



POLITECNICO
MILANO 1863

CONFIGURABLE MARKOV DECISION PROCESSES

A. M. METELLI, M. MUTTI AND M. RESTELLI

{albertomaria.metelli, marcello.restelli}@polimi.it
{mirco.mutti}@mail.polimi.it



MOTIVATIONS AND PROBLEM

- In many real-world problems it is possible to **configure the environment** during the learning process:
 - Driving on a track*: the driver/engineer configures the vehicle settings (e.g., seasonal tires, stability and vehicle attitude, engine model, ...)
 - Student-Teacher interaction*: the teacher adapts its teaching style to the student's needs (e.g., simplify the questions, slow down the presentation of concepts, ...)
- Configuring the environment can bring some **benefits**:
 - learn policy with **higher performances** (e.g., configure the car to learn a better driving policy)
 - speed up** the learning process (e.g., configure the car to learn faster a good driving policy)
- How to choose an **environment configuration** so that the performance of the **optimal policy** is maximized?

CONTRIBUTIONS

- We propose a new framework, **Configurable Markov Decision Process** (Conf-MDP), to model sequential decision-making problems in which the **transition model** is configurable.
- We present a safe learning algorithm, **Safe Policy-Model Iteration** (SPMI), able to *jointly* and *adaptively* optimize the policy and the environment configuration.
- We provide an **empirical evaluation** to highlight the benefits of the environment configuration and the effectiveness of SPMI.

CONFIGURABLE MARKOV DECISION PROCESSES

Markov Decision Process (MDP)

$$(S, A, R, P, \gamma, \mu)$$

Configurable Markov Decision Process (Conf-MDP)

$$(S, A, R, \gamma, \mu, \mathcal{P}, \Pi)$$

- Learning in an MDP means to recover an **optimal policy** π^* under a **fixed environment** P
- The environment P is fixed and **cannot** be controlled
- The policy π belongs to a policy space Π

- Learning in a Conf-MDP means to recover an **optimal policy** π^* and an **optimal environment configuration** P^*
- The environment P belongs to a the model space \mathcal{P} and **can** be configured
- The policy π belongs to a policy space Π

SAFE POLICY-MODEL ITERATION

- SPMI finds the optimal model-policy pair:

$$P^*, \pi^* = \arg \max_{P \in \mathcal{P}, \pi \in \Pi} J_\mu^{P, \pi}$$

- At each iteration SPMI decides whether to update the **model** P , the **policy** π or **both** (*joint* and *adaptive*)
- SPMI is **safe** as it ensures a **monotonic performance improvement** [Pirota et al., 2013]
- SPMI optimizes a **lower-bound** of the performance improvement:

$$J_\mu^{P', \pi'} - J_\mu^{P, \pi} \geq B(P', \pi', P, \pi)$$

BOUND ON PERFORMANCE IMPROVEMENT

$$B(P', \pi', P, \pi) = \frac{\mathbb{A}_{P, \pi, \mu}^{P', \pi} + \mathbb{A}_{P, \pi, \mu}^{P, \pi'}}{1 - \gamma} - \frac{\gamma \Delta Q^{P, \pi} D}{2(1 - \gamma)^2}$$

Labels: *dissimilarity*, *model advantage*, *policy advantage*

POLICY AND MODEL UPDATE

$$P' = (1 - \beta)P + \beta \bar{P} \quad \pi' = (1 - \alpha)\pi + \alpha \bar{\pi}$$

ALGORITHM

```

initialize  $\pi_0, P_0$ 
for  $i = 0, 1, 2, \dots$  until convergence do
   $\bar{P}_i = \text{ModelChooser}(P_i)$ 
   $\bar{\pi}_i = \text{PolicyChooser}(\pi_i)$ 
   $\alpha_i, \beta_i = \arg \max_{\alpha, \beta \in [0, 1]} B(P', \pi', P, \pi)$ 
   $\pi_{i+1} = \alpha_i \bar{\pi}_i + (1 - \alpha_i) \pi_i$ 
   $P_{i+1} = \beta_i \bar{P}_i + (1 - \beta_i) P_i$ 
end for
  
```

MORE ON CONF-MDPs

MODEL AND POLICY CHOOSER

How to select the **target** \bar{P}_i and $\bar{\pi}_i$?

- greedy** chooser
 - select the policy/model that maximizes the advantage function
 - might generate oscillations
- persistent** chooser
 - select between the *greedy* and the old target
 - avoids oscillations

MODEL AND POLICY SPACES

How to represent the model space \mathcal{P} ?

- unconstrained** model space: any model is valid
- parametric** model space: the model P_ω depends on parameters ω
 - e.g., convex hull of *vertex models*

P-GRADIENT THEOREM

Extension of the *Policy Gradient Theorem* [Sutton et al., 2000] to model learning:

$$\nabla_\omega J_\mu^{P_\omega} = \int \delta_\mu^{P_\omega}(s, a) \nabla_\omega P_\omega(s'|s, a) U^{P_\omega}(s, a, s') ds' da ds.$$

$$U^{P_\omega}(s, a, s') = R(s, a) + \gamma V^{P_\omega}(s') - V^{P_\omega}(s).$$

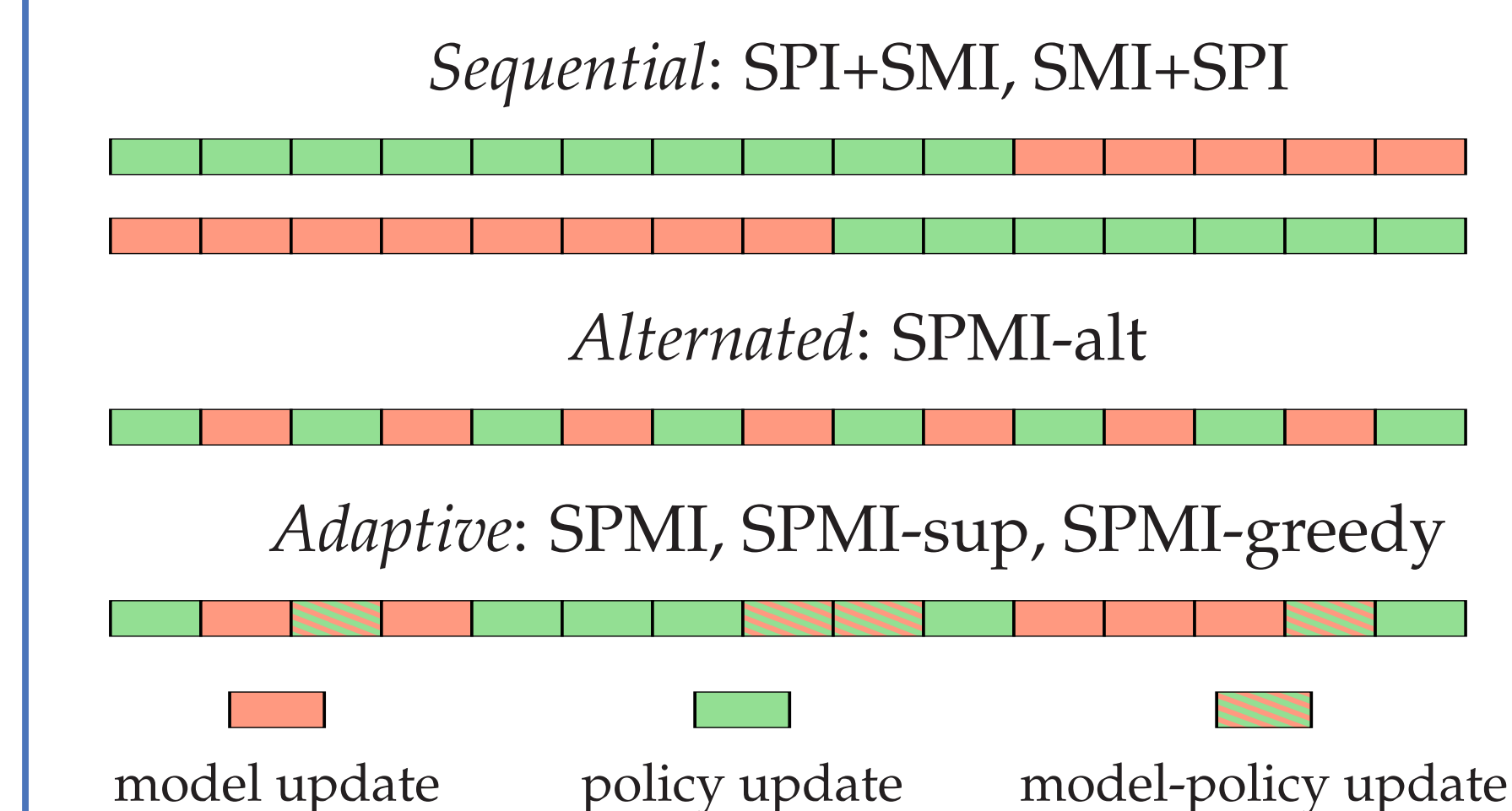
REFERENCES

Matteo Pirota, Marcello Restelli, Alessio Pecorino, and Daniele Calandriello. Safe policy iteration. In *Proceedings of the 30th International Conference on International Conference on Machine Learning*, volume 28 of ICML'13, pages 307–315, 2013.

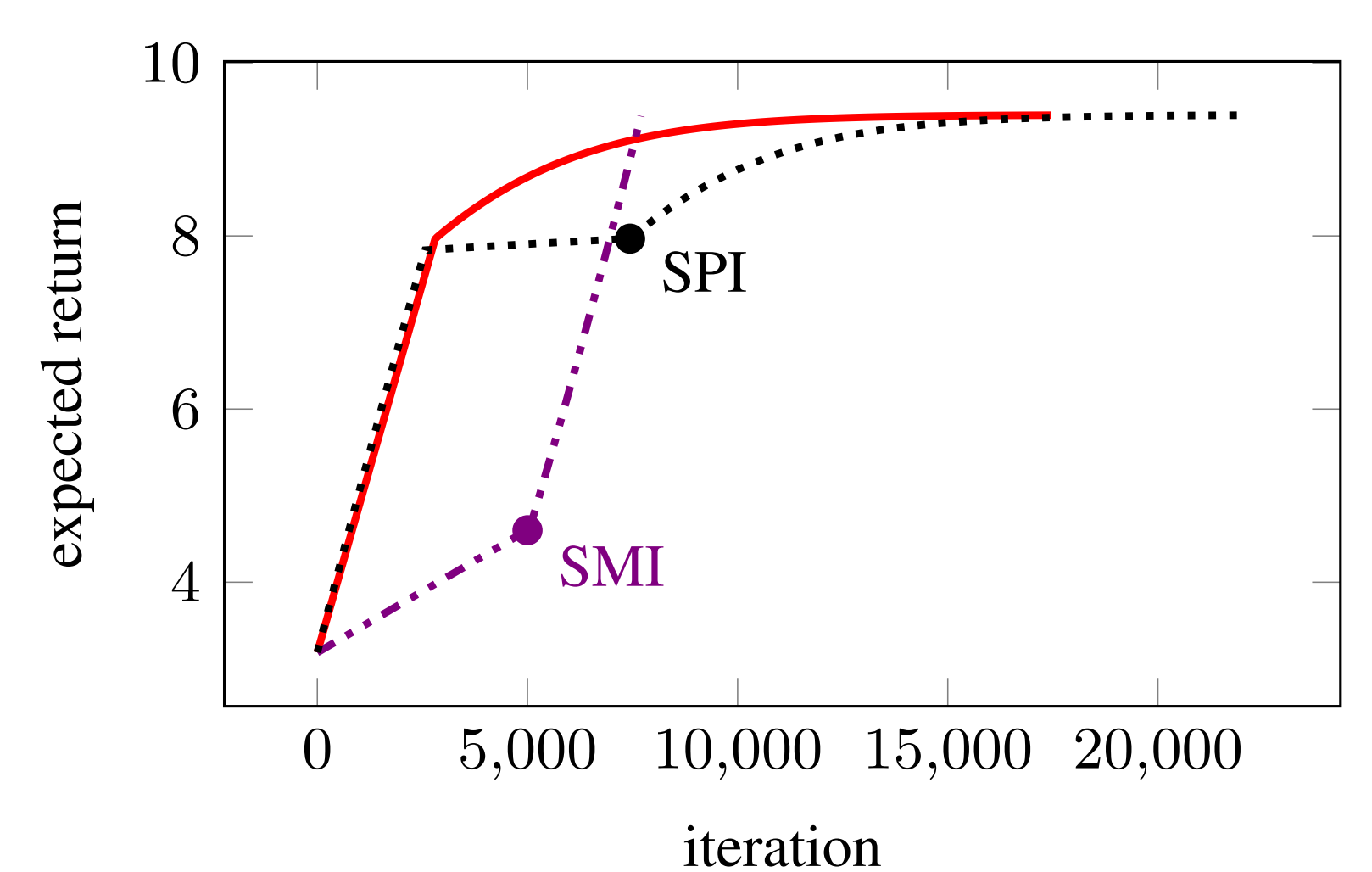
Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, pages 1057–1063, 2000.

EXPERIMENTAL EVALUATION

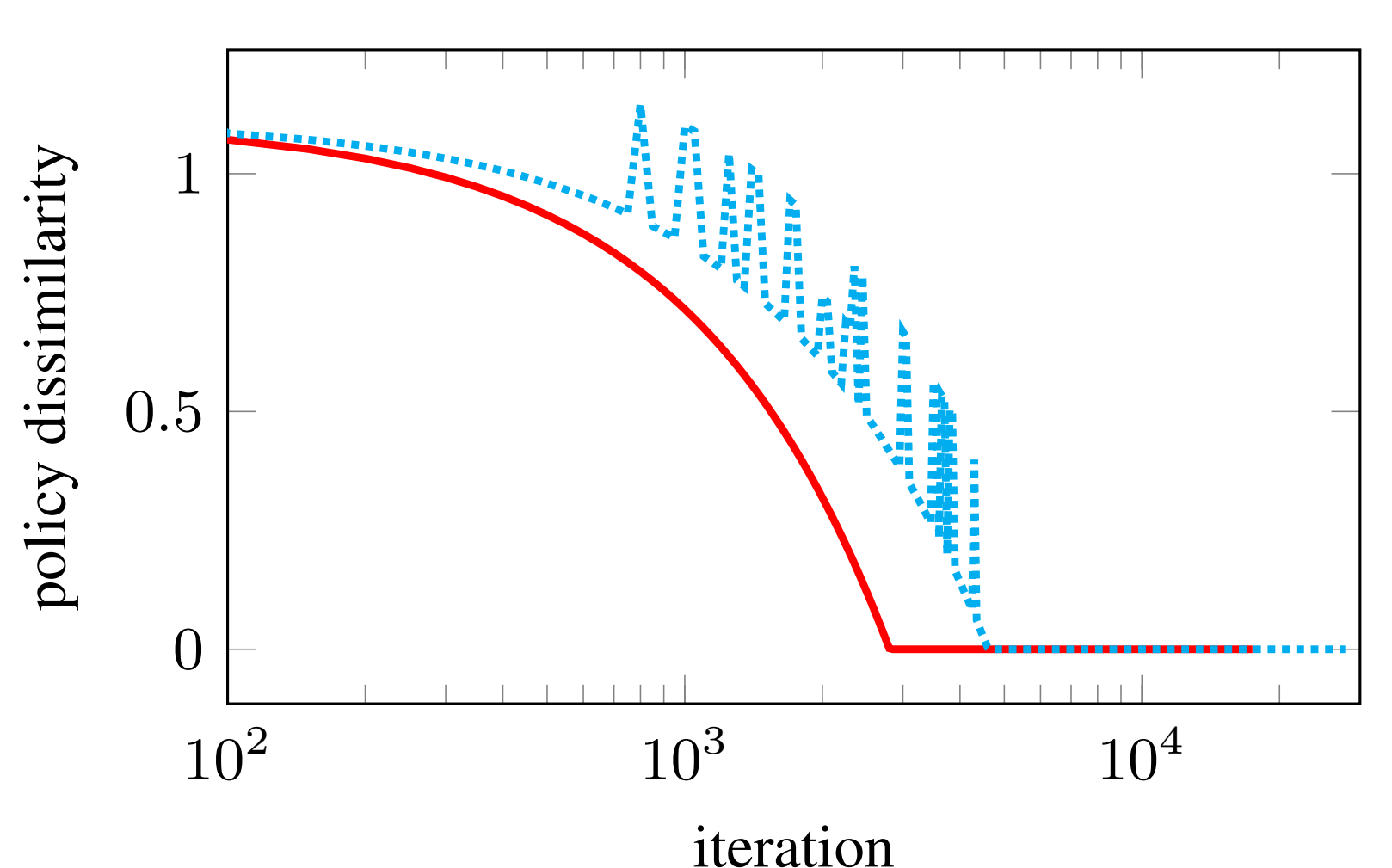
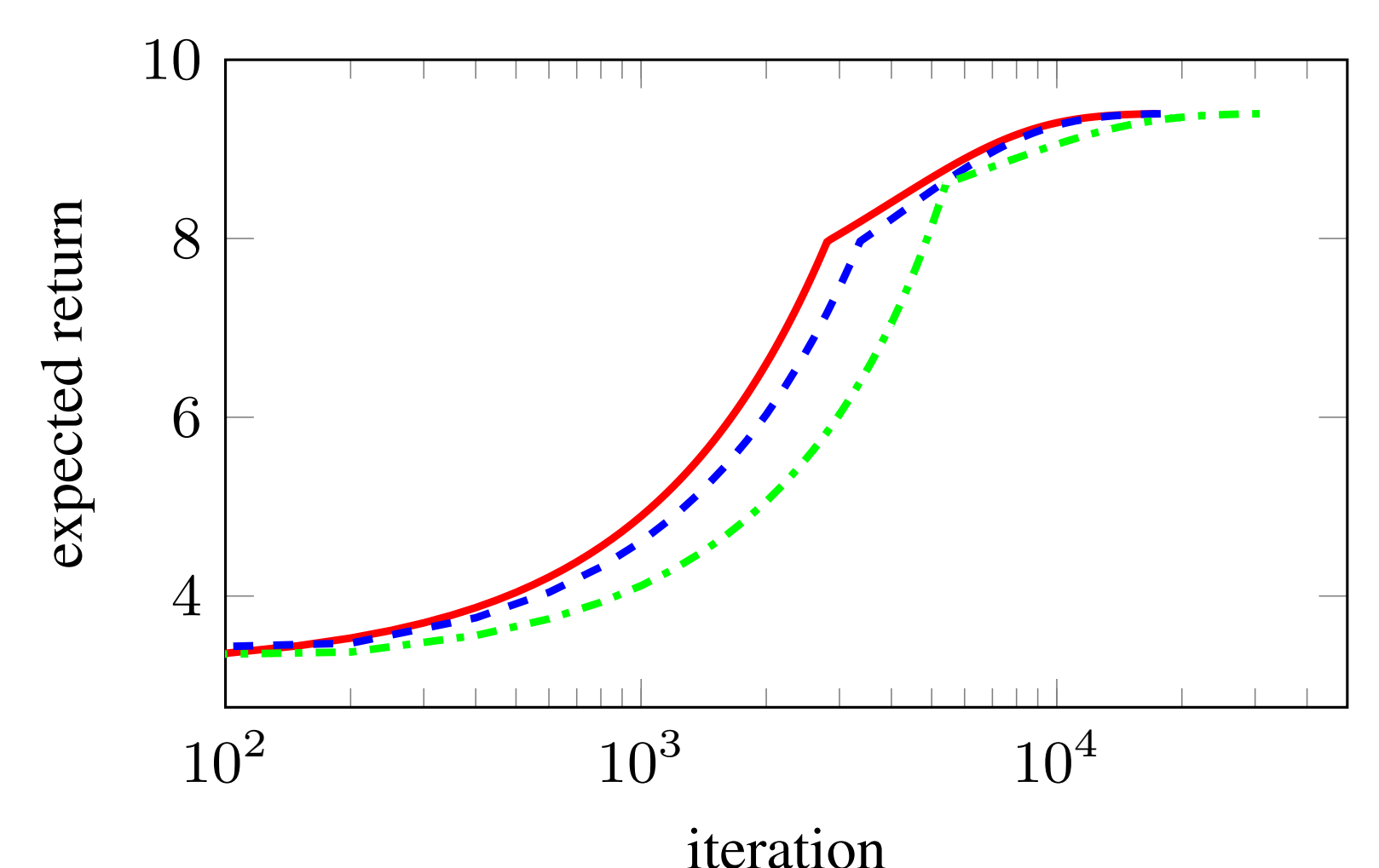
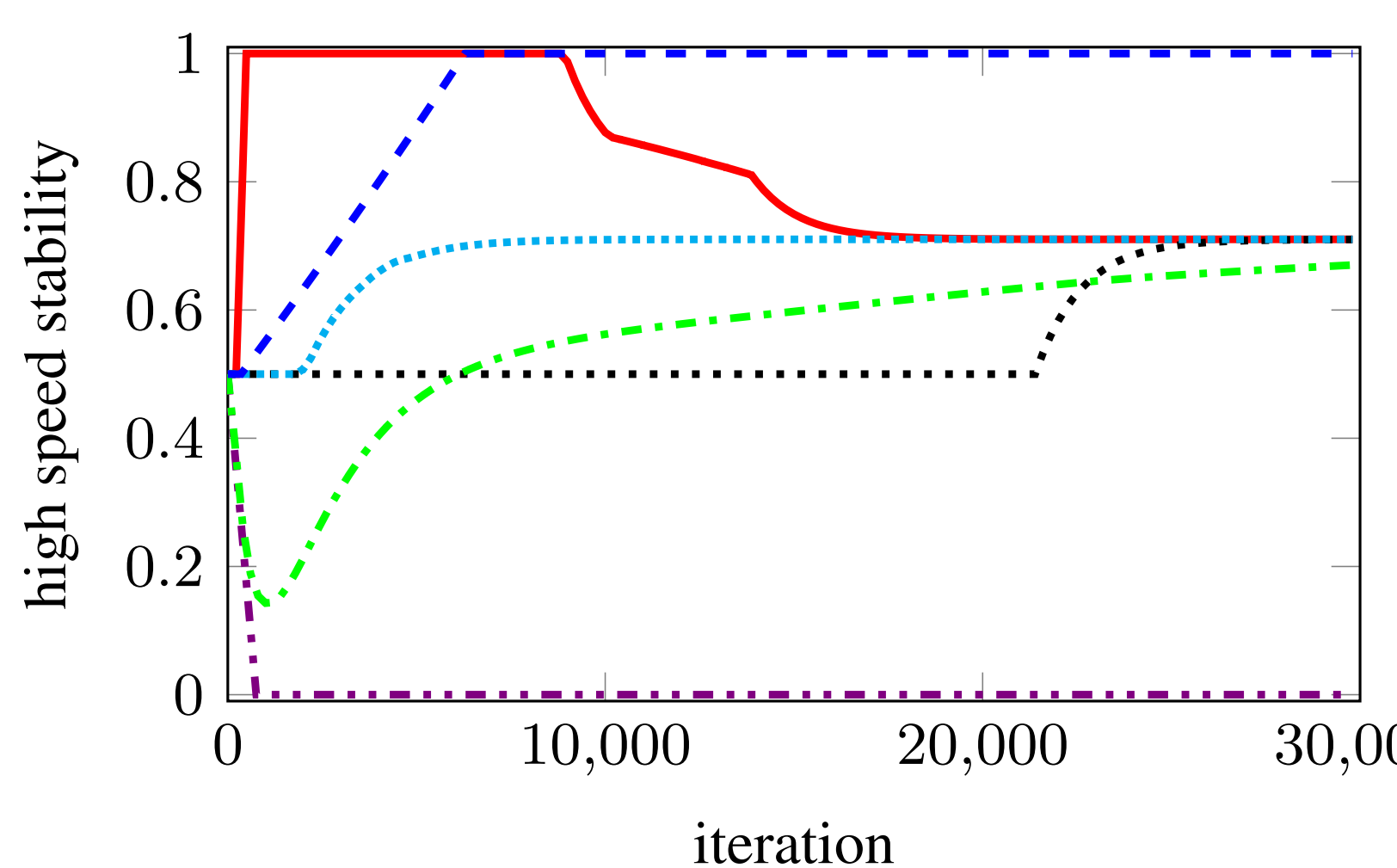
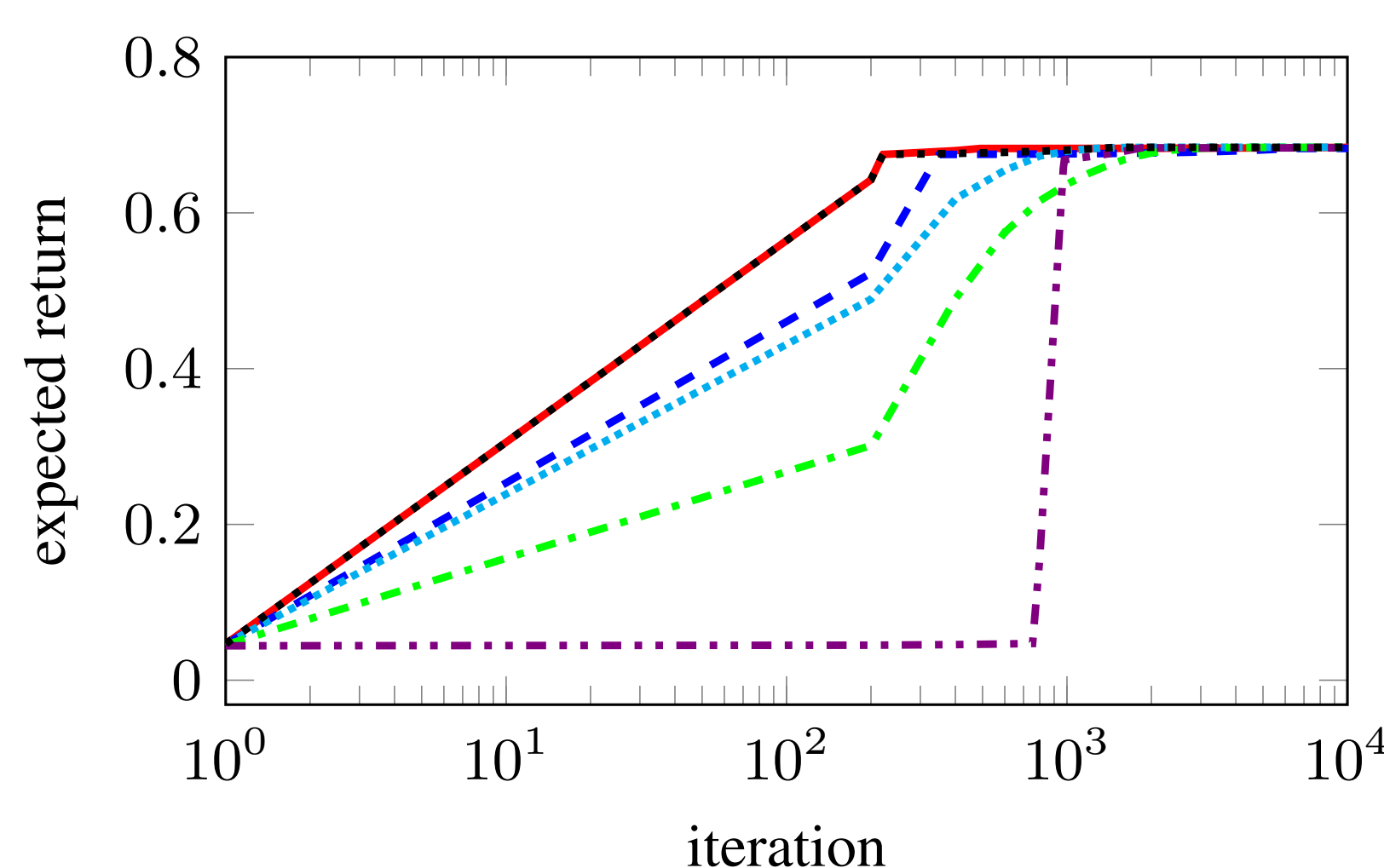
COMPARED ALGORITHMS



STUDENT-TEACHER



RACETRACK



..... SPI+SMI - - - SMI+SPI — SPMI - - - SPMI-sup - - - SPMI-alt - - - SPMI-greedy