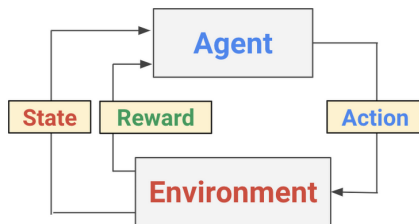


# On Optimization Formulations and Algorithms of Markov Decision Problems

Lexing Ying

work with Haoya Li, Yuhua Zhu, Samarth Gupta, Hsiangfu Yu, Inderjit Dhillon  
MATH-IMS Joint Applied Mathematics Colloquium  
Chinese University of Hong Kong  
Feb 4, 2022

## Markov decision problems



- ▶ Reinforcement learning (RL), online decision making
- ▶ Game playing: AlphaGo, AlphaZero
- ▶ Robotics
- ▶ Science: AlphaFold

This talk: optimization formulations and algorithms

# MDP

Consider a Markov decision process (MDP)  $M = (S, A, P, r, \gamma)$

- ▶  $S$ : state space
- ▶  $A$ : action space
- ▶  $P$ : transition tensor,  $P_{st}^a$  probability from  $s$  to  $t$  under action  $a$
- ▶  $r$ : reward,  $r_s^a$  reward at  $s$  under action  $a$

Policy  $\pi$

- ▶  $\pi_s^a$ : probability of taking action  $a$  at state  $s$
- ▶  $\sum_{a \in A} \pi_s^a = 1$

Under a fixed  $\pi$ , the behavior follows  $P^\pi$  and  $r^\pi$

$$P_{st}^\pi = \sum_{a \in A} P_{st}^a \pi_s^a, \quad r_s^\pi = \sum_{a \in A} r_s^a \pi_s^a$$

The value function  $v_s^\pi \equiv \mathbb{E} \left[ \sum_{m=0}^{\infty} \gamma^m r_{s_m}^{a_m} \mid s_0 = s \right]$  where

- ▶  $(s_0, a_0, s_1, a_1, s_2, a_2, \dots)$ ,  $a_m \sim \pi_{s_m}^A, s_{m+1} \sim P_{s_m, S}^{a_m}$

$v^\pi \in \mathbb{R}^{|S|}$  satisfies

$$v^\pi = r^\pi + \gamma P^\pi v^\pi$$

Goal: maximize the value function!

## Optimization formulations

- ▶ Primal: value (policy) iteration
- ▶ Primal-dual: actor-critic
- ▶ Dual: policy gradient

## New optimization algorithms

- ▶ Quasi-Newton policy gradient (for dual form)
- ▶ Accelerated primal-dual algorithm (for primal-dual form)

## Background: linear programming forms

Primal ( $e$  positive)

$$\min_{v: b - Av \leq 0} e^T v$$

Primal dual

$$\min_v \max_{u \geq 0} e^T v + u^T (b - Av)$$

(by minimax theorem)

$$\max_{u \geq 0} \min_v e^T v + u^T (b - Av)$$

Dual

$$\max_{u \geq 0, A^T u = e} b^T u$$

## Background: entropy regularized version

(Generalized) Shannon entropy, convex

$$h(u) = \sum_i u_i \log \frac{u_i}{\sum_j u_j} = \sum_i u_i \log u_i - \left( \sum_i u_i \right) \log \left( \sum_i u_i \right)$$

Primal dual

$$\min_v \max_{u>0} e^T v + u^T (b - Av) - h(u)$$

$$\max_{u>0} \min_v e^T v + u^T (b - Av) - h(u)$$

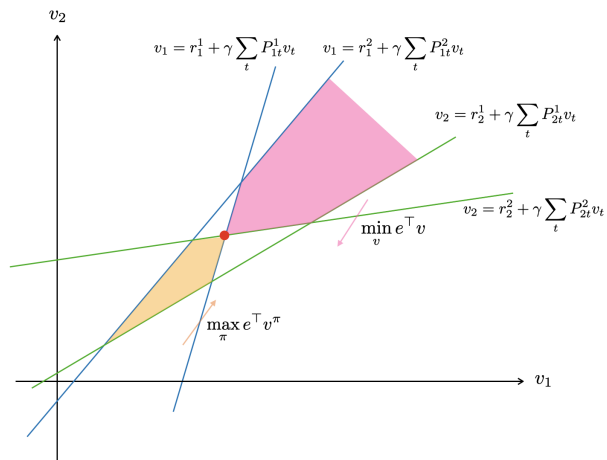
Primal

$$\min_{1^T \exp(b - Av) \leq 1} e^T v$$

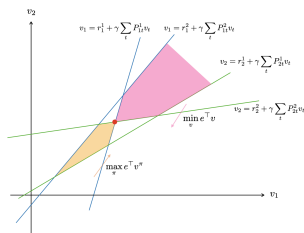
Dual

$$\max_{A^T u = e} b^T u - h(u)$$

## Going back to MDP



# Non-regularized MDP



Primal

$$\min_v \sum_{s \in S} e_s v_s, \text{ s.t. } \forall a, s, r_s^a - \left( v_s - \gamma \sum_{t \in S} P_{st}^a v_t \right) \leq 0$$

Primal-dual ( $u_s^a$  Lagrange multipliers)

$$\min_{v_s} \max_{u_s^a \geq 0} \sum_{s \in S} e_s v_s + \sum_{s, a} (r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t - v_s) u_s^a$$

$$\min_v \max_{u^a \geq 0} e^T v + \sum_{a \in A} (u^a)^T (r^a + \gamma P^a v - v)$$

Dual

$$\max_{u^a \geq 0} \sum_{a \in A} (r^a)^T u^a, \text{ s.t. } \sum_{a \in A} (I - \gamma (P^a)^T) u^a = e$$



## Connections with RL algorithms

Primal:

$$\min_v \sum_{s \in S} e_s v_s, \text{ s.t. } \forall a, s, r_s^a - \left( v_s - \gamma \sum_{t \in S} P_{st}^a v_t \right) \leq 0$$

- ▶ KKT condition is equivalent to Bellman equation

$$\begin{cases} r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t^* \leq v_s^*, & \text{for } \forall s, a \\ \sum_{a \in A} (u_s^a - \gamma \sum_{t \in S} P_{ts}^a u_t^a) = e_s, & \text{for } \forall s \\ u_s^a (r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t^* - v_s^*) = 0 & \text{for } \forall s, a \end{cases} \Leftrightarrow v_s^* = \max_{a \in A} \left( r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t^* \right)$$

- ▶ Algorithm: value/policy iteration

Dual:

$$\max_{u^a \geq 0} \sum_{a \in A} (r^a)^\top u^a, \text{ s.t. } \sum_{a \in A} (I - \gamma(P^a)^\top) u^a = e$$

- ▶ Let  $u_s^a = w_s \pi_s^a$ . Then  $(I - \gamma(P^\pi)^\top) w = e$ ,  $\sum_{a \in A} (r^a)^\top u^a = (r^\pi)^\top w^\pi$

$$\max_{\pi} e^\top (I - \gamma P^\pi)^{-1} r^\pi$$

- ▶ Algorithm: policy gradient

## Entropy regularized MDP

$$h(u_s) = \sum_{a \in A} u_s^a \log \frac{u_s^a}{\sum_{b \in A} u_s^b}$$

Primal-dual

$$\min_{v_s} \sup_{u_s^a > 0} \sum_{s \in S} e_s v_s + \sum_{s,a} \left( r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t - v_s \right) u_s^a - \tau \sum_{s \in S} h(u_s)$$

Primal

$$\min_{v_s} e^\top v, \text{ s.t. } \forall s, v_s \geq \tau \log \left( \sum_{a \in A} \exp \left( \frac{r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t}{\tau} \right) \right)$$

Dual

$$\sup_{u^a > 0} \sum_{a \in A} (r^a)^\top u^a - \sum_{s \in S} h(u_s), \text{ s.t. } \sum_{a \in A} (I - \gamma(P^a)^\top) u^a = e$$

## Connections with RL algorithms

### Primal

$$\min_{v_s} e^\top v, \text{ s.t. } \forall s, v_s \geq \tau \log \left( \sum_{a \in A} \exp \left( \frac{r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t}{\tau} \right) \right)$$

- ▶ KKT condition is equivalent to Bellman equation

$$v_s^* = \tau \log \left( \sum_{a \in A} \exp \left( \frac{r_s^a + \gamma \sum_{t \in S} P_{st}^a v_t^*}{\tau} \right) \right)$$

- ▶ Algorithm: value/policy iteration

### Dual

$$\sup_{u^a > 0} \sum_{a \in A} (r^a)^\top u^a - \sum_{s \in S} h(u_s), \text{ s.t. } \sum_{a \in A} (I - \gamma(P^a)^\top) u^a = e$$

- ▶ Let  $u_s^a = w_s \pi_s^a$ . Then  $(I - \gamma(P^\pi)^\top) w = e$   $\sum_{a \in A} (r^a)^\top u^a = (r^\pi)^\top w^\pi$

$$\max_{\pi} e^\top (I - \gamma P^\pi)^{-1} (r^\pi - \tau h^\pi)$$

where  $h_s^\pi = \sum_{a \in A} \pi_s^a \log \pi_s^a$

- ▶ Algorithm: policy gradient

## Optimization formulations

- ▶ Primal: value (policy) iteration
- ▶ Primal-dual: actor-critic
- ▶ Dual: policy gradient

## New optimization algorithms

- ▶ 1. Quasi-Newton policy gradient (for dual form)
- ▶ 2. Accelerated natural gradient ascent descent (for primal-dual form)

## New algorithm 1

$$\max_{\pi} e^{\top} (I - \gamma P^{\pi})^{-1} (r^{\pi} - \tau h^{\pi})$$

Consider gradient-type dynamics

- ▶ 1st order gradient descent: slow
- ▶ 2nd order Newton method: Hessian too expensive to estimate

Our contribution: quasi-Newton

- ▶ Diagonal approximation of Hessian
- ▶ As cheap as a 1st-order method
- ▶ Converges like 2nd-order methods: super exponential convergence

## Hessian approximation

Gradient

$$\frac{\partial E}{\partial \pi_s^a} = (r_s^a - \tau(\log \pi_s^a + 1) - [(I - \gamma P^a)v^\pi]_s + c_s)(w^\pi)_s,$$

Hessian

$$\frac{\partial^2 E}{\partial \pi_s^a \partial \pi_t^b} = H_{st}^{ab} + \text{remainder}$$

Theorem

For the choice

$$H_{st}^{ab}(\pi) = -\tau \delta_{st}^{ab} \frac{(w^\pi)_s}{\pi_s^a},$$

$$\|\partial^2 E(\pi) - H(\pi)\| = O(\|\pi - \pi^*\|), \quad \pi^* \text{ exact soln}$$

Preconditioned gradient descent

$$\frac{d\pi_s^a}{dt} = -(H^{-1} \nabla_\pi E)_s^a = \pi_s^a (r_s^a - \tau(\log \pi_s^a + 1) - [(I - \gamma P^a)v^\pi]_s + c_s) / \tau.$$

In terms of  $\theta_s^a = \log \pi_s^a$

$$\frac{d\theta_s^a}{dt} = (r_s^a - \tau(\theta_s^a + 1) - [(I - \gamma P^a)v^\pi]_s + c_s) / \tau.$$

## Quasi-Newton algorithm

$$\frac{d\theta_s^a}{dt} = (r_s^a - \tau(\theta_s^a + 1) - [(I - \gamma P^a)v^\pi]_s + c_s)/\tau.$$

Taking explicit Euler with step size  $\eta$

$$\theta_s^a \leftarrow \eta(r_s^a - \tau - [(I - \gamma P^a)v^\pi]_s + c_s)/\tau + (1 - \eta)\theta_s^a$$

Back in terms of  $\pi$

$$\pi_s^a \leftarrow (\pi_s^a)^{1-\eta} \exp(\eta(r_s^a + (\gamma P^a v^\pi)_s)/\tau)$$

When  $\eta = 1$

$$\pi_s^a \leftarrow \exp((r_s^a + (\gamma P^a v^\pi)_s)/\tau)$$

and this is entropy regularized natural policy gradient method

# Super exponential convergence

## Theorem

With step size  $\eta = 1$ ,

$$\left\| \pi^{(k+1)} - \pi^* \right\| \leq C \left\| \pi^k - \pi^* \right\|^2,$$

Our approach generalizes to **entropies of form**

$$(h^\pi)_s = \sum_a \phi \left( \frac{\pi_s^a}{\gamma_a} \right) \gamma_a,$$

where

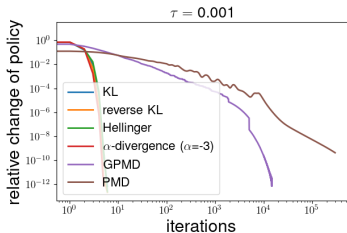
- ▶  $\phi$  is convex on  $(0, +\infty)$  and  $\phi(1) = 0$ ,
- ▶  $\gamma$  is a prior distribution over  $A$ .

Examples

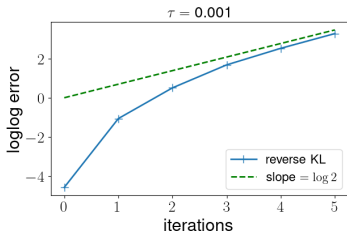
- ▶  $\phi(x) = x \log x$ , KL-divergence
- ▶ Reverse KL-divergence
- ▶ Hellinger divergences



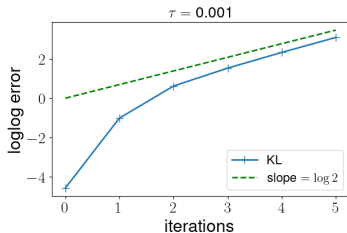
# Example 1



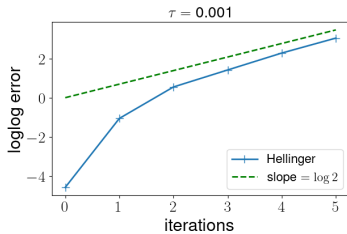
Relative policy error



loglog rKL



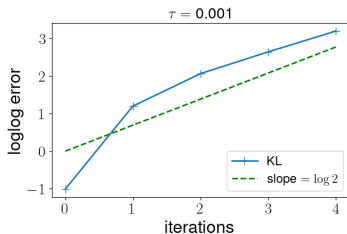
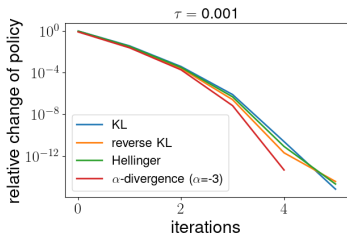
loglog KL



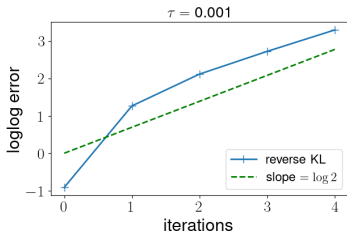
loglog Hellinger

Green line slope = log 2 (Newton quadratic convergence)

## Example 2

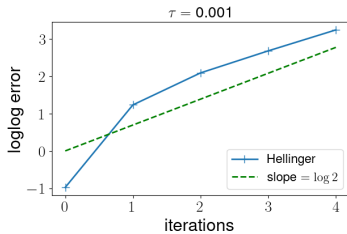


Relative policy error



loglog rKL

loglog KL



loglog Hellinger

Green line slope =  $\log 2$  (Newton quadratic convergence)

## New algorithm 2

$$\min_v \max_u \sum_{s \in S} e_s v_s + \sum_{s \in S, a \in A} u_s^a (r_s^a - ((I - \gamma P^a)v)_s) - \tau \sum_{s \in S, a \in A} u_s^a \log(u_s^a / \tilde{u}_s)$$

where  $\tilde{u}_s := \sum_{a \in A} u_s^a$ . By introducing  $K^a = I - \gamma P^a$ , in short

$$\min_v \max_u \sum_s e_s v_s + \sum_s^a u_s^a (r_s^a - (K^a v)_s) - \tau \sum_{sa} u_s^a \log(u_s^a / \tilde{u}_s).$$

Consider 1st order methods

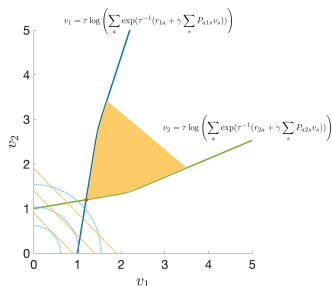
- ▶ Simple gradient descent (GD) in  $v$  and gradient ascent (GA) in  $u$  fail to converge due to lack of strict convexity/concavity
- ▶ Example:  $\min_v \max_u v \cdot u \Rightarrow \begin{cases} \frac{dv}{dt} = -u \\ \frac{du}{dt} = v \end{cases}$
- ▶ Improved versions gives  $1/T^c$  convergence

Our contribution

- ▶ Quadratically convexified in  $v$
- ▶ Improved preconditioning in  $u$
- ▶ Fast exponential convergence

## Quadratic convexification

Main idea in the primal form



For primal-dual, the following are equivalent

$$\min_v \max_u E_0(v, u) = \sum_s e_s v_s + \sum_{sa} u_s^a \left( r_s^a - \sum_{s'} K_{ss'}^a v_{s'} \right) - \tau \sum_{sa} u_s^a \log \frac{u_s^a}{\tilde{u}_s}$$

$$\min_v \max_u E(v, u) = \frac{\alpha}{2} \sum_s v_s^2 + \sum_{sa} u_s^a \left( r_s^a - \sum_{s'} K_{ss'}^a v_{s'} \right) - \tau \sum_{sa} u_s^a \log \frac{u_s^a}{\tilde{u}_s}$$

$E(v, u)$  is quadratically convexified in  $v$

## Natural GAD

For  $E(v, u)$ , apply natural GAD

$$\frac{\partial E}{\partial v_{s'}} = \alpha v_{s'} - \sum_{sa} K_{ss'}^a u_s^a,$$

$$\frac{\partial E}{\partial u_s^a} = \left( r_s^a - \sum_{s'} K_{ss'}^a v_{s'} \right) - \tau \log \frac{u_s^a}{\tilde{u}_s}$$

$$\frac{\partial^2 E}{\partial v_s \partial v_{s'}} = \alpha \delta_{ss'},$$

$$\frac{\partial^2 E}{\partial u_s^a \partial u_{s'}^{a'}} = -\tau \delta_{ss'} \left( \frac{\delta^{aa'}}{u_s^a} - \frac{1}{\tilde{u}_s} \right) \approx -\tau \delta_{ss'} \frac{\delta^{aa'}}{u_s^a}$$

Natural GAD (NGAD)

$$\frac{dv_{s'}}{dt} = - \left( v_{s'} - \frac{1}{\alpha} \sum_{sa} K_{ss'}^a u_s^a \right)$$

$$\frac{du_s^a}{dt} = -u_s^a \left( \log \frac{u_s^a}{\tilde{u}_s} - \frac{1}{\tau} \left( r_s^a - \sum_{s'} K_{ss'}^a v_{s'} \right) \right)$$

# Convergence

Let  $(v^*, u^*)$  be the exact solution.

## Theorem

$$L(v, u) = \frac{\alpha}{2} \sum_{s \in S} |v_s - v_s^*|^2 + \tau \sum_{s \in S, a \in A} \left( u_s^{*a} \log \frac{u_s^{*a}}{u_s^a} + u_s^a - u_s^{*a} \right)$$

*is a Lyapunov function of the NGAD dynamics*

## Theorem

*The NGAD dynamics converges globally to  $(v^*, u^*)$ .*

## Theorem

*The NGAD dynamics converges to  $(v^*, u^*)$  at rate  $O(e^{-ct})$ .*

Q: In practice, the convergence is still slow. Can this be improved?

## Another look at the Hessian

$$\frac{\partial^2 E}{\partial u^2} = \begin{bmatrix} H_1 & & \\ & \ddots & \\ & & H_{|S|} \end{bmatrix}, \quad H_s = \text{diag}((u_{s\cdot})^{-1}) - \frac{1}{\tilde{u}_s} \mathbf{1}_{|A|} \mathbf{1}_{|A|}^\top, \quad s \in S.$$

Natural grad:  $\frac{\partial^2 E}{\partial u^2} \approx \begin{bmatrix} \text{diag}((u_{1\cdot})^{-1}) & & \\ & \ddots & \\ & & \text{diag}((u_{|S|\cdot})^{-1}) \end{bmatrix}$  with inverse

$$\begin{bmatrix} \text{diag}(u_{1\cdot}) & & \\ & \ddots & \\ & & \text{diag}(u_{|S|\cdot}) \end{bmatrix} \equiv \begin{bmatrix} \tilde{u}_1(\text{diag}(\pi_{1\cdot})) & & \\ & \ddots & \\ & & \tilde{u}_{|S|}(\text{diag}(\pi_{|S|\cdot})) \end{bmatrix}$$

But the pseudo-inverse of  $\frac{\partial^2 E}{\partial u^2}$  is

$$\begin{bmatrix} \tilde{u}_1(\text{diag}(\pi_{1\cdot}) - \pi_{1\cdot} \pi_{1\cdot}^\top) & & \\ & \ddots & \\ & & \tilde{u}_{|S|}(\text{diag}(\pi_{|S|\cdot}) - \pi_{|S|\cdot} \pi_{|S|\cdot}^\top) \end{bmatrix}$$

# Interpolate

An interpolation of the two gives

$$\begin{bmatrix} \tilde{u}_1(\text{diag}(\pi_{1\cdot}) - c\pi_{1\cdot}\pi_{1\cdot}^\top) \\ \vdots \\ \tilde{u}_{|S|}(\text{diag}(\pi_{|S|\cdot}) - c\pi_{|S|\cdot}\pi_{|S|\cdot}^\top) \end{bmatrix}$$

where  $c \in (0, 1)$

This gives interpolated NGAD (I-NGAD)

$$\begin{aligned} \frac{dv_{s'}}{dt} &= - \left( v_{s'} - \frac{1}{\alpha} \sum_{sa} K_{ss'}^a u_s^a \right) \\ \frac{du_s}{dt} &= -\tilde{u}_s \left( \text{diag}(\pi_s) - c\pi_s(\pi_s)^\top \right) \left( \log \frac{u_s}{\tilde{u}_s} - \frac{1}{\tau} \left( r_s - \sum_{s'} K_{ss'} v_{s'} \right) \right) \end{aligned}$$



# Convergence

Let  $(v^*, u^*)$  be the exact solution.

## Theorem

$$L_c(v, u) = \frac{\alpha}{2} \sum_s |v_s - v_s^*|^2 + \tau \left( \sum_{sa} \left( u_s^{*a} \log \frac{u_s^{*a}}{u_s^a} + u_s^a - u_s^{*a} \right) + \frac{c}{1-c} \sum_s \left( \tilde{u}_s^* \log \frac{\tilde{u}_s^*}{\tilde{u}_s} + \tilde{u}_s - \tilde{u}_s^* \right) \right)$$

*is a Lyapunov function of the I-NGAD dynamics*

## Theorem

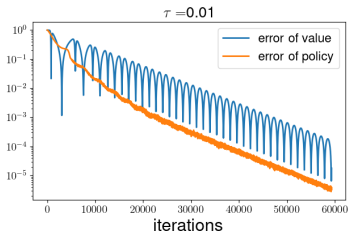
*The I-NGAD dynamics converges globally to  $(v^*, u^*)$ .*

## Theorem

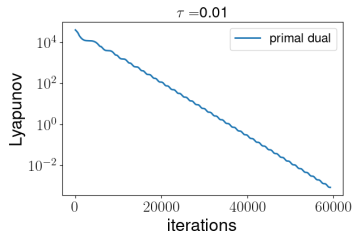
*The I-NGAD dynamics converges to  $(v^*, u^*)$  at rate  $O(e^{-ct})$ .*

In practice, the convergence is FAST

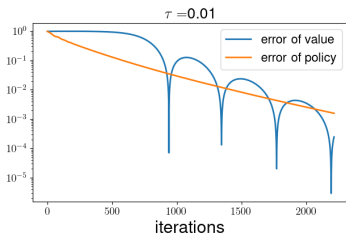
# Example 1 (no noise)



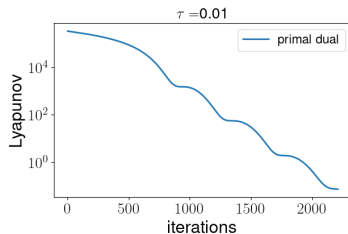
NGAD: Errors



NGAD: Lyapunov function



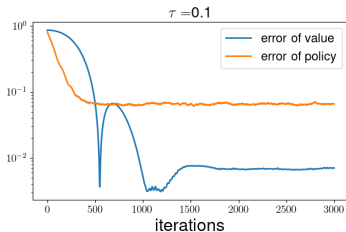
I-NGAD: Errors



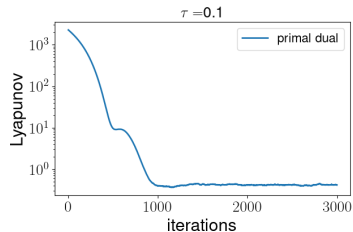
I-NGAD: Lyapunov function

## Example 2 (Gaussian noise)

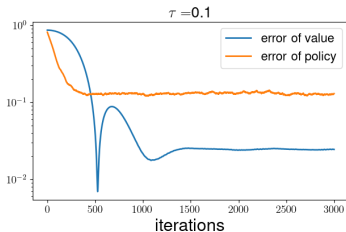
NGAD converges slowly. I-NGAD works well ( $\sigma$  noise level)



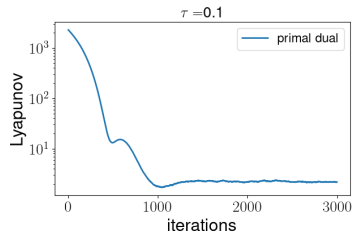
I-NGAD: Errors ( $\sigma = 0.1$ )



I-NGAD: Lyapunov function ( $\sigma = 0.1$ )



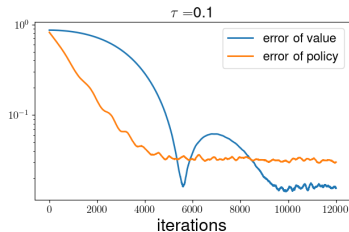
I-NGAD: Errors ( $\sigma = 0.2$ )



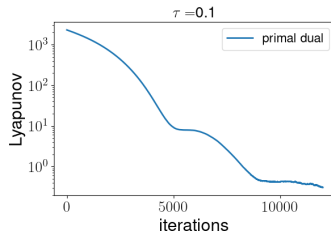
I-NGAD: Lyapunov function ( $\sigma = 0.2$ )

## Example 3 (Sampling noise)

NGAD fails to converge. I-NGAD works well



I-NGAD: Errors



I-NGAD: Lyapunov function

# Summary

## Linear/convex programming formulations for MDP

- ▶ Non-regularized and entropy-regularized
- ▶ Primal form, Bellman equation, value (policy) iteration
- ▶ Primal-dual form (actor-critic algorithms)
- ▶ Dual, policy gradient

## New algorithms

- ▶ Quasi-Newton policy gradient algorithm, dual form, super exponential convergence
- ▶ Interpolated natural gradient ascent descent (I-NGAD), primal-dual form, exponential convergence

# Thank you

## References

- ▶ Lexing Ying, Yuhua Zhu, A Note on Optimization Formulations of Markov Decision Processes. To appear in Communications in Mathematical Sciences.
- ▶ Haoya Li, Samarth Gupta, Hsiangfu Yu, Lexing Ying, Inderjit Dhillon. Quasi-Newton policy gradient algorithms. arXiv:2110.0239
- ▶ Haoya Li, Hsiangfu Yu, Lexing Ying, Inderjit Dhillon. Accelerating Primal-dual Methods for Regularized Markov Decision Processes.