

Reinforcement Learning

Emma Brunskill
Stanford University
Winter 2018
Midterm Review

Reinforcement Learning Involves

- Optimization
- Delayed consequences
- Generalization
- Exploration

Learning Objectives

- Define the key features of reinforcement learning that distinguishes it from AI and non-interactive machine learning (as assessed by exams).
- Given an application problem (e.g. from computer vision, robotics, etc), decide if it should be formulated as a RL problem; if yes be able to define it formally (in terms of the state space, action space, dynamics and reward model), state what algorithm (from class) is best suited for addressing it and justify your answer (as assessed by the project and exams).
- Implement in code common RL algorithms such as a deep RL algorithm, including imitation learning (as assessed by the homeworks).
- Describe (list and define) multiple criteria for analyzing RL algorithms and evaluate algorithms on these metrics: e.g. regret, sample complexity, computational complexity, empirical performance, convergence, etc (as assessed by homeworks and exams).
- Describe the exploration vs exploitation challenge and compare and contrast at least two approaches for addressing this challenge (in terms of performance, scalability, complexity of implementation, and theoretical guarantees) (as assessed by an assignment and exams).

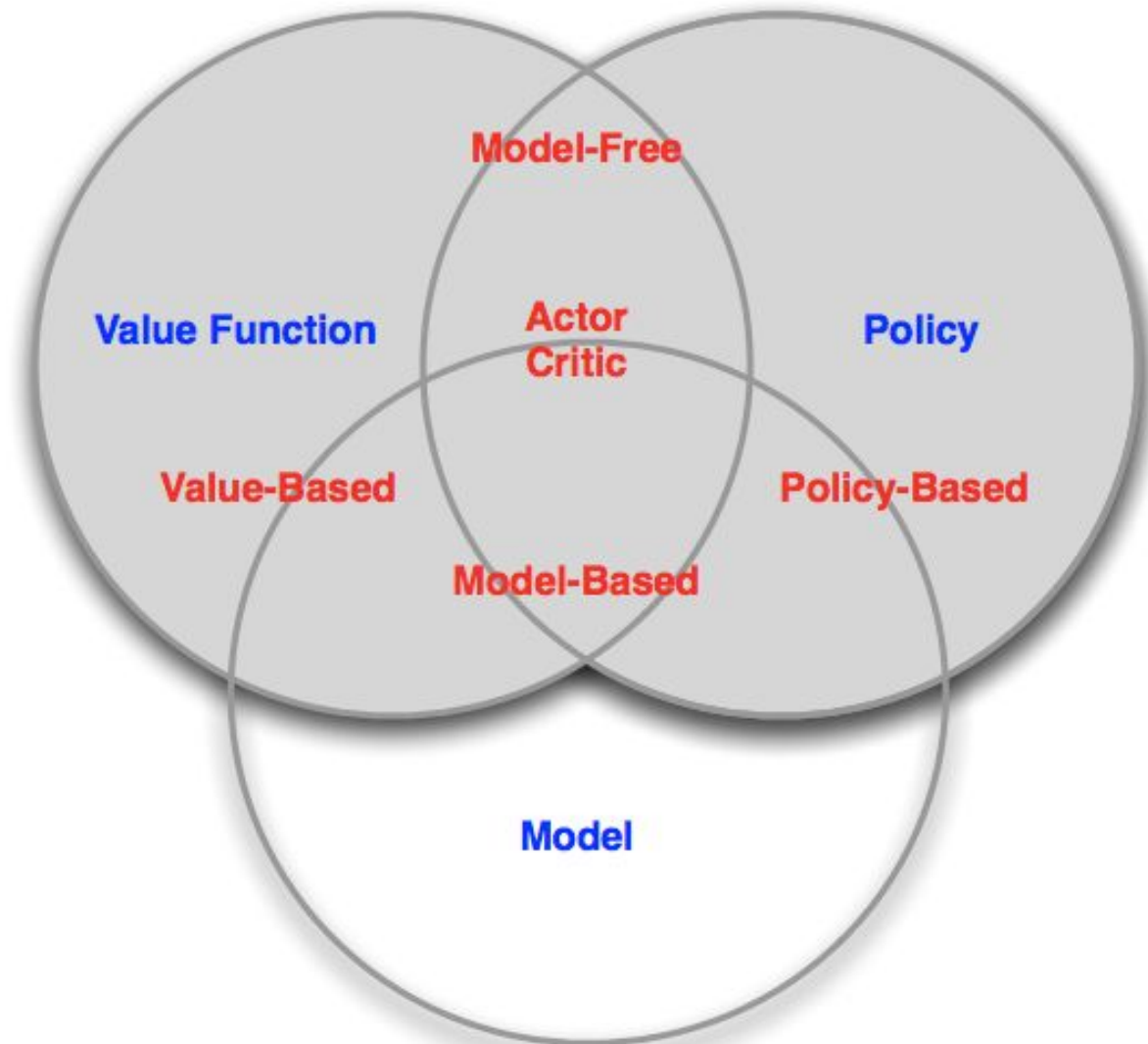
Learning Objectives

- Define the key features of reinforcement learning that distinguishes it from AI and non-interactive machine learning (as assessed by exams).
- Given an application problem (e.g. from computer vision, robotics, etc), decide if it should be formulated as a RL problem; if yes be able to define it formally (in terms of the state space, action space, dynamics and reward model), state what algorithm (from class) is best suited for addressing it and justify your answer (as assessed by the project and exams).
- Describe (list and define) multiple criteria for analyzing RL algorithms and evaluate algorithms on these metrics: e.g. regret, sample complexity, computational complexity, empirical performance, convergence, etc (as assessed by homeworks and exams).

What We've Covered So Far

- Markov decision process planning
- Model free policy evaluation
- Model-free learning to make good decisions
- Value function approximation, focus on model-free methods
- Imitation learning
- Policy search

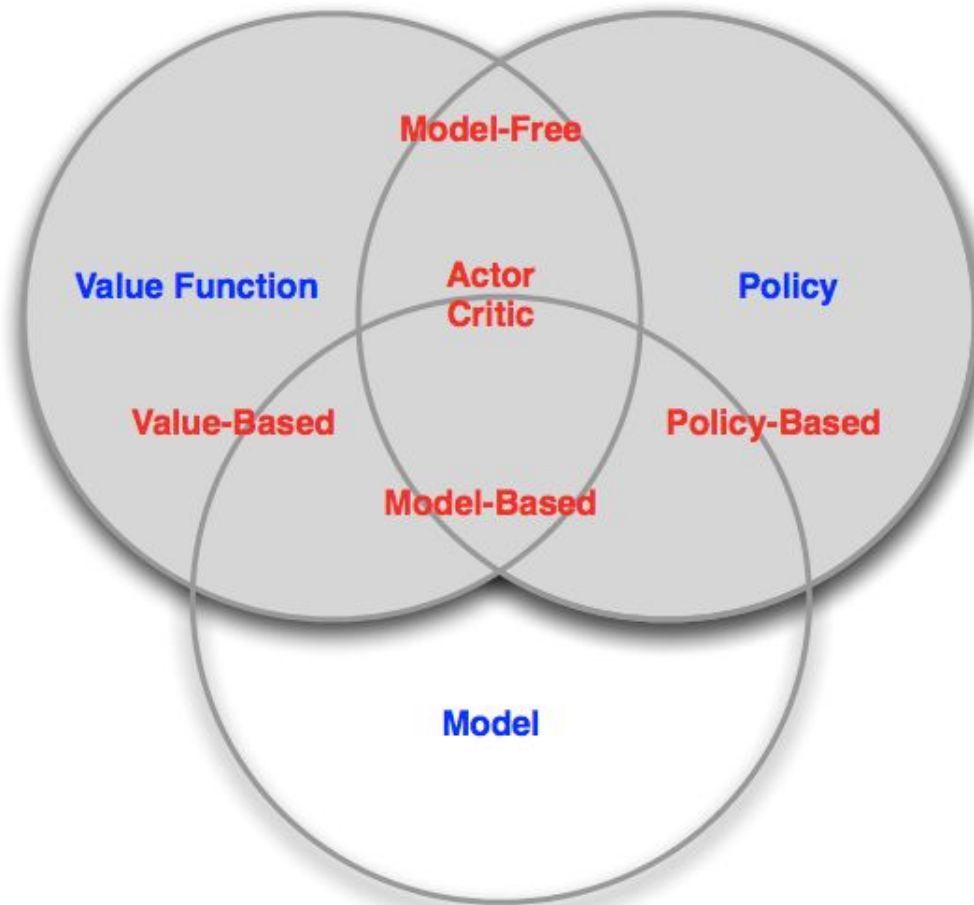
Reinforcement Learning



Reinforcement Learning

model \rightarrow value \rightarrow policy

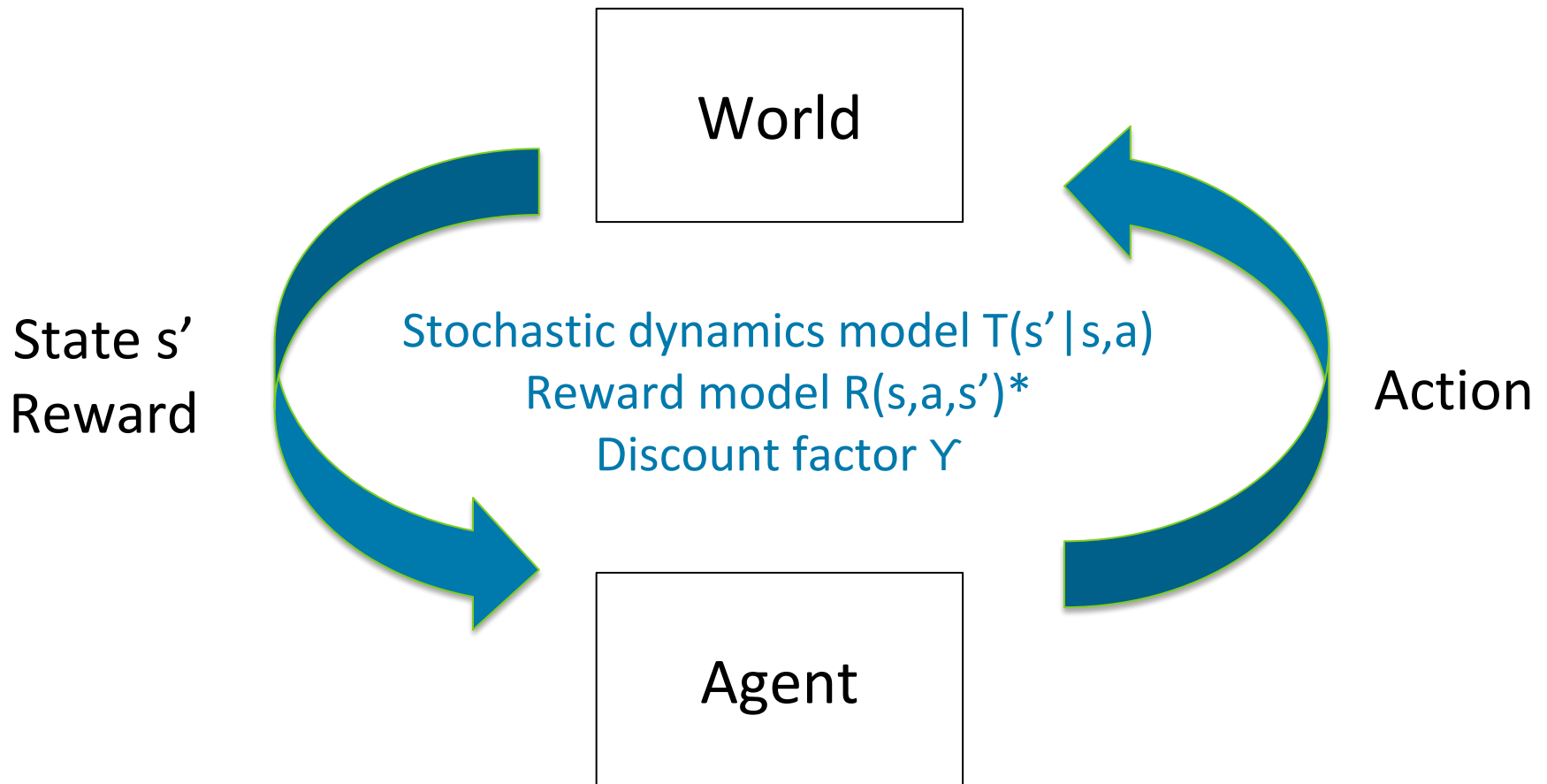
(ordering sufficient but not necessary, e.g. having a model is not required to learn a value)



What We've Covered So Far

- Markov decision process planning
- Model free policy evaluation
- Model-free learning to make good decisions
- Value function approximation, focus on model-free methods
- Imitation learning
- Policy search

Model: Frequently model as a Markov Decision Process, $\langle S, A, R, T, \gamma \rangle$



Policy mapping from state \rightarrow action

MDPs

- Define a MDP $\langle S, A, R, T, \gamma \rangle$
- Markov property
 - What is this, why is it important
- What are the MDP models / values V / state-action values Q / policy
- What is MDP planning? What is difference from reinforcement learning?
 - Planning = know the reward & dynamics
 - Learning = don't know reward & dynamics

Bellman Backup Operator

$$V_{k+1}(s) = \max_a \left[r(s, a) + \gamma \sum_{s' \in S} p(s'|a, s) V_k(s') \right]$$

 Bellman backup

- Bellman backup is a contraction if discount factor, $\gamma < 1$
- Bellman contraction operator: with repeated applications, guaranteed to converge to a single fixed point (the optimal value)

Value vs Policy Iteration

- Value iteration:
 - Compute optimal value if horizon= k
 - Note this can be used to compute optimal policy if horizon = k
 - Increment k
- Policy iteration:
 - Compute infinite horizon value of a policy
 - Use to select another (better) policy
 - Closely related to a very popular method in RL: policy gradient

Policy Iteration (PI)

1. $i=0$; Initialize $\pi_0(s)$ randomly for all states s
2. Converged = 0;
3. While $i \neq 0$ or $|\pi_i - \pi_{i-1}| > 0$
 - $i=i+1$
 - Policy **evaluation**: Compute V^π
 - Policy **improvement**:

$$Q^{\pi_i}(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) V^{\pi_i}(s')$$

$$\pi_{i+1}(s) = \arg \max_a Q^{\pi_i}(s, a)$$

Check Your Understanding

Consider finite state and action MDP and use a lookup table representation, $\gamma < 1$, infinite H :

- Does the initial setting of the value function in value iteration impact the final computed value? Why/why not?
- Do value iteration and policy iteration always yield the same solution?
- Is the number of iterations needed for PI on a tabular MDP with $|A|$ actions and $|S|$ states bounded?

What We've Covered So Far

- Markov decision process planning
- Model free policy evaluation
- Model-free learning to make good decisions
- Value function approximation, focus on model-free methods
- Imitation learning
- Policy search

Model-free Passive RL

- Directly estimate Q or V of a policy from data
- The Q function for a particular policy is the **expected discounted sum of future rewards** obtained by following policy starting with (s,a)
- For Markov decision processes,

$$Q^{\pi_i}(s, a) = r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi_i}(s')$$

Model-free Passive RL

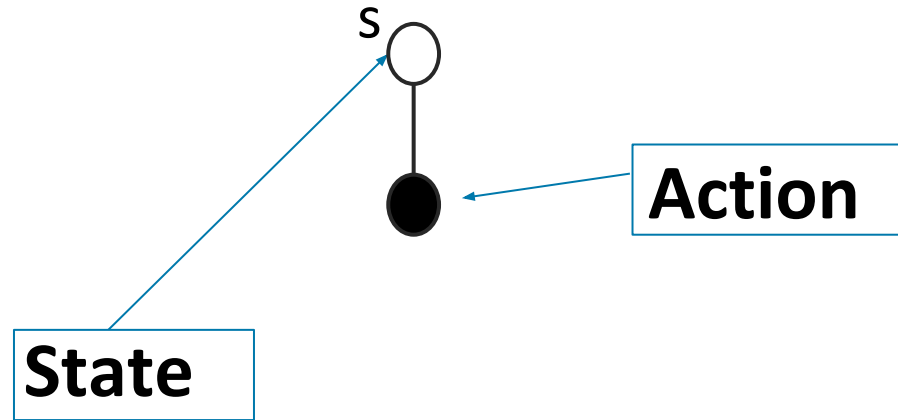
- Directly estimate Q or V of a policy from data
- The Q function for a particular policy is the **expected discounted sum of future rewards** obtained by following policy starting with (s,a)
- For Markov decision processes,

$$Q^{\pi_i}(s, a) = r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi_i}(s')$$

- Consider episodic domains
 - Act in world for H steps, then reset back to state sampled from starting distribution
- MC: directly average episodic rewards
- TD/Q-learning: use a “target” to bootstrap

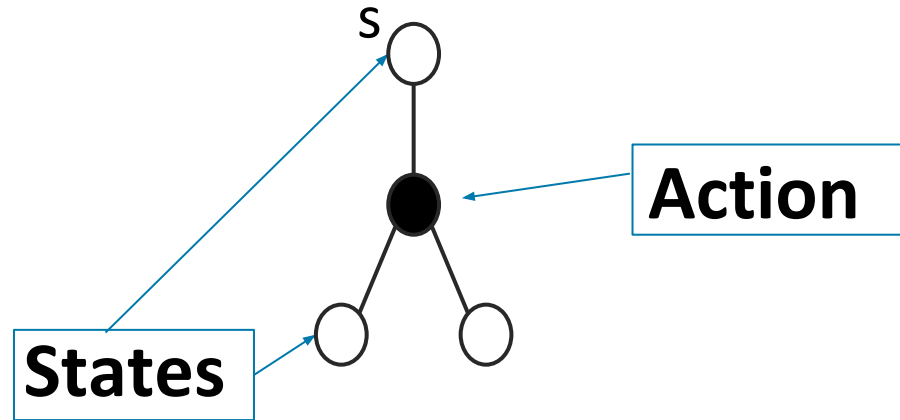
Dynamic Programming Policy Evaluation

$$V^{\pi}(s) \leftarrow \mathbb{E}_{\pi}[r_t + \gamma V_{i-1} | s_t = s]$$



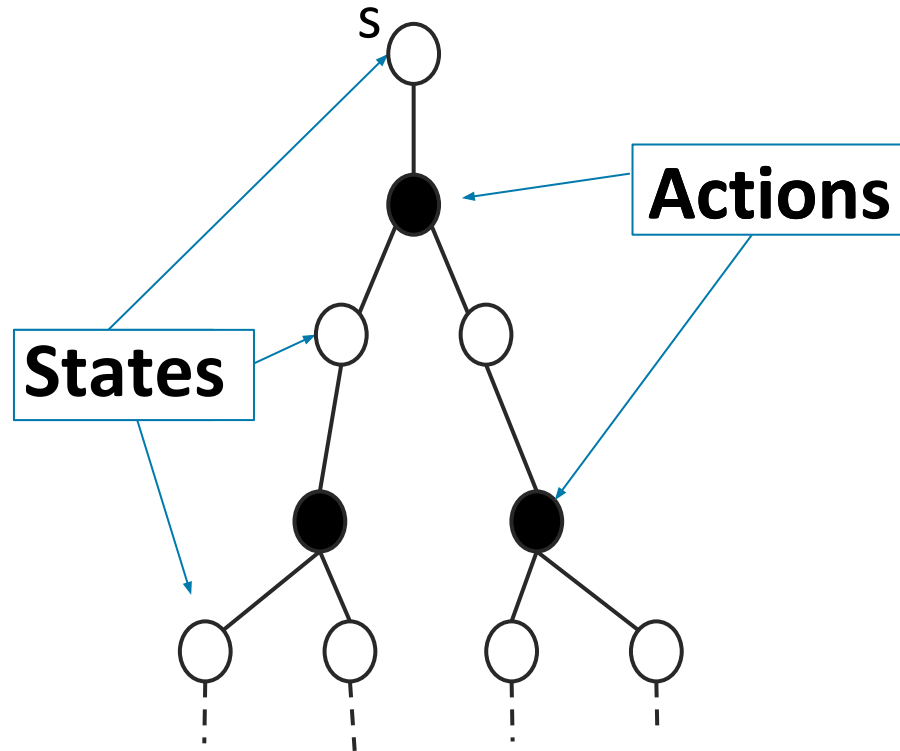
Dynamic Programming Policy Evaluation

$$V^{\pi}(s) \leftarrow \mathbb{E}_{\pi}[r_t + \gamma V_{i-1} | s_t = s]$$



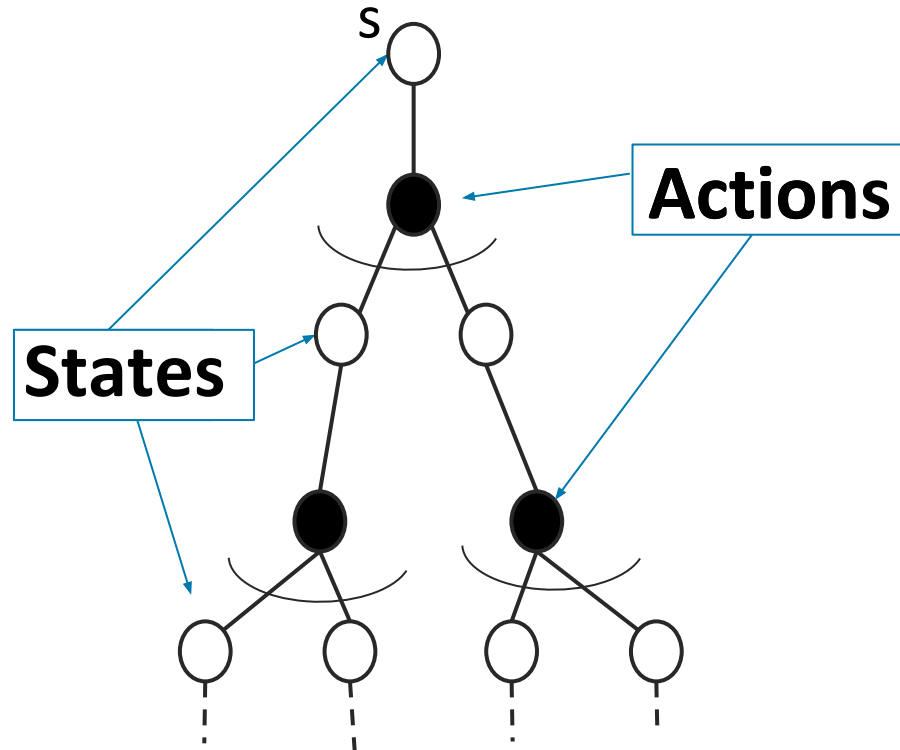
Dynamic Programming Policy Evaluation

$$V^\pi(s) \leftarrow \mathbb{E}_\pi[r_t + \gamma V_{i-1} | s_t = s]$$



Dynamic Programming Policy Evaluation

$$V^\pi(s) \leftarrow \mathbb{E}_\pi[r_t + \gamma V_{i-1} | s_t = s]$$

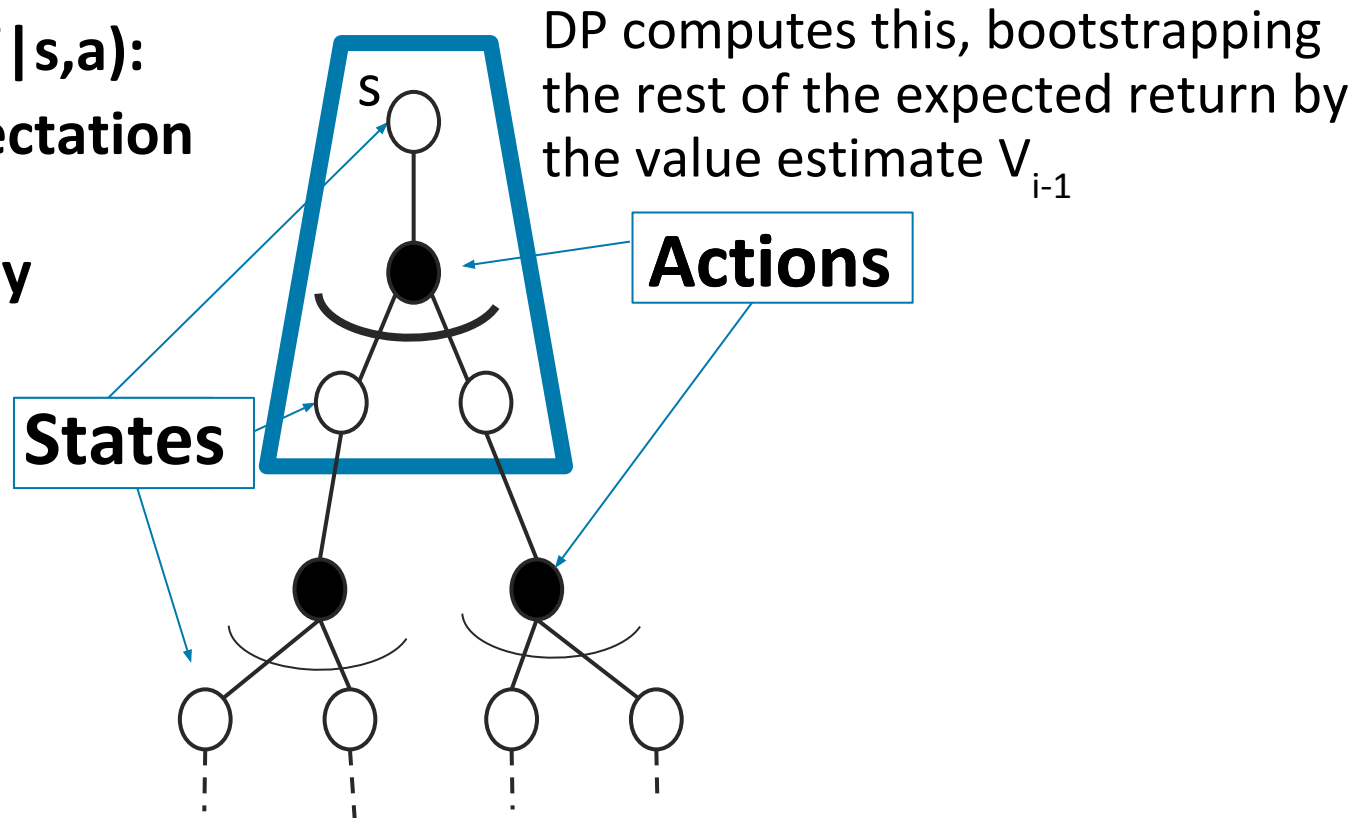


 = **Expectation**

Dynamic Programming Policy Evaluation

$$V^\pi(s) \leftarrow \mathbb{E}_\pi[r_t + \gamma V_{i-1} | s_t = s]$$

Know model $P(s' | s, a)$:
reward and expectation
over next states
computed exactly



 = **Expectation**

- Bootstrapping: Update for V uses an estimate

MC Policy Evaluation

$$V^{\pi}(s) = V^{\pi}(s) + \alpha(G_{it} - V^{\pi}(s))$$

$$V^\pi(s) = V^\pi(s) + \alpha(G_{it} - V^\pi(s))$$

The diagram illustrates a Markov Decision Process (MDP) with the following components:

- States:** Represented by circles. A box labeled "States" points to a white circle. A blue circle at the top is labeled "S".
- Actions:** Represented by curved lines branching from a state. A box labeled "Actions" points to one of these lines.
- Transitions:** Represented by straight lines connecting states. A thick blue line connects the top blue state to a blue state below it, which then connects to a terminal state "T" (a blue square). A black line connects the top blue state to a black state below it, which then branches into two white states.
- Terminal State:** A blue square labeled "T" at the end of a thick blue path.
- Annotations:** Text at the top left reads "The estimate of the return to expectation", with a blue arrow pointing to the top blue state "S".

 = Expectation
 = **Terminal state**

Temporal Difference Policy Evaluation

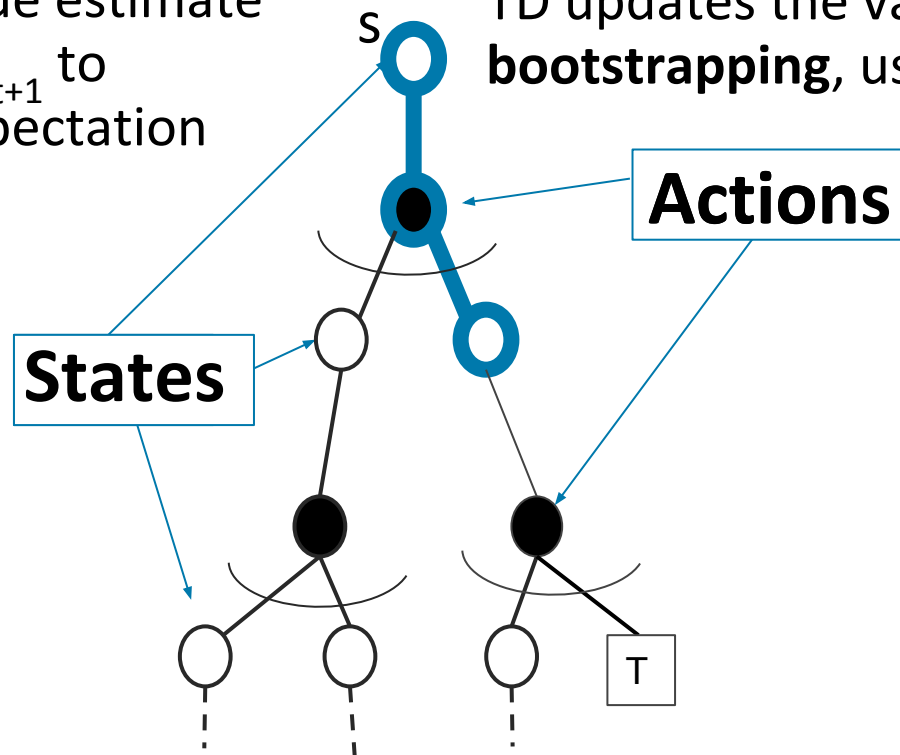
$$V^{\pi}(s_t) = V^{\pi}(s_t) + \alpha ([r_t + \gamma V^{\pi}(s_{t+1})] - V^{\pi}(s_t))$$

Temporal Difference Policy Evaluation

$$V^\pi(s_t) = V^\pi(s_t) + \alpha ([r_t + \gamma V^\pi(s_{t+1})] - V^\pi(s_t))$$

TD updates the value estimate using a **sample** of s_{t+1} to approximate an expectation

TD updates the value estimate by **bootstrapping**, uses estimate of $V(s_{t+1})$



\cup = Expectation

T = Terminal state

Check Your Understanding?

(Answer Yes/No/NA to Each Algorithm for Each Part)

- Usable when no models of current domain
 - DP: MC: TD:
- Handles continuing (non-episodic) domains
 - DP: MC: TD:
- Handles Non-Markovian domains
 - DP: MC: TD:
- Converges to true value of policy in limit of updates*
 - DP: MC: TD:
- Unbiased estimate of value
 - DP: MC: TD:

* For tabular representations of value function.

Some Important Properties to Evaluate Model-free Policy Evaluation Algorithms

- Bias/variance characteristics
- Data efficiency
- Computational efficiency

What We've Covered So Far

- Markov decision process planning
- Model free policy evaluation
- **Model-free learning to make good decisions**
- Value function approximation, focus on model-free methods
- Imitation learning
- Policy search

Q-Learning

- Update $Q(s,a)$ every time experience (s,a,s',r)
 - Create new target / sample estimate

$$\begin{aligned}Q_{samp}(s, a) &= r + \gamma V(s') \\ &= r + \gamma \max_{a'} Q(s', a')\end{aligned}$$

- Update estimate of $Q(s,a)$

$$Q(s, a) = (1 - \alpha)Q(s, a) + \alpha Q_{samp}(s, a)$$

Q-Learning Properties

- If acting randomly*, Q-learning converges Q^*
 - Optimal Q values
 - Finds optimal policy
- Off-policy learning
 - Can act in one way
 - But learning values of another π (the optimal one!)

*Again, under mild reachability assumptions

Check Your Understanding: T/F (or T/F and under what conditions)

- In an MDP with finite state- and action spaces using a lookup table, Q-learning with a ϵ -greedy policy converges to the optimal policy in the limit of infinite data.
- Monte-Carlo estimation cannot be used in MDPs with large state-spaces.
- Model-based reinforcement learning is always more data efficient than model-free RL

What We've Covered So Far

- Markov decision process planning
- Model free policy evaluation
- Model-free learning to make good decisions
- **Value function approximation, focus on model-free methods**
- Imitation learning
- Policy search

Monte Carlo vs TD Learning: Convergence in On Policy Case

- Evaluating value of a single policy

$$MSVE(w) = \sum_{s \in \mathcal{S}} d(s) \left(V^\pi(s) - \tilde{V}^\pi(s, w) \right)^2$$

- where
 - $d(s)$ is generally the on-policy π stationary distrib
 - $\tilde{V}(s, w)$ is the value function approximation

Convergence given infinite amount of data?

Monte Carlo Convergence: Linear VFA

- Evaluating value of a single policy

$$MSVE(w) = \sum_{s \in S} d(s) \left(V^\pi(s) - \tilde{V}^\pi(s, w) \right)^2$$

- where
 - $d(s)$ is generally the on-policy π stationary distrib
 - $\tilde{V}(s, w)$ is the value function approximation
- Linear VFA: $V(s) = \sum w_i f_i(s)$
- **Monte Carlo converges to min MSE possible!**

$$MSVE(w_{MC}) = \min_w \sum_{s \in S} d(s) \left(V^\pi(s) - \tilde{V}^\pi(s, w) \right)^2$$

TD Learning Convergence: Linear VFA

- Evaluating value of a single policy

$$MSVE(w) = \sum_{s \in S} d(s) \left(V^\pi(s) - \tilde{V}^\pi(s, w) \right)^2$$

- where
 - $d(s)$ is generally the on-policy π stationary distrib
 - $\tilde{V}(s, w)$ is the value function approximation
- Linear VFA: $V(s) = \sum w_i f_i(s)$
- **TD converges to constant factor of best MSE**

$$\begin{aligned} MSVE(w_{TD}) &= \frac{1}{1 - \gamma} \min_w \sum_{s \in S} d(s) \left(V^\pi(s) - \tilde{V}^\pi(s, w) \right)^2 \\ &= \frac{1}{1 - \gamma} MSVE(w_{MC}) \end{aligned}$$

Off Policy Learning

Can use importance sampling with MC to estimate the value of a policy didn't try (but that has support in behavior policy)

Q-learning with function approximation can diverge (not converge even given infinite data)

Deep Learning & Model Free Q learning

$$\mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[\left(\underbrace{r + \gamma \max_{a'} Q(s', a'; w_i^-)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right)^2 \right]$$

- Running stochastic gradient descent
- Now use a deep network to approximate Q

Challenges

- Challenge of using function approximation
 - Local updates (s, a, r, s') highly correlated
 - “Target” (approximation to true value of s') can change quickly and lead to instabilities

DQN: Q-learning with DL

- Experience replay of mix of prior (s_i, a_i, r_i, s_{i+1}) tuples to update $Q(w)$
- Fix target $Q(w^-)$ for number of steps, then update
- Optimize MSE between current Q and Q target

$$\mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[\left(\underbrace{r + \gamma \max_{a'} Q(s', a'; w_i^-)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right)^2 \right]$$

- Use stochastic gradient descent

Deep RL

- Experience replay is hugely helpful
- Target stablization is also helpful
- No guarantees on convergence (yet)
- Some other influential ideas
 - Double Q (two separate networks, each act as a “target” for each other)
 - Dueling: separate value and advantage
 - Many advances in deep RL build on prior ideas for RL with look up table representations

Check Your Understanding: T/F (or T/F and under what conditions)

- In finite state spaces with features that can represent the true value function, TD learning with value function approximation always finds the true value function of the policy given sufficient data.

What We've Covered So Far

- Markov decision process planning
 - Model free policy evaluation
 - Model-free learning to make good decisions
 - Value function approximation, focus on model-free methods
 - **Imitation learning**
 - **Policy search**
- **These will only be tested at a lighter level, since the homework will be post the midterm**

Imitation Learning

- Behavioral cloning
 - Definition
 - What can go wrong

Imitation Learning

- Inverse reinforcement learning
 - Formulation
 - How many reward models are compatible with a demonstration of the state-action sequence assuming that sequence comes from the optimal policy?
 - What's a popular way to choose among these?

Policy Search

- Why are stochastic parameterized policies useful?
- Are policy gradient methods the only form of policy search? If not, are they the best type?
- Does the likelihood ratio policy gradient require us to know the dynamics model?
- Give 2 ideas that are used to reduce the variance of the default likelihood ratio policy gradient estimator

Midterm review

- Markov decision process planning
- Model free policy evaluation
- Model-free learning to make good decisions
- Value function approximation, focus on model-free methods
- Imitation learning
- Policy search

Midterm review

- To study: go through lecture notes, do all of assignments 1 and 2
- Do the practice midterm
 - This is only the second offering of this class, so unfortunately we only have 1 past midterm
 - Ignore parts on R-max and other topics we did not cover before the midterm
- Reach out to us on piazza or during office hours with any questions
- You can bring a 1 sided 1 page of notes.
- Good luck!