

Comprehensive Question Bank on the Theoretical Guarantees of Value-Based Reinforcement Learning

By Taha Majlesi

July 17, 2025

Abstract

This document contains a comprehensive set of questions based on the provided text, "Theoretical Guarantees of Value-Based Reinforcement Learning: From Bellman Equations to Algorithmic Convergence." It is divided into two parts: a set of 80 multiple-choice questions designed to test factual recall and conceptual understanding, and a set of 30 descriptive questions requiring detailed explanations, derivations, and comparisons. A complete answer key is provided for all questions.

Contents

I	Multiple-Choice Questions	2
1	Questions	2
II	Descriptive Questions	13
2	Questions	13
III	Answer Key	16
3	Answers to Multiple-Choice Questions	16
4	Answers to Descriptive Questions	16

Part I

Multiple-Choice Questions

1 Questions

1. What is the primary mathematical framework used to formalize sequential decision-making problems in reinforcement learning?
 - (a) **A) Markov Decision Process (MDP)**
 - (b) B) Bellman Equation System
 - (c) C) Banach Fixed-Point Theorem
 - (d) D) Generalized Policy Iteration (GPI)
2. In the MDP tuple $(\mathcal{S}, \mathcal{A}, P, \gamma)$, what does \mathcal{S} represent?
 - (a) A) The set of all possible actions.
 - (b) **B) A finite set of all possible states.**
 - (c) C) The probability of transitioning between states.
 - (d) D) The reward function.
3. The Markov Property states that the future is independent of the past given the:
 - (a) A) entire history of actions.
 - (b) B) initial state.
 - (c) **C) present state.**
 - (d) D) reward received at the previous step.
4. What is the main purpose of the discount factor, γ ?
 - (a) A) To increase the value of future rewards.
 - (b) B) To ensure the agent always chooses the action with the highest immediate reward.
 - (c) **C) To determine the present value of future rewards and ensure returns are finite.**
 - (d) D) To define the number of states in the MDP.
5. An agent that is "myopic" or short-sighted would have a discount factor γ close to:
 - (a) **A) 0**
 - (b) B) 1
 - (c) C) 0.5
 - (d) D) -1
6. What does a policy, π , define in reinforcement learning?
 - (a) A) The value of being in a state.
 - (b) B) The environment's dynamics.
 - (c) **C) The agent's strategy or blueprint for behavior.**
 - (d) D) The total accumulated reward.
7. A stochastic policy, $\pi(a|s)$, provides:

- (a) A) a single action for each state.
 - (b) B) a probability distribution over actions for each state.
 - (c) C) the expected return for each state.
 - (d) D) the transition probability to the next state.
8. The ultimate goal of a reinforcement learning agent is to find:
- (a) A) the value function with the smallest values.
 - (b) B) a policy that visits every state.
 - (c) C) the transition model of the environment.
 - (d) D) an optimal policy, π^* .
9. What does the state-value function, $v_\pi(s)$, quantify?
- (a) A) The immediate reward for being in state s .
 - (b) B) The long-term value of being in state s while following policy π .
 - (c) C) The probability of reaching state s .
 - (d) D) The best action to take in state s .
10. The action-value function, $q_\pi(s, a)$, is also commonly known as the:
- (a) A) Advantage function.
 - (b) B) State function.
 - (c) C) Q-function.
 - (d) D) Reward function.
11. Why is the action-value function (q_π) often more useful for control than the state-value function (v_π)?
- (a) A) It is easier to compute.
 - (b) B) It does not depend on the policy.
 - (c) C) It allows for direct comparison of different actions within a state.
 - (d) D) It only considers immediate rewards.
12. The Bellman expectation equation provides a recursive relationship for:
- (a) A) optimal value functions only.
 - (b) B) value functions under a specific policy π .
 - (c) C) the policy itself.
 - (d) D) the environment's transition probabilities.
13. The Bellman expectation equation for $v_\pi(s)$ decomposes the value of a state into the immediate reward and the:
- (a) A) value of the previous state.
 - (b) B) maximum possible reward in the MDP.
 - (c) C) discounted value of successor states.
 - (d) D) probability of the most likely action.
14. For a finite MDP, the Bellman expectation equation for v_π defines a system of:

- (a) A) non-linear equations.
 - (b) B) $|\mathcal{S}|$ linear equations in $|\mathcal{S}|$ variables.
 - (c) C) $|\mathcal{A}|$ quadratic equations.
 - (d) D) inequalities that cannot be solved directly.
15. The process of solving for the value function of a fixed policy is known as:
- (a) A) Policy Improvement.
 - (b) B) Value Iteration.
 - (c) C) Policy Evaluation.
 - (d) D) Bootstrapping.
16. The optimal state-value function, $v^*(s)$, is defined as the:
- (a) A) average of $v_\pi(s)$ over all possible policies.
 - (b) B) minimum of $v_\pi(s)$ over all possible policies.
 - (c) C) maximum of $v_\pi(s)$ over all possible policies.
 - (d) D) value function for a random policy.
17. What is the relationship between $v^*(s)$ and $q^*(s, a)$?
- (a) A) $v^*(s) = \sum_a q^*(s, a)$
 - (b) B) $v^*(s) = \max_a q^*(s, a)$
 - (c) C) $v^*(s) = \min_a q^*(s, a)$
 - (d) D) $v^*(s) = \mathbb{E}[q^*(s, a)]$
18. The Bellman optimality equation differs from the Bellman expectation equation due to the presence of the:
- (a) A) summation operator.
 - (b) B) discount factor.
 - (c) C) max operator.
 - (d) D) policy term $\pi(a|s)$.
19. The non-linearity of the Bellman optimality equation means it must be solved using:
- (a) A) direct matrix inversion.
 - (b) B) a single analytical step.
 - (c) C) iterative solution methods.
 - (d) D) linear programming.
20. The Bellman optimality operator, \mathcal{T}^* , maps a q-function to:
- (a) A) a policy.
 - (b) B) a scalar reward.
 - (c) C) a new q-function that is one step closer to optimal.
 - (d) D) the probability of the best action.
21. The optimal action-value function, q^* , is a \mathcal{T}^* of the Bellman optimality operator.

- (a) A) derivative
 - (b) B) integral
 - (c) C) fixed point
 - (d) D) starting point
22. What mathematical theorem is the cornerstone for proving the existence and uniqueness of an optimal value function?
- (a) A) The Central Limit Theorem
 - (b) B) The Banach Fixed-Point Theorem
 - (c) C) The Policy Improvement Theorem
 - (d) D) Bayes' Theorem
23. The standard metric used to measure the distance between two value functions is the:
- (a) A) L1 norm (Manhattan distance).
 - (b) B) L2 norm (Euclidean distance).
 - (c) C) L-infinity norm (maximum absolute difference).
 - (d) D) Cosine similarity.
24. A mapping \mathcal{T} is a contraction if it brings any two points in a metric space:
- (a) A) farther apart.
 - (b) B) to the same point in one step.
 - (c) C) closer together by at least a certain factor.
 - (d) D) to an orthogonal position.
25. The proof in the text shows that the Bellman optimality operator \mathcal{T}^* is a contraction with which contraction constant?
- (a) A) 1
 - (b) B) α
 - (c) C) γ
 - (d) D) $|\mathcal{S}|$
26. The fact that \mathcal{T}^* is a contraction mapping guarantees all of the following EXCEPT:
- (a) A) A unique optimal solution q^* exists.
 - (b) B) Iterative application of \mathcal{T}^* will converge to q^* .
 - (c) C) The convergence will happen in a single iteration.
 - (d) D) The convergence is guaranteed regardless of the starting q-function.
27. A higher discount factor γ (closer to 1) corresponds to:
- (a) A) a stronger contraction and faster convergence.
 - (b) B) a weaker contraction and slower guaranteed convergence.
 - (c) C) no change in the contraction rate.
 - (d) D) a non-contractive operator.
28. Dynamic Programming (DP) algorithms require what crucial piece of information?
- (a) A) A set of sample trajectories.
 - (b) B) An initial optimal policy.

- (c) C) A perfect model of the environment (P and R).
 - (d) D) A pre-computed value function.
29. The Value Iteration algorithm finds the optimal value function by iteratively applying the:
- (a) A) Bellman expectation equation.
 - (b) B) Bellman optimality equation as an update rule.
 - (c) C) Policy Improvement Theorem.
 - (d) D) L-infinity norm.
30. In Value Iteration, how is the optimal policy found?
- (a) A) It is iterated upon in each step of the algorithm.
 - (b) B) It is initialized randomly and never updated.
 - (c) C) It is extracted at the end, after the value function has converged.
 - (d) D) It is the average of all policies considered.
31. The update rule $V_{k+1} = \mathcal{T}^*V_k$ is the concise representation of which algorithm?
- (a) A) Policy Iteration
 - (b) B) Q-Learning
 - (c) C) Value Iteration
 - (d) D) Monte Carlo Evaluation
32. The Policy Iteration algorithm alternates between which two steps?
- (a) A) Initialization and Termination.
 - (b) B) Value Iteration and Policy Extraction.
 - (c) C) Policy Evaluation and Policy Improvement.
 - (d) D) Exploration and Exploitation.
33. The "Policy Evaluation" step in Policy Iteration involves:
- (a) A) finding the best action for each state.
 - (b) B) computing the value function for the current policy.
 - (c) C) updating the policy to be greedy.
 - (d) D) checking for convergence of the policy.
34. The "Policy Improvement" step in Policy Iteration involves:
- (a) A) solving a system of linear equations.
 - (b) B) calculating the L-infinity norm between two policies.
 - (c) C) forming a new, better policy by acting greedily with respect to the current value function.
 - (d) D) running more episodes to gather data.
35. What theorem guarantees that the greedy update in Policy Iteration leads to a better or equal policy?
- (a) A) The Banach Fixed-Point Theorem
 - (b) B) The Bellman Equation Theorem
 - (c) C) The Policy Improvement Theorem
 - (d) D) The Law of Large Numbers

36. Policy Iteration is guaranteed to converge for a finite MDP because the number of possible deterministic policies is:
- (a) A) infinite.
 - (b) B) dependent on the discount factor.
 - (c) C) finite.
 - (d) D) always equal to the number of states.
37. What is Generalized Policy Iteration (GPI)?
- (a) A) A specific algorithm that is faster than Value Iteration.
 - (b) B) A general template describing the interaction between policy evaluation and improvement.
 - (c) C) Another name for the Bellman optimality equation.
 - (d) D) A method for handling continuous action spaces.
38. How can Value Iteration be viewed within the GPI framework?
- (a) A) As a method with no policy evaluation.
 - (b) B) As a version where policy evaluation is truncated after one Bellman backup.
 - (c) C) As a method with no policy improvement.
 - (d) D) As a method that only works for deterministic policies.
39. What is the main trade-off between Value Iteration and Policy Iteration?
- (a) A) Simplicity vs. correctness.
 - (b) B) Memory usage vs. CPU usage.
 - (c) C) Number of iterations vs. computational cost per iteration.
 - (d) D) Applicability to episodic vs. continuing tasks.
40. Which algorithm typically has a higher computational cost per iteration?
- (a) A) Value Iteration
 - (b) B) Policy Iteration
 - (c) C) Both have the same cost per iteration.
 - (d) D) It depends on the discount factor.
41. Which algorithm often converges in fewer iterations?
- (a) A) Value Iteration
 - (b) B) Policy Iteration
 - (c) C) Both require the same number of iterations.
 - (d) D) It is impossible to say.
42. The complexity of one iteration of Value Iteration is typically:
- (a) A) $O(|\mathcal{S}|^3)$
 - (b) B) $O(|\mathcal{A}||\mathcal{S}|^2)$
 - (c) C) $O(|\mathcal{S}|)$
 - (d) D) $O(\log |\mathcal{S}|)$
43. The expensive step in Policy Iteration is policy evaluation, which can have a complexity of
- if solved via matrix inversion.*

44. A) $O(|\mathcal{S}|^2)$
45. B) $O(|\mathcal{A}||\mathcal{S}|)$
46. C) $O(|\mathcal{S}|^3)$
47. D) $O(|\mathcal{S}| \log |\mathcal{S}|)$

The return G_t for an episodic task is defined as:

1. A) $\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
2. B) $\sum_{k=1}^{T-t} R_{t+k}$
3. C) R_{t+1}
4. D) $\max_k R_{t+k}$

The probabilistic interpretation of the discount factor γ relates it to:

1. A) the probability of choosing a random action.
2. B) a constant probability of the episode continuing.
3. C) the probability of receiving a positive reward.
4. D) the certainty of the transition model.

The effective planning horizon of an agent is approximately:

1. A) γ
2. B) $1 - \gamma$
3. C) $1/(1 - \gamma)$
4. D) $\gamma/(1 - \gamma)$

The state s_t being a sufficient statistic of the history is a direct consequence of:

1. A) the Bellman equation.
2. B) the discount factor being less than 1.
3. C) the Markov Property.
4. D) having a finite action space.

The relationship $v_{\pi}(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s) q_{\pi}(s, a)$ shows that the state-value is the π .
of the action-values under policy

1. A) maximum
2. B) minimum
3. C) expected value
4. D) sum

Richard Bellman's Principle of Optimality is directly embodied in the:

1. A) definition of the return G_t .
2. B) policy evaluation step.
3. C) Bellman optimality equation.
4. D) L-infinity norm.

A complete metric space is one where every *sequence converges to a limit within the space.*

- A) monotonic
- B) arithmetic
- C) Cauchy
- D) geometric

The proof that \mathcal{T}^* is a contraction relies on a key property of which operator?

1. A) The summation operator
2. B) The maximum operator
3. C) The expectation operator
4. D) The argmax operator

The final step of the contraction proof uses the fact that $\sum_{s',r} p(s',r|s,a)$ is equal to:

1. A) 0
2. B) 1
3. C) γ
4. D) $|\mathcal{S}|$

The convergence of Value Iteration is a "beautiful consequence" of which theorem?

1. A) The Policy Improvement Theorem
2. B) The Law of Total Probability
3. C) The Banach Fixed-Point Theorem
4. D) The Central Limit Theorem

In Policy Iteration, if the policy does not change after the improvement step, it means:

1. A) the algorithm is stuck in a local minimum.
2. B) the discount factor is too low.
3. C) the current policy is already optimal.
4. D) the evaluation step was not run for long enough.

The total number of possible deterministic policies in a finite MDP is:

1. A) $|\mathcal{S}| \times |\mathcal{A}|$
2. B) $|\mathcal{S}| + |\mathcal{A}|$
3. C) $|\mathcal{A}|^{|\mathcal{S}|}$
4. D) $|\mathcal{S}|^{|\mathcal{A}|}$

Modern RL algorithms like Q-learning can be understood as forms of:

1. A) direct matrix inversion.
2. B) exhaustive search.
3. C) Generalized Policy Iteration (GPI).
4. D) supervised learning.

The term $p(s',r|s,a)$ represents the:

1. A) probability of being in state s and taking action a .
2. B) expected reward for taking action a .
3. C) joint probability of transitioning to state s' and receiving reward r , given state s and action a .
4. D) policy's probability of choosing action a .

An optimal policy π^* is one that achieves a expected return than any other policy from every state.

- A) lower or equal
- B) strictly lower
- C) higher or equal

D) strictly different

The Bellman optimality operator \mathcal{T}^* is applied to what kind of function?

1. A) A policy function $\pi(a|s)$
2. B) An action-value function $q(s, a)$
3. C) A transition function $p(s'|s, a)$
4. D) A reward function $r(s, a)$

The theoretical framework discussed in the text primarily applies to which type of MDPs?

1. A) Infinite and continuous
2. B) Partially observable
3. C) Finite
4. D) Non-stationary

The identity $v^*(s) = \max_a q^*(s, a)$ is crucial for deriving the Bellman optimality equation for:

1. A) $q^*(s, a)$
2. B) $v^*(s)$
3. C) $\pi^*(s)$
4. D) G_t

The "dance of evaluation and improvement" is a metaphor for:

1. A) Value Iteration
2. B) Policy Iteration / GPI
3. C) The Bellman equation
4. D) The Markov property

If $\gamma = 0$, the agent's objective is to maximize the:

1. A) total sum of all future rewards.
2. B) final reward at the end of the episode.
3. C) immediate reward R_{t+1} .
4. D) variance of the rewards.

The existence of at least one optimal deterministic policy is guaranteed for:

1. A) all reinforcement learning problems.
2. B) continuous state spaces only.
3. C) any finite MDP.
4. D) problems with a discount factor of 1.

The equation $q_\pi(s, a) = \sum_{s', r} p(s', r|s, a)[r + \gamma v_\pi(s')]$ expresses the action-value in terms of the:

1. A) action-value of the next state.
2. B) state-value of the next state.
3. C) current policy.
4. D) immediate reward only.

The core reason Value Iteration converges is that its update rule is a:

1. A) linear operator.
2. B) policy improvement step.
3. C) contraction mapping.
4. D) stochastic approximation.

In the context of the Banach Fixed-Point Theorem, the "space" is the set of all:

1. A) possible policies.
2. B) states and actions.
3. C) bounded real-valued functions over the state-action space.
4. D) sample trajectories.

The Policy Improvement Theorem requires that for the new policy π' to be better, $q_{\pi}(s, \pi'(s))$ must be $\geq v_{\pi}(s)$ for all states.

1. A) less than
2. B) strictly equal to
3. C) greater than or equal to
4. D) unrelated to

Which algorithm is generally simpler to implement?

1. A) Value Iteration
2. B) Policy Iteration
3. C) Both are equally complex.
4. D) Generalized Policy Iteration

The final conclusion of the text is that value-based RL methods are:

1. A) effective heuristics with no theoretical backing.
2. B) only applicable to game theory.
3. C) principled algorithms grounded in solid mathematical analysis.
4. D) less effective than evolutionary algorithms.

The state-transition probability $P(s'|s, a)$ is derived from the full dynamics $p(s', r|s, a)$ by:

1. A) maximizing over r .
2. B) taking the expectation over r .
3. C) summing over all possible rewards r .
4. D) ignoring the reward r .

The expected immediate reward $r(s, a)$ is an expectation over:

1. A) future policies.
2. B) all possible actions.
3. C) the next state s' and reward r .
4. D) the discount factor γ .

The Bellman expectation operator (for policy evaluation) is also a γ -contraction. This guarantees that:

1. A) the policy is optimal.
2. B) the iterative policy evaluation process converges to the true value function of that policy.

3. C) the policy will improve in the next step.
4. D) the MDP has a unique solution.

If Value Iteration is terminated early using a threshold ϵ , the resulting value function is:

1. A) guaranteed to be the exact optimal value function V^* .
2. B) completely random and useless.
3. C) an approximation of the optimal value function V^* .
4. D) the value function for a random policy.

The "unifying theme" that connects Value Iteration and Policy Iteration is:

1. A) The Markov Property.
2. B) The Banach Fixed-Point Theorem.
3. C) Generalized Policy Iteration (GPI).
4. D) The concept of return G_t .

The proof of convergence for Policy Iteration relies on monotonic improvement and:

1. A) the contraction mapping property.
2. B) the algorithm being simple to implement.
3. C) the finite number of possible deterministic policies.
4. D) the discount factor being close to 1.

The theoretical guarantees discussed in the text primarily concern the *of algorithms.*

- A) sample efficiency
- B) memory complexity
- C) convergence
- D) parallelizability

Part II

Descriptive Questions

2 Questions

1. **The Markov Property:** Explain the Markov Property in the context of an MDP. Why is this property a "direct causal link that enables the recursive structure of the Bellman equations"?
2. **The Role of the Discount Factor (γ):** Describe the three intertwined roles of the discount factor γ as presented in the text: mathematical necessity, behavioral preference, and probabilistic interpretation.
3. **Value Functions:** Compare and contrast the state-value function (v_π) and the action-value function (q_π). Why is the action-value function generally more critical for policy improvement?
4. **Bellman Expectation Equation Derivation:** Provide a step-by-step derivation of the Bellman expectation equation for the state-value function, $v_\pi(s)$. Explain the key step where the recursive nature becomes apparent.
5. **Bellman Optimality vs. Expectation:** What is the crucial difference in the mathematical form of the Bellman optimality equation for q^* compared to the Bellman expectation equation for q_π ? What is the implication of this difference for solving these equations?
6. **Principle of Optimality:** State Richard Bellman's Principle of Optimality and explain how it is mathematically captured in the Bellman optimality equation for $v^*(s)$.
7. **Fixed-Point Problem:** Explain what it means for the optimal value function q^* to be a "fixed point" of the Bellman optimality operator \mathcal{T}^* . Why is this reframing conceptually important?
8. **Contraction Mapping:** What is a contraction mapping? What two properties must be proven to apply the Banach Fixed-Point Theorem to the Bellman optimality operator?
9. **Contraction Proof Sketch:** Briefly outline the key steps used to prove that the Bellman optimality operator \mathcal{T}^* is a γ -contraction. What key inequality involving the 'max' operator is used?
10. **Banach Fixed-Point Theorem's Guarantees:** What are the two main guarantees provided by the Banach Fixed-Point Theorem once we have proven that \mathcal{T}^* is a contraction? Why are these guarantees so fundamental to value-based RL?
11. **Value Iteration Algorithm:** Describe the Value Iteration algorithm step-by-step. How does its main update rule directly implement the concept of applying the Bellman optimality operator?
12. **Convergence of Value Iteration:** Explain precisely why Value Iteration is guaranteed to converge to the optimal value function. Your answer should explicitly reference the concepts of contraction mappings and fixed points.
13. **Policy Iteration Algorithm:** Describe the two main, alternating steps of the Policy Iteration algorithm. What is computed in each step?

14. **Policy Evaluation Step:** In Policy Iteration, how is the policy evaluation step typically performed? What equation must be solved, and what does its solution represent?
15. **Policy Improvement Theorem:** State the Policy Improvement Theorem. How does this theorem guarantee that the policy improvement step in Policy Iteration always leads to a better (or equally good) policy?
16. **Convergence of Policy Iteration:** Explain why Policy Iteration is guaranteed to converge to an optimal policy in a finite number of iterations for a finite MDP.
17. **Generalized Policy Iteration (GPI):** What is the core idea behind Generalized Policy Iteration (GPI)? How do Value Iteration and Policy Iteration represent two different points on the spectrum of GPI?
18. **Computational Trade-off:** Compare Value Iteration and Policy Iteration in terms of their computational complexity per iteration and the typical number of iterations required for convergence.
19. **Choosing Between VI and PI:** Based on their trade-offs, describe a scenario where you would prefer to use Value Iteration and a scenario where you would prefer Policy Iteration.
20. **Policy Extraction:** In Value Iteration, the policy is extracted only after the value function has converged. Write down the mathematical formula for this policy extraction step and explain what it means in words.
21. **Linearity vs. Non-linearity:** Why is the Bellman expectation equation considered a system of linear equations, while the Bellman optimality equation is non-linear?
22. **The L-infinity Norm:** Define the L-infinity norm, $\|\cdot\|_\infty$, for value functions. Why is this specific metric suitable for the convergence proofs in RL?
23. **Intuition of the γ -Contraction:** Provide an intuitive explanation for why the Bellman optimality operator is a γ -contraction. If two value functions differ by at most ϵ , by how much will they differ after one application of the operator?
24. **Relationship between v^* and q^* :** Derive the Bellman optimality equation for $v^*(s)$ starting from the identity $v^*(s) = \max_a q^*(s, a)$.
25. **The Role of the Model:** The DP methods discussed (VI and PI) assume a "perfect model of the environment." What does this mean, and which components of the MDP tuple must be known?
26. **Termination Condition:** What is the practical termination condition for the Value Iteration algorithm? How does this relate to the L-infinity norm?
27. **The "Dance of Policy and Value":** Elaborate on the metaphor of the "dance of policy and value" used to describe GPI. What happens when this dance stabilizes?
28. **From q^* to π^* :** Once you have found the optimal action-value function q^* , how do you determine the optimal policy π^* ? Is this policy guaranteed to be deterministic?
29. **Proof of Policy Improvement Theorem:** Sketch the proof of the Policy Improvement Theorem. What is the key operation that is repeated to show $v_{\pi'}(s) \geq v_\pi(s)$?

30. **Unified Role of γ :** The text concludes by highlighting the unified role of the discount factor γ . Summarize how this single parameter connects the mathematical, behavioral, and algorithmic aspects of reinforcement learning.

Part III

Answer Key

3 Answers to Multiple-Choice Questions

- | | | | | | | | | | |
|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. B | 3. C | 4. C | 5. A | 6. C | 7. B | 8. D | 9. B | 10. C |
| 2. C | 12. B | 13. C | 14. B | 15. C | 16. C | 17. B | 18. C | 19. C | |
| | 20. C | | | | | | | | |
| 3. C | 22. B | 23. C | 24. C | 25. C | 26. C | 27. B | 28. C | 29. B | |
| | 30. C | | | | | | | | |
| 4. C | 32. C | 33. B | 34. C | 35. C | 36. C | 37. B | 38. B | 39. C | |
| | 40. B | | | | | | | | |
| 5. B | 42. B | 43. C | 44. C | 45. B | 46. C | 47. C | 48. C | 49. B | |
| | 50. C | | | | | | | | |
| 6. B | 52. C | 53. B | 54. C | 55. C | 56. C | 57. B | 58. C | 59. B | |
| | 60. C | | | | | | | | |
| 7. C | 62. B | 63. C | 64. C | 65. C | 66. B | 67. C | 68. C | 69. A | |
| | 70. C | | | | | | | | |
| 8. C | 72. B | 73. C | 74. B | 75. C | 76. C | 77. C | 78. C | 79. C | |
| | 80. C | | | | | | | | |

4 Answers to Descriptive Questions

1. **The Markov Property:** The Markov Property states that the future is independent of the past given the present. In an MDP, this means the probability of the next state S_{t+1} and reward R_{t+1} depends only on the current state S_t and current action A_t , not on any prior states, actions, or rewards. The state S_t is a "sufficient statistic" of the history. This property is the causal link because it allows the value of a state to be expressed recursively. The value of being in state s can be determined by the immediate rewards from s and the values of the possible next states s' , without needing to know how state s was reached. This enables the one-step lookahead structure of the Bellman equations.
2. **The Role of the Discount Factor (γ):**
 - **Mathematical Necessity:** For continuing (infinite-horizon) tasks, the sum of rewards could diverge to infinity. A discount factor $\gamma < 1$ ensures that the infinite geometric series of discounted rewards converges to a finite value, making the problem mathematically well-posed.
 - **Behavioral Preference:** It models the agent's preference for immediate versus future rewards. A γ near 0 creates a "myopic" agent that prioritizes short-term gains. A γ near 1 creates a "farsighted" agent that can learn complex, long-term strategies.
 - **Probabilistic Interpretation:** It can be interpreted as the probability of the process continuing. At each step, there is a $(1 - \gamma)$ probability of termination. This implies the agent's effective planning horizon is about $1/(1 - \gamma)$ steps.

3. Value Functions:

- The **state-value function**, $v_\pi(s)$, gives the expected return from starting in state s and following policy π . It answers, "How good is it to be in this state?"
- The **action-value function**, $q_\pi(s, a)$, gives the expected return from starting in state s , taking action a , and *then* following policy π . It answers, "How good is it to take this action in this state?"

The action-value function is more critical for policy improvement because to improve a policy, an agent must compare the consequences of different actions. $v_\pi(s)$ averages over all actions the policy might take, hiding this information. $q_\pi(s, a)$ provides the value for each specific action, allowing the agent to identify a better action than what the current policy prescribes.

4. Bellman Expectation Equation Derivation:

- Start with the definition: $v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$.
- Expand the return: $G_t = R_{t+1} + \gamma G_{t+1}$. So, $v_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]$.
- Use linearity of expectation: $v_\pi(s) = \mathbb{E}_\pi[R_{t+1} | S_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s]$.
- Expand the expectation over the policy's actions and the environment's dynamics: $v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']]$.
- Recursive Step:** Recognize that $\mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']$ is, by definition, the value of the successor state, $v_\pi(s')$.
- Substitute back to get the final form: $v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_\pi(s')]$.

5. Bellman Optimality vs. Expectation: The crucial difference is the presence of a **maximization operator** ('max') in the optimality equation.

- **Expectation for q_π :** $q_\pi(s, a) = \sum_{s', r} p(s', r|s, a) [r + \gamma \sum_{a'} \pi(a'|s') q_\pi(s', a')]$. It averages over the next actions according to policy π .
- **Optimality for q^* :** $q^*(s, a) = \sum_{s', r} p(s', r|s, a) [r + \gamma \max_{a'} q^*(s', a')]$. It takes the maximum value over the next actions.

Implication: The expectation equation is linear in the value functions. For a finite MDP, this defines a system of linear equations that can be solved directly. The 'max' operator makes the optimality equation non-linear, meaning it cannot be solved with a single matrix inversion and requires iterative methods like Value Iteration.

- Principle of Optimality:** The principle states: "An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision." This is captured in the Bellman optimality equation for $v^*(s)$ by the structure: $v^*(s) = \max_{a \in \mathcal{A}} \sum_{s', r} p(s', r|s, a) [r + \gamma v^*(s')]$. The equation breaks the problem into two parts: the initial decision ('max_a') and the value of the state resulting from that decision ($v^*(s')$), which is assumed to be the value from following an optimal policy thereafter.
- Fixed-Point Problem:** For q^* to be a "fixed point" of \mathcal{T}^* means that applying the operator to the function returns the exact same function:

$q^* = \mathcal{T}^*q^*$. This reframing is conceptually important because it allows us to move from the specific domain of reinforcement learning to the general and powerful mathematical field of fixed-point theory. By showing the operator has certain properties (i.e., it's a contraction), we can use established theorems (like the Banach Fixed-Point Theorem) to prove that a unique solution exists and that an iterative algorithm will find it.

8. **Contraction Mapping:** A contraction mapping is a function that, when applied to any two points in a metric space, is guaranteed to bring those points closer together by at least a certain factor $\alpha \in [0, 1)$. To apply the Banach Fixed-Point Theorem, we must prove:

- (a) The space of value functions is a **non-empty complete metric space**.
- (b) The Bellman optimality operator \mathcal{T}^* is a **contraction mapping** on that space.

9. **Contraction Proof Sketch:**

- (a) Start with the L-infinity distance between the operator applied to two q-functions: $\|\mathcal{T}^*q_1 - \mathcal{T}^*q_2\|_\infty$.
- (b) Substitute the definition of the operator. The reward terms r cancel out.
- (c) Factor out the discount factor γ and use the triangle inequality.
- (d) Use the key inequality for the maximum operator: $|\max_x f(x) - \max_x g(x)| \leq \max_x |f(x) - g(x)|$.
- (e) Recognize that $\max_{a'} |q_1(s', a') - q_2(s', a')|$ is less than or equal to the L-infinity norm over all states, $\|q_1 - q_2\|_\infty$.
- (f) Factor out the constant norm term. The remaining sum of probabilities $\sum p(s', r|s, a)$ equals 1, leaving the final result: $\|\mathcal{T}^*q_1 - \mathcal{T}^*q_2\|_\infty \leq \gamma \|q_1 - q_2\|_\infty$.

10. **Banach Fixed-Point Theorem's Guarantees:**

- (a) **Existence and Uniqueness:** It guarantees that there exists one and only one fixed point, which in our case is the unique optimal value function q^* .
- (b) **Convergence:** It guarantees that the sequence generated by iteratively applying the operator from any arbitrary starting point will converge to this unique fixed point.

These guarantees are fundamental because they assure us that the problem we are trying to solve has a single, well-defined answer, and that the iterative algorithms we design to find it are not just heuristics but are provably correct and will eventually succeed.

11. **Value Iteration Algorithm:**

- (a) **Initialization:** Start with an arbitrary value function V_0 (e.g., all zeros).
- (b) **Iteration:** For each step k , compute the next value function V_{k+1} for all states s by applying the Bellman optimality backup: $V_{k+1}(s) \leftarrow \max_a \sum_{s', r} p(s', r|s, a)[r + \gamma V_k(s')]$.
- (c) **Termination:** Stop when the change in the value function is small, i.e., $\|V_{k+1} - V_k\|_\infty < \epsilon$.
- (d) **Policy Extraction:** Extract the optimal policy by acting greedily with respect to the final value function V^* .

The update rule is a direct implementation because it takes the current value function (V_k) and applies the one-step lookahead and maximization (' \max_a ') defined by the Bellman optimality operator to produce the new value function (V_{k+1}).

12. **Convergence of Value Iteration:** Value Iteration is guaranteed to converge because its update rule, $V_{k+1} = \mathcal{T}^*V_k$, is a repeated application of the Bellman optimality operator. As proven, \mathcal{T}^* is a γ -contraction mapping on the space of value functions. The Banach Fixed-Point Theorem states that for any contraction mapping, iteratively applying the operator from any starting point is guaranteed to converge to the unique fixed point of that operator. Therefore, the sequence of value functions $\{V_k\}$ must converge to the unique optimal value function V^* .
13. **Policy Iteration Algorithm:**
 - (a) **Policy Evaluation:** Given the current policy π_k , compute its state-value function v_{π_k} . This means finding the value function that satisfies the Bellman expectation equation for π_k .
 - (b) **Policy Improvement:** Using the computed value function v_{π_k} , form a new, improved policy π_{k+1} by acting greedily with respect to v_{π_k} for every state.

These two steps are repeated until the policy no longer changes.

14. **Policy Evaluation Step:** The policy evaluation step is performed by solving the system of linear equations defined by the Bellman expectation equation for the current policy π_k : $v(s) = \sum_a \pi_k(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v(s')]$. This can be solved directly (matrix inversion) or, more commonly, iteratively by repeatedly applying the equation as an update rule until the value function converges. The solution, v_{π_k} , represents the true expected long-term return for following policy π_k from any state.
15. **Policy Improvement Theorem:** The theorem states that if we have two deterministic policies, π and π' , such that for all states s , the value of acting according to the new policy for one step is better than the value of the old policy ($q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$), then the new policy π' must be as good as, or better than, the old policy π everywhere ($v_{\pi'}(s) \geq v_{\pi}(s)$). This guarantees improvement because the greedy update in Policy Iteration constructs the new policy π_{k+1} to explicitly satisfy the condition $q_{\pi_k}(s, \pi_{k+1}(s)) \geq v_{\pi_k}(s)$, thus ensuring the new policy is an improvement.
16. **Convergence of Policy Iteration:** Convergence is guaranteed for two reasons:
 - (a) The Policy Improvement Theorem guarantees that each new policy is strictly better than the previous one, unless the policy is already optimal. This means the value function is monotonically increasing, and the algorithm cannot visit the same policy twice.
 - (b) For a finite MDP, there is a finite number of possible deterministic policies ($|\mathcal{A}|^{|\mathcal{S}|}$).

Since each step generates a strictly better policy from a finite set, the process must terminate in a finite number of steps. It terminates when the policy no longer improves, which means it has reached the optimal policy.

17. **Generalized Policy Iteration (GPI):** GPI is a general template or idea, not a specific algorithm. It describes the fundamental interaction

in RL where two processes, policy evaluation and policy improvement, compete and cooperate. The policy is improved with respect to the value function, and the value function is driven towards the value of the current policy.

- **Policy Iteration** performs a full, complete policy evaluation before doing any improvement.
- **Value Iteration** performs only a single step of evaluation (one Bellman backup) within its improvement loop. It represents an extreme point on the GPI spectrum where evaluation is heavily truncated.

18. **Computational Trade-off:**

- **Value Iteration (VI):** Has a low computational cost per iteration ($O(|\mathcal{A}||\mathcal{S}|^2)$), but may require many iterations to converge to the precise optimal value function.
- **Policy Iteration (PI):** Has a very high computational cost per iteration, as the policy evaluation step requires solving a full system of linear equations (e.g., $O(|\mathcal{S}|^3)$ or many iterative sweeps). However, it often converges in a surprisingly small number of iterations.

The trade-off is (many, cheap iterations) vs. (few, expensive iterations).

19. **Choosing Between VI and PI:**

- **Prefer Value Iteration** when the state space $|\mathcal{S}|$ is very large, making the policy evaluation step of PI computationally prohibitive. Its simpler implementation is also an advantage.
- **Prefer Policy Iteration** when the number of policies is not excessively large and the policy evaluation step is manageable. It often converges much faster in terms of total wall-clock time because it requires so few iterations.

20. **Policy Extraction:** The formula is: $\pi^*(s) \leftarrow \arg \max_{a \in \mathcal{A}} \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')]$. In words, this means: "For each state s , look one step ahead for every possible action a . Calculate the expected return for taking that action, which is the immediate reward plus the discounted value of the resulting state (using the converged optimal value function V^*). The optimal policy is to choose the action a that maximizes this one-step lookahead value."

21. **Linearity vs. Non-linearity:**

- The **Bellman expectation equation** is linear because the value function terms ($v_\pi(s)$ and $v_\pi(s')$) appear as linear terms. The equation is a weighted average, and there are no products or non-linear functions (like max, log, square) of the value function variables.
- The **Bellman optimality equation** is non-linear because of the 'max' operator. The maximum of a set of variables is a non-linear function, preventing the system from being solved with standard linear algebra.

22. **The L-infinity Norm:** The L-infinity norm distance between two value functions q_1 and q_2 is defined as: $\|q_1 - q_2\|_\infty = \max_{s, a} |q_1(s, a) - q_2(s, a)|$. It is the largest absolute difference between the functions at any single state-action pair. This metric is suitable because the Bellman backup is a "worst-case" update over all states. The L-infinity norm captures this worst-case error, and showing that this maximum error shrinks at each step is a very strong guarantee of convergence for the entire function.

23. **Intuition of the γ -Contraction:** The Bellman operator performs a one-step backup, updating a state's value based on the values of its successors. The key is that the successor values are discounted by γ . If two value functions q_1 and q_2 have a maximum difference of ϵ , when we back up their values, this difference ϵ is multiplied by $\gamma < 1$. Therefore, the new maximum difference between the updated value functions, \mathcal{T}^*q_1 and \mathcal{T}^*q_2 , can be at most $\gamma\epsilon$. Each iteration "shrinks" the distance between any two value functions.
24. **Relationship between v^* and q^* :**
- (a) Start with the identity: $v^*(s) = \max_a q^*(s, a)$.
 - (b) The optimal action-value $q^*(s, a)$ is the expected return for taking action a and then following the optimal policy. The value of the next state s' will therefore be $v^*(s')$. So, $q^*(s, a) = \sum_{s', r} p(s', r|s, a)[r + \gamma v^*(s')]$.
 - (c) Substitute this expression for $q^*(s, a)$ back into the identity: $v^*(s) = \max_a \left\{ \sum_{s', r} p(s', r|s, a)[r + \gamma v^*(s')] \right\}$. This is the Bellman optimality equation for v^* .
25. **The Role of the Model:** A "perfect model of the environment" means that the agent has full knowledge of the environment's dynamics. Specifically, it must know the transition dynamics function, P , and the reward function, r . In the tuple $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$, the functions $p(s', r|s, a)$ (or the derived $P(s'|s, a)$ and $r(s, a)$) must be known to perform the summations and expectations in the Bellman updates.
26. **Termination Condition:** The practical termination condition for Value Iteration is when the maximum change across all state values in an iteration is less than some small positive threshold ϵ . This is written as: $\|V_{k+1} - V_k\|_\infty < \epsilon$. This directly uses the L-infinity norm, as it checks the maximum absolute difference between the value function vectors from one iteration to the next.
27. **The "Dance of Policy and Value":** This metaphor describes the interplay in GPI. The "policy" leads the dance by becoming greedy with respect to the current "value function." Then, the "value function" follows, being updated to be consistent with the new policy. They take turns leading. This dance stabilizes when they are in perfect harmony: the policy is greedy with respect to the value function, and the value function is the true value function for that policy. At this point, the Bellman optimality equation is satisfied, and the optimal solution has been found.
28. **From q^* to π^* :** Once q^* is known, the optimal policy π^* is found by simply acting greedily at every state: $\pi^*(s) = \arg \max_{a \in \mathcal{A}} q^*(s, a)$. For any finite MDP, there is always at least one such optimal policy, and it is guaranteed that at least one of them is deterministic. This greedy extraction will find one such deterministic optimal policy.
29. **Proof of Policy Improvement Theorem:** The proof sketch involves unrolling the value function definitions. Start with the inequality from the theorem's premise: $v_\pi(s) \leq q_\pi(s, \pi'(s))$. Expand $q_\pi(s, \pi'(s))$: $v_\pi(s) \leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$. The key operation is to repeatedly apply this inequality. Since $v_\pi(S_{t+1}) \leq q_\pi(S_{t+1}, \pi'(S_{t+1}))$, we can substitute this inside the expectation. By repeatedly unrolling the value function

and applying the inequality at each step, we show that the value under the old policy π is a lower bound that is pushed up at each time step, ultimately showing that $v_\pi(s) \leq v_{\pi'}(s)$.

30. **Unified Role of γ :** The discount factor γ elegantly unifies three core aspects of RL:
- (a) **Mathematical:** As a tool with $\gamma < 1$, it ensures the infinite sum of rewards in continuing tasks converges, making the problem well-posed.
 - (b) **Behavioral:** As a modeling parameter, it defines the agent's character. A low γ creates a myopic agent, while a high γ creates a farsighted one, controlling its planning horizon.
 - (c) **Algorithmic:** It emerges as the contraction constant of the Bellman optimality operator. This directly ties the agent's behavioral farsightedness to the guaranteed rate of convergence of algorithms like Value Iteration. A higher γ (more farsighted) means a weaker contraction and slower guaranteed convergence.