Model-based Reinforcement Learning

Course: Reinforcement Learning, Bachelor AI, Leiden University

Lecturer: Thomas Moerland

Content

I: Integrated view of planning and learning

II: Combination of planning and learning

Content

I: Integrated view of planning and learning

- 1. Access to the MDP dynamics
- 2. Planning versus learning
- 3. Types of back-ups

II: Combination of planning and learning

Content

I: Integrated view of planning and learning

- 1. Access to the MDP dynamics
- 2. Planning versus learning
- 3. Types of back-ups

II: Combination of planning and learning

- 4. Model-based reinforcement learning (MBRL)
- 5. Learning a model
- 6. Model-based RL algorithms

Dynamics function: p(s'|s,a)

Which next state s can we observe after taking action a in state s

Dynamics function: p(s'|s,a)

Which next state s can we observe after taking action a in state s

Two important considerations:

1) Is our access *reversible* or *irreversible*?

Access to the environment dynamics

- <u>Reversible access</u>
 - We can take try any action in any state we want
 - Similar to planning in our head
 - Call such access a **model**



Access to the environment dynamics

- <u>Reversible access</u>
 - We can take try any action in any state we want
 - Similar to planning in our head
 - Call such access a **model**



- <u>Irreversible access</u>
 - When we take an action, we have to continue from the next state
 - Similar to the real world
 - Call such access an **environment**



Dynamics function: p(s'|s,a)

Which next state s can we observe after taking action a in state s

Two important considerations:

1) Is our access *reversible* or *irreversible*? (model versus environment)

Dynamics function: p(s'|s,a)

Which next state s can we observe after taking action a in state s

Two important considerations:

- 1) Is our access *reversible* or *irreversible*? (model versus environment)
- 2) Do we get the *full distribution* of p(s'|s,a), or a *sample*?

Dynamics function: p(s'|s,a)

Which next state s can we observe after taking action a in state s

Two important considerations:

- 1) Is our access *reversible* or *irreversible*? (model versus environment)
- 2) Do we get the *full distribution* of p(s'|s,a), or a *sample*?

$$p(s'=1) = 0.1$$
 $s' = 3!$ (this time)
 $p(s'=2) = 0.2$
 $p(s'=3) = 0.7$

Dynamics function: p(s'|s,a)

Which next state s can we observe after taking action a in state s

Two important considerations:

- 1) Is our access *reversible* or *irreversible*? (model versus environment)
- 2) Do we get the *full distribution* of p(s'|s,a), or a *sample*?

$$p(s'=1) = 0.1$$
 $s' = 3!$ (this time)
 $p(s'=2) = 0.2$
 $p(s'=3) = 0.7$

Next page: 2-by-2 overview of these considerations

















Summary 1: Access to MDP dynamics

1. Reversible versus irreversible access

2. Distributional versus sample models

2) Difference between planning and learning

Planning

Reinforcement learning

Planning

Reinforcement learning





Planning

Reinforcement learning









Planning

Reinforcement learning



1/3

(0/1

1/1

(2/3)

(0/1)

(0/1



Both can solve the same MDP optimization problem





So what discriminates planning from reinforcement learning?

So what discriminates planning from reinforcement learning?

Two factors:

1)	Access to the MDP dynamics:	reversible	or	irreversible
2)	Storage of the solution:	local	or	global

So what discriminates planning from reinforcement learning?

Two factors:

1)	Access to the MDP dynamics:	reversible	or	irreversible
2)	Storage of the solution:	local	or	global

Local versus global solution

Question: can you guess what the difference between a local and global solution might be?

Local versus global solution

- A **local** solution temporarily stores a solution for a subset of all states
 - focus on current state
 - discarded after execution
 - e.g., a planning tree





Local versus global solution

- A **local** solution temporarily stores a solution for a subset of all states
 - focus on current state
 - discarded after execution
 - e.g., a planning tree

- A **global** solution permanently stores estimates for all states
 - e.g., a value table
 - only option in (model-free) RL, because we have to move forward and do not know when we get back to a state







Problem: we have two possible distinctions

Problem: we have two possible distinctions

	<u>Local solution</u>	<u>Global solution</u>
<u>Reversible MDP</u> access		
<u>Irreversible</u> MDP access		

Problem: we have two possible distinctions

	Local solution	<u>Global solution</u>
<u>Reversible MDP</u> <u>access</u>	Planning e.g., MCTS	
<u>Irreversible</u> MDP access		
Problem: we have two possible distinctions

	<u>Local solution</u>	<u>Global solution</u>
<u>Reversible MDP</u> <u>access</u>	Planning e.g., MCTS	
<u>Irreversible</u> MDP access		Model-free RL e.g., Q-learning

Problem: we have two possible distinctions

	<u>Local solution</u>	<u>Global solution</u>
<u>Reversible MDP</u> <u>access</u>	Planning e.g., MCTS	Borderline/Model-based RL e.g., Dynamic Programming
<u>Irreversible</u> MDP access		Model-free RL e.g., Q-learning

Problem: we have two possible distinctions

	<u>Local solution</u>	<u>Global solution</u>
<u>Reversible MDP</u> access	Planning e.g., MCTS	Borderline/Model-based RL e.g., Dynamic Programming
<u>Irreversible</u> MDP access	(impossible)	Model-free RL e.g., Q-learning

Q: Why is it impossible to use a local solution when we have irreversible access?

Problem: we have two possible distinctions

	<u>Local solution</u>	<u>Global solution</u>
<u>Reversible MDP</u> access	Planning e.g., MCTS	Borderline/Model-based RL e.g., Dynamic Programming
<u>Irreversible</u> MDP access	(impossible)	Model-free RL e.g., Q-learning

Q: Why is it impossible to use a local solution when we have irreversible access? A: Local solutions get discarded after execution of the real action, but if the environment is irreversible, we directly throw away our new solution after the first sample.

Summary 2: Planning versus learning

- 1. Planning:reversible model +local solution
- 2. Model-free RL: irreversible model + global solution
- 3. Borderline/model-based RL: reversible model + global solution

3) Comparison of back-ups

Expected versus sample back-ups

Our access to the MDP dynamics also influences the way we can back-up information!

Expected versus sample back-ups

Our access to the MDP dynamics also influences the way we can back-up information!

Need to distinguish:

- **Expected** back-ups (mostly planning)
- **Sample** back-ups (mostly RL)

Expected versus sample back-ups

Our access to the MDP dynamics also influences the way we can back-up information!

Need to distinguish:

- **Expected** back-ups (mostly planning)
- **Sample** back-ups (mostly RL)

We will illustrate this difference with back-up diagrams

A back-up diagram graphically represents a one-step back-up equation

A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the Bellman equation



A back-up diagram graphically represents a one-step back-up equation



A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the Bellman equation



To compute a new estimate of V(s)

A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the Bellman equation



We sum over all actions

A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the Bellman equation



A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the Bellman equation



A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the Bellman equation

= **full expectation** over the actions (policy) and next states (dynamics)



policy evaluation

A back-up diagram graphically represents a one-step back-up equation



A back-up diagram graphically represents a one-step back-up equation



A back-up diagram graphically represents a one-step back-up equation



A back-up diagram graphically represents a one-step back-up equation



A back-up diagram graphically represents a one-step back-up equation



A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the <u>Temporal Difference Learning</u>

This is a **sample** back-up over the actions (policy) and next states (dynamics)



A back-up diagram graphically represents a one-step back-up equation

Back-up diagram for the <u>Temporal Difference Learning</u>

This is a **sample** back-up over the actions (policy) and next states (dynamics)



Sample estimates will converge to the correct equation over multiple updates/samples









State-action values Q(s,a)

Q: Categorize Q-learning.

Spectrum of possible 1-step back-ups



State-action values Q(s,a)



State-action values Q(s,a)



Advice: carefully study these 1-step back-up diagrams, they provide a lot of insight/overview

Can also extend the back-up over multiple-steps
















Back-up diagrams can provide much intuition about the space of possible back-ups

Back-up diagrams can provide much intuition about the space of possible back-ups

Essentially four considerations:

- 1) The width over the actions: sample (shallow) or expected (broad)
- 2) The back-up policy: on-policy or off-policy
- 3) The width over the dynamics: sample (shallow) or expected (broad)
- 4) The depth of the back-up: 1-step (shallow) or full depth (deep)

Back-up diagrams can provide much intuition about the space of possible back-ups

Essentially four considerations:

- 1) The width over the actions: sample (shallow) or expected (broad)
- 2) The back-up policy: on-policy or off-policy
- 3) The width over the dynamics: sample (shallow) or expected (broad)
- 4) The depth of the back-up: 1-step (shallow) or full depth (deep)



Question: can you indicate each decision for plain SARSA?

Back-up diagrams can provide much intuition about the space of possible back-ups

Essentially four considerations:

- 1) The width over the actions: **sample** (shallow) or expected (broad)
- 2) The back-up policy: **on-policy** or off-policy
- 3) The width over the dynamics: **sample** (shallow) or expected (broad)
- 4) The depth of the back-up: **1-step** (shallow) or full depth (deep)



Question: can you indicate each decision for plain SARSA?

Break

Interaction with the real world is typically irreversible.

Q: How could we still obtain a model of how the world works?

When the environment is irreversible, maybe we can still

learn a reversible model from data



Inspired by the way humans acquire (reversible) dynamics models:

we learn them from real-world (irreversible) experience











Two important steps:

1. How to learn a model from data? (Sec. 5)

2. How to integrate planning updates and learning updates? (Sec. 6)

5) Learning a model

Learning a model

Main question:

Given a dataset of observed transitions (s,a,r,s'), how can we estimate

1) the dynamics p(s'|s,a) and

2) The reward r(s,a,s') function?

Learning a model

Main question:

Given a dataset of observed transitions (s,a,r,s'), how can we estimate

1) the dynamics p(s'|s,a) and

2) The reward r(s,a,s') function?

- 1. Track counts **n(s,a,s')**
 - a. Number of times we observed s' after taking a in s.

- 1. Track counts **n(s,a,s')**
 - a. Number of times we observed s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')

- 1. Track counts **n(s,a,s')**
 - a. Number of times we observed s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \ge |A| \ge |S|$

- 1. Track counts **n(s,a,s')**
 - a. Number of times we observed s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \ge |A| \ge |S|$
- 2. Estimate p(s'|s,a) by normalizing the observed counts:

$$\hat{p}(s'|s,a) = \frac{n(s,a,s')}{n(s,a)}$$

- 1. Track counts **n(s,a,s')**
 - a. Number of times we observed s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \times |A| \times |S|$
- 2. Estimate p(s'|s,a) by normalizing the observed counts:



- 1. Track counts **n(s,a,s')**
 - a. Number of times we observed s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \times |A| \times |S|$
- 2. Estimate p(s'|s,a) by normalizing the observed counts:



Total number of trials at s,a

Example:

After taking a=1 in s=1, we have stored the following counts:

- n(s=1,a=1,s'=1) = 4
- n(s=1,a=1,s'=2) = 2
- n(s=1,a=1,s'=3) = 6

Example:

After taking a=1 in s=1, we have stored the following counts:

- n(s=1,a=1,s'=1) = 4
- n(s=1,a=1,s'=2) = 2
- n(s=1,a=1,s'=3) = 6

Q: Compute
$$p(s'|s=1,a=1)$$

$$\hat{p}(s'|s,a) = \frac{n(s,a,s')}{n(s,a)}$$

Example:

After taking a=1 in s=1, we have stored the following counts:

- n(s=1,a=1,s'=1) = 4
- n(s=1,a=1,s'=2) = 2
- n(s=1,a=1,s'=3) = 6

Q: Compute
$$p(s'|s=1,a=1)$$

$$\hat{p}(s'|s,a) = \frac{n(s,a,s')}{n(s,a)}$$

A:

- $p(s'=1|s=1,a=1) = 4/12 = \frac{1}{3}$
- $p(s'=2|s=1,a=1) = 2/12 = \frac{1}{6}$
- $p(s'=3|s=1,a=1) = 6/12 = \frac{1}{2}$

Learning a model

Main question:

Given a dataset of observed transitions (s,a,r,s'), how can we estimate

- 1) the dynamics p(s'|s,a) and
- 2) the reward r(s,a,s') function?

- 1.
- Also track total transition rewards **R**_{sum}(**s**,**a**,**s**') a. Sum of all observed rewards when reaching s' after taking a in s.

- 1. Also track total transition rewards **R**_{sum}(s,a,s')
 - a. Sum of all observed rewards when reaching s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \ge |A| \ge |S|$

- 1. Also track total transition rewards **R**_{sum}(**s**,**a**,**s**')
 - a. Sum of all observed rewards when reaching s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \times |A| \times |S|$
- 2. Estimate r(s,a,s') by computing the average transition reward:

$$\hat{r}(s, a, s') = \frac{R_{\text{sum}}(s, a, s')}{n(s, a, s')}$$

- 1. Also track total transition rewards **R**_{sum}(s,a,s')
 - a. Sum of all observed rewards when reaching s' after taking a in s.
 - b. Can estimate from transition data (s,a,r,s')
 - c. Array of size $|S| \times |A| \times |S|$
- 2. Estimate r(s,a,s') by computing the average transition reward:


Model estimation pseudocode

Full pseudocode in the lecture notes

Model estimation pseudocode

Full pseudocode in the lecture notes

Algorithm 1: Tabular model update pseudo-code. PS = prioritized sweeping.

Input: Maximum number of timesteps T. **Initialization**: Initialize n(s, a, s') = 0 and $R_{sum}(s, a, s') = 0$ $\forall s \in S, a \in A$ **repeat** T **times**

Observe
$$\langle s, a, r, s' \rangle$$

 $n(s, a, s') \leftarrow n(s, a, s') + 1$
 $R_{sum}(s, a, s') \leftarrow R_{sum}(s, a, s') + r$
 $\hat{p}(s'|s, a) = \frac{n(s, a, s')}{\sum_{s'} n(s, a, s')}$
 $\hat{r}(s, a, s') = \frac{R_{sum}(s, a, s')}{n(s, a, s')}$
 $\hat{p}(s, a|s') = \frac{n(s, a, s')}{\sum_{s, a} n(s, a, s')}$

/* Observe new transition */
/* Update transition counts */
 /* Update reward sums */
/* Estimate transition function */
 /* Estimate reward function */
 /* Reverse model (only for PS) */

end

Estimate tabular p(s'|s,a) and r(s,a,s') from observed data tuples <s,a,r,s'>

Estimate tabular p(s'|s,a) and r(s,a,s') from observed data tuples <s,a,r,s'>

1. Maintain counts n(s,a,s') and total rewards R_{sum}(s,a,s').

Estimate tabular p(s'|s,a) and r(s,a,s') from observed data tuples <s,a,r,s'>

- 1. Maintain counts n(s,a,s') and total rewards R_{sum}(s,a,s').
- 2. Compute transition model p(s'|s,a) from normalizing counts.

Estimate tabular p(s'|s,a) and r(s,a,s') from observed data tuples <s,a,r,s'>

- 1. Maintain counts n(s,a,s') and total rewards R_{sum}(s,a,s').
- 2. Compute transition model p(s'|s,a) from normalizing counts.
- 3. Compute reward model r(s,a,s') from averaging total rewards.

Estimate tabular p(s'|s,a) and r(s,a,s') from observed data tuples <s,a,r,s'>

- 1. Maintain counts n(s,a,s') and total rewards R_{sum}(s,a,s').
- 2. Compute transition model p(s'|s,a) from normalizing counts.
- 3. Compute reward model r(s,a,s') from averaging total rewards.

Next section: how may this model be useful?

We want to make use of the learned model!

(combine planning and learning)

Discuss three algorithms:

- a) Real-time Dynamic Programming (RTDP)
- b) Dyna
- c) Prioritized sweeping

Discuss three algorithms:

- a) Real-time Dynamic Programming (RTDP)
- b) Dyna
- c) Prioritized sweeping

We already discussed **Dynamic Programming** (DP) in an earlier lecture

We already discussed <u>Dynamic Programming</u> (DP) in an earlier lecture

- Classic bridging algorithm between planning and learning
- Sweep through the entire state space
- At each state update with Bellman optimality equation
- Guaranteed convergence

We already discussed Dynamic Programming (DP) in an earlier lecture

- Classic bridging algorithm between planning and learning
- Sweep through the entire state space
- At each state update with Bellman optimality equation
- Guaranteed convergence

Main problem = <u>curse of dimensionality</u>









Solution of **real-time DP**: apply DP updates on traces sampled from the start

(All model-free RL approaches do this by definition, since the access to the MDP is then irreversible)







Discuss three algorithms:

- a) Real-time Dynamic Programming (RTDP)
- b) Dyna
- c) Prioritized sweeping

Main idea:

Learn a model to generate additional transition data,

apply standard update to these simulated transitions.









learn model





learn model

Algorithm 2: Dyna Q-learning with ϵ -greedy exploration.

Input: Number of planning updates K, exploration parameter $\epsilon \in (0, 1]$, learning rate $\alpha \in (0, 1]$, discount parameter $\gamma \in [0, 1]$, maximum number of timesteps T. Initialization: Initialize $\hat{Q}(s, a) = 0$, n(s, a, s') = 0, $R_{sum}(s, a, s') = 0$ $\forall s \in S, a \in A$. for t = 1...T do $s \leftarrow \text{current state}$ /* Reset when environment terminates */ $a \sim \pi_{\epsilon\text{-greedy}}(a|s)$ /* Sample action */ $r, s' \sim p(r, s'|s, a)$ /* Simulate environment */ $\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha \cdot [r + \gamma \cdot \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$ /* Update Q-table */ end

Algorithm 2: Dyna Q-learning with ϵ -greedy exploration.

Input: Number of planning updates K, exploration parameter $\epsilon \in (0, 1]$, learning rate $\alpha \in (0, 1]$, discount parameter $\gamma \in [0, 1]$, maximum number of timesteps T. **Initialization**: Initialize $\hat{Q}(s, a) = 0$, n(s, a, s') = 0, $R_{sum}(s, a, s') = 0$ $\forall s \in \mathcal{S}, a \in \mathcal{A}$. for t = 1...T do /* Reset when environment terminates */ $s \leftarrow \text{current state}$ $a \sim \pi_{\epsilon\text{-greedy}}(a|s)$ /* Sample action */ $r, s' \sim p(r, s'|s, a)$ /* Simulate environment */ /* Update model (Alg.1) */ $\hat{p}(s', r | s, a) \leftarrow \text{Update}(s, a, r, s')$ Learn model $\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha \cdot [r + \gamma \cdot \max_{a'} \hat{Q}(s',a') - \hat{Q}(s,a)]$ /* Update Q-table */ repeat K times Make k /* State to plan on */ $s \leftarrow$ random previously observed state planning /* Planning action */ $a \leftarrow$ previously taken action in state s $\begin{array}{ll} s',r \sim \hat{p}(s',r|s,a) & /* \text{ Simulate model }*/\\ \hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha \cdot [r + \gamma \cdot \max_{a'} \hat{Q}(s',a') - \hat{Q}(s,a)] & /* \text{ Update Q-table }*/ \end{array}$ updates in between every real end \mathbf{end} step

Algorithm 2: Dyna Q-learning with ϵ -greedy exploration.

Algorithm you will implement in the assignment!



WITHOUT PLANNING (n=0)

				G
				1
S				

WITH PLANNING (n=50)



Discuss three algorithms:

- a) Real-time Dynamic Programming (RTDP)
- b) Dyna
- c) Prioritized sweeping
If the value estimate of a state changes a lot,

then the states that precede it should probably also be updated.

If the value estimate of a state changes a lot,

then the states that precede it should probably also be updated.

Main idea:

Use a backward/reverse model to identify states that likely need updating (backwards search to spread information faster)







Prioritized sweeping Multistep updates will propagate further along the trace





Main idea: *prioritize* states that deserve an update & additional backward search

Main idea: *prioritize* states that deserve an update & additional backward search

- Maintain priority queue, with priority as the <u>absolute TD error</u>

$$\mathbf{p} \leftarrow |r + \gamma \cdot \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)|$$

$$\overset{\text{New Q-learning}}{\underset{\text{estimate}}{}}$$
Current estimate

Pseudocode in assignment/lecture notes

Algorithm 3: Prioritized sweeping (Q-learning with ϵ -greedy exploration).			
Input: Number of planning updates K , exploration parameter $\epsilon \in (0, 1]$, learning rate $\alpha \in (0, 1]$, discount parameter $\gamma \in [0, 1]$, maximum number of timesteps T ,			
priority threshold θ .			
Initialization: Initialize $Q(s, a) = 0$, $n(s, a, s) = 0$, $R_{sum}(s, a, s) = 0$, $\forall s \in S, a \in A$, and prioritized queue PO			
for $t = 1$. T do			
$s \leftarrow \text{current state}$ /* Reset	when environment terminates */		
$a \sim \pi_{\epsilon \text{-greedy}}(a s)$	/* Sample action */		
$r, s' \sim p(r, s' s, a)$	/* Simulate environment */		
$\hat{p}(s', r s, a) \leftarrow \text{Update}(s, a, r, s')$	/* Update model (Alg.1) */		
$\mathbf{p} \leftarrow r + \gamma \cdot \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) $	/* Compute priority \mathbf{p} */		
$\mathbf{if} \ \mathbf{p} > \theta \ \mathbf{then}$			
Insert (s, a) into PQ with priority p /	<pre>/* State-action needs update */</pre>		
end			
/// Start sampling from PQ to perform updates			
repeat K times			
$s, a \leftarrow \text{pop highest priority from PQ}$ /* S	Sample PQ , break when empty */		
$s', r \sim \hat{p}(s', r s, a)$	/* Simulate model */		
$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha \cdot [r + \gamma \cdot \max_{a'} \hat{Q}(s',a') - \hat{Q}(s,a)]$ /* Update Q-table */			
/// Loop over all state action that may lead to state s			
for all (\bar{s}, \bar{a}) with $\hat{p}(\bar{s}, \bar{a} s) > 0$ do			
$ar{r}=\hat{r}(ar{s},ar{a},s)$	<pre>/* Get reward from model */</pre>		
$\mathbf{p} \leftarrow \bar{r} + \gamma \cdot \max_a \hat{Q}(s, a) - \hat{Q}(\bar{s}, \bar{a}) $ if $\mathbf{p} > \theta$ then	/* Compute priority \mathbf{p} */		
Insert (s, a) into PQ with priority p / end	<pre>/* State-action needs update */</pre>		
end			
end			
end			

Algorithm 3: Prioritized sweeping (Q-learning with ϵ -greedy exploration).	
Input: Number of planning updates K, exploration parameter $\epsilon \in (0, 1]$, learning rate	
$\alpha \in (0, 1]$, discount parameter $\gamma \in [0, 1]$, maximum number of timesteps T,	
priority threshold $\hat{\theta}$.	
Initialization : Initialize $Q(s, a) = 0$, $n(s, a, s') = 0$, $R_{sum}(s, a, s') = 0$ $\forall s \in \mathcal{S}, a \in \mathcal{A}$,	— Maintain a quaua
and prioritized queue PQ.	Maintain a queue
for $t = 1T$ do	
$s \leftarrow \text{current state}$ /* Reset when environment terminates */	
$a \sim \pi_{\epsilon\text{-greedy}}(a s)$ /* Sample action */	
$r, s' \sim p(r, s' s, a)$ /* Simulate environment */	
$p(s', r s, a) \leftarrow \text{Update}(s, a, r, s')$ /* Update model (Alg.1) */	
$\mathbf{p} \leftarrow r + \gamma \cdot \max_{a'} Q(s', a') - Q(s, a) $ /* Compute priority \mathbf{p} */	
Insert (s, a) into PQ with priority p /* State-action needs update */	
end	
/// Start sampling from PQ to perform updates	
repeat K times	
$s, a \leftarrow \text{pop highest priority from PQ}$ /* Sample PQ, break when empty */	
$s', r \sim \hat{p}(s', r s, a)$ /* Simulate model */	
$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [r + \gamma \cdot \max_{a'} Q(s',a') - Q(s,a)]$ /* Update Q-table */	
/// Loop over all state action that may lead to state s	
for all (\bar{s}, \bar{a}) with $\hat{p}(\bar{s}, \bar{a} s) > 0$ do	
$ar{r}=\hat{r}(ar{s},ar{a},s)$ /* Get reward from model */	
$\mathbf{p} \leftarrow ar{r} + \gamma \cdot \max_a \hat{Q}(s, a) - \hat{Q}(ar{s}, ar{a}) $ /* Compute priority \mathbf{p} */	
$\mathbf{if} \ \mathbf{p} \! > \! \theta \ \mathbf{then}$	
Insert (s, a) into PQ with priority p /* State-action needs update */	
\mathbf{end}	
\mathbf{end}	
\mathbf{end}	
end	

Algorithm 3: Prioritized sweeping (Q-learning	with ϵ -greedy exploration).	
Input: Number of planning updates K , explo	ration parameter $\epsilon \in (0, 1]$, learning rate	
$\alpha \in (0, 1]$, discount parameter $\gamma \in [0, 1]$], maximum number of timesteps T ,	
priority threshold θ .		
Initialization : Initialize $\hat{Q}(s, a) = 0$, $n(s, a, s)$	$(s) = 0, R_{sum}(s, a, s') = 0 \forall s \in \mathcal{S}, a \in \mathcal{A},$	Maintain a guava
and prioritized queue PQ.		Maintain a queue
for $t = 1T$ do		
$s \leftarrow \text{current state}$ /* I	Reset when environment terminates */	
$a \sim \pi_{\epsilon ext{-greedy}}(a s)$	/* Sample action */	Loom a forward and
$r, s' \sim p(r, s' s, a)$	/* Simulate environment */	
$\hat{p}(s', r s, a) \leftarrow \text{Update}(s, a, r, s')$	/* Update model (Aig.1) */	hackward model
$\mathbf{p} \leftarrow r + \gamma \cdot \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) $	/* Compute priority \mathbf{p} */	Dackwaru mouer
$\mathbf{if} \ \mathbf{p} > \theta \ \mathbf{then}$		
Insert (s, a) into PQ with priority p	<pre>/* State-action needs update */</pre>	
\mathbf{end}		
/// Start sampling from PQ to perfo	rm updates	
repeat K times		
$s, a \leftarrow \text{pop highest priority from PQ}$	/* Sample PQ , break when empty */	
$s', r \sim \hat{p}(s', r s, a)$	/* Simulate model */	
$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha \cdot [r + \gamma \cdot \max_{a'} \hat{Q}]$	$(s',a') - \hat{Q}(s,a)$] /* Update Q-table */	
/// Loop over all state action that	t may lead to state s	
for all (\bar{s}, \bar{a}) with $\hat{p}(\bar{s}, \bar{a} s) > 0$ do		
$ar{r} = \hat{r}(ar{s},ar{a},s)$	<pre>/* Get reward from model */</pre>	
$\mathbf{p} \leftarrow \bar{r} + \gamma \cdot \max_a \hat{Q}(s, a) - \hat{Q}(\bar{s}, \bar{a}) $	/* Compute priority p */	
if $\mathbf{p} > \theta$ then		
Insert (s, a) into PQ with priori	ty \mathbf{p} /* State-action needs update */	
end		









Algorithm 3: Prioritized sweeping (Q-learning with ϵ -greedy exploration).

Input: Number of planning updates K, exploration parameter $\epsilon \in (0, 1]$, learning rate $\alpha \in (0,1]$, discount parameter $\gamma \in [0,1]$, maximum number of timesteps T, priority threshold θ . **Initialization**: Initialize $\hat{Q}(s, a) = 0$, n(s, a, s') = 0, $R_{sum}(s, a, s') = 0$ $\forall s \in \mathcal{S}, a \in \mathcal{A}$, and prioritized queue PQ. for t = 1...T do $s \leftarrow \text{current state}$ /* Reset when environment terminates */ $a \sim \pi_{\epsilon\text{-greedy}}(a|s)$ /* Sample action */ $r, s' \sim p(r, s'|s, a)$ /* Simulate environment */ $\hat{p}(s', r|s, a) \leftarrow \text{Update}(s, a, r, s')$ /* Update model (Alg.1) */ $\mathbf{p} \leftarrow |r + \gamma \cdot \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)|$ /* Compute priority p */ if $p > \theta$ then Insert (s, a) into PQ with priority **p** /* State-action needs update */ end /// Start sampling from PQ to perform updates repeat K times $s, a \leftarrow \text{pop highest priority from PQ}$ /* Sample PQ, break when empty */ $s', r \sim \hat{p}(s', r|s, a)$ /* Simulate model */ $\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha \cdot [r + \gamma \cdot \max_{a'} \hat{Q}(s',a') - \hat{Q}(s,a)]$ /* Update Q-table */ /// Loop over all state action that may lead to state s for all (\bar{s}, \bar{a}) with $\hat{p}(\bar{s}, \bar{a}|s) > 0$ do $\bar{r} = \hat{r}(\bar{s}, \bar{a}, s)$ /* Get reward from model */ $\mathbf{p} \leftarrow |\bar{r} + \gamma \cdot \max_a \hat{Q}(s, a) - \hat{Q}(\bar{s}, \bar{a})|$ /* Compute priority p */ if $p > \theta$ then Insert (s, a) into PQ with priority **p** /* State-action needs update */ end end end end

You will also implement this algorithm for the assignment!



Summary 6: Model-based RL algorithms

Planning and learning can be combined in a variety of ways

Summary 6: Model-based RL algorithms

Planning and learning can be combined in a variety of ways

- E.g., <u>Dyna</u>
 - Learn forward model
 - Use model to sample additional transition data
 - Apply standard model-free RL update to simulated experience as well.

Summary 6: Model-based RL algorithms

Planning and learning can be combined in a variety of ways

- E.g., <u>Dyna</u>
 - Learn forward model
 - Use model to sample additional transition data
 - Apply standard model-free RL update to simulated experience as well.
- E.g., Prioritized sweeping
 - Learn a backward model
 - Use backward model to identify states that will likely change on the next update
 - Prioritize these states for updating in a separate queue

To Do



Read

- Sutton and Barto, Chapter 8
- Lecture notes
- Take care to study the back-up diagrams and associated equations!

To Do



Read

- Sutton and Barto, Chapter 8
- Lecture notes
- Take care to study the back-up diagrams and associated equations!

Assignment:

- 1. Implement two MBRL algorithms:
 - a. Dyna
 - b. Prioritized sweeping
- 2. Investigate their performance
- 3. Write a report



Questions?