



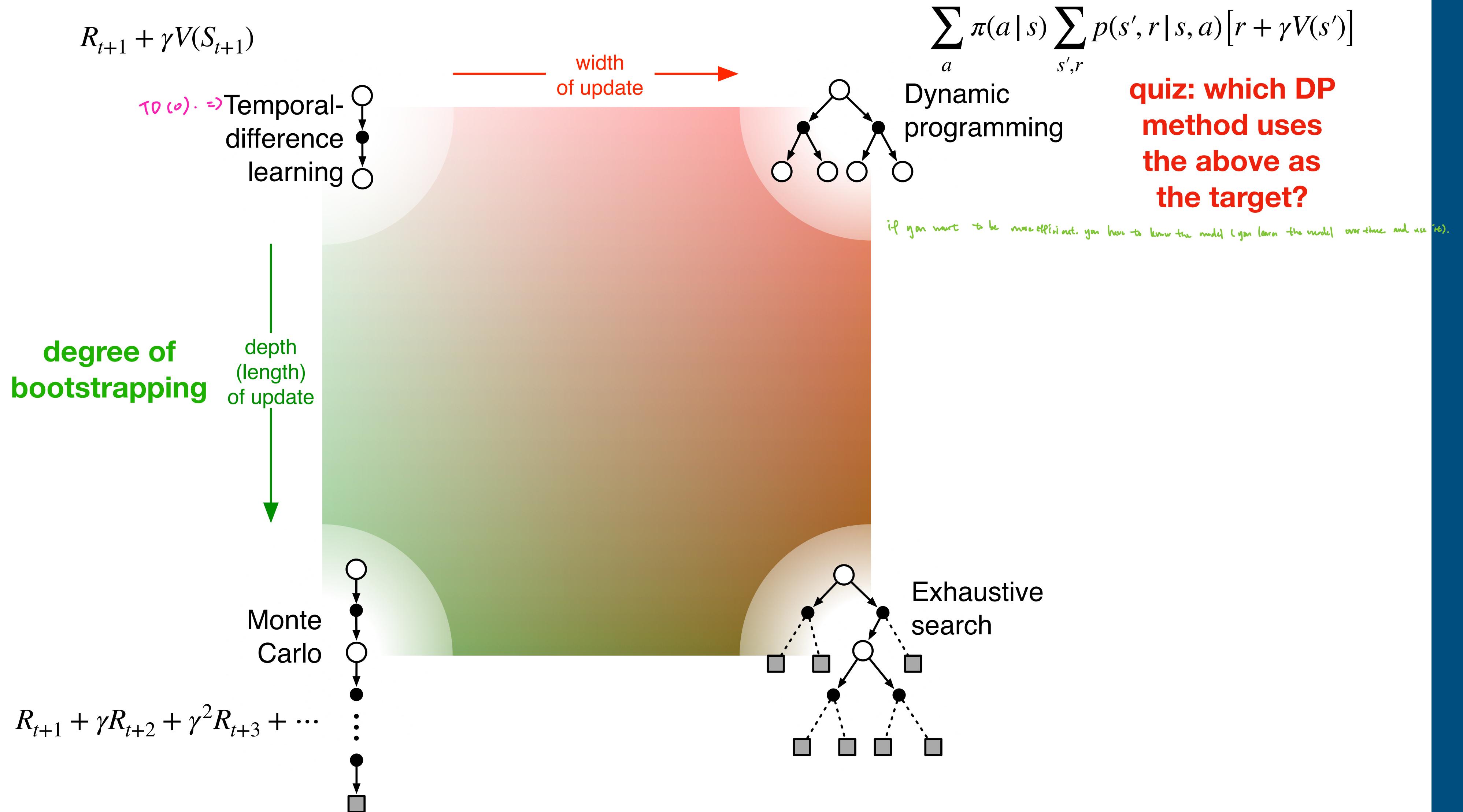
Planning, Learning and Acting

Rupam Mahmood

March 9, 2020



First: Summary of Part 1



Learning and planning

Model of the environment: used by agent to predict environment's response to actions

For example, if the agent has the transition probabilities $p(s',r|s,a)$ as in dynamic programming

Planning: a process that takes model as an input and produces/improves a policy

Planning produces a policy: decision time planning.

Dynamic programming methods are planning methods

A planning method uses the model to simulate the environment to produce simulated experience which is used to produce/improve a policy



Deepak Ranganatha Sastry Mamillapalli 下午 1:17
对所有人说

DR so all model-based learning uses planning? => Yes, model plays the role to plan

but if the model isn't accurate, wouldn't the simulated experience be different and biased from real experience? Yes.

Model learning: the process of using real experience to improve the model

→ 没办法直接用经验来更新模型，但可以估计它们。

For example, if we are estimating the transition probabilities based on real experience

Deepak Ranganatha Sastry Mamillapalli 下午 1:24
对所有人说

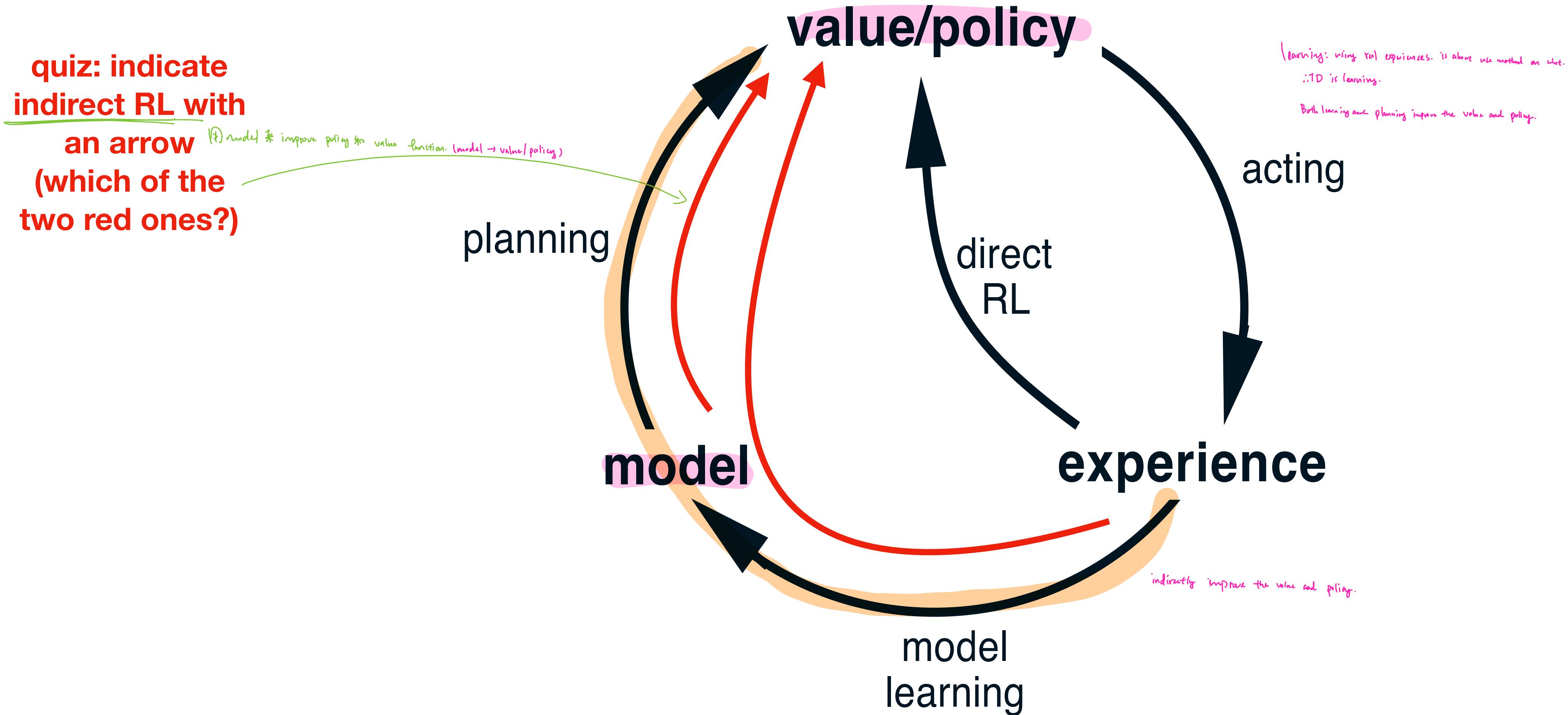
DR if we are making updates to both the model and value functions, then they should both converge to their true values right? => TD does not. 有更新值函数的梯度，但 TD 不一样。
TD (c) converge to the solution & TD 不一样。
me due. 因为它捕获到的奖励是不一样的。
会更快

Direct reinforcement learning: the process of using real experience to directly improve the value function and policy

For example, TD or MC methods for prediction and control

Indirect reinforcement learning: the process of using real experience to improve the value function and policy using the model through planning

An architecture for planning, learning and acting at the same time



Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ $\Rightarrow Q$ learning update
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment) \Rightarrow update the model after the Q learning update.
- (f) Loop repeat n times:

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

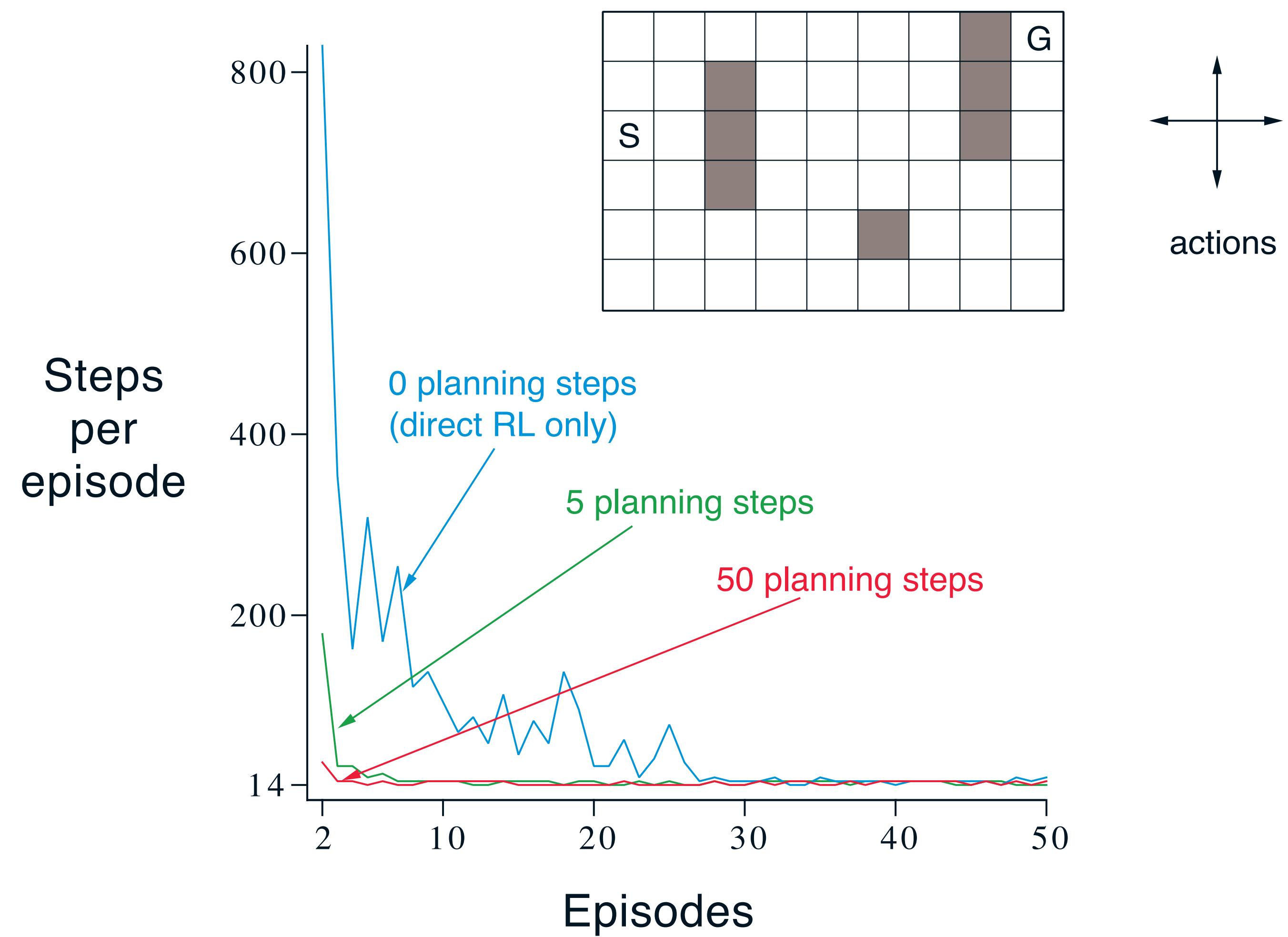
$R, S' \leftarrow Model(S, A)$ \Rightarrow sampling from the model.

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ \Rightarrow use the sampling to do another update.
(\Rightarrow after n loops sampling for update)

Which part is direct RL? What is this method called? (a) \sim (d): Q learning via real experience to update.

Which part is model learning? (e)

Which part is planning? (f)



How many planning steps should we take?

Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
 \Rightarrow storing the sample.
- (f) Loop repeat n times:

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

Dyna-Q+

1. Adds a bonus $\kappa\sqrt{\tau(s, a)}$ to reward in planning

$\tau(s, a)$ denotes the number of time steps (s, a) has not been tried

2. Actions that have not been tried from a previously visited state are allowed to be considered in planning

Where would you put these steps in Dyna-Q to get Dyna-Q+?

Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- (f) Loop repeat n times:
 - $S \leftarrow$ random previously observed state
 - $A \leftarrow$ random action previously taken in S
 - $R, S' \leftarrow Model(S, A)$
 - $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

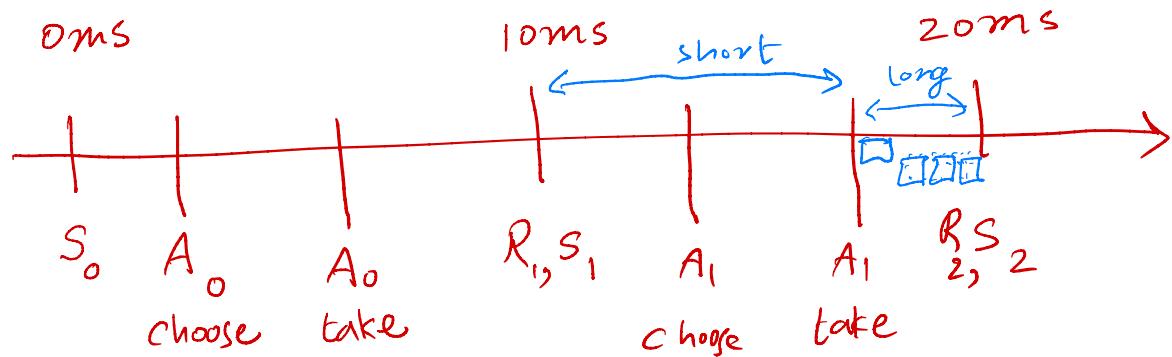
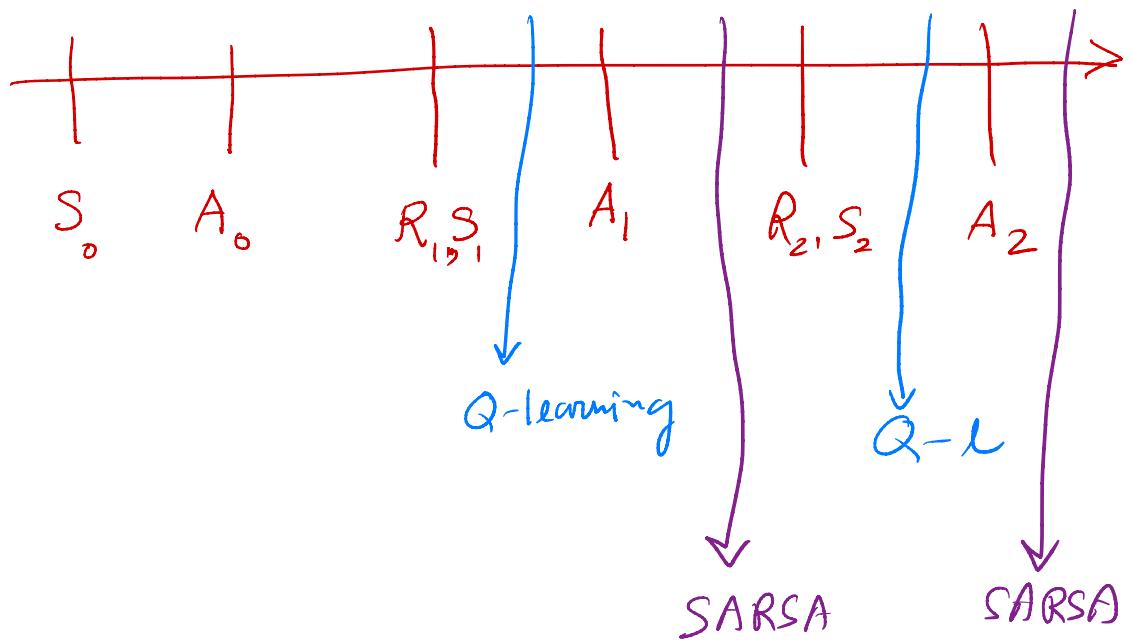
$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

Dyna-Q+: calculating visitation counts

Consider an MDP with one actions (L) and two states with the following episode

S_0	A_0	S_1	A_1	S_2	A_2	S_3	A_3	S_4	A_4
y	L	x	L	x	L	y	L	x	L

Calculate $\tau(s, a)$ for all state-action pairs at each step



- ① model is incorrect because environment is ~~perfect~~ not perfect.
(more and more test, more and more sample \rightarrow converge to true model. If it's ~~perfect~~ you are using the wrong model).

 2. model learn using function approximation that generalized imperfectly.
 3. suboptimal policy compared by planning \rightarrow allowing the modeling error.

If the model is optimistic, you can quickly discover the error. Although it's wrong.

② wrong model will lead you to discover the error in the model.

③ Dyna Q⁺: add bonus at update. (# of time step you try to make the transition : the bonus). $R + k\frac{1}{t}$
5 Dyna Q \Rightarrow $R - k\frac{1}{t}$ is a small value is added to reward. \Rightarrow ~~optimistic~~ - ~~greedy~~.
is Dyna Q with an exploration bonus that encourages exploration.
LL-optimistic \Rightarrow Dyna Q Mr. - \Rightarrow ~~greedy~~ to \Rightarrow ~~greedy~~.

Worksheet 9: Planning, Learning & Acting

CMPUT 397
March 19, 2021

6. (Exercise 8.2 S&B) Why did the Dyna agent with exploration bonus, Dyna-Q+, perform better in the first phase as well as in the second phase of the blocking experiment in Figure 8.4?

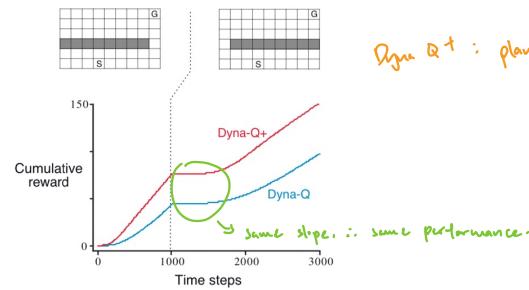


Figure 8.4: Average performance of Dyna agents on a blocking task. The left environment was used for the first 1000 steps, the right environment for the rest. Dyna-Q+ is Dyna-Q with an exploration bonus that encourages exploration. ■

As the first episode, (see step 10), Dyna Q+ agent has the exploration bonus, γ is 0.95, ϵ is 0.05, α is 0.1. It visits 50 pairs that are relatively unvisited \Rightarrow it maintains the Dyna Q+ Q-table goal state. To Dyna Q+ it takes a random policy, $\pi = 1/5$ (mean of the visitable 50 pairs).

Worksheet 9: Planning, Learning & Acting

CMPUT 397
March 19, 2021

7. (Exercise 8.3 S&B) **Challenge Question:** Careful inspection of Figure 8.5 reveals that the difference between Dyna-Q+ and Dyna-Q narrowed slightly over the first part of the experiment. What is the reason for this?

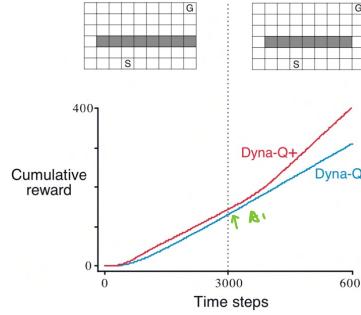


Figure 8.5: Average performance of Dyna agents on a shortcut task. The left environment was used for the first 3000 steps, the right environment for the rest. ■

Ans: ∵ Dyna Q+ keep explore. is it doesn't always follow the optimal path. ∴ It narrowed slightly in performance.

Worksheet 9: Planning, Learning & Acting

CMPUT 397
March 19, 2021

with ϵ -greedy target policy

3. Modify the Tabular Dyna-Q algorithm so that it uses Expected Sarsa instead of Q-learning. Assume that the target policy is ϵ -greedy. What should we call this algorithm?

Tabular Dyna-Q

```
Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ 
Loop forever:
  (a)  $S \leftarrow$  current (nonterminal) state
  (b)  $A \leftarrow \epsilon\text{-greedy}(S, Q)$ 
  (c) Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$ 
  (d)  $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
  (e)  $Model(S, A) \leftarrow R, S'$  (assuming deterministic environment)
  (f) Loop repeat  $n$  times:
     $S \leftarrow$  random previously observed state
     $A \leftarrow$  random action previously taken in  $S$ 
     $R, S' \leftarrow Model(S, A)$ 
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
```

Ans: ∵ Tabular Dyna Q ≈ Expected Sarsa:

∴ d) to f):

Tabular Dyna Q ≈ Expected Sarsa update: $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \sum_{a'} \pi(a'|S) Q(S', a) - Q(S, A)]$

With ϵ -greedy target policy ∴ Before each update, calculate π : $\pi(a|S) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(S)|} & \text{if } a = \text{argmax}_a Q(S, a) \\ \frac{\epsilon}{|\mathcal{A}(S)|}, & \text{otherwise} \end{cases}$

∴ Tabular Dyna - Expected Sarsa.

④ How many planning steps should we take?

As many as possible within the time / budget constraint.

Takes time away from actually interacting with the environment.

Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- $S \leftarrow$ current (nonterminal) state
- $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- Take action A ; observe resultant reward, R , and state, S'
- $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- Loop repeat n times:

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

①

