1.	Which approach ensures continual (never-ending) exploration? ( <b>Select all that apply</b> )	1 / 1 point
	<ul> <li>Exploring starts</li> <li>Correct</li> <li>Correct! Exploring starts guarantee that all state-action pairs are visited an infinite number of times in the limit of an infinite number of episodes.</li> </ul>	
	On-policy learning with a <b>deterministic</b> policy	
	$ ightharpoonup$ On-policy learning with an $\epsilon$ -soft policy	
	$\bigcirc$ Correct! $\epsilon$ -soft policies assign non-zero probabilities to all state-action pairs.	
	$lacksquare$ Off-Policy learning with an $\epsilon$ -soft behavior policy and a <b>deterministic</b> target policy	
	$\bigcirc$ Correct Correct! $\epsilon$ -soft policies have non-zero probabilities for all actions in all states. The behavior policy is used to generate samples and should be exploratory.	
	$lacksquare$ Off-Policy learning with an $\epsilon$ -soft target policy and a <b>deterministic</b> behavior policy	
2.	When can Monte Carlo methods, as defined in the course, be applied? (Select all that apply)	1/1 point
	When the problem is <b>continuing</b> and given a batch of data containing sequences of states, actions, and rewards	
	☐ When the problem is <b>continuing</b> and there is a model that produces samples of the next state and reward	
	When the problem is <b>episodic</b> and given a batch of data containing sample episodes (sequences of states, actions, and rewards)	

**⊘** Correct

Correct! Well-defined returns are available in episodic tasks.

- When the problem is **episodic** and there is a model that produces samples of the next state and reward
  - **⊘** Correct

Correct! Well-defined returns are available in episodic tasks.

**3.** Which of the following learning settings are examples of off-policy learning? (Select all that apply)

1/1 point

- Learning the optimal policy while continuing to explore
  - **⊘** Correct

Correct! An off-policy method with an exploratory behavior policy can assure continual exploration.

- Learning from data generated by a human expert
  - ✓ Correct

Correct! Applications of off-policy learning include learning from data generated by a non-learning agent or human expert. The policy that is being learned (the target policy) can be different from the human expert's policy (the behavior policy).

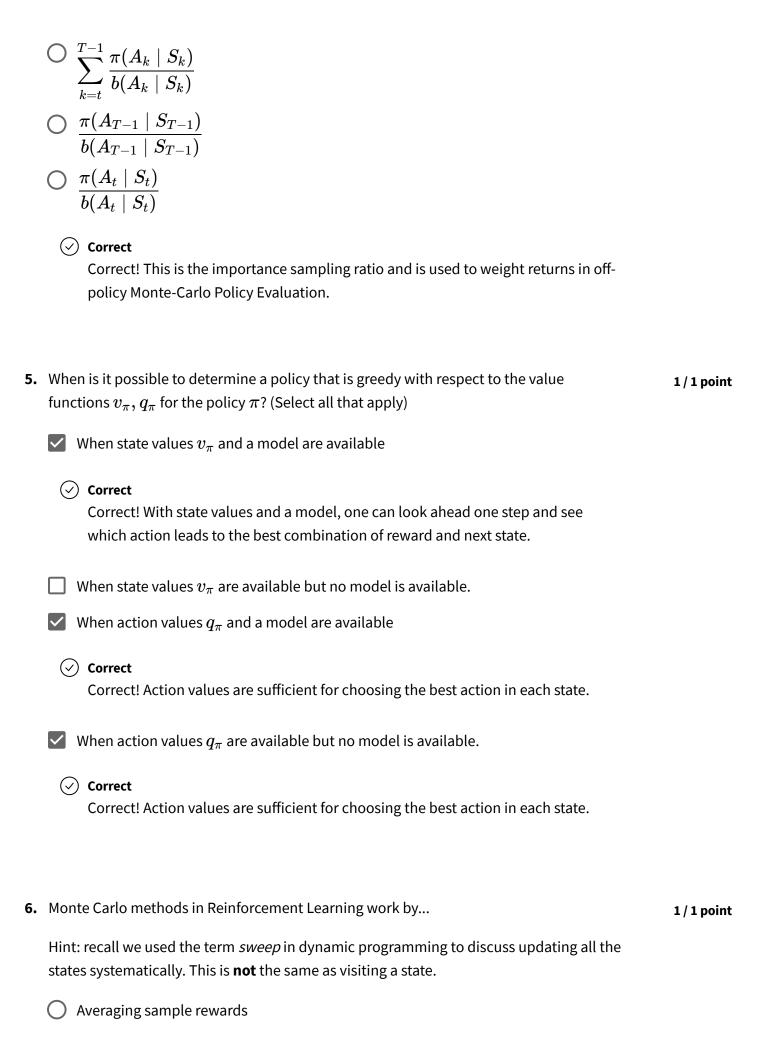
**4.** If a trajectory starts at time t and ends at time T, what is its relative probability under the target policy  $\pi$  and the behavior policy b?

1 / 1 point

Hint: pay attention to the time subscripts of A and S in the answers below.

Hint: Sums and products are not the same things!

 $igotimes \prod_{k=t}^{T-1} rac{\pi(A_k \mid S_k)}{b(A_k \mid S_k)}$ 



	Averaging sample returns	
	Performing <b>sweeps</b> through the state set	
	O Planning with a model of the environment	
	Correct  Correct! Monte Carlo methods in Reinforcement Learning sample and average returns much like bandit methods sample and average rewards.	
7.	Suppose the state $s$ has been visited three times, with corresponding returns $8,4$ , and $3$ . What is the current Monte Carlo estimate for the value of $s$ ?	1 / 1 point
	$\bigcirc$ 3	
	O 15	
	5	
	$\bigcirc$ 3.5	
	○ Correct     Correct! The Monte Carlo estimate for the state value is the average of sample returns observed from that state.	
8.	When does Monte Carlo prediction perform its first update?	1/1 point
	After the first time step	
	After every state is visited at least once	
	At the end of the first episode	
	Correct Correct! Monte Carlo Prediction updates value estimates at the end of an episode.	

**9.** In Monte Carlo prediction of state-values, **memory** requirements depend on (Select all

that apply).

1/1 point

Hint: think of the two data structures used in the algorithm	
The number of states	
<ul> <li>Correct! Monte Carlo Prediction needs to store the estimated value for each state.</li> </ul>	
☐ The number of possible actions in each state	
✓ The length of episodes	
<ul> <li>Correct         Correct! Monte Carlo Prediction needs to store the sequence of states and rewards. during an episode     </li> </ul>	
<b>10.</b> In an $\epsilon$ -greedy policy over ${\cal A}$ actions, what is the probability of the highest valued action if there are no other actions with the same value?	1 / 1 point
$\bigcirc \ 1-\epsilon$	
$\bigcirc$ $\epsilon$	

Correct! The highest valued action still has a chance of being selected as an

 $\bigcirc$   $1 - \epsilon + \frac{\epsilon}{A}$ 

**⊘** Correct

exploratory action.

 $\bigcirc \frac{\epsilon}{\mathcal{A}}$