



EXPLOITING STRUCTURE AND UNCERTAINTY OF BELLMAN UPDATES IN MARKOV DECISION PROCESSES



POLITECNICO
MILANO 1863

D. TATEO, C. D'ERAMO, A. NUARA, M. RESTELLI, A. BONARINI

{davide.tateo, carlo.deramo, alessandro.nuara, marcello.restelli, andrea.bonarini}@polimi.it

PROBLEM

- Learning is difficult in highly stochastic environments
- Uncertainty in action-value function estimates propagates
- Some algorithms face this problem focusing on the bias of the estimate

CONTRIBUTIONS

1. Split the estimate in two components:
 - The expected reward $\tilde{R}(x, u)$
 - The expected next state value function $\tilde{Q}(x, u)$
2. Use different learning rates for the two components
3. We provide empirical results showing the effectiveness of our approach

RQ-LEARNING ALGORITHM

IDEA

Split the action-value function in two components:

- $\tilde{R}(x, u) = \mathbb{E}_{x' \sim \mathcal{P}(x'|x, u)} [r(x, u, x')]$
- $\tilde{Q}(x, u) = \mathbb{E}_{x' \sim \mathcal{P}(x'|x, u)} \left[\max_{u'} Q^*(x', u') \right]$
- $Q^*(x, u) = \tilde{R}(x, u) + \gamma \tilde{Q}(x, u)$

EMPIRICAL RESULTS