

Autonomous and Mobile Robotics M

26 January 2023 - Theory

Some questions may have more than one correct answers: for each question, indicate all the correct answers.

1. A non-holonomic constraint:
 - ☒ cannot be fully integrated
 - ☐ can be written in the configuration space
 - ☒ does not restrict the space of configurations but the instant robot mobility
2. In a trajectory following problem, the robot must asymptotically perform a desired trajectory
 - ☒ that depends on a free parameter s
 - ☐ that depends on time t
 - ☐ that must be represented by a polynomial function
3. In map-based navigation, the robot:
 - ☒ plans the trajectory using a map of the environment
 - ☐ must update the planned path on the based of sensor information
 - ☐ navigates the environment on the base of the sensor information only
4. Sequential Monte-Carlo Localization:
 - ☒ is based on the simultaneous evaluation of multiple potential robot configurations called particles;
 - ☒ resamples the particles on the based of the weights evaluated after prediction, innovation and normalization;
 - ☐ works only in case the robot configuration can be described by a Gaussian distribution.
5. The future discounted reward is defined as:
 - ☐ $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots = \sum_{i=0}^{\infty} \gamma^i R_{t+i}$
 - ☒ $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{i=0}^{\infty} \gamma^i R_{t+i+1}$
 - ☐ $G_t = R_{t-1} + \gamma R_{t-2} + \gamma^2 R_{t-3} + \dots = \sum_{i=0}^{\infty} \gamma^i R_{t-i-1}$
6. The Bellman optimality equation for the action value function can be written as
 - ☐ $q_*(s, a) = \max v_{\pi}(s)$
 - ☒ $q_*(s, a) = \mathbb{E}_*[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a]$
 - ☐ $q_*(s, a) = p(s', r | s, a) [r + \gamma \max_{a' \in \mathcal{A}} q_*(s', a')]$
7. Monte-Carlo reinforcement learning:
 - ☐ can be applied to non-episodic tasks
 - ☐ requires the knowledge of the reward model
 - ☒ does not require the knowledge of the transition model
8. Under which hypotheses Monte-Carlo reinforcement learning methods converge to optimal value function?
 - ☒ infinite episodes
 - ☐ deterministic policy
 - ☒ exploring starts
9. In backward view Sarsa(λ), the eligibility trace update can be defined as
 - ☐ $E_t(s) = \gamma E_{t-1}(s) + \lambda(S_t = s)$
 - ☐ $E_t(s) = \gamma \lambda E_{t-1}(s) + 1(S_t = s)$
 - ☒ $E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + 1(S_t = s, A_t = a)$
10. TD(λ) with Value Function Approximation is defined as
 - ☐ $\Delta \mathbf{w} = \alpha(G_t - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
 - ☐ $\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$
 - ☒ $\Delta \mathbf{w} = \alpha(G_t^{\lambda} - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$

Autonomous and Mobile Robotics M

26 January 2023 - Exercise

The student is asked to solve the following problem.

Let us consider a fully observable and deterministic environment with 5 states $s_{\{1,\dots,5\}}$.

s_1	s_2	s_3	s_4	s_5
-------	-------	-------	-------	-------

- Action set : $\{\text{TryLeft}, \text{TryRight}\}$
- Rewards:
 - +1 in state s_1
 - -2 in state s_3
 - +4 in state s_5
 - 0 in all other states
- Transition model:
 - $p(s_1|s_1, \text{TryLeft}) = p(s_5|s_5, \text{TryRight}) = 1$
 - $p(s_1|s_1, \text{TryRight}) = p(s_2|s_1, \text{TryRight}) = 0.5$
 - $p(s_1|s_2, \text{TryLeft}) = p(s_2|s_2, \text{TryLeft}) = 0.5$
 - $p(s_2|s_2, \text{TryRight}) = p(s_3|s_2, \text{TryRight}) = 0.5$
 - ...
- Policy: $\pi(\text{TryLeft}|s_{\{1,\dots,5\}}) = \pi(\text{TryRight}|s_{\{1,\dots,5\}}) = 0.5$
- Discount factor $\gamma = 1$

Starting from an arbitrary initialisation of the state value function, compute the first iteration of the state value function evaluation provided by a Dynamic Programming algorithm with asynchronous backup assuming the random policy π .

$v_\pi(s_1)$	$v_\pi(s_2)$	$v_\pi(s_3)$	$v_\pi(s_4)$	$v_\pi(s_5)$

Solution:

The state value function is initialized to 0 for all the states.

$v_\pi(s_1)$	$v_\pi(s_2)$	$v_\pi(s_3)$	$v_\pi(s_4)$	$v_\pi(s_5)$
0.75	-0.0625	-1.0156	0.2461	3.0615