla$ -TD learning:** (new)

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|\nabla h_\alpha^\theta(X) - \nabla h_\alpha(X)\|^2], \quad X \sim \pi$$



## Discounted cost value function

Cost function:  $c : \mathbb{R}^d \rightarrow \mathbb{R}$

Discount factor:  $\alpha < 1$

Value function: 
$$h_\alpha(x) = \sum_{t=0}^{\infty} \alpha^t \mathbb{E}[c(X(t)) \mid X(0) = x]$$

**Goal of TD learning:** Approximate  $h_\alpha$  in parameterized class  $\{h_\alpha^\theta : \theta \in \mathbb{R}^\ell\}$

**Goal of  $\nabla$ -TD learning:** (new)

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|\nabla h_\alpha^\theta(X) - \nabla h_\alpha(X)\|^2], \quad X \sim \pi$$

Example: affine parameterization,

$$h_\alpha^\theta(x) = \kappa(\theta) + \sum_{j=1}^{\ell} \theta_j \psi_j(x), \quad \psi : \mathbb{R}^d \rightarrow \mathbb{R}^\ell$$

## Discounted cost value function

Value function: 
$$h_\alpha(x) = \sum_{t=0}^{\infty} \alpha^t \mathbf{E}[c(X(t)) \mid X(0) = x]$$

Goal of  $\nabla$ -TD learning: 
$$\theta^* = \arg \min_{\theta} \mathbf{E}_{\pi} [\|\nabla h_{\alpha}^{\theta}(X) - \nabla h_{\alpha}(X)\|^2]$$

Recover missing constant:

$$h_{\alpha}^{\theta}(x) = \theta^T \psi(x) + \kappa(\theta)$$

$$\kappa(\theta) = -\theta^T \pi(\psi) + \pi(c)/(1 - \alpha)$$

$$\implies \pi(h_{\alpha}^{\theta}) = \pi(h_{\alpha}) = \pi(c)/(1 - \alpha)$$

# Gradient Representation

Goal of  $\nabla$ -TD learning:

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|\nabla h_{\alpha}^{\theta}(X) - \nabla h_{\alpha}(X)\|^2], \quad X \sim \pi$$

Representation:

$$\nabla h_{\alpha}(x) = \sum_{t=0}^{\infty} \alpha^t \nabla P^t c(x)$$

$$\begin{aligned} \implies \nabla h_{\alpha} &= \Omega_{\alpha} \nabla c := \sum_{t=0}^{\infty} \alpha^t Q^t \nabla c(x) \\ &= \sum_{t=0}^{\infty} \alpha^t \mathbb{E}[\mathcal{S}^T(t) \nabla c(X(t)) \mid X(0) = x] \end{aligned}$$

Adjoint Representation of  $\theta^*$ Goal of  $\nabla$ -TD learning:

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|\nabla h_{\alpha}^{\theta}(X) - \nabla h_{\alpha}(X)\|^2], \quad X \sim \pi$$

Solution:

$$\theta^* = M^{-1}b$$

$$M = \mathbb{E}_{\pi}[(\nabla\psi(X))^T \nabla\psi(X)]$$

$$b = \mathbb{E}_{\pi}[(\nabla\psi(X))^T \nabla h_{\alpha}(X)]$$

$$\nabla h_{\alpha} = \Omega_{\alpha} \nabla c$$

Adjoint Representation of  $\theta^*$ Goal of  $\nabla$ -TD learning:

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|\nabla h_{\alpha}^{\theta}(X) - \nabla h_{\alpha}(X)\|^2], \quad X \sim \pi$$

Solution:

$$\theta^* = M^{-1}b$$

$$M = \mathbb{E}_{\pi}[(\nabla\psi(X))^T \nabla\psi(X)]$$

$$b = \mathbb{E}_{\pi}[(\nabla\psi(X))^T \nabla h_{\alpha}(X)]$$

$$\nabla h_{\alpha} = \Omega_{\alpha} \nabla c$$

Adjoint gives causal representation:

$$b_i = \langle \partial^i \psi, \Omega_{\alpha} \nabla c \rangle$$

$$= \langle \Omega_{\alpha}^{\dagger} \partial^i \psi, \nabla c \rangle$$

Adjoint Representation of  $\theta^*$ Goal of  $\nabla$ -TD learning:

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|\nabla h_{\alpha}^{\theta}(X) - \nabla h_{\alpha}(X)\|^2], \quad X \sim \pi$$

Solution:

$$\theta^* = M^{-1}b$$

$$M = \mathbb{E}_{\pi}[(\nabla\psi(X))^T \nabla\psi(X)]$$

$$b = \mathbb{E}_{\pi}[(\nabla\psi(X))^T \nabla h_{\alpha}(X)]$$

$$\nabla h_{\alpha} = \Omega_{\alpha} \nabla c$$

Adjoint gives causal representation:

$$b_i = \langle \partial^i \psi, \Omega_{\alpha} \nabla c \rangle$$

$$= \langle \Omega_{\alpha}^{\dagger} \partial^i \psi, \nabla c \rangle = \mathbb{E}_{\pi}[\varphi(t)^T \nabla c(X(t))],$$

$$\varphi(t) = \sum_{k=0}^{\infty} \alpha^k [\mathcal{A}(1+t-k) \mathcal{A}(2+t-k) \cdots \mathcal{A}(t)]^T \nabla \psi(X(t-k)), \quad t \in \mathbb{Z}$$

## Differential Least Squares Temporal Difference Algorithm

 $\nabla$ -LSTD algorithm

$$\varphi(t) = \alpha \mathcal{A}^T(t) \varphi(t-1) + \nabla \psi(X(t))$$

$$b(t) = (1 - \gamma_t) b(t-1) + \gamma_t \varphi(t)^T \nabla c(X(t))$$

$$M(t) = (1 - \gamma_t) M(t-1) + \gamma_t \nabla \psi(X(t)) \nabla \psi(X(t))^T$$

$$\theta(t) = M^{-1}(t) b(t)$$

# Differential Least Squares Temporal Difference Algorithm

## $\nabla$ -LSTD algorithm

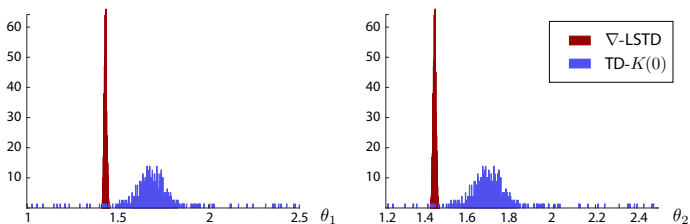
$$\varphi(t) = \alpha \mathcal{A}^T(t) \varphi(t-1) + \nabla \psi(X(t))$$

$$b(t) = (1 - \gamma_t) b(t-1) + \gamma_t \varphi(t)^T \nabla c(X(t))$$

$$M(t) = (1 - \gamma_t) M(t-1) + \gamma_t \nabla \psi(X(t)) \nabla \psi(X(t))^T$$

$$\theta(t) = M^{-1}(t) b(t)$$

Algorithm is amazing:





# Conclusions

New Banach space for Markov processes is just right for our goals:

$$\mathcal{D}h = \mathcal{G}(h, \nabla h)$$

# Conclusions

New Banach space for Markov processes is just right for our goals:

$$\mathcal{D}h = \mathcal{G}(h, \nabla h)$$

Gaps:

Relationship with Lyapunov exponents remains a mystery.

Needed for a firmer theory for the  $\nabla$ -LSTD algorithm.

# Conclusions

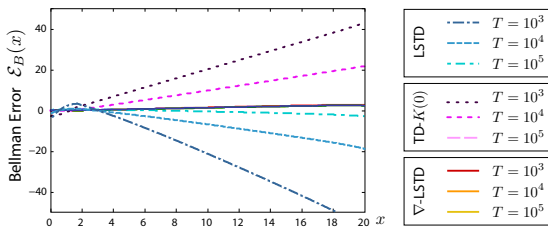
New Banach space for Markov processes is just right for our goals:

$$\mathcal{D}h = \mathcal{G}(h, \nabla h)$$

Gaps:







Relationship with Lyapunov exponents remains a mystery.

Needed for a firmer theory for the  $\nabla$ -LSTD algorithm.



Thank you!

# Selected References

-  A. Devraj, I. Kontoyiannis, and S. P. Meyn, “Exponential ergodicity and Lyapunov exponents Part I: Markov chains in discrete time,” In preparation, 2016.
-  I. Kontoyiannis and S. P. Meyn. “Computable exponential bounds for screened estimation and simulation,” *Ann. Appl. Probab.*, 18(4):1491–1518, 2008.
-  I. Kontoyiannis and S. P. Meyn. “Approximating a diffusion by a finite-state hidden Markov model,” *Stochastic Proc. and their Appl.*, 2016.
-  S. P. Meyn and R. L. Tweedie, *Markov chains and stochastic stability*, 2nd ed. Cambridge Mathematical Library, 2009.
-  A. M. Devraj and S. P. Meyn. Differential TD learning for value function approximation. In *55th Conference on Decision and Control*, pages 6347–6354, Dec 2016.
-  S. P. Meyn, *Control Techniques for Complex Networks*. Cambridge University Press, 2007, pre-publication edition available online.