

# On Model Accuracy and Planning Horizon

Nan Jiang, Alex Kulesza, Satinder Singh  
 Computer Science & Engineering, University of Michigan  
 Richard Lewis  
 Department of Psychology, University of Michigan



## Abstract

For MDPs with a large discount factor, it is common to use reduced horizon in planning to speed computation. However, when the model available to the agent is inaccurate, the policy found using a shorter planning horizon can be actually better than the one learned with the true horizon. The reason, we argue, is that planning horizon is a complexity control parameter for the policy class to be learned. We propose two complexity measures controlled by planning horizon. The implications provided by the loss bounds based on the proposed complexity measures are empirically verified.

## Outline

Given MDP  $M = \langle S, A, T, R, \gamma_{\text{eval}} \rangle$ , consider an agent that

- plans with an inaccurate model  $\widehat{M}$ , estimated from dataset of size  $n$ ,
- uses a discount factor  $\gamma \leq \gamma_{\text{eval}}$  in planning,
- outputs  $\pi_{\widehat{M}, \gamma}^*$  (i.e. the **Certainty-Equivalent** policy)

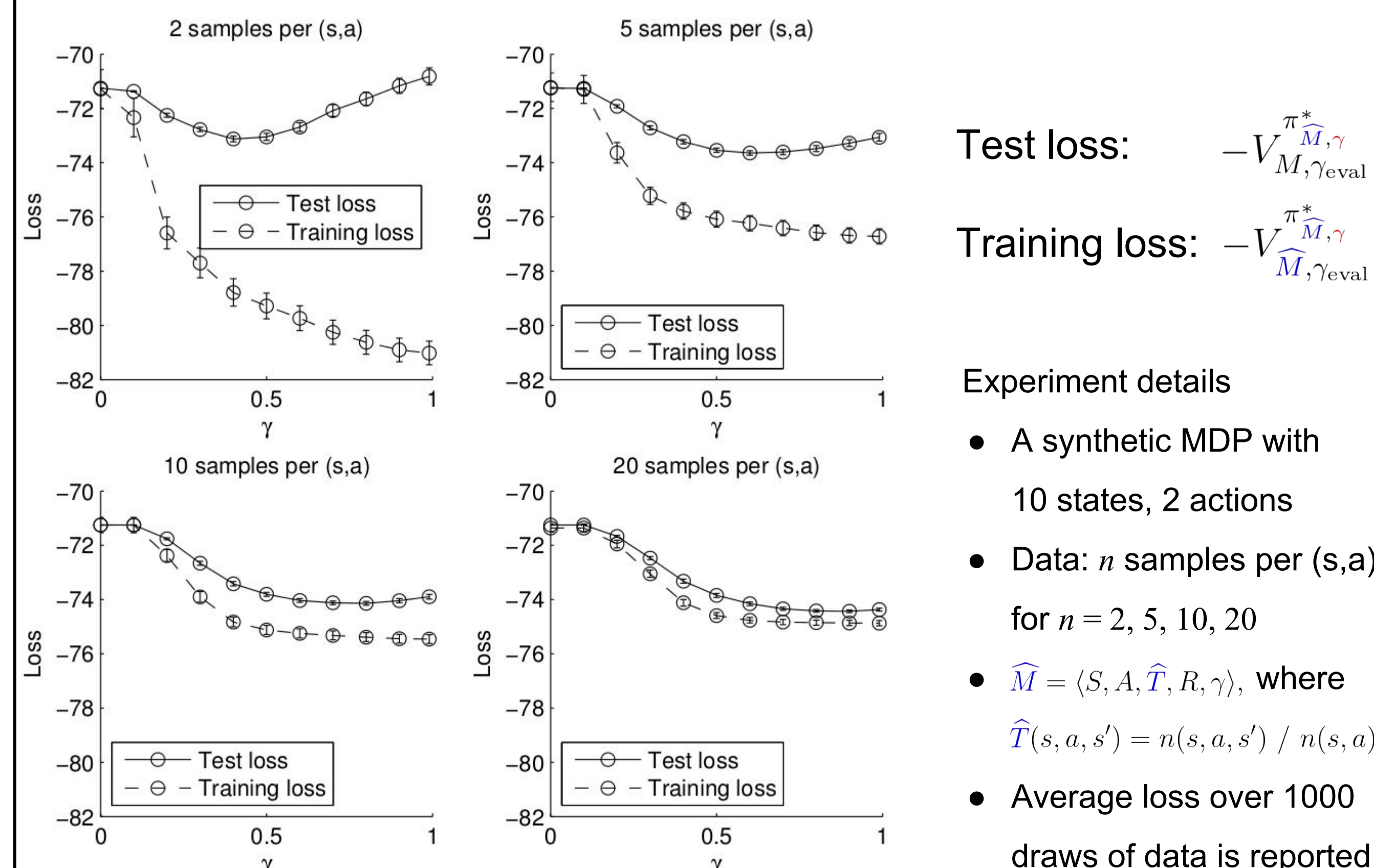
Question: what is the optimal  $\gamma$  that minimizes planning loss, namely

$$\gamma^* = \arg \min_{\gamma \in [0, \gamma_{\text{eval}}]} \left\| V_{M, \gamma_{\text{eval}}}^{\pi_{\widehat{M}, \gamma}^*} - V_{M, \gamma_{\text{eval}}}^{\pi_{\widehat{M}, \gamma}^*} \right\|_{\infty}$$

Our work:

- Empirical observation:  $\gamma^* < \gamma_{\text{eval}}$ ;  $n \uparrow \Rightarrow \gamma^* \downarrow$
- Theoretical explanation:  $\gamma$  is a complexity control parameter

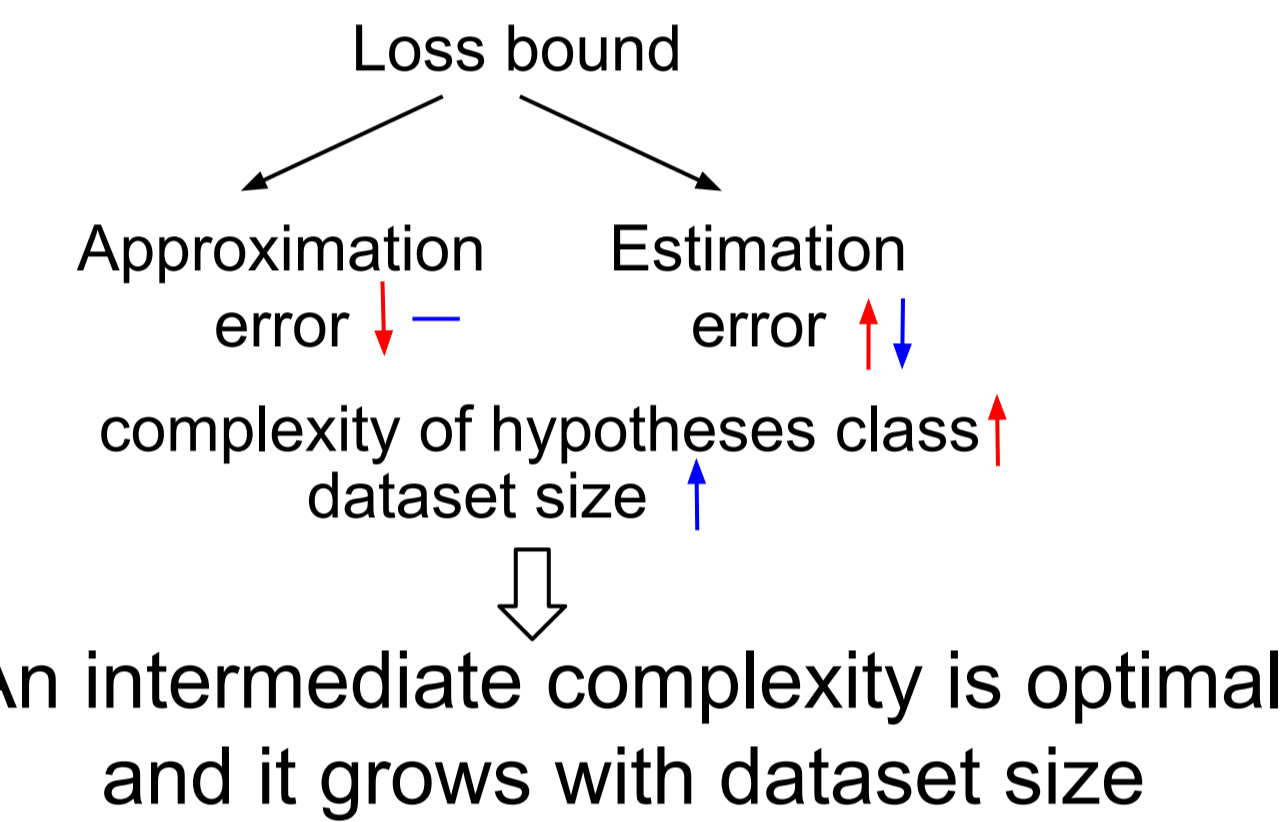
## Empirical Illustration



## Theoretical framework

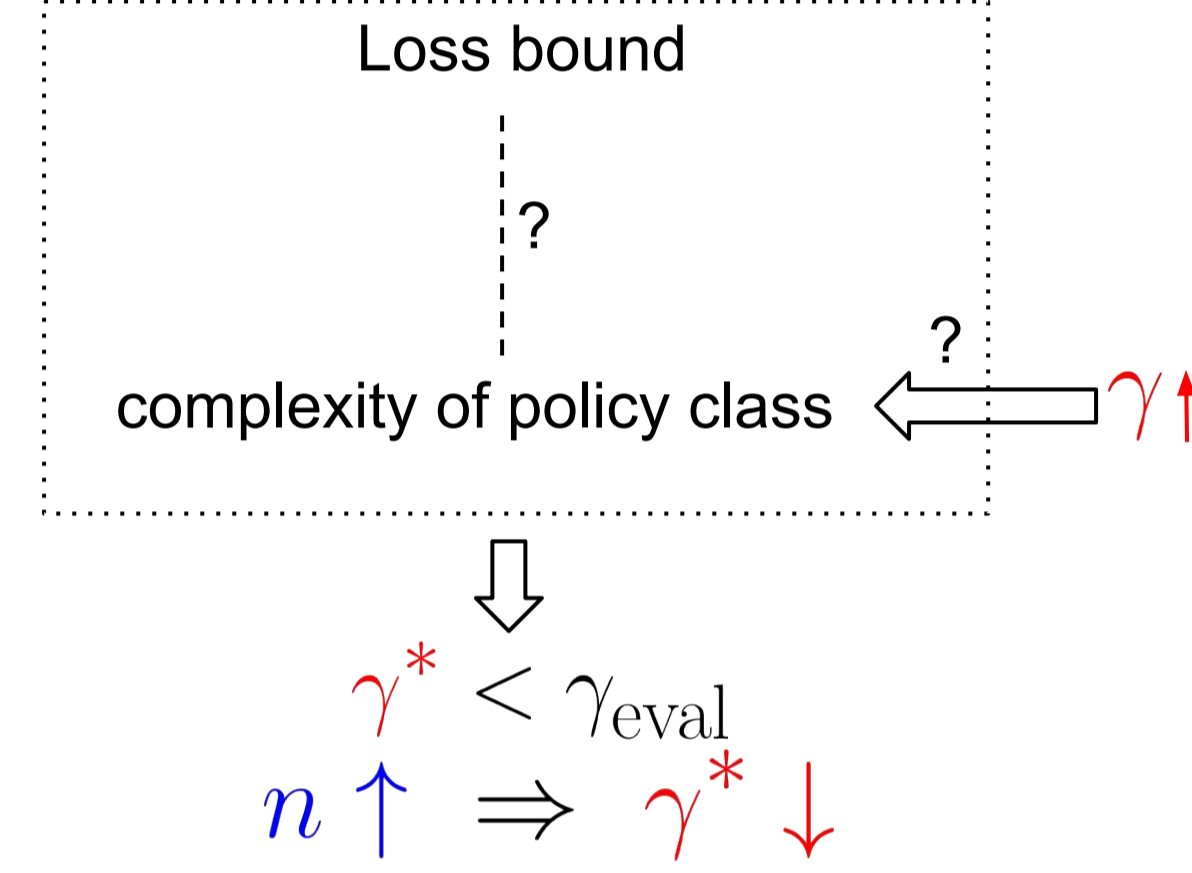
Empirical risk minimization

- A dataset
- A class of hypotheses
- Select the hypothesis that minimizes training error



Certainty-equivalent planning

- Empirical model  $\widehat{M}$
- A class of policies
- Select the policy that maximizes  $V_{\widehat{M}, \gamma}^{\pi}$



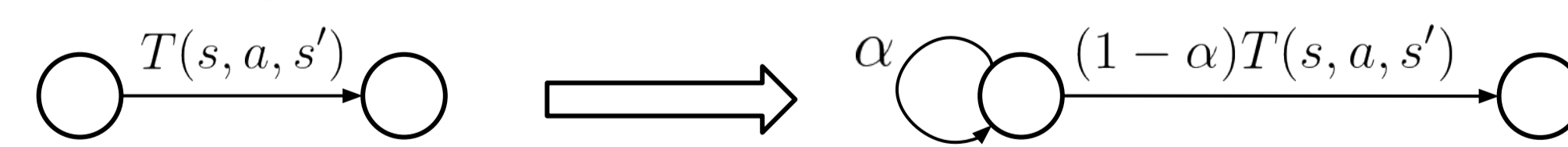
## A Counting Complexity Measure

- Observation: fixing  $R$  and  $\gamma$ , some  $\pi$  can never be optimal
- Define  $\Pi_{R, \gamma}$  as the set of policies that *can* be optimal under  $R$  and  $\gamma$ , i.e.  $\Pi_{R, \gamma} = \{ \pi : \exists T, \text{ s.t. } \pi \text{ is optimal in } \langle S, A, T, R, \gamma \rangle \}$

- Theorem 1 ( $|\Pi_{R, \gamma}|$  increasing with  $\gamma$ )

- $\gamma \leq \gamma' \Rightarrow \Pi_{R, \gamma} \subseteq \Pi_{R, \gamma'}$
- $|\Pi_{R, 0}| = 1$ , if  $\forall s \in S, \arg \max_{a \in A} R(s, a)$  is unique.
- $\exists \gamma$ , s.t.  $|\Pi_{R, \gamma}| \geq |A|^{|S|-2}$ , if  $\exists s, s' \in S, \max_{a \in A} R(s, a) > \max_{a' \in A} R(s', a')$ .

Idea for proving Claim 1: if  $T$  makes  $\pi$  optimal under  $\gamma$ , construct  $T'$ :



Any policy's value function w.r.t.  $(T, \gamma)$  is proportional to that w.r.t.  $(T', \gamma')$ .

- Theorem 2 (Planning loss bound). w.h.p. planning loss is bounded by

$$\frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} R_{\max} + \frac{2R_{\max}}{(1 - \gamma)^2} \sqrt{\frac{1}{2n} \log \frac{2|S||A||\Pi_{R, \gamma}|}{\delta}}$$

$$\max_{\pi: S \rightarrow A} \left\| V_{M, \gamma_{\text{eval}}}^{\pi} - V_{M, \gamma}^{\pi} \right\|_{\infty}$$

## Rademacher Complexity Measure

- More general: does not require fixed reward

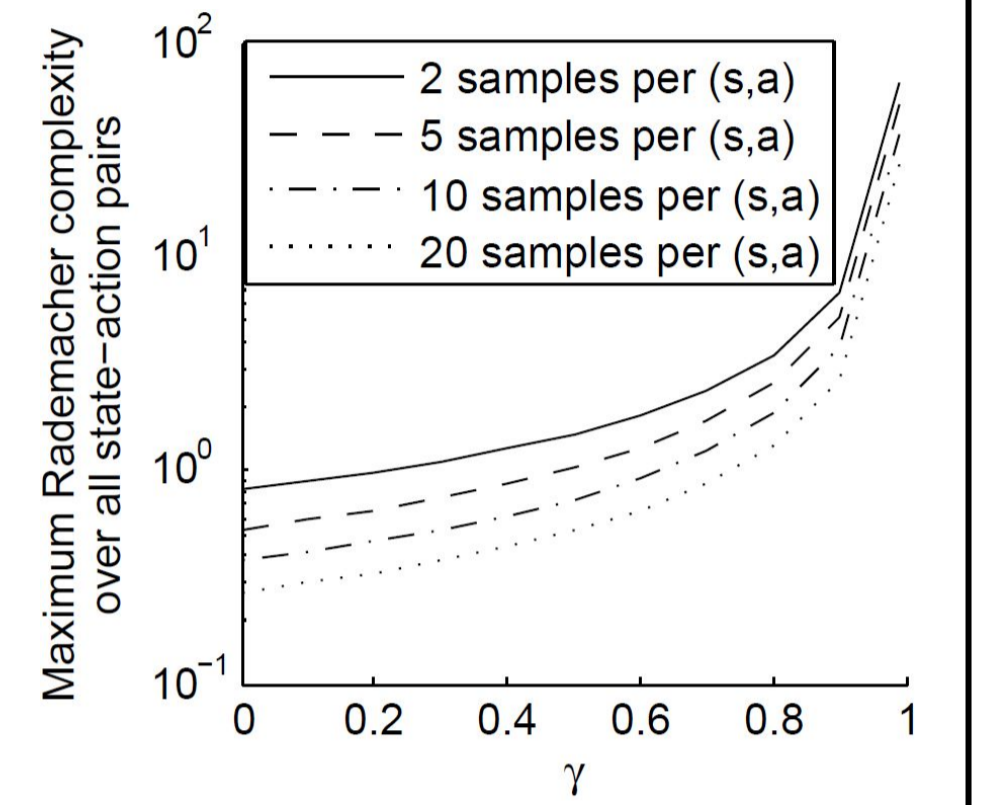
- Define function class  $\mathcal{F}_{M, \gamma} = \{ f_{M, \gamma}^{\pi} : \pi \in S \rightarrow A \}$ , where

$$f_{M, \gamma}^{\pi}(r, s') = r + \gamma V_{M, \gamma}^{\pi}(s')$$

- Complexity measure:  $\max_{s \in S, a \in A} \widehat{\mathfrak{R}}_{D_{s, a}}(\mathcal{F}_{M, \gamma})$

- $D_{s, a}$ : reward & next-state pairs sampled from  $(s, a)$

$$\widehat{\mathfrak{R}}_{D_{s, a}}(\mathcal{F}_{M, \gamma}) = \mathbb{E}_{\sigma_i \sim \text{Unif}\{-1, 1\}} \left\{ \sup_{\pi: S \rightarrow A} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i f_{M, \gamma}^{\pi}(r_i, s'_i) \right) \right\}$$



- Increasing with  $\gamma$ : shown empirically

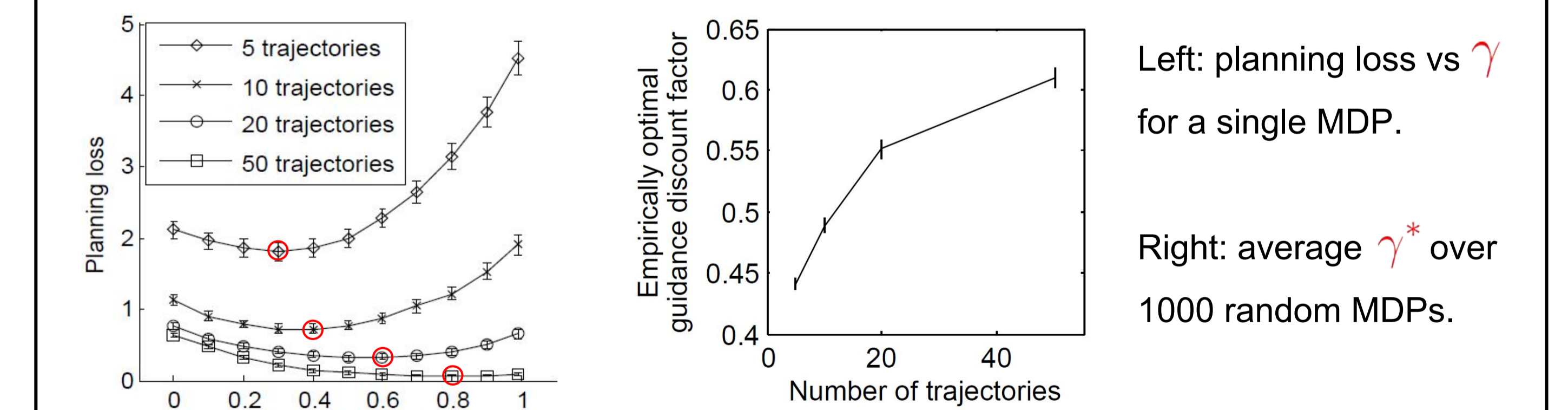
- Planning loss bound:

$$\frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} R_{\max} + \frac{1}{1 - \gamma} \left( 2 \max_{s \in S, a \in A} \widehat{\mathfrak{R}}_{D_{s, a}}(\mathcal{F}_{M, \gamma}) + \frac{3R_{\max}}{1 - \gamma} \sqrt{\frac{1}{2n} \log \frac{4|S||A|}{\delta}} \right)$$

$$\max_{s \in S, a \in A} \max_{\pi: S \rightarrow A} \left| \frac{1}{n} \sum_{(r, s') \in D_{s, a}} f_{M, \gamma}^{\pi}(r, s') - \mathbb{E}_{(r, s') \sim \mathbb{P}_{s, a}} \{ f_{M, \gamma}^{\pi}(r, s') \} \right|$$

## Experiments in more realistic settings

- Synthetic MDPs, data collected via trajectories



Experiment details: MDPs with 10 states, 2 actions, 5 next-states, T and R generated randomly. Trajectories start from random initial state followed by purely random exploration.

- Optimal planning depth in UCT
  - UCT has a perfect simulator, but still needs to translate sample trajectories to action choice
  - #trajectories  $\uparrow \Rightarrow$  optimal planning depth  $\uparrow$

Experiment details: RockSample (7x8) domain, exploration parameter separately optimized for each (planning depth, #trajectories) pair.

