

Simulated TD RPEs under the optimal policy

Stationary environments with varying reward probability p

