

Model-based and Model-free DRL

Yu Hong
MSc D-INFK
30 Mar 2021, DNN Seminar

Outline

1. Model based vs. model free DRL
2. Guarantee and usage for Model based method
3. Successor features

Outline

1. Model based vs. model free DRL
2. Guarantee and usage for Model based method
3. Successor features

How to get more 'likes'?



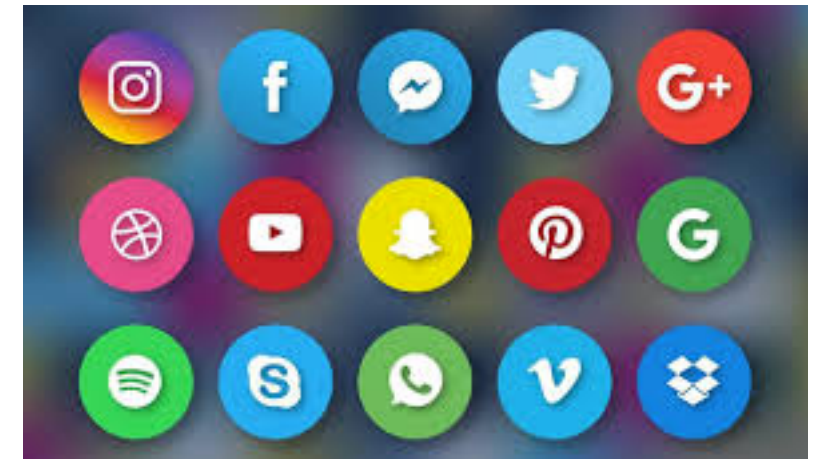
Model-based

ETH zürich Department of Humanities, Social and Political Sciences

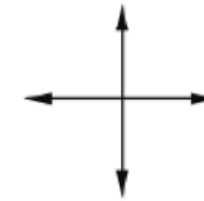
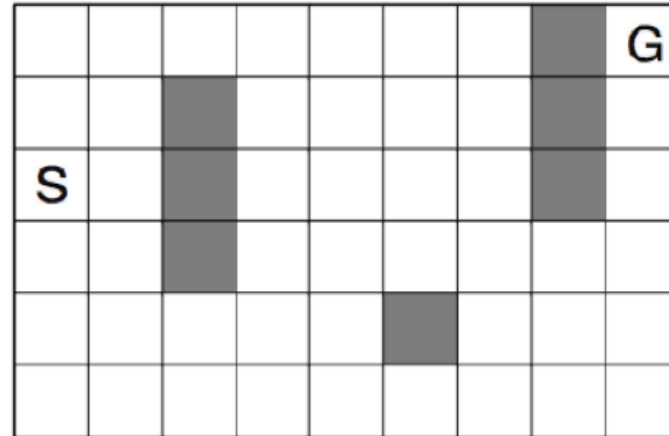
News & Events The Department Research Studies Doctoral Studies Continuing Education Information & Documents

Studies overview	Science in Perspective Teaching Committee Science in Perspective SiP FAQs	BA in Public Policy Teaching Commission BA Public Policy	MA in History and Philosophy of Knowledge Teaching Committee MA HPK	MA in Comparative and International Studies Teaching Committee MA CIS
MA in Science Didactics	MSc in Science, Technology and Policy	Teacher Training		

Model-free

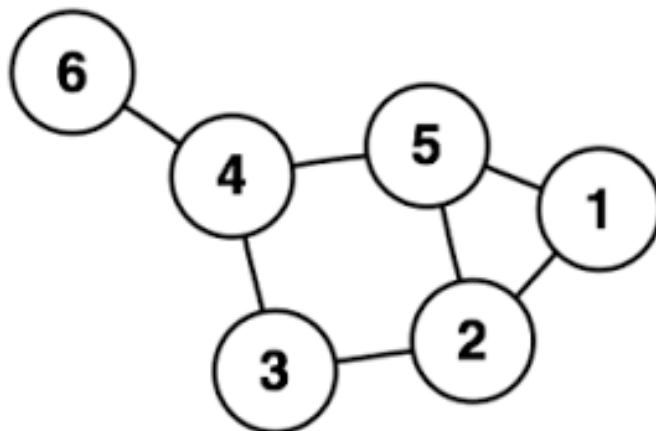


Simple maze



actions

Model-based



Model-free

	U	D	L	R
1	[0	0	0	0
2	[0	0	0	0
3	[0	0	0	0
4	[0	0	0	0

Start	State - 1	State - 2
Snake	State - 3	Treasure
	State - 4	State - 4

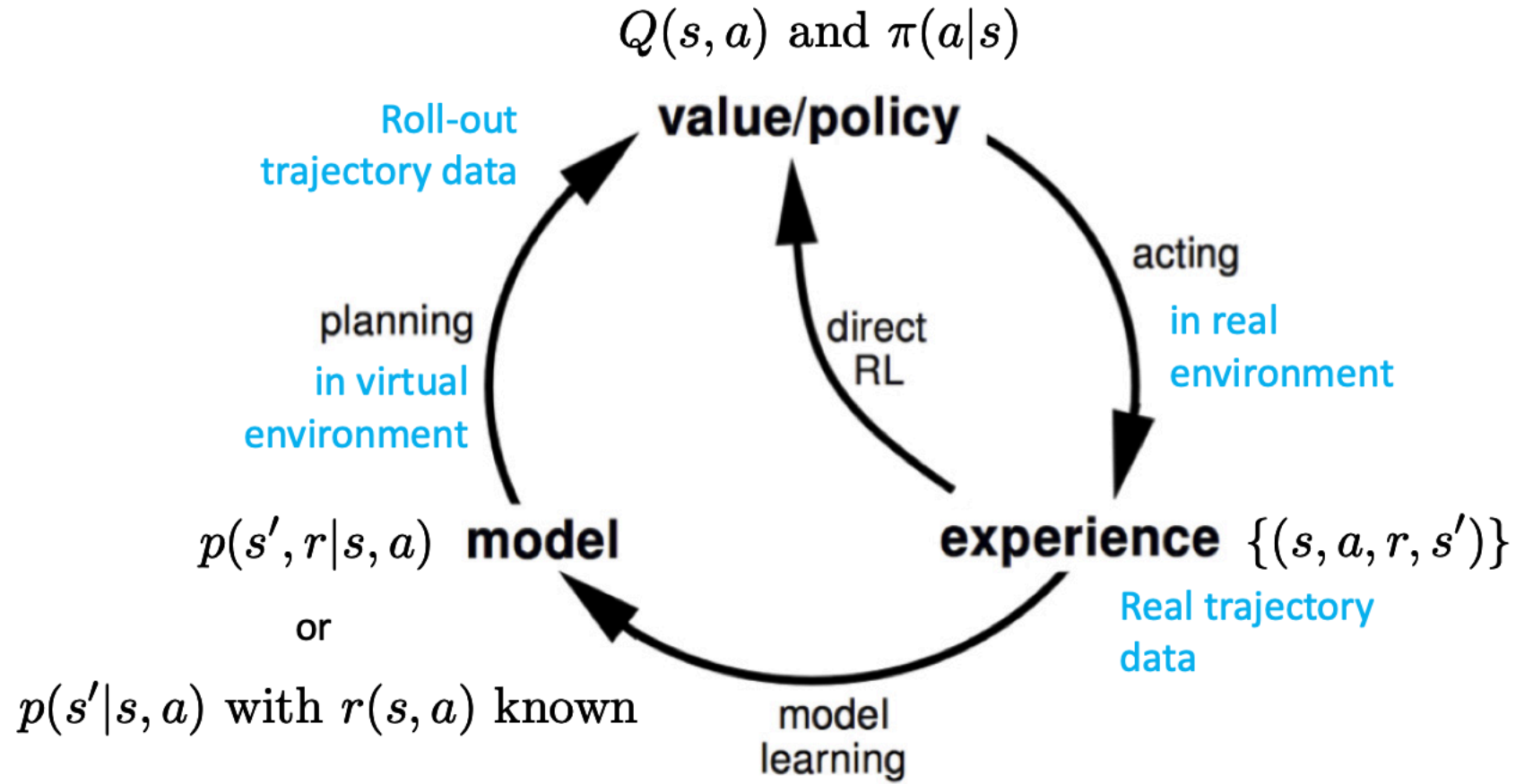
Model-based and Model-free DRL

	Model-based	Model-free
Model	supervised learning	No model
Environment	Less interaction	Interaction
Value function	Based on model	Based on samples

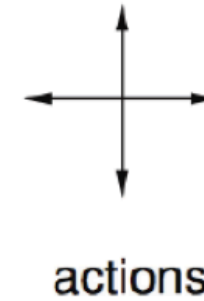
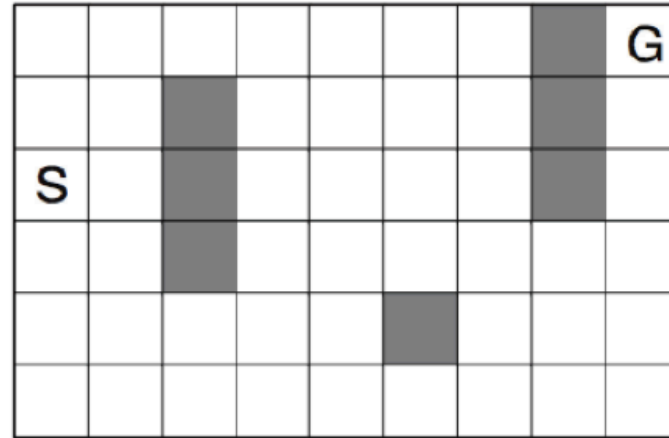
Model-based and Model-free DRL

	Model-based	Model-free
Model	On policy	On-policy/Off-policy
Environment	Less interaction	Interaction
Value function	Compounded error	Best asymptotic performance

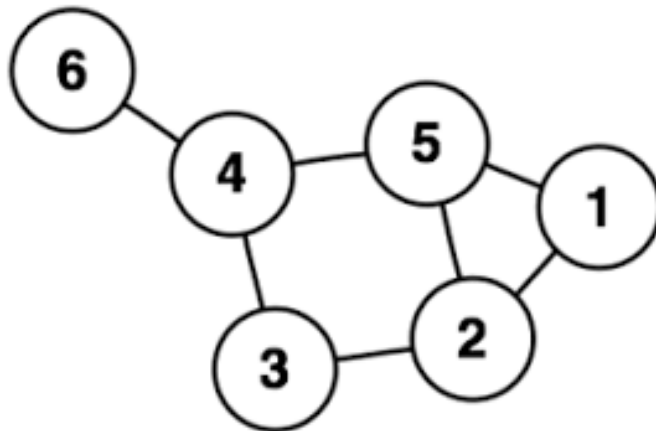
Dyna-Q



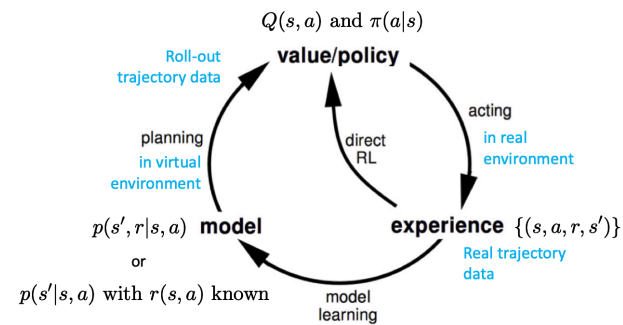
Simple maze



Model-based



Combination

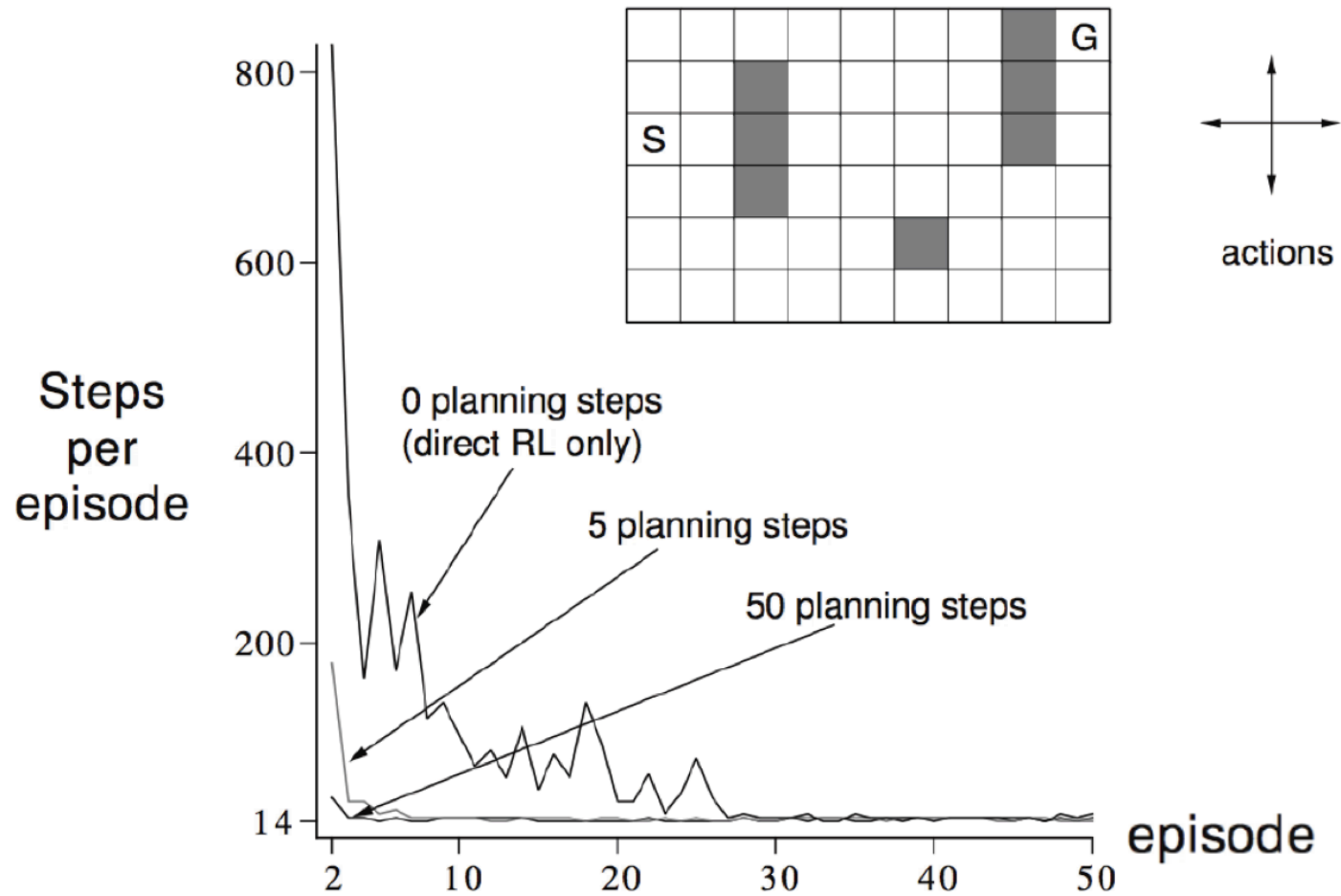


Model-free

	U	D	L	R
1	[0 0 0 0			
2	[0 0 0 0			
3	[0 0 0 0			
4	[0 0 0 0			

Start	State - 1	State - 2
Snake	State - 3	Treasure
	State - 4	State - 4

Simple maze



Outline

1. Model based vs. model free DRL
2. Guarantee and usage for Model based method
3. Successor features

Monotonic improvement guarantee

$$V^{\pi, M^*} \geq V^{\pi, \hat{M}} - D(\hat{M}, \pi)$$

assumption: (R₁) $\forall \pi, d(\pi, \pi_D) \leq \delta$

(R₂) $\hat{M} = M^* \Rightarrow D(\hat{M}, \pi) = 0$

(R₃) $D(\hat{M}, \pi)$ has form $\mathbb{E}_{\tau \sim \pi_D, M^*} [f(\hat{M}, \pi, \tau)]$

Monotonic improvement guarantee

Algorithm 1 Meta-Algorithm for Model-based RL

Inputs: Initial policy π_0 . Discrepancy bound D and distance function d that satisfy equation (R1) and (R2).

For $k = 0$ to T :

$$\pi_{k+1}, M_{k+1} = \underset{\pi \in \Pi, M \in \mathcal{M}}{\operatorname{argmax}} \quad V^{\pi, M} - D_{\pi_k, \delta}(M, \pi) \quad (3.3)$$

$$\text{s.t. } d(\pi, \pi_k) \leq \delta \quad (3.4)$$

Monotonic improvement guarantee

Theorem 3.1. Suppose that $M^* \in \mathcal{M}$, that D and d satisfy equation (R1) and (R2), and the optimization problem in equation (3.3) is solvable at each iteration. Then, Algorithm 1 produces a sequence of policies π_0, \dots, π_T with monotonically increasing values:

$$V^{\pi_0, M^*} \leq V^{\pi_1, M^*} \leq \dots \leq V^{\pi_T, M^*} \quad (3.5)$$

Moreover, as $k \rightarrow \infty$, the value V^{π_k, M^*} converges to some $V^{\bar{\pi}, M^*}$, where $\bar{\pi}$ is a local maximum of V^{π, M^*} in domain Π .

assumption: (R1) $\forall \pi, \pi_D, d(\pi, \pi_D) \leq \delta$

(R2) $\hat{M} = M^* \Rightarrow D(\hat{M}, \pi) = 0$

(R3) $D(\hat{M}, \pi)$ has form $\mathbb{E}_{\tau \sim \pi_D, M^*} [f(\hat{M}, \pi, \tau)]$

$$\pi_{k+1}, M_{k+1} = \underset{\pi \in \Pi, M \in \mathcal{M}}{\operatorname{argmax}} V^{\pi, M} - D_{\pi_k, \delta}(M, \pi) \\ \text{s.t. } d(\pi, \pi_k) \leq \delta$$

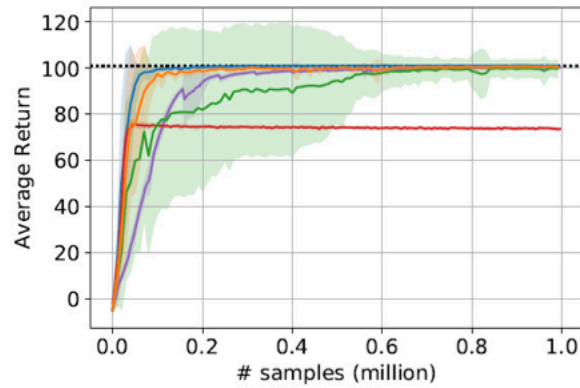
Proof: $\forall k,$

$$\begin{aligned} & V^{\pi_{k+1}, M_{k+1}^*} \stackrel{(R1)}{\geq} V^{\pi_{k+1}, M_{k+1}} - D(M_{k+1}, \pi_{k+1}) \\ & \stackrel{\text{meta algorithm}}{\geq} V^{\pi_k, M^*} - \underbrace{D(M^*, \pi_k)}_{(R2) = 0} \\ & = V^{\pi_k, M^*} \end{aligned}$$

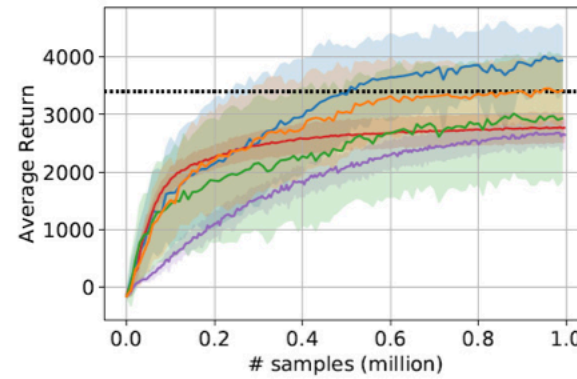
Algorithmic Framework for Model-based Deep Reinforcement Learning with Theoretical Guarantees

Meta algorithm-stochastic lower bound optimisation

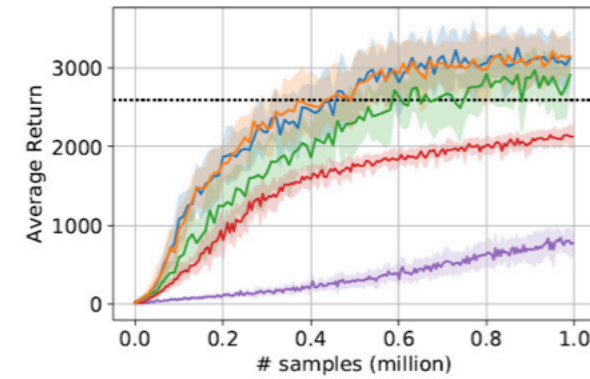
All envs with maximum horizon of 500



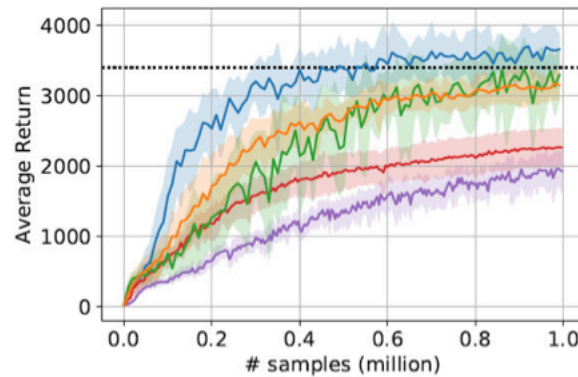
(a) Swimmer



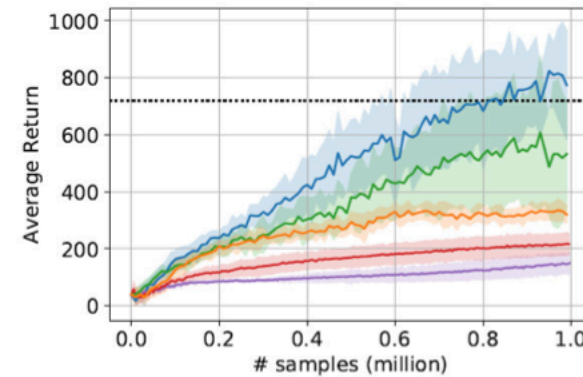
(b) Half Cheetah



(c) Ant



(d) Walker



(e) Humanoid

— SLBO — SLBO-MSE — MB-TRPO — SAC — MF-TRPO

Algorithmic Framework for Model-based Deep Reinforcement Learning with Theoretical Guarantees

How to perform rollout?

policy shift $\epsilon_{TV} = \max_S D_{TV}(\pi \parallel \pi_D)$

model shift $\epsilon_m = \max_t \mathbb{E}_{s \sim \pi_D} \left[D_{TV}(p(s', r | s, a) \parallel p_\theta(s', r | s, a)) \right]$

$$V^{\pi, M^*} \geq V^{\pi, \hat{M}} - \left[\frac{2\gamma r_{\max}(\epsilon_m + 2\epsilon_{TV})}{(1-\gamma)^2} + \frac{4\gamma r_{\max}\epsilon_{\pi}}{(1-\gamma)} \right]$$

When to Trust Your Model: Model-Based Policy Optimization

How to perform rollout?

Theorem 4.2. $V^{\text{Trk}, M^*} \geq V^{\text{Trk}, \hat{M}} - 2\gamma_{\max} \left[\frac{\gamma^{k+1} \epsilon_{\pi}}{1-\gamma} + \frac{\gamma^k + 2}{1-\gamma} \epsilon_{\pi} + \frac{k}{1-\gamma} (\epsilon_m + 2\epsilon_{\pi}) \right]$

\Rightarrow optimal $k=0$!!!

Do not use the model!

When to Trust Your Model: Model-Based Policy Optimization

How to perform rollout?

$$\text{model shift } \epsilon_m = \max_t \mathbb{E}_{s \sim \pi_t} \left[D_{TV} (p(s', r | s, a) \parallel p_\theta(s', r | s, a)) \right]$$

a new model error (first-order) approximation

$$\epsilon_{m'} = \max_t \mathbb{E}_{s \sim \pi_t} \left[D_{TV} (p(s', r | s, a) \parallel p_\theta(s', r | s, a)) \right]$$

$$\text{approximation: } \hat{\epsilon}_{m'}(\epsilon_\pi) \approx \epsilon_m + \epsilon_\pi \frac{d\epsilon_{m'}}{d\epsilon_\pi}$$

When to Trust Your Model: Model-Based Policy Optimization

How to perform rollout?

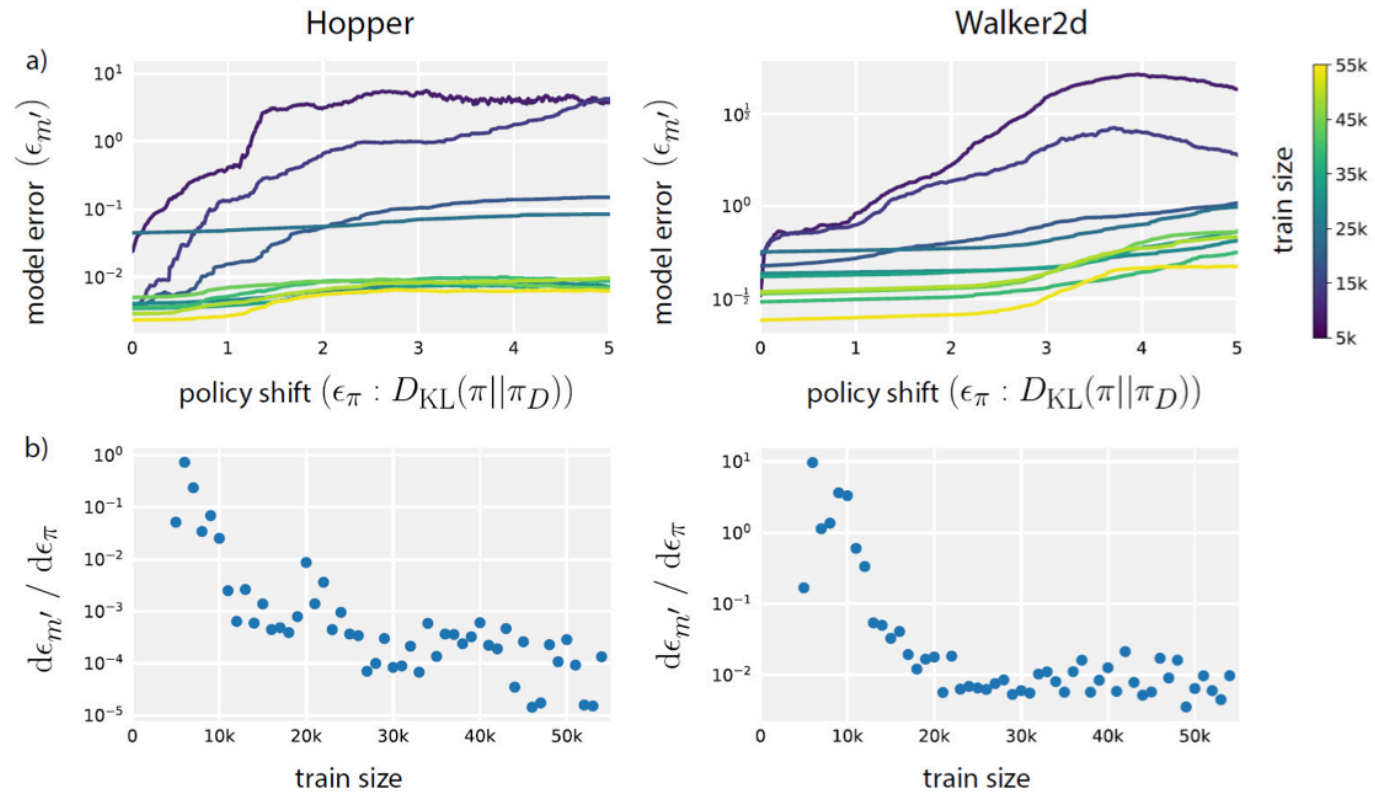
The bound is rewritten as

$$V^{\pi_k, M^*} \geq V^{\pi_k, \hat{M}} - 2r_{\max} \left[\frac{\gamma^{k+1} \epsilon_{\pi}}{(1-\gamma)^2} + \frac{\gamma^k \epsilon_{\pi}}{1-\gamma} + \frac{k}{1-\gamma} \epsilon_{M'} \right]$$

\Rightarrow if $\frac{d\epsilon_{M'}}{d\epsilon_{\pi}}$ small enough, $k \geq 0$

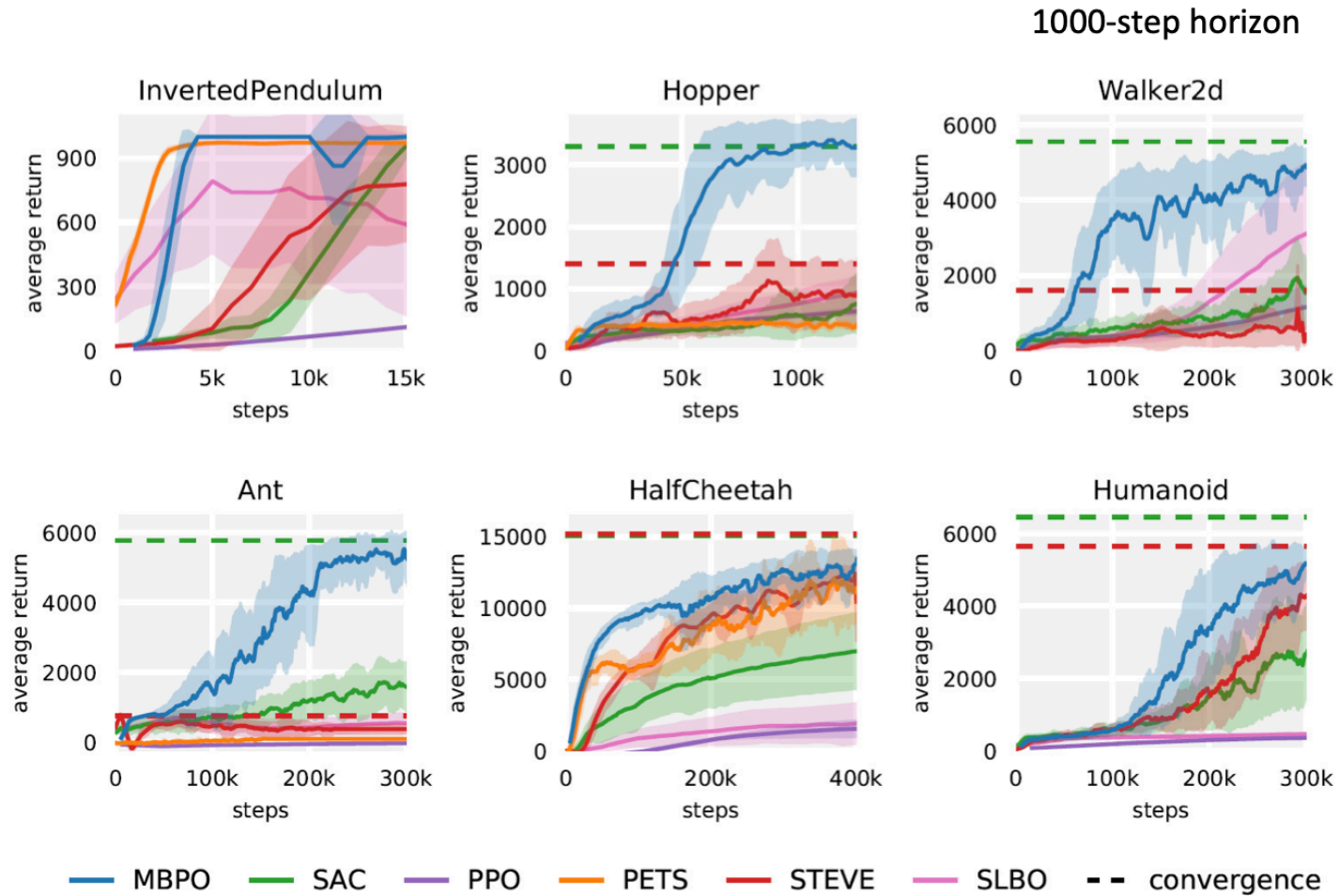
When to Trust Your Model: Model-Based Policy Optimization

Empirical analysis of $\frac{d\epsilon_{m'}}{d\epsilon_{\pi}}$



When to Trust Your Model: Model-Based Policy Optimization

Model-based policy optimisation

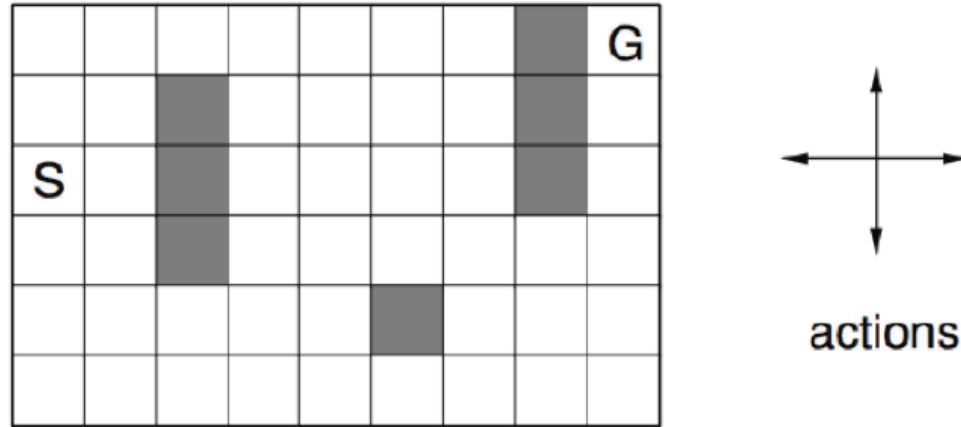


When to Trust Your Model: Model-Based Policy Optimization

Outline

1. Model based vs. model free DRL
2. Guarantee and usage for Model based method
3. Successor features

Successor representation



$$M^\pi(s_i, s_j) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k \mathbb{I}[s_{t+k} = s_j] \mid s_t = s_i \right]$$

$$V^\pi(s_i) = \sum_j M^\pi(s_i, s_j) \underbrace{R(s_j)}$$

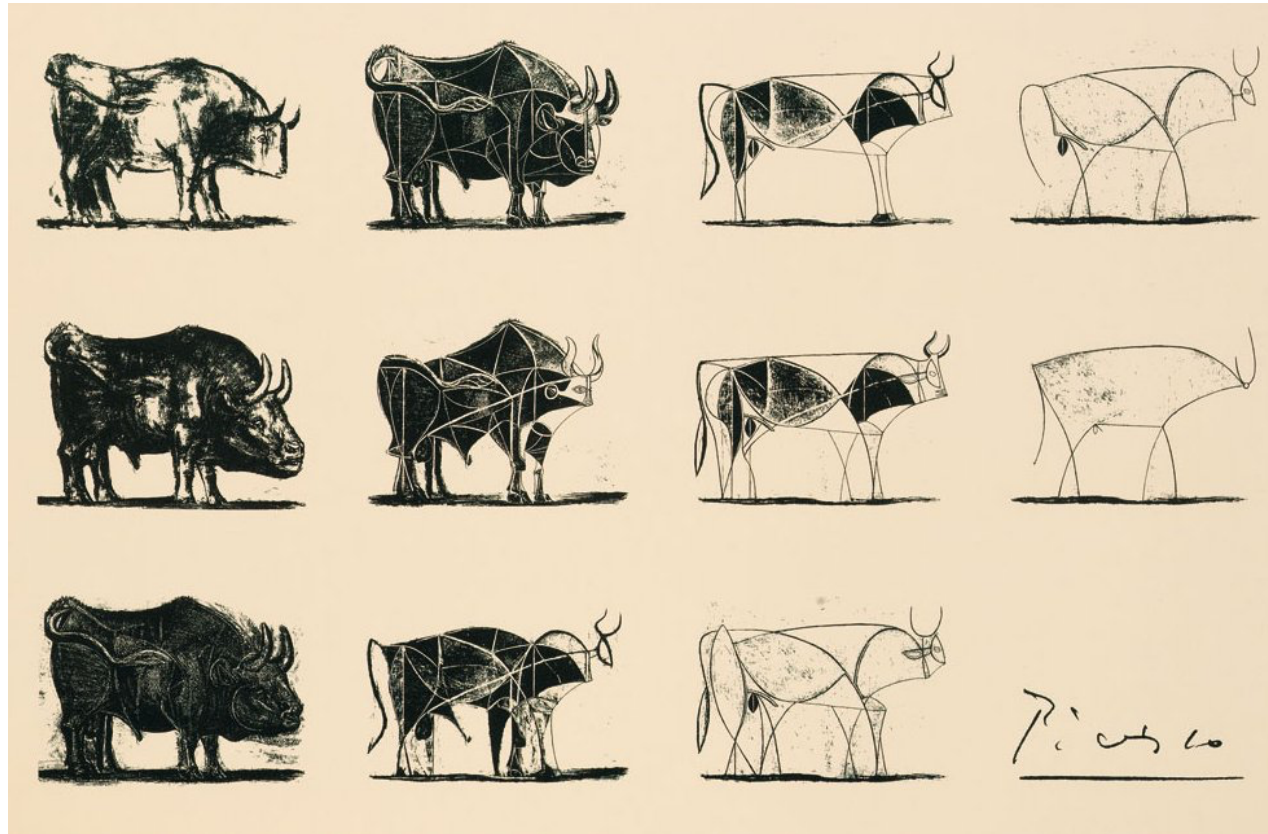
assumption: $R(s_j) = \phi(s_j) \cdot \omega$

$$= \sum_j M^\pi(s_i, s_j) \phi(s_j) \omega = \underbrace{\left[\sum_j M^\pi(s_i, s_j) \phi(s_j) \right]}_{\psi(s_i)} \cdot \omega$$

summary

Model based/ model free/ successor features

Justify model usage and how to use model



“I never made a painting as a work of art. It's all research.” – Pablo Picasso

Thanks



Dayan, Peter. "Improving generalization for temporal difference learning: The successor representation." *Neural Computation* 5.4 (1993): 613-624.

Luo, Yuping, et al. "Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees." *arXiv preprint arXiv:1807.03858* (2018).

Barreto, André, et al. "Successor features for transfer in reinforcement learning." *arXiv preprint arXiv:1606.05312* (2016).

Barreto, Andre, et al. "Transfer in deep reinforcement learning using successor features and generalised policy improvement." *International Conference on Machine Learning*. PMLR, 2018.

Janner, Michael, et al. "When to trust your model: Model-based policy optimization." *arXiv preprint arXiv:1906.08253* (2019).

