

Overestimation in Q-Learning

[Deep Reinforcement Learning with Double Q-learning](#)

Hado van Hasselt, Arthur Guez, David Silver. AAI 2016

[Non-delusional Q-learning and value-iteration](#)

Tyler Lu, Dale Schuurmans, Craig Boutilier. NeurIPS 2018

Yang Liu

19.03.2019

TD Method

“Learn a guess from a guess”

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[Y_t - Q(s_t, a_t)]$$

Y_t : TD target combining **Sample Reward** and **Current Estimate**

Q-Learning: TD Control

“Learn a guess from a guess”

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [Y_t - Q(s_t, a_t)]$$

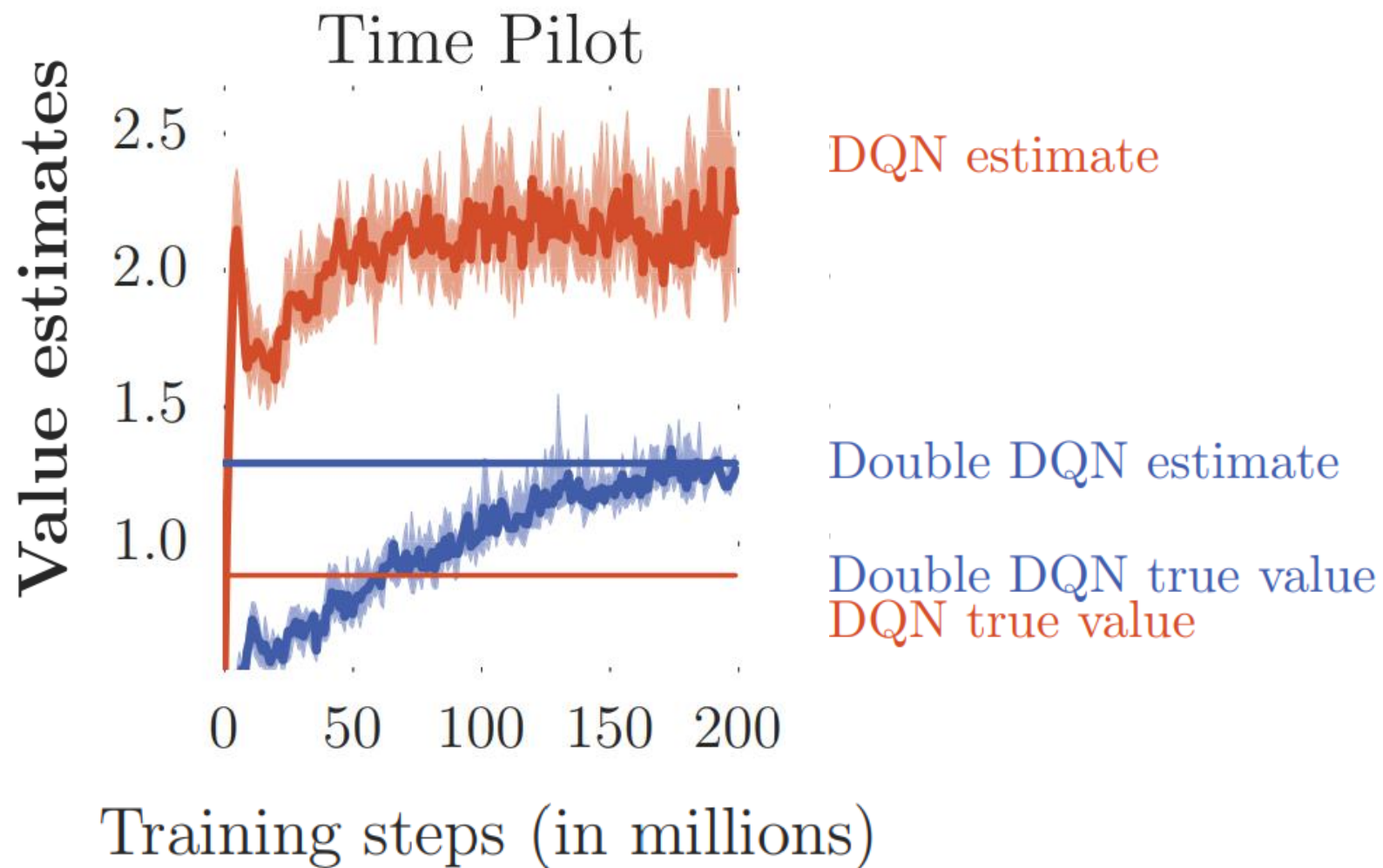
$$Y_t = r_{t+1} + \gamma \max_a Q(s_{t+1}, a)$$

Sample Reward Current Estimate

Converge to (Watkins 1989)

$$q_*(s, a) := \max_{\pi} q_{\pi}(s, a) = \mathbb{E}[R_t + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

Overestimation in Q-Learning



What's Wrong?

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[Y_t - Q(s_t, a_t)]$$

$$Y_t = R_{t+1} + \gamma \max_{a \in \mathcal{A}} Q(s_{t+1}, a)$$

$\max_{a \in \mathcal{A}}$ \Rightarrow two biases \Rightarrow Pathological Behavior

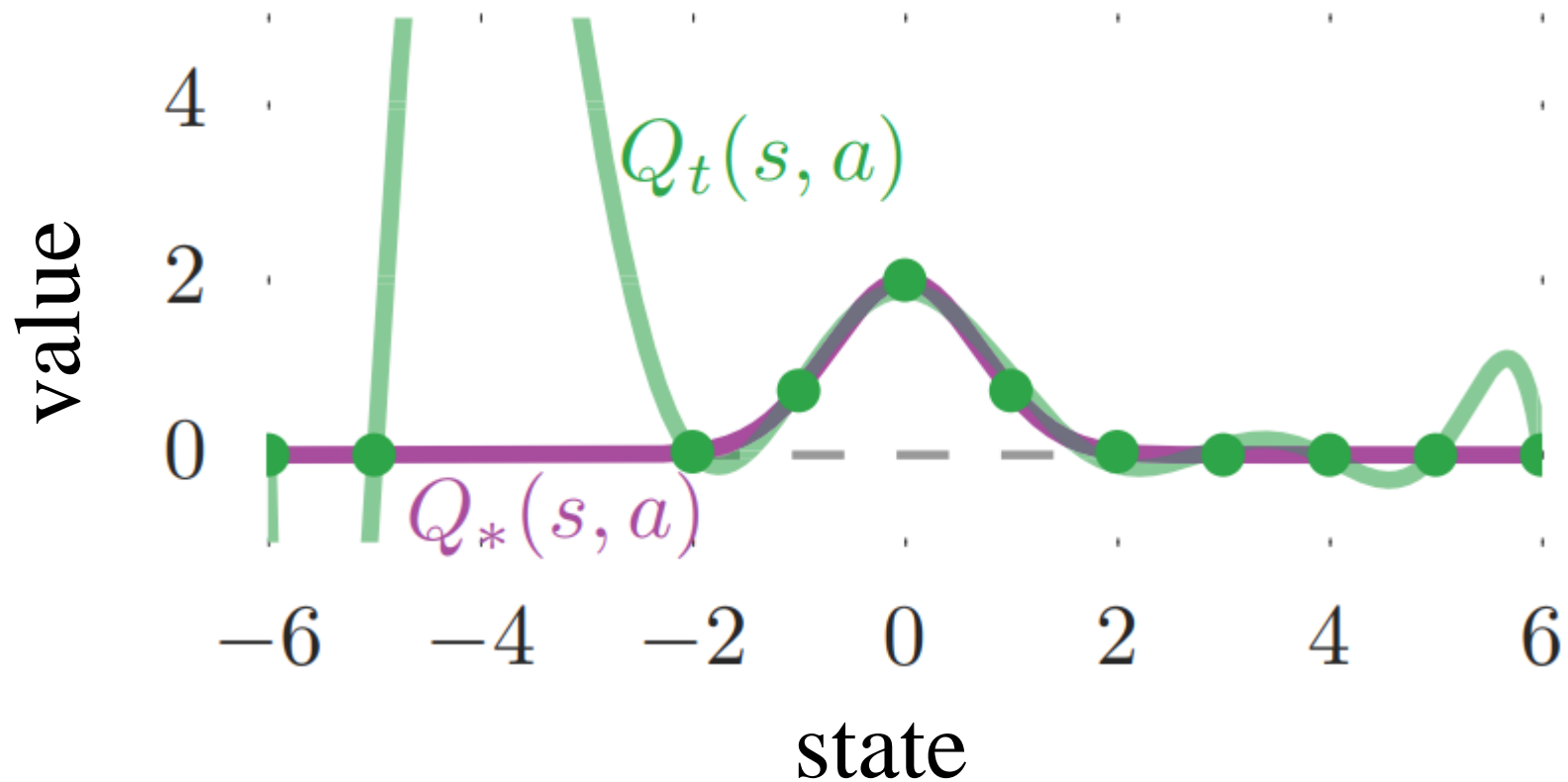
Maximization Bias

Inferior Policy

Delusional Bias

Divergence

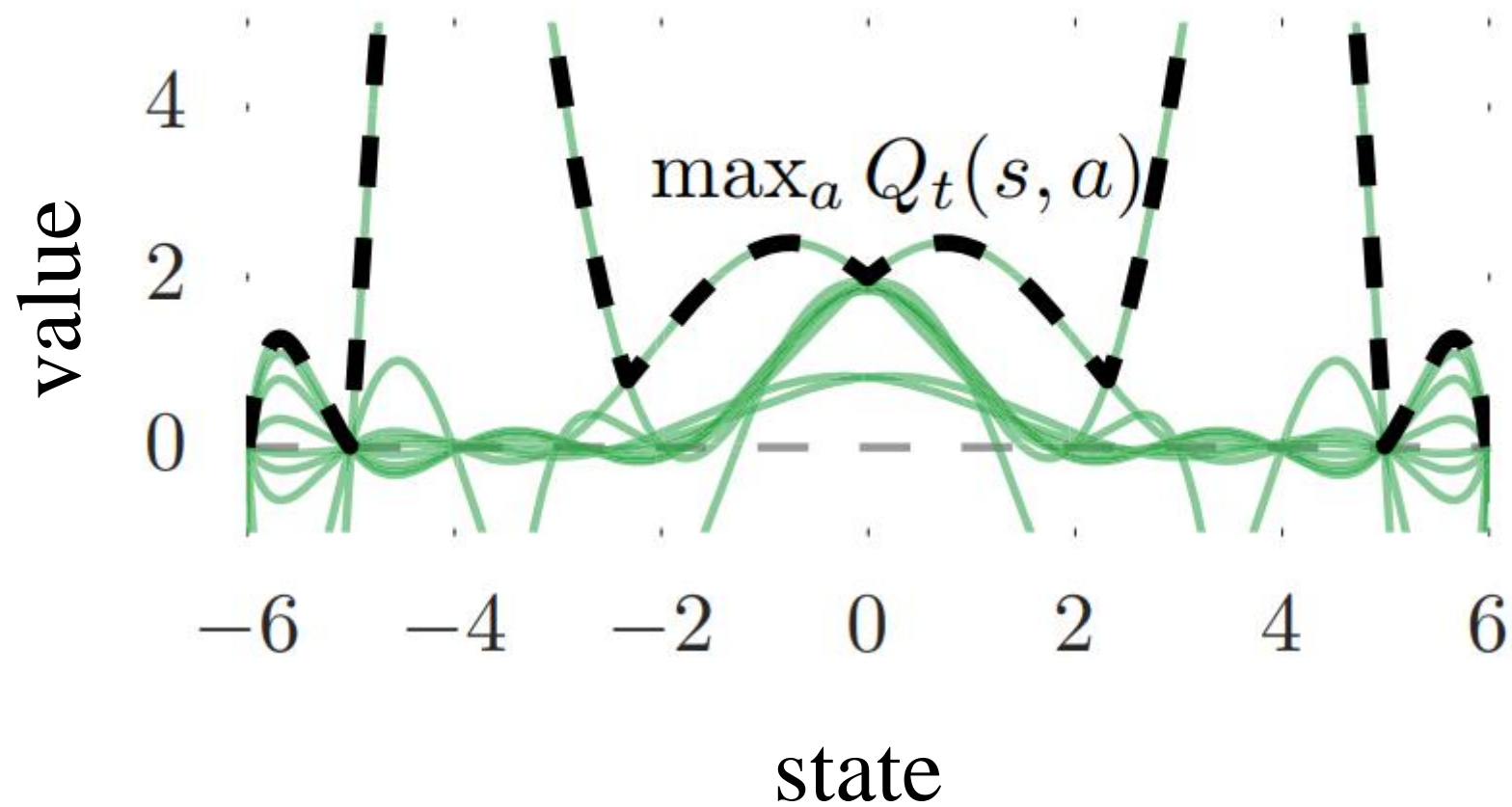
Maximization Bias



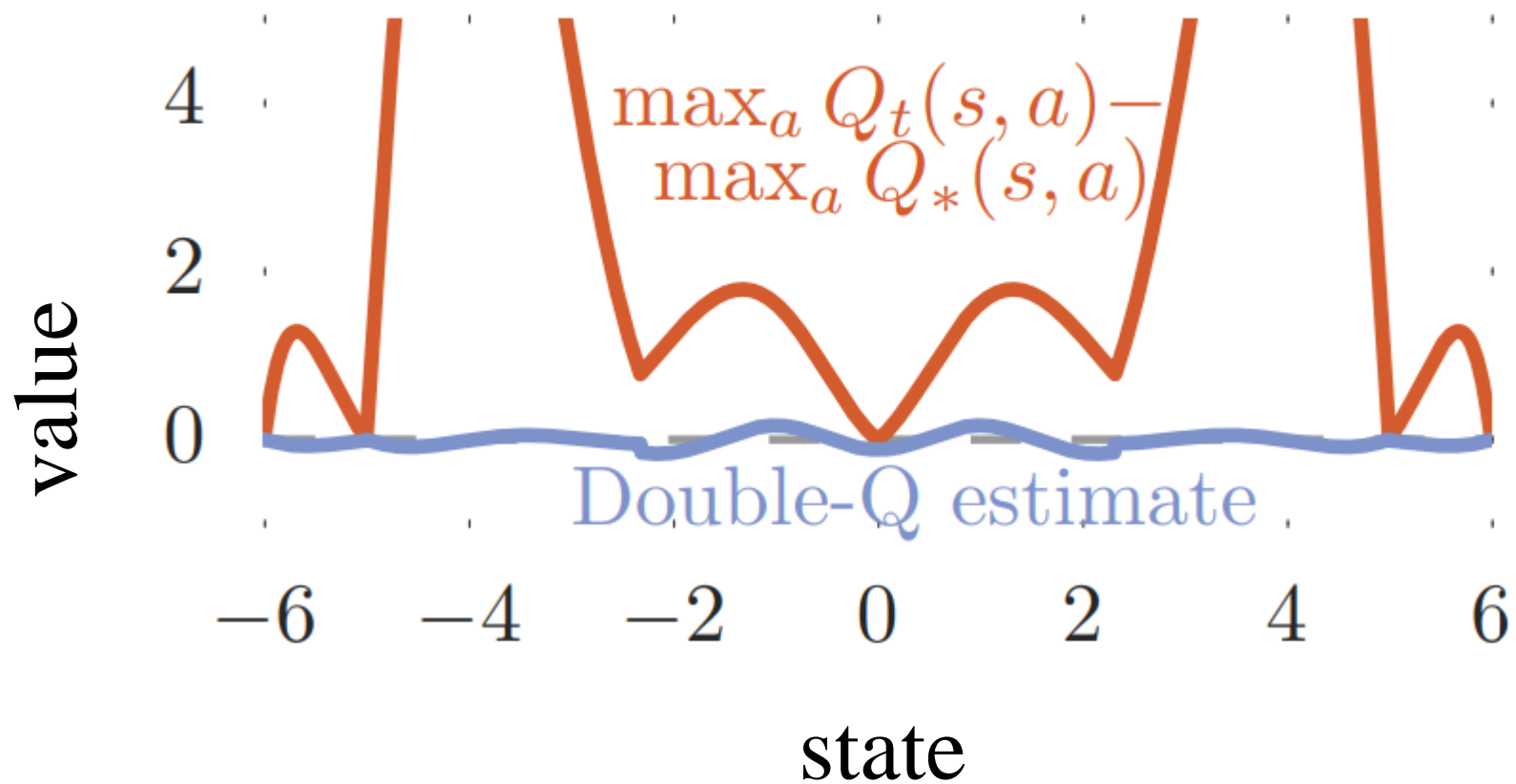
$$Q_*(s, a) = V_*(s) = 2e^{-s^2}, \forall a$$

$Q_t(s, a)$ 9-degree polynomial

Maximization Bias



Maximization Bias



Maximization Bias

By convexity of max function, for R.V. X_1, \dots, X_n and their sample mean μ_1, \dots, μ_n and their realization $\widehat{\mu}_1, \dots, \widehat{\mu}_n$

$$\max_i \mathbb{E}[X_i] = \max_i \mathbb{E}[\mu_i] \leq \mathbb{E}[\max_i \mu_i] \approx \max_i \widehat{\mu}_i$$

$$\begin{aligned} q_*(s, a) &= \mathbb{E}[R_{t+1} | S_t = s, A_t = a] + \gamma \mathbb{E}[\max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a] \\ &\approx r_{t+1} + \gamma \mathbb{E} \left[\max_{a'} q_*(s_{t+1}, a') \right] \\ &\leq r_{t+1} + \gamma \mathbb{E} \left[\max_{a'} Q(s_{t+1}, a') \right] \\ &= r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') = Y_t \end{aligned}$$

Double Q-Learning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[Y_t - Q(s_t, a_t)]$$

Single Q-Learning

$$Y_t = r_{t+1} + \gamma \max_a Q(s_{t+1}, a)$$

Double Q-Learning

$$Y_t = r_{t+1} + \gamma Q_2(s_{t+1}, \operatorname{argmax}_a Q_1)$$

Theorem: Double Estimator is bounded by $\max_i \mathbb{E}[X_i]$

Double DQN

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha [Y_t - Q(s_t, a_t; \boldsymbol{\theta}_t)] \nabla_{\boldsymbol{\theta}} Q(s_t, a_t; \boldsymbol{\theta}_t)$$

DQN

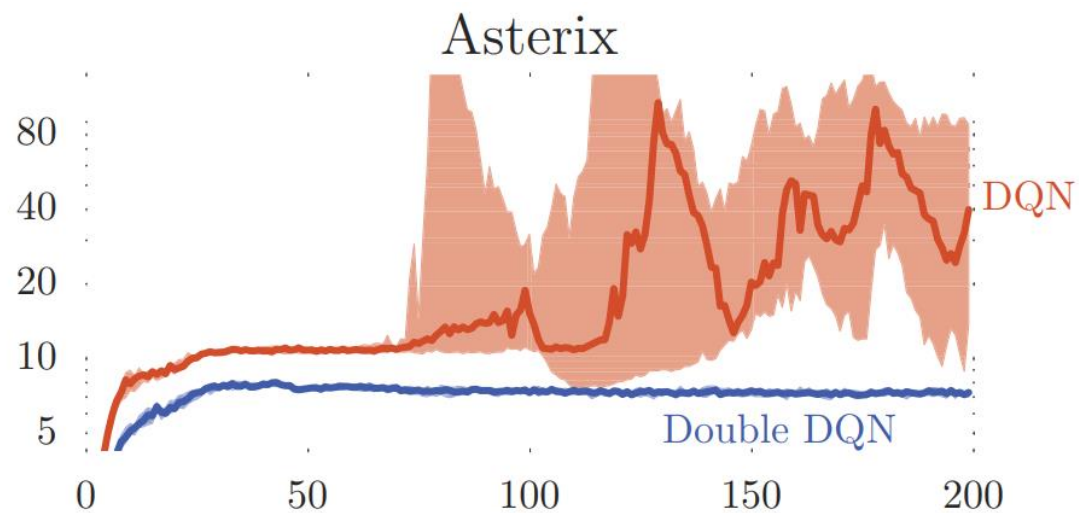
$$Y_t = r_{t+1} + \gamma \max_a Q(s_{t+1}, a; \boldsymbol{\theta}_t^-)$$

Double DQN

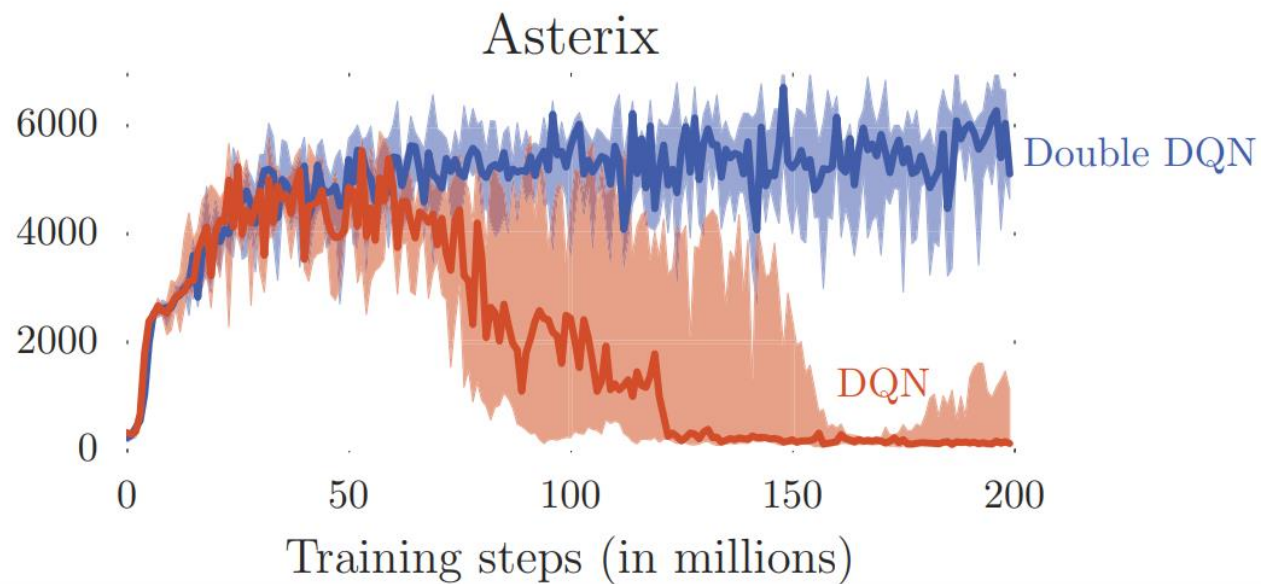
$$Y_t = r_{t+1} + \gamma Q(s_{t+1}, \operatorname{argmax}_a Q(s_{t+1}, a; \boldsymbol{\theta}_t); \boldsymbol{\theta}_t^-)$$

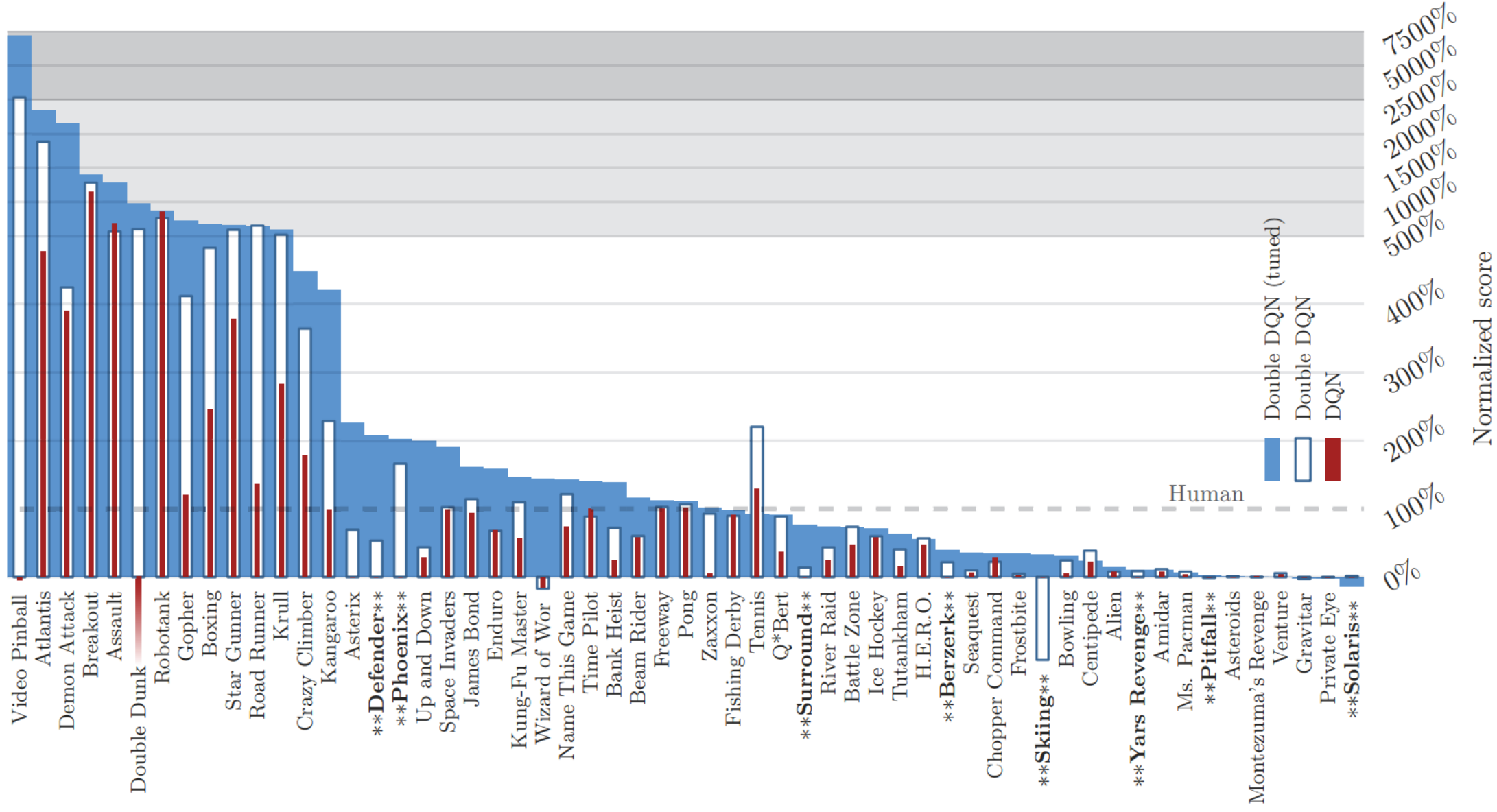
Result

Value Estimate



Score





Comments

Pros

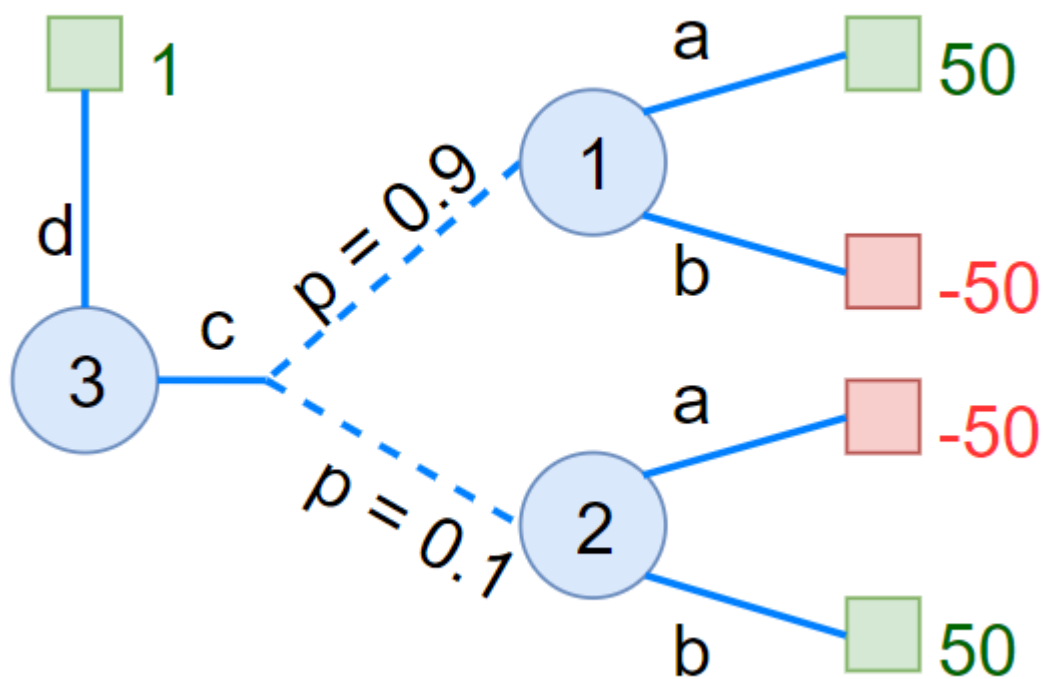
- Cheap modification
- Improvement on **stability** and **score**
- Widely adopted by following Deepmind's work

Cons

- Still biased estimate(proven **underestimate**)
- Experience Replay Buffer, the two estimates are **not independent**

Delusional Bias

A Simple MDP

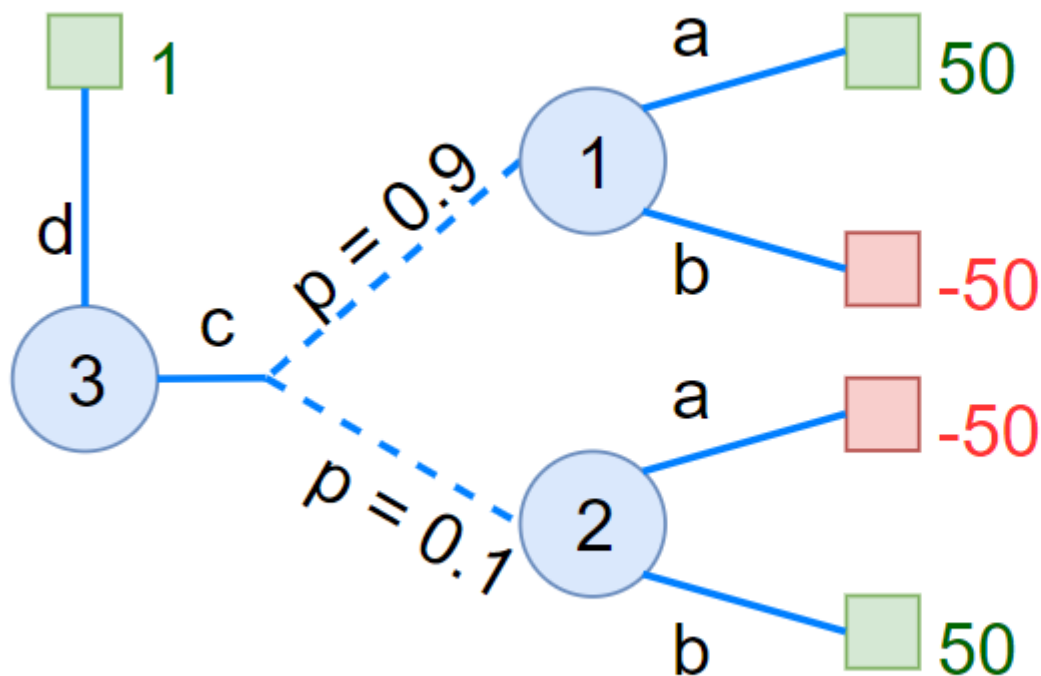


Optimal Policy: 50

3-c 1-a 2-b

Delusional Bias

Consider Linear Approximation $Q(s, a) = \theta\phi(s, a), \theta_0 > 0$



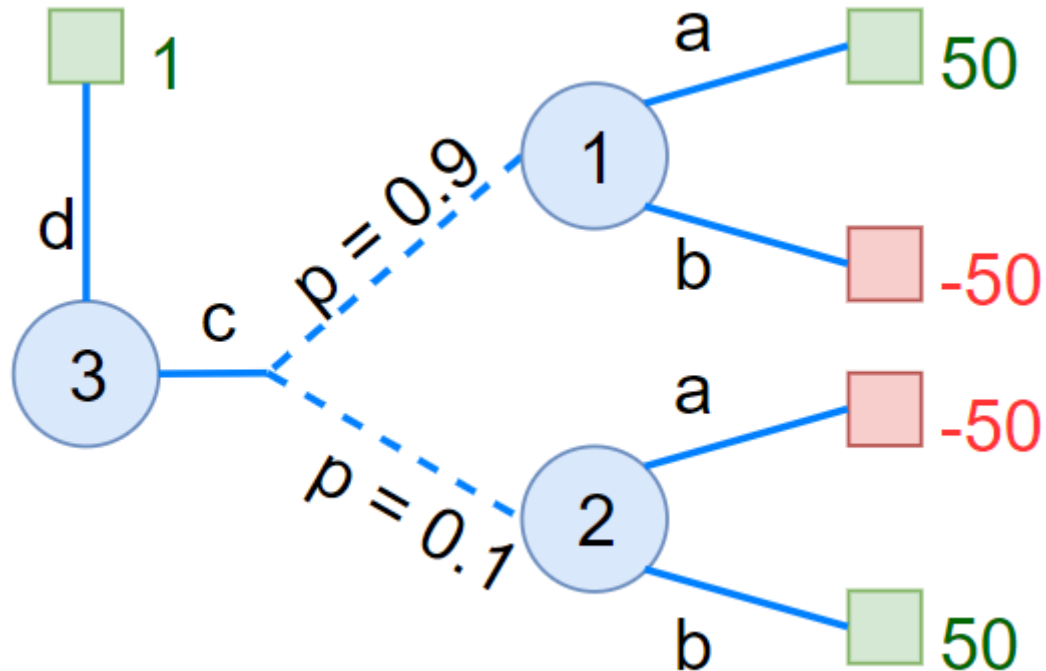
$\phi(1, a) = 1$	$\phi(1, a) = -1$
$\phi(2, a) = 1$	$\phi(2, b) = -1$

Optimal Policy: 40

3-c 1-a 2-a

Consider Linear Approximation $Q(s, a) = \theta \phi(s, a)$

$\phi(1, a) = 1$	$\phi(1, b) = -1$
$\phi(2, a) = 1$	$\phi(2, b) = -1$



Inconsistent Back up

Back up at (1, a): $\theta \uparrow$

Back up at (2, b): $\theta \downarrow \downarrow$

$\Rightarrow \theta_* \approx 0 \Rightarrow Q(3, c) \approx 0$

$\Rightarrow \pi_*(3) = d$ with value **1**

Delusional Bias

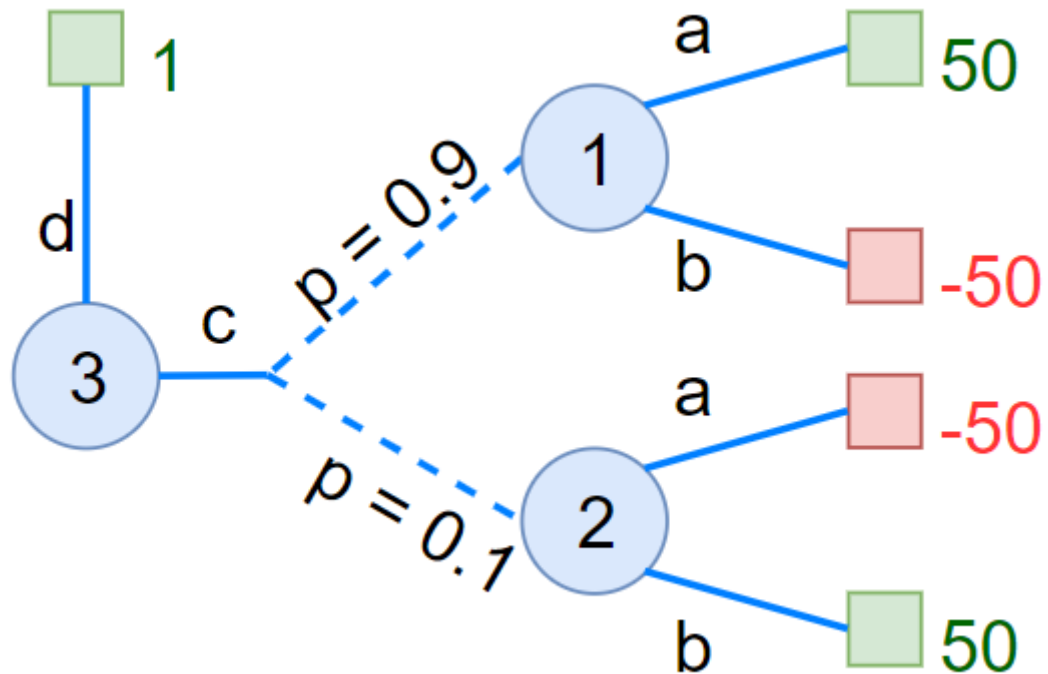
Delusional Bias occurs whenever a backed-up value estimate is derived from action choices that are **not realizable** in the underlying *class(Greedy Policy Class)*.

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t + \alpha [Y_t - Q(s_t, a_t; \boldsymbol{\theta}_t)] \nabla_{\boldsymbol{\theta}} Q(s_t, a_t; \boldsymbol{\theta}_t) \\ Y_t &= r_{t+1} + \gamma \max_{a \in \mathcal{A}} Q(s_{t+1}, a; \boldsymbol{\theta}_t)\end{aligned}$$

Unconstrained maximization \Rightarrow Overestimation in target Y_t

Delusional Bias

Delusional Bias occurs whenever a backed-up value estimate is derived from action choices that are **not realizable** in the underlying *class*.



Inconsistent Backup at (1,a) and (2,b)
⇒ Delusional Bias ⇒ Poor policy

Delusional Bias

Explains following pathological behavior:

- Poor policy
- Divergence
- Cyclic Behavior
- Discounting Paradox

γ -score of model trained with γ' higher than model trained with γ

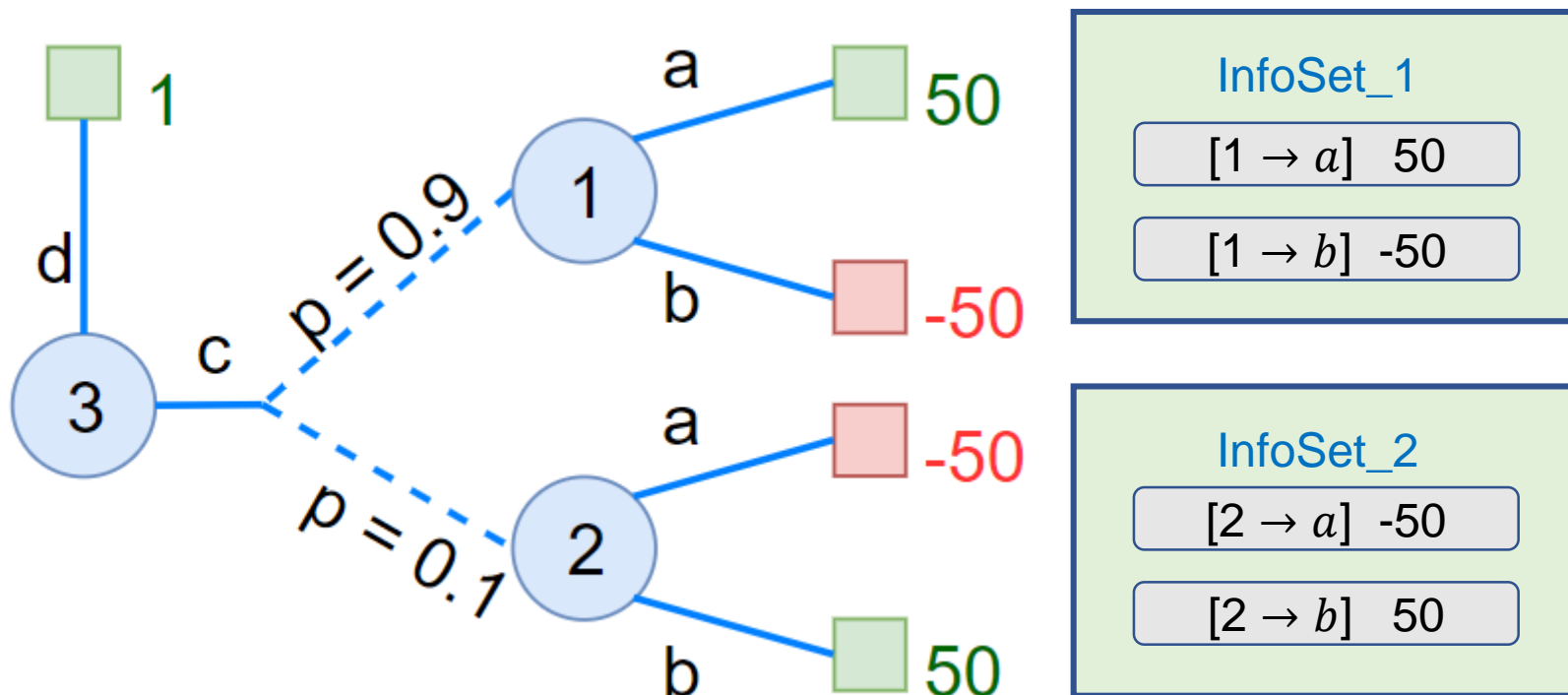
Also Arises in Value Iteration

$$V(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s'} p(s'|s, a) [r_{s,a} + \gamma V(s')]$$

Full state Bellman Backup

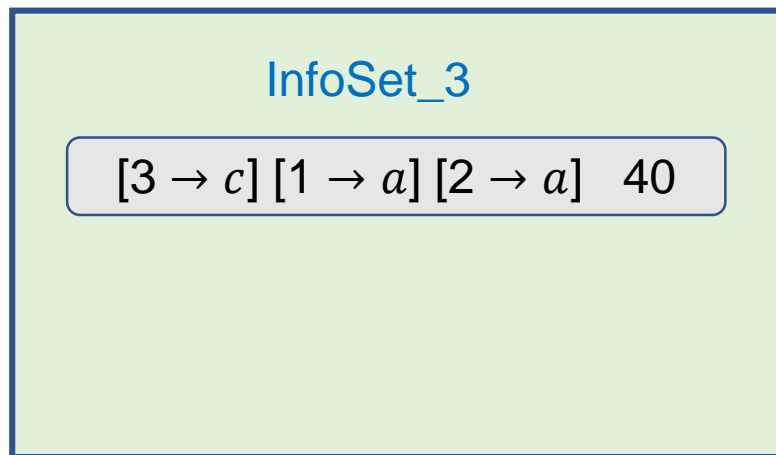
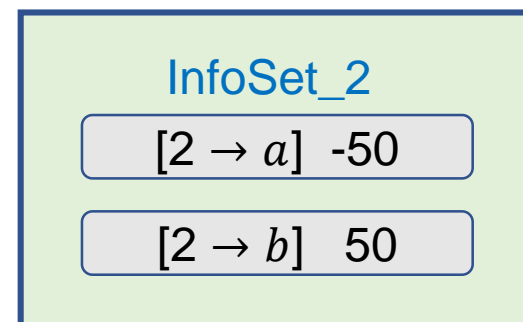
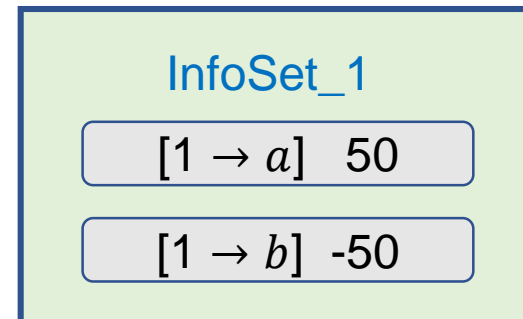
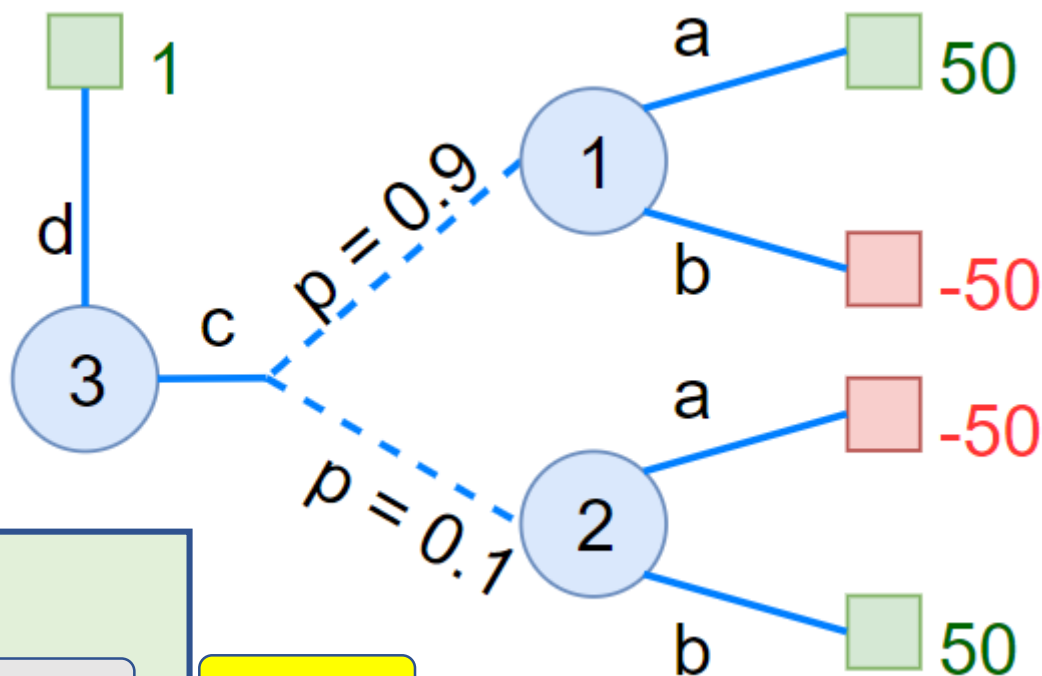
Policy-Class Value Iteration

Idea: just track every feasible paths using *information set*



Policy-Class Value Iteration

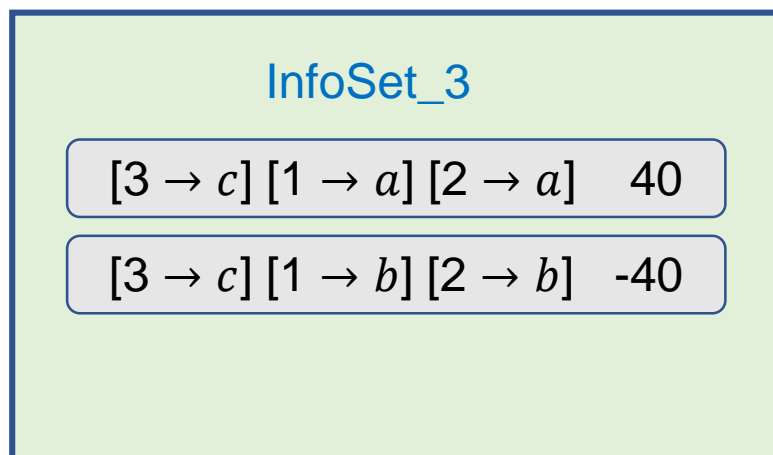
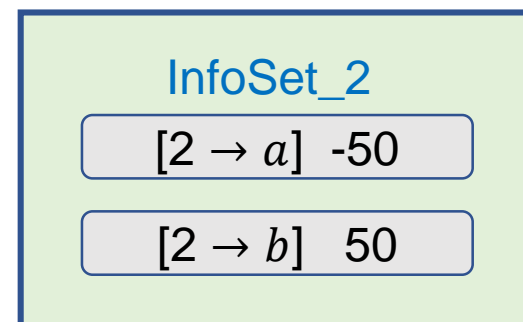
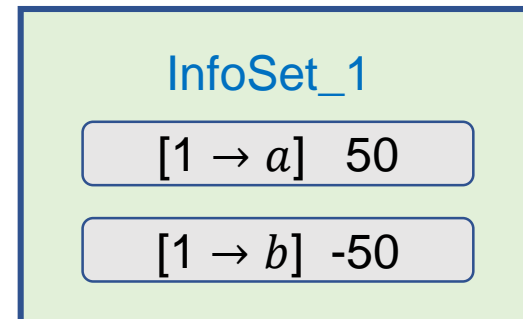
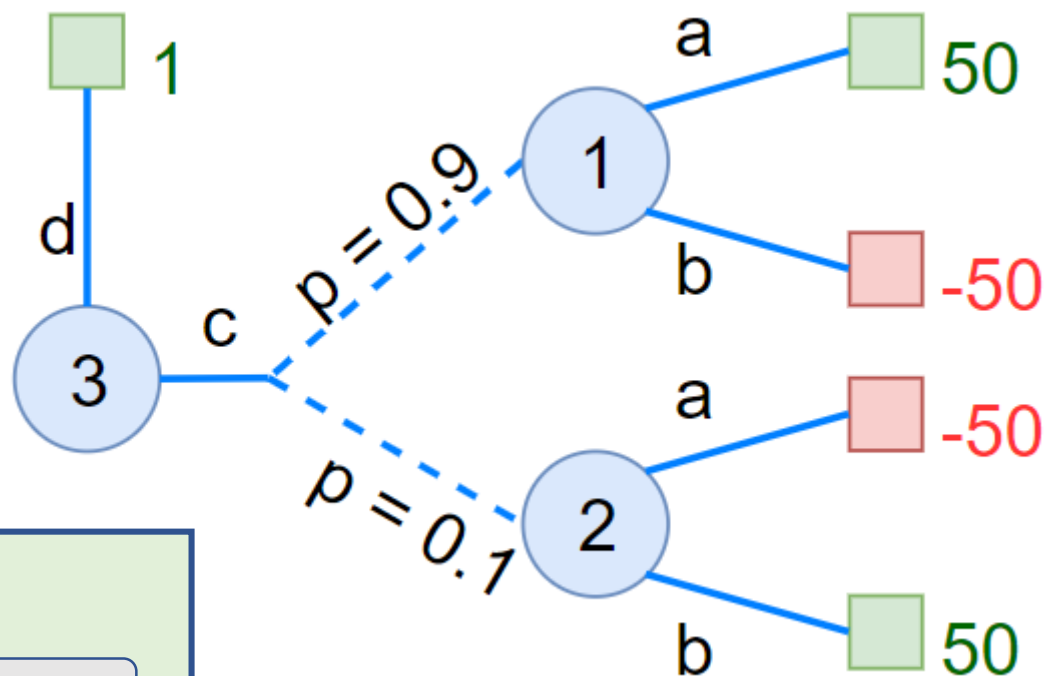
Idea: just track every feasible paths using *information set*



feasible

Policy-Class Value Iteration

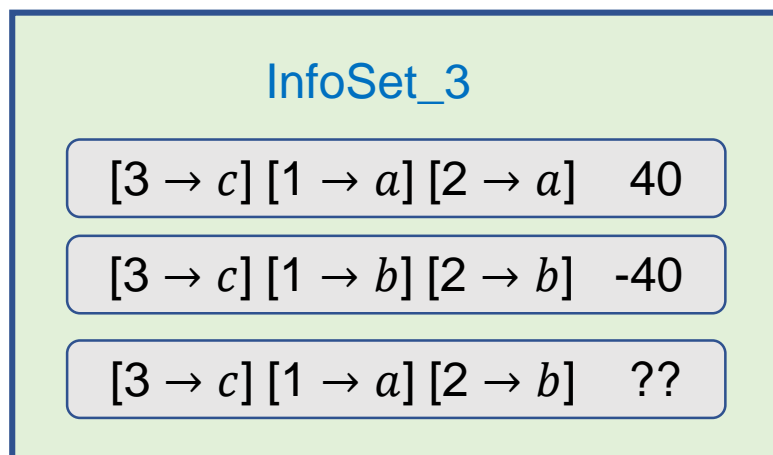
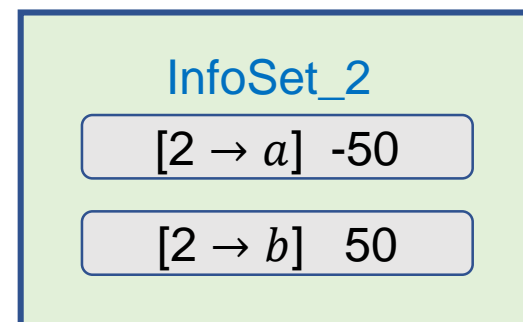
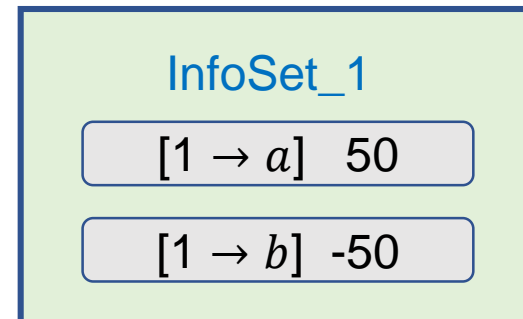
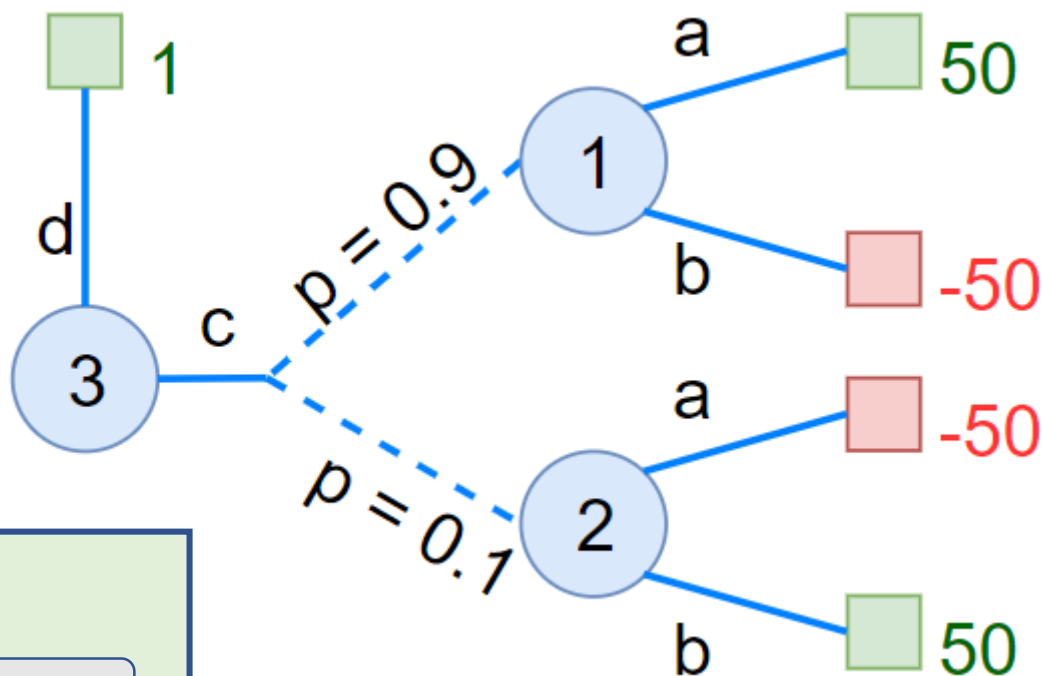
Idea: just track every feasible paths using *information set*



feasible

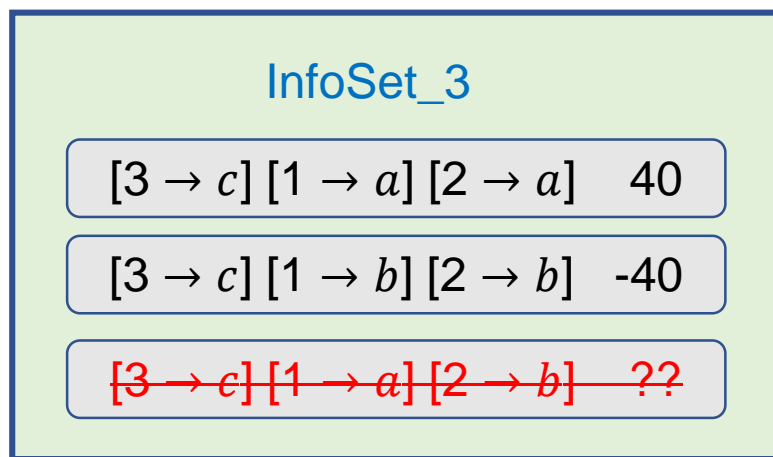
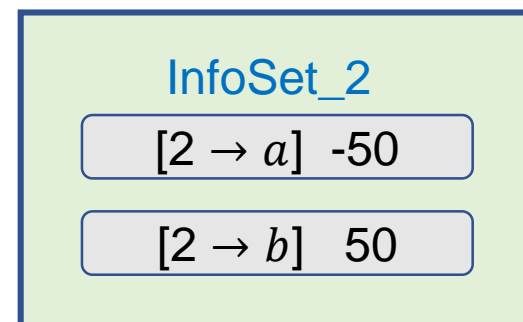
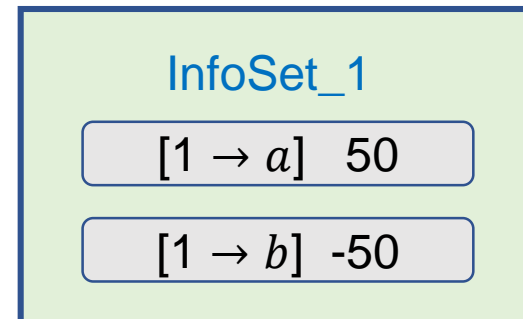
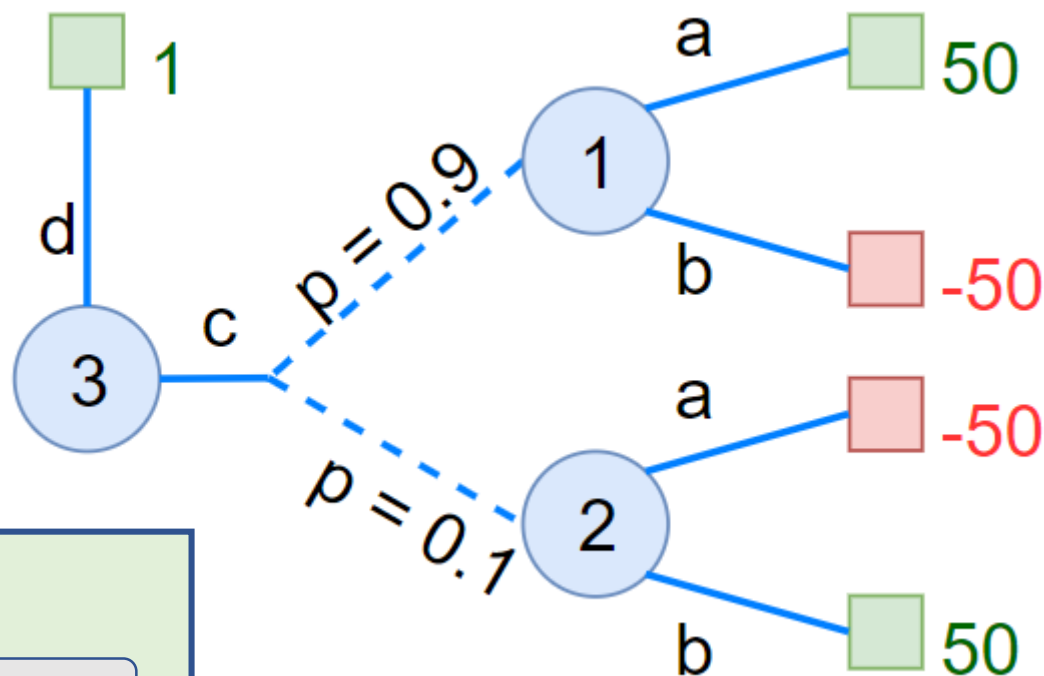
Policy-Class Value Iteration

Idea: just track every feasible paths using *information set*



Policy-Class Value Iteration

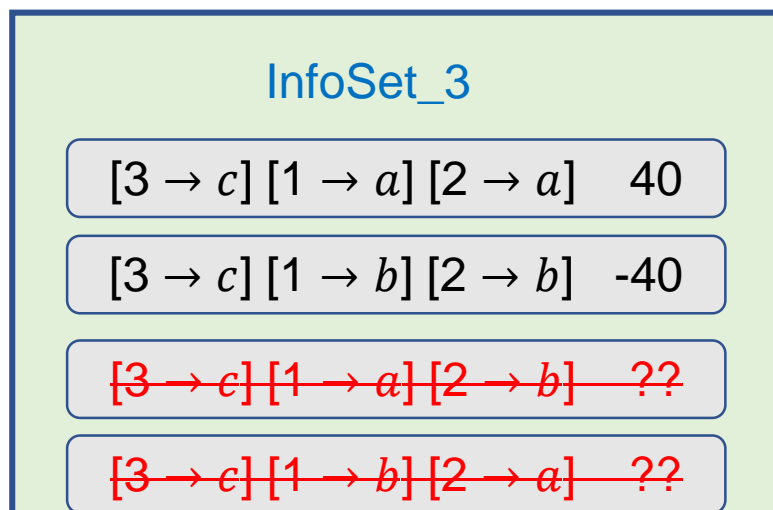
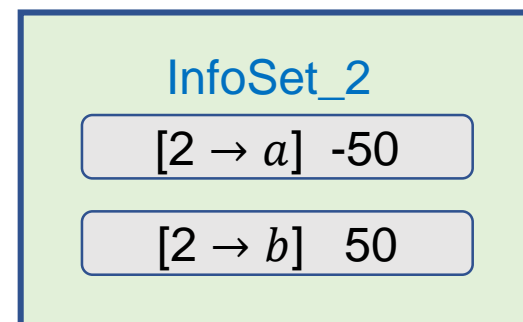
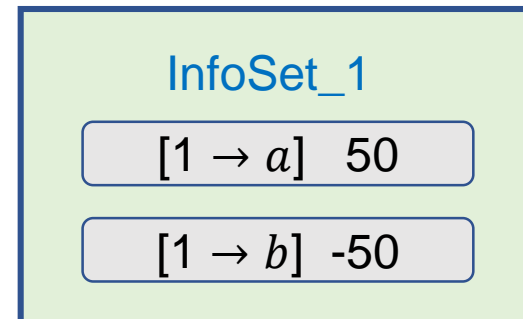
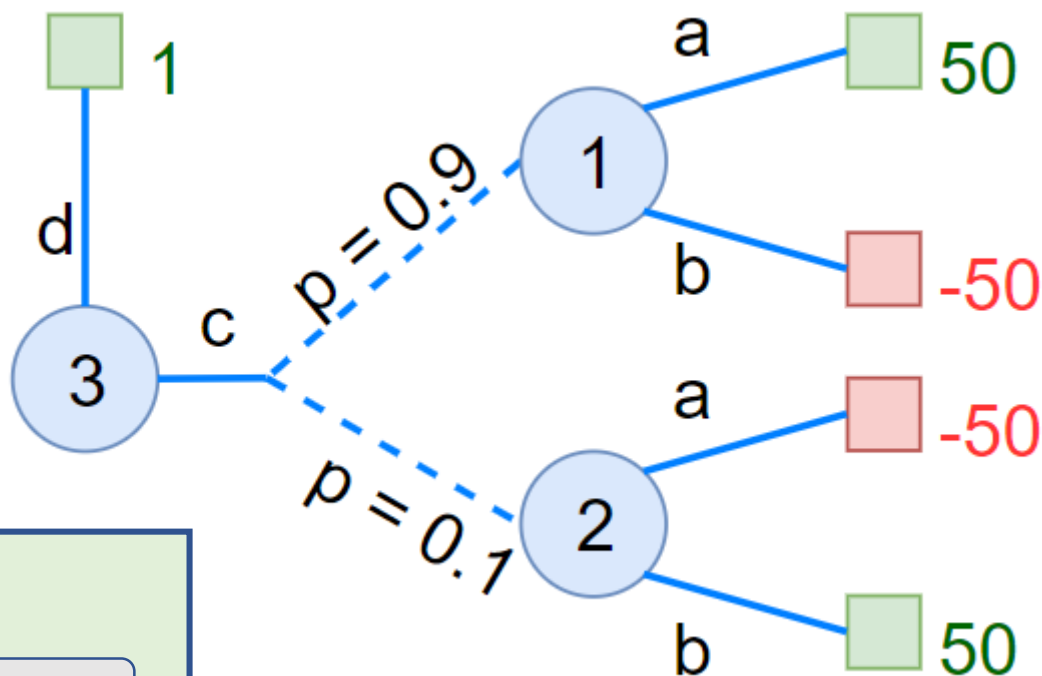
Idea: just track every feasible paths using *information set*



infeasible

Policy-Class Value Iteration

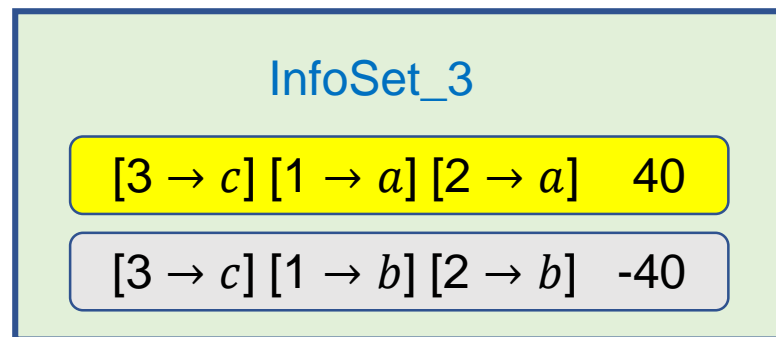
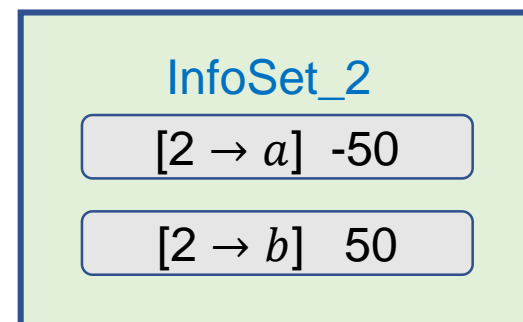
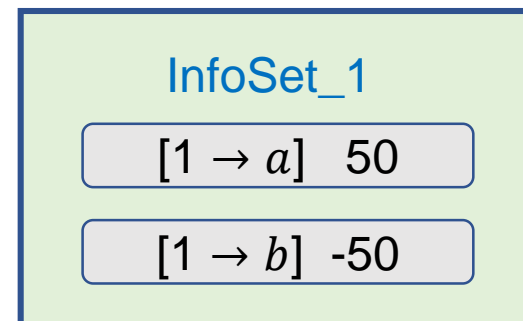
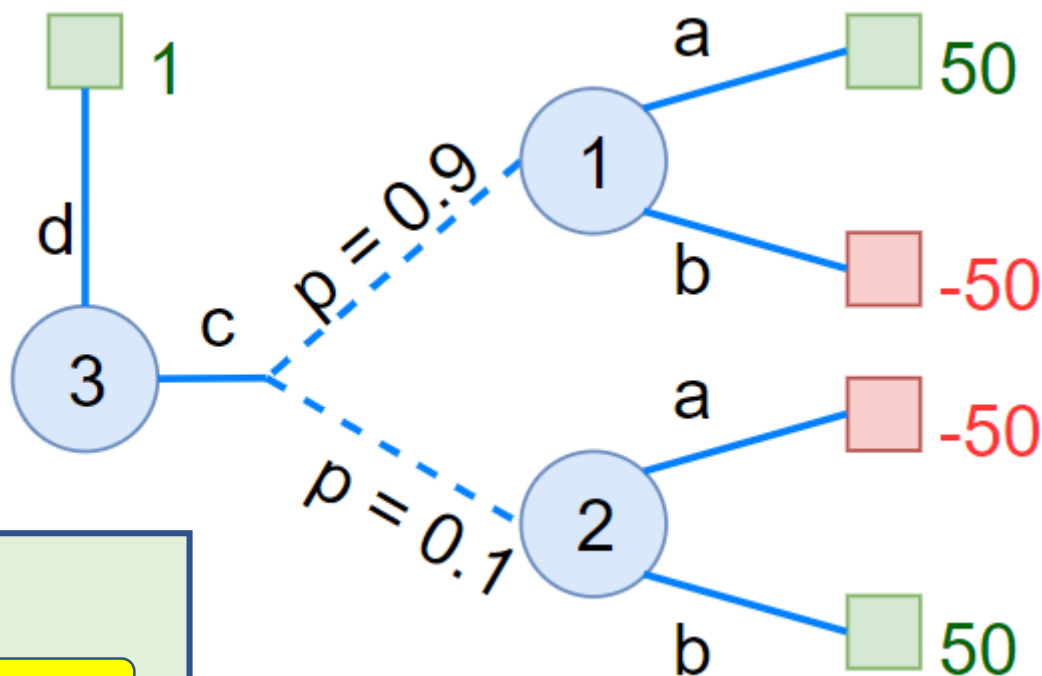
Idea: just track every feasible paths using *information set*



infeasible

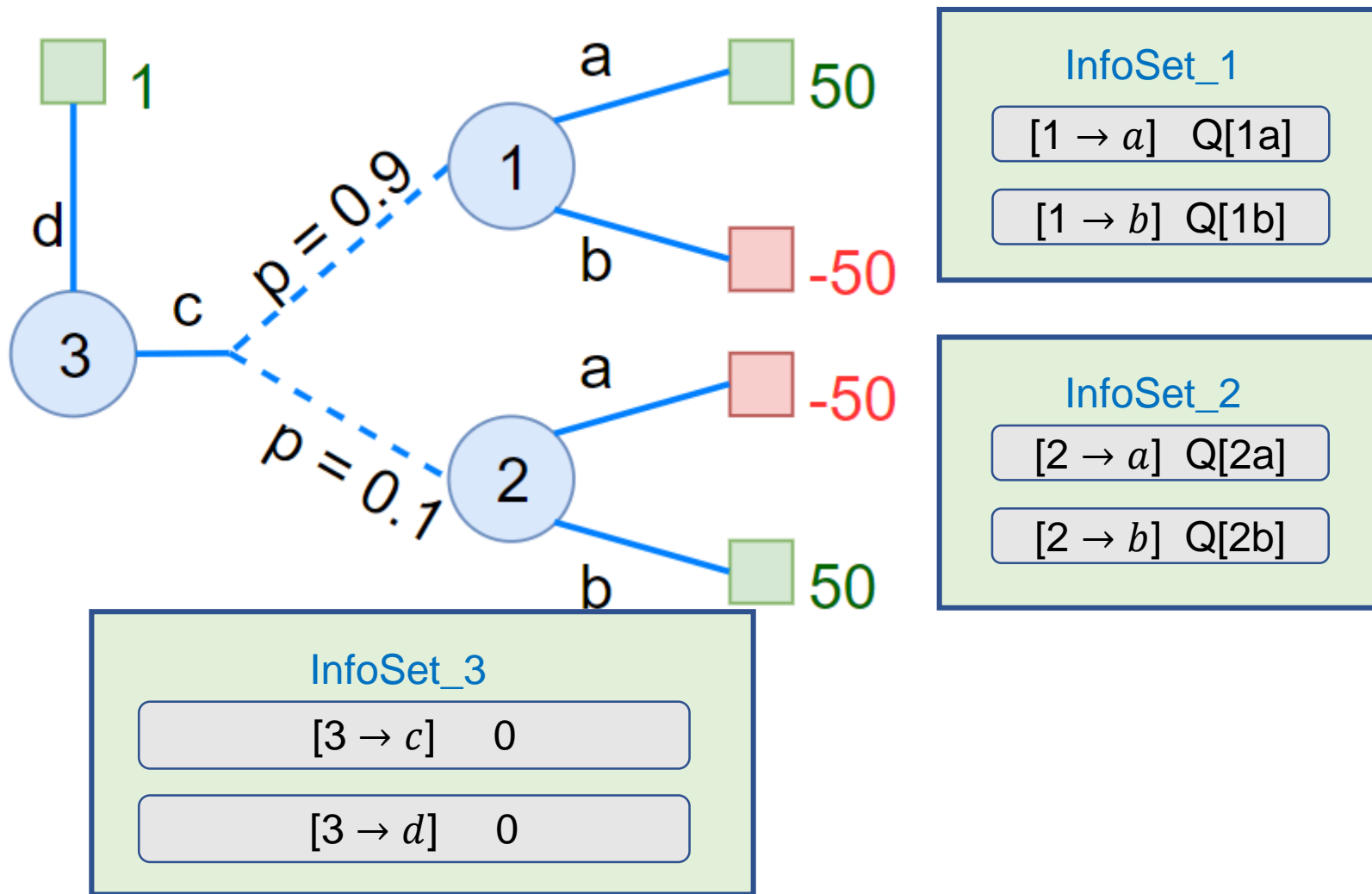
Policy-Class Value Iteration

Idea: just track every feasible paths using *information set*



Policy-Class Q-Learning

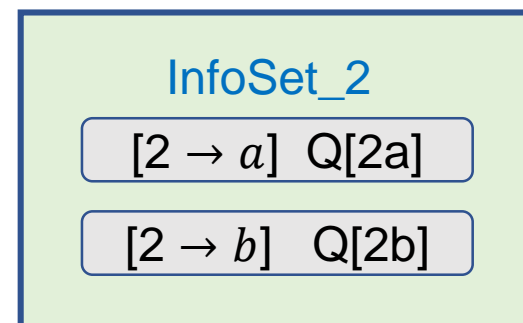
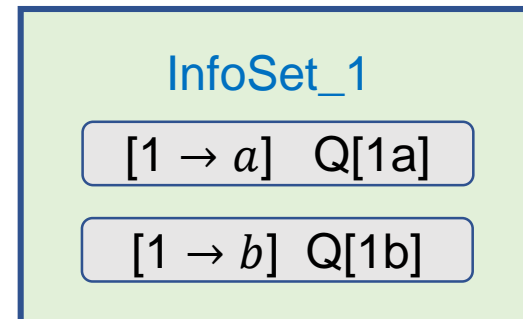
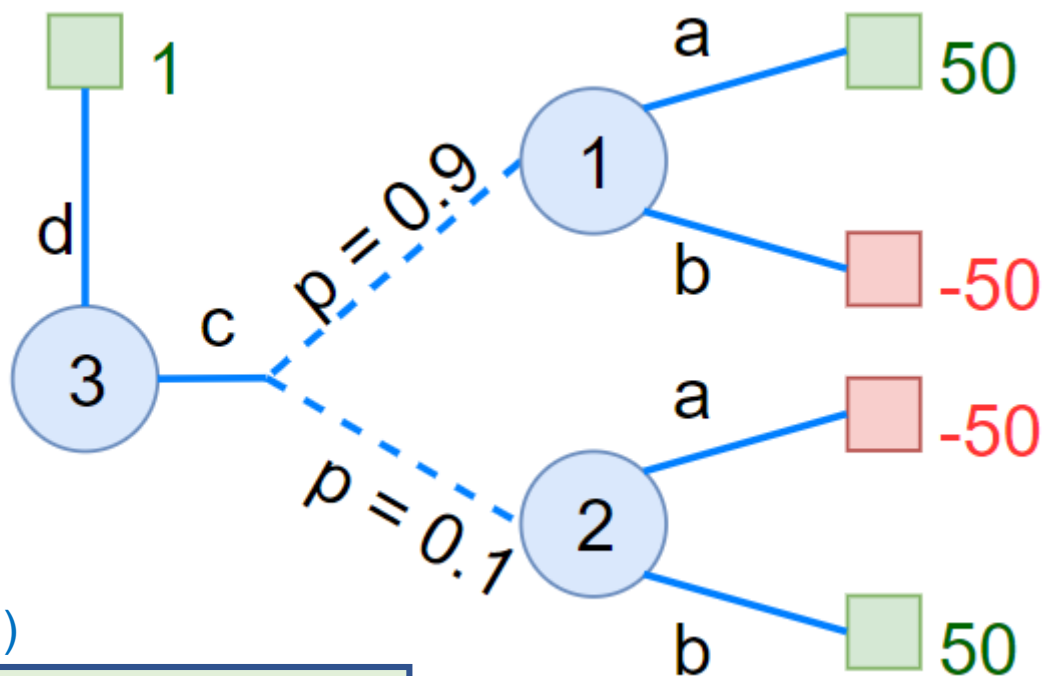
Sample Backup instead of full state Backup



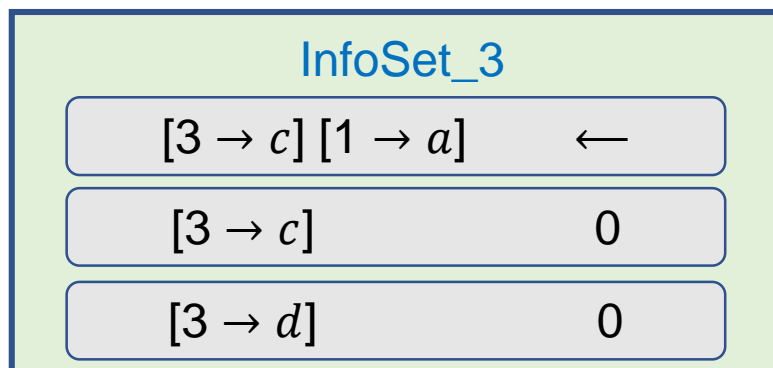
Initialization

Policy-Class Q-Learning

Sample Backup instead of full state Backup



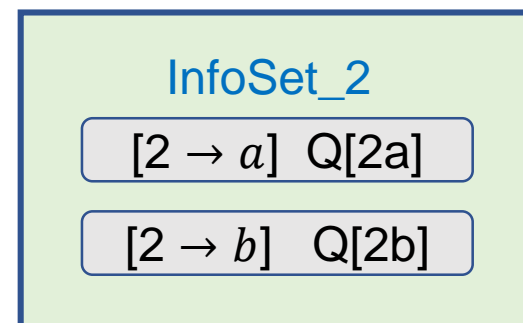
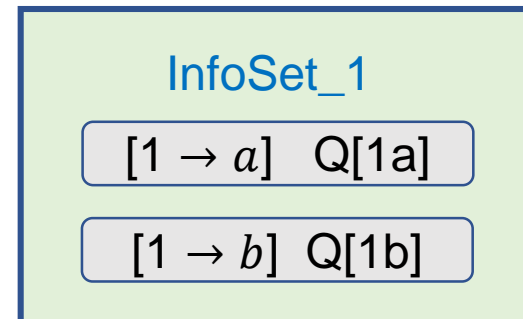
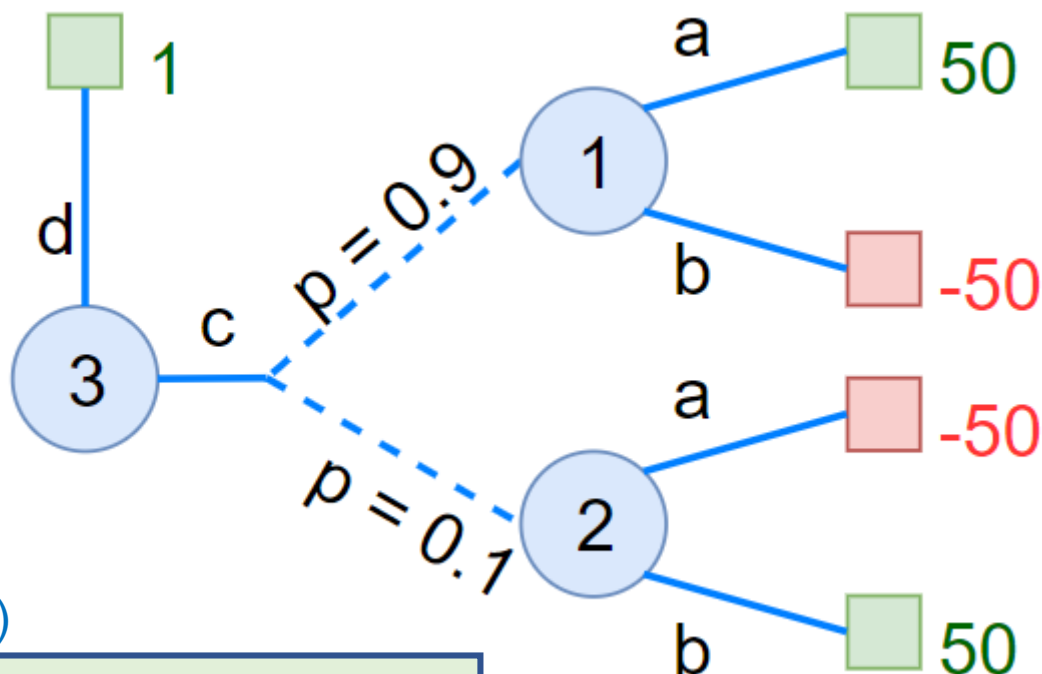
Sample experience (3, c, 0, 1)



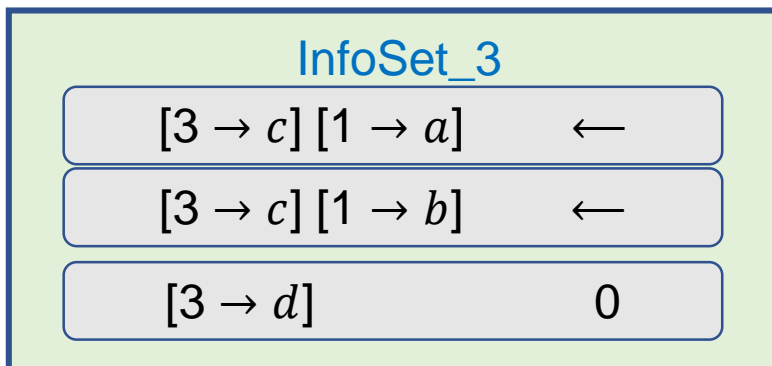
feasible $Q[3c][1a] \leftarrow Q[3c] + \alpha_t^{sa}(r + \gamma Q[1a] - Q[3c])$

Policy-Class Q-Learning

Sample Backup instead of full state Backup



Sample experience (3, c, 0, 1)



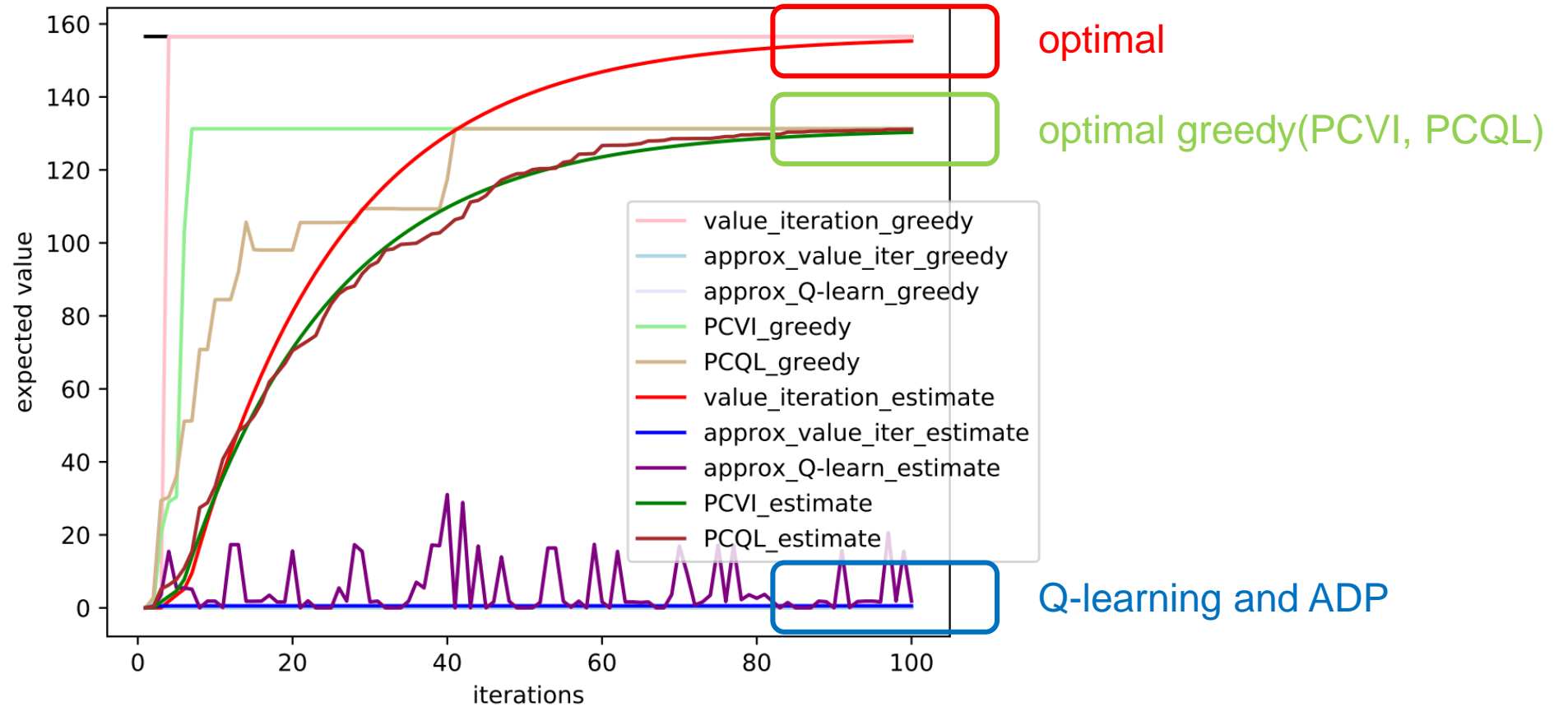
feasible

$$Q[3c][1b] \leftarrow Q[3c] + \alpha_t^{sa} (r + \gamma Q[1b] - Q[3c])$$

Theorems

- PCVI and PCQL converge
- Converge to optimal policy in given class/No delusion
- Bounded runtime

Result



4x4 Grid world using 4 features

Comment:

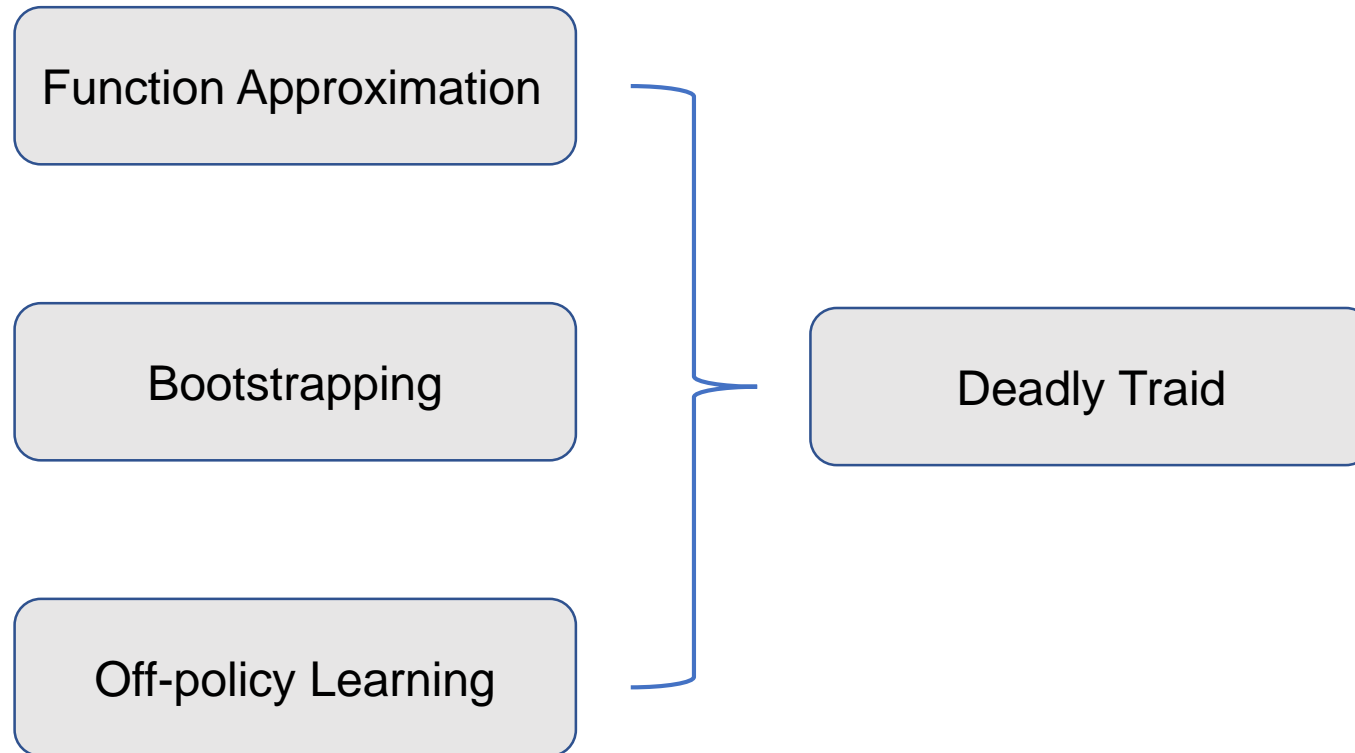
Pros:

- Identified Delusional Bias and its consequence
- Come up with **proven** algorithm PCVI and PCQL
- Heuristic methods for large models

Cons:

- No **scalable** solution
- Provide no result on DNN

More on Deadly Traid



Reference

[The Value Function Polytope in Reinforcement Learning](#) Robert Dadashi, Adrien Ali Taïga, Nicolas Le Roux, Dale Schuurmans, Marc G. Bellemare. Arxiv 2019

[Deep Reinforcement Learning and the Deadly Triad](#) Hado Van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, Joseph Modayil. Arxiv 2018

[Non-delusional Q-learning and value-iteration](#) Tyler Lu, Dale Schuurmans, Craig Boutilier. NeurIPS 2018

[Deep Reinforcement Learning with Double Q-learning](#) Hado van Hasselt, Arthur Guez, David Silver. AAI 2016

[Reinforcement Learning: An Introduction](#) Richard Sutton and Andrew Barto. 2016

[Double Q-learning](#) Hado van Hasselt. NeurIPS 2010

[The Optimizer's Curse: Skepticism and Postdecision Surprise in Decision Analysis](#) James Smith and Robert Winkler. Management Science 2006

[Q-learning](#) Christopher Watkins and Peter Dayan. Machine Learning 1992