

Comprehensive Question Bank on the Theory of Value-Based Reinforcement Learning

By Taha Majlesi

July 17, 2025

Contents

I	Multiple-Choice Questions (80 Questions)	1
1	Part I: Foundations of Reinforcement Learning	1
2	Part II: The Bellman Equations	5
3	Part III: Fixed-Point Theory	8
4	Part IV: Synthesis and Proofs	12
5	Part V: Conclusion and Implications	14
II	Explanatory Questions (30 Questions)	16
	Answers to Multiple-Choice Questions	19
	Answers to Explanatory Questions	22

Part I

Multiple-Choice Questions (80 Questions)

1 Part I: Foundations of Reinforcement Learning

1. What is the primary goal of a Reinforcement Learning agent?
 - a. To classify data accurately.
 - b. To predict a static value from a dataset.
 - c. To learn a policy that maximizes a cumulative reward signal.
 - d. To minimize the interaction time with the environment.
2. A value function in RL serves to:
 - a. Define the agent's possible actions.
 - b. Quantify the long-term desirability of states or actions.
 - c. Store the history of all rewards received.
 - d. Model the environment's transition probabilities.
3. Which of the following is NOT one of the three fundamental theoretical questions about the optimal value function mentioned in the introduction?
 - a. Existence
 - b. Uniqueness
 - c. Stability
 - d. Computability
4. The theoretical framework for guaranteeing the properties of the optimal value function relies on:
 - a. Supervised Learning Theory
 - b. The Central Limit Theorem
 - c. The Banach Fixed-Point Theorem

- d. Bayesian Inference
5. What does the symbol \mathcal{S} represent in an MDP?
- a. The set of all possible strategies.
 - b. The set of all possible states.
 - c. The set of all possible scalar rewards.
 - d. The set of all possible successor states.
6. The Markov Property states that the future is:
- a. Dependent on the entire history of states and actions.
 - b. Completely random and unpredictable.
 - c. Independent of the past given the present state.
 - d. Determined solely by the agent's policy.
7. The transition probability function $p(s', r|s, a)$ defines:
- a. The probability of being in state s' and receiving reward r after taking action a .
 - b. The probability of transitioning to state s' and receiving reward r , given the current state s and action a .
 - c. The reward for transitioning from state s to s' .
 - d. The policy's probability of choosing action a .
8. What is the valid range for the discount factor γ ?
- a. $\gamma > 1$
 - b. $\gamma \in [0, 1]$
 - c. $\gamma \in [0, 1)$
 - d. $\gamma \in (-1, 1)$
9. A discount factor γ close to 0 makes the agent:
- a. "Farsighted," valuing future rewards highly.
 - b. "Myopic," focusing on immediate rewards.
 - c. Risk-averse.
 - d. Indifferent to the timing of rewards.

10. Besides modeling a preference for immediate rewards, what is another interpretation of the discount factor γ ?
 - a. The learning rate of the agent.
 - b. The probability that the agent "survives" to the next time step.
 - c. The complexity of the environment.
 - d. The number of states in the MDP.
11. A deterministic policy $\pi(s)$ maps:
 - a. Each state to a probability distribution over actions.
 - b. Each state to a single, specific action.
 - c. Each state-action pair to a reward.
 - d. Each state to a value.
12. The return, G_t , is defined as:
 - a. The immediate reward R_{t+1} .
 - b. The average of all future rewards.
 - c. The sum of discounted future rewards.
 - d. The maximum possible reward in the episode.
13. The state-value function $V^\pi(s)$ represents:
 - a. The immediate reward for being in state s .
 - b. The probability of reaching state s .
 - c. The expected return when starting in state s and following policy π .
 - d. The best possible return from state s .
14. The action-value function $q^\pi(s, a)$ answers the question:
 - a. "What is the best action in state s ?"
 - b. "How good is it to be in state s ?"
 - c. "What is the probability of taking action a in state s ?"
 - d. "How good is it to take action a in state s , and then commit to policy π ?"
15. How is $V^\pi(s)$ related to $q^\pi(s, a)$?

- a. $V^\pi(s) = \max_a q^\pi(s, a)$
 - b. $V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q^\pi(s, a)$
 - c. $V^\pi(s) = \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} q^\pi(s, a)$
 - d. $V^\pi(s) = q^\pi(s, \pi(s))$
16. The 5-tuple defining an MDP is:
- a. $(\mathcal{S}, \mathcal{A}, \pi, R, \gamma)$
 - b. $(\mathcal{S}, \mathcal{A}, P, R, G_t)$
 - c. $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$
 - d. $(\mathcal{S}, V, q, \pi, \gamma)$
17. The Markov Property is crucial because it:
- a. Guarantees rewards are always positive.
 - b. Ensures the number of states is finite.
 - c. Allows future predictions based only on the current state, not the full history.
 - d. Makes all policies deterministic.
18. The sum $\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r|s, a)$ must equal:
- a. 0
 - b. 1
 - c. γ
 - d. $|\mathcal{S}|$
19. An agent with $\gamma = 1$ in an infinite-horizon problem might face what issue?
- a. The value functions could diverge to infinity.
 - b. The agent would only care about the next reward.
 - c. The policy would become random.
 - d. The transition probabilities would become invalid.
20. Learning an accurate action-value function allows an agent to:
- a. Formulate an effective policy without a model of the environment.
 - b. Predict the exact sequence of future states.
 - c. Guarantee a positive reward at every step.
 - d. Eliminate all uncertainty from the problem.

2 Part II: The Bellman Equations

21. The Bellman equations provide what kind of structure for value functions?
 - a. A linear structure.
 - b. A recursive structure.
 - c. A chaotic structure.
 - d. A historical structure.
22. The Bellman expectation equation is used for:
 - a. Policy improvement.
 - b. Policy evaluation.
 - c. Finding the optimal policy directly.
 - d. Initializing the value function.
23. The Bellman expectation equation for $V^\pi(s)$ expresses the value of a state as:
 - a. The maximum possible value of its successor states.
 - b. The immediate reward plus the discounted value of the most likely successor state.
 - c. A weighted average of the values of its potential successor states.
 - d. The sum of all possible future rewards.
24. For a finite MDP, the Bellman expectation equation for V^π gives a system of:
 - a. $|\mathcal{S}|$ non-linear equations.
 - b. $|\mathcal{S}| \times |\mathcal{A}|$ linear equations.
 - c. $|\mathcal{S}|$ linear equations in $|\mathcal{S}|$ variables.
 - d. $|\mathcal{A}|$ non-linear equations.
25. Which algorithms are based on the Bellman expectation equation for q^π ?
 - a. Value Iteration and Policy Iteration.
 - b. Gradient Descent and Newton's Method.
 - c. K-Means and PCA.
 - d. Temporal-Difference (TD) learning and SARSA.

26. The ultimate goal of reinforcement learning, moving from evaluation to control, requires using the:
 - a. Bellman expectation equation.
 - b. Bellman optimality equation.
 - c. Markov chain equation.
 - d. Return calculation formula.
27. An optimal policy π^* is defined as a policy where:
 - a. $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all $s \in \mathcal{S}$ and all policies π .
 - b. It yields the highest immediate reward from every state.
 - c. It is deterministic.
 - d. It explores every state-action pair.
28. The optimal state-value function, $V^*(s)$, is defined as:
 - a. $\mathbb{E}[V^{\pi}(s)]$ over all policies π .
 - b. $\arg \max_{\pi} V^{\pi}(s)$.
 - c. $\max_{\pi} V^{\pi}(s)$.
 - d. The value function of a random policy.
29. The Principle of Optimality states that an optimal policy's remaining decisions must:
 - a. Revisit the initial state.
 - b. Constitute an optimal policy with regard to the state resulting from the first decision.
 - c. Be random to ensure exploration.
 - d. Follow the same action as the first decision.
30. The key difference between the Bellman expectation and optimality equations is the replacement of a policy-weighted average with a:
 - a. Minimization operator.
 - b. Maximization operator.
 - c. Summation over all policies.

- d. Direct matrix inversion.
31. The 'max' operator in the Bellman optimality equation implicitly defines:
- a. The discount factor.
 - b. The set of states.
 - c. The optimal policy itself.
 - d. The reward function.
32. The Bellman optimality equation for q^* turns the problem of finding q^* into:
- a. Solving a system of linear equations.
 - b. Solving a system of non-linear equations.
 - c. A simple table lookup problem.
 - d. A supervised learning problem.
33. The Bellman expectation equation is considered "passive" because:
- a. It does not require computation.
 - b. It describes the value of a fixed, given strategy.
 - c. It only works for environments with no rewards.
 - d. It cannot be solved iteratively.
34. The Bellman optimality equation is considered "active" because:
- a. It defines the value of the best possible strategy.
 - b. It requires an agent to be physically moving.
 - c. It changes the environment's dynamics.
 - d. It is only applicable to episodic tasks.
35. If you know $q^*(s, a)$, how do you find the optimal action in state s ?
- a. Choose an action randomly.
 - b. Choose the action a that minimizes $q^*(s, a)$.
 - c. Choose the action a that maximizes $q^*(s, a)$.
 - d. Average the values of $q^*(s, a)$ over all a .

3 Part III: Fixed-Point Theory

36. A fixed point x of an operator \mathcal{T} satisfies which equation?
- a. $\mathcal{T}(x) = 0$
 - b. $\mathcal{T}(x) = x$
 - c. $\mathcal{T}(x) = 1$
 - d. $\mathcal{T}(x) = \mathcal{T}(0)$
37. In a dynamical system, a fixed point corresponds to:
- a. A state of maximum change.
 - b. An initial state.
 - c. An equilibrium state.
 - d. A terminal state.
38. The iterative method $x_{k+1} = f(x_k)$ is known as:
- a. Newton's method.
 - b. Gradient descent.
 - c. Fixed-point iteration.
 - d. Bisection method.
39. Which of the following is NOT a required property of a metric $d(x, y)$?
- a. Non-negativity
 - b. Symmetry
 - c. Linearity
 - d. Triangle Inequality
40. The property $d(x, y) = 0 \iff x = y$ is called:
- a. Identity of Indiscernibles
 - b. Symmetry
 - c. Non-negativity
 - d. Completeness

41. A mapping \mathcal{T} is a contraction if $d(\mathcal{T}(x), \mathcal{T}(y)) \leq \alpha d(x, y)$ for:
- $\alpha = 1$
 - $\alpha > 1$
 - $0 \leq \alpha < 1$
 - any $\alpha \in \mathbb{R}$
42. What is a mapping called if the contraction constant α is allowed to be 1?
- Expansive
 - Non-expansive
 - Isometric
 - Injective
43. Why is the condition $\alpha < 1$ absolutely critical for a contraction mapping?
- It ensures the mapping is linear.
 - It guarantees the existence and uniqueness of a fixed point.
 - It makes the metric space finite.
 - It ensures the operator is invertible.
44. A complete metric space is one where:
- Every sequence has a limit.
 - The space is finite.
 - Every Cauchy sequence converges to a limit within the space.
 - The distance between any two points is 1.
45. Which of these is a famous example of a space that is NOT complete?
- The set of real numbers, \mathbb{R} .
 - The set of rational numbers, \mathbb{Q} .
 - The set of all bounded functions.
 - A high-dimensional real space, \mathbb{R}^n .
46. The Banach Fixed-Point Theorem guarantees that a contraction mapping on a non-empty complete metric space has:

- a. At least one fixed point.
 - b. Exactly one fixed point.
 - c. A finite number of fixed points.
 - d. No fixed points.
47. The proof of the Banach Fixed-Point Theorem is described as "constructive" because it:
- a. Assumes the fixed point exists and derives a contradiction.
 - b. Provides the algorithm for finding the fixed point.
 - c. Is very short and elegant.
 - d. Was constructed by Stefan Banach himself.
48. In the proof of existence, the sequence $x_{n+1} = \mathcal{T}(x_n)$ is first shown to be:
- a. A monotonic sequence.
 - b. A bounded sequence.
 - c. A Cauchy sequence.
 - d. A divergent sequence.
49. The proof of uniqueness for the Banach Fixed-Point Theorem uses what common technique?
- a. Proof by induction.
 - b. Proof by construction.
 - c. Proof by contradiction.
 - d. Proof by exhaustion.
50. The inequality $d(x_m, x_n) \leq \frac{\alpha^n}{1-\alpha} d(x_1, x_0)$ is used to show that:
- a. The sequence is Cauchy.
 - b. The limit is a fixed point.
 - c. The fixed point is unique.
 - d. The operator is a contraction.
51. The space of all bounded functions, where value functions reside, is:

- a. Not a metric space.
 - b. A complete metric space.
 - c. An incomplete metric space.
 - d. Not a vector space.
52. The function $f(x) = x + 1$ on the real line is an example of a function that is:
- a. A contraction with no fixed point.
 - b. Non-expansive with no fixed point.
 - c. A contraction with one fixed point.
 - d. Non-expansive with one fixed point.
53. The geometric interpretation of a fixed point for a real function $f(x)$ is where the graph $y = f(x)$ intersects:
- a. The x-axis.
 - b. The y-axis.
 - c. The line $y = x$.
 - d. The line $y = -x$.
54. The Banach Fixed-Point Theorem is also known as the:
- a. Contraction Mapping Theorem.
 - b. Intermediate Value Theorem.
 - c. Mean Value Theorem.
 - d. Brouwer Fixed-Point Theorem.
55. To show the limit x^* of the sequence x_n is a fixed point, one proves that:
- a. $d(x^*, x_n) \rightarrow 0$.
 - b. $d(x^*, \mathcal{T}(x^*)) = 0$.
 - c. $d(x_n, x_{n+1}) \rightarrow 0$.
 - d. $\mathcal{T}(x^*) = x_0$.

4 Part IV: Synthesis and Proofs

56. To apply fixed-point theory to RL, we reframe the problem of solving the Bellman optimality equation as:
- a. A matrix inversion problem.
 - b. A search for a fixed point of an operator.
 - c. A linear programming problem.
 - d. A supervised classification task.
57. The Bellman Optimality Operator, \mathcal{T}^* , takes a(n) _____ as input and produces a new one as output.
- a. Policy
 - b. Action-value function
 - c. State
 - d. Reward
58. The space of action-value functions for a finite MDP can be represented as a vector in:
- a. $\mathbb{R}^{|S|}$
 - b. $\mathbb{R}^{|A|}$
 - c. $\mathbb{R}^{|S| \times |A|}$
 - d. \mathbb{R}
59. The metric used to prove the Bellman operator is a contraction is the:
- a. Euclidean norm (L_2 -norm).
 - b. Manhattan norm (L_1 -norm).
 - c. Infinity norm (L_∞ -norm).
 - d. Hamming distance.
60. The infinity norm, $\|q_1 - q_2\|_\infty$, measures the:
- a. Average difference between q_1 and q_2 .
 - b. Sum of absolute differences between q_1 and q_2 .
 - c. Largest absolute difference between q_1 and q_2 across all state-action pairs.

- d. Smallest difference between q_1 and q_2 .
61. The proof of contraction for the Bellman operator shows that $\|\mathcal{T}^*q_1 - \mathcal{T}^*q_2\|_\infty \leq \alpha\|q_1 - q_2\|_\infty$, where α is equal to:
- a. 1
 - b. γ
 - c. $1 - \gamma$
 - d. $p(s', r|s, a)$
62. In the proof of contraction, which terms cancel out when taking the difference $|(\mathcal{T}^*q_1)(s, a) - (\mathcal{T}^*q_2)(s, a)|$?
- a. The discount factor γ .
 - b. The expected immediate reward terms.
 - c. The maximization operators.
 - d. The transition probabilities.
63. The key property of the max function used in the proof is:
- a. $|\max f(x) - \max g(x)| \leq \max |f(x) - g(x)|$
 - b. $\max(f + g) = \max f + \max g$
 - c. $\max(c \cdot f) = c \cdot \max f$
 - d. $\max(-f) = -\min(f)$
64. Why is the choice of the infinity norm fundamental to the strength of the result?
- a. It is the only norm that makes the space complete.
 - b. It is the easiest to compute.
 - c. It guarantees uniform convergence across the entire state-action space.
 - d. It is the only norm that is a valid metric.
65. The Value Iteration algorithm's update rule is equivalent to:
- a. One step of gradient descent.
 - b. A random walk in the value function space.
 - c. The repeated application of the Bellman Optimality Operator, $q_{k+1} \leftarrow \mathcal{T}^*q_k$.
 - d. The application of the Bellman Expectation Operator.
66. The convergence of Value Iteration is a direct consequence of:

- a. The Law of Large Numbers.
 - b. The Banach Fixed-Point Theorem applied to the Bellman Optimality Operator.
 - c. The fact that rewards are bounded.
 - d. The finiteness of the state space.
67. The final implication of the synthesis is that Value Iteration is guaranteed to converge to:
- a. A locally optimal value function.
 - b. A good approximation of the optimal value function.
 - c. The unique optimal action-value function, q^* .
 - d. One of many possible optimal value functions.
68. The convergence of Value Iteration holds regardless of:
- a. The discount factor γ .
 - b. The initial value function estimate q_0 .
 - c. The reward function.
 - d. The transition probabilities.
69. Uniform convergence ensures that:
- a. The algorithm converges in a single step.
 - b. The error improves for some states, but may worsen for others.
 - c. The value estimates improve across all states and actions simultaneously.
 - d. The policy converges, but the value function does not.
70. The logical syllogism for the final proof relies on which two premises?
- a. The space is complete and the operator is a contraction.
 - b. The state space is finite and the action space is finite.
 - c. The rewards are bounded and the policy is deterministic.
 - d. The problem is an MDP and the discount factor is less than 1.

5 Part V: Conclusion and Implications

71. The theoretical guarantees provide a foundation of _____ for a large class of RL algorithms.

- a. Speed
 - b. Simplicity
 - c. Trust
 - d. Heuristics
72. The uniqueness of q^* ensures that the policy derived by acting greedily with respect to it is:
- a. Truly optimal.
 - b. Near-optimal.
 - c. One of many optimal policies.
 - d. A stochastic policy.
73. The greedy policy extraction rule is:
- a. $\pi^*(s) = \text{random action}$
 - b. $\pi^*(s) = \arg \min_a q^*(s, a)$
 - c. $\pi^*(s) = \arg \max_a q^*(s, a)$
 - d. $\pi^*(s) = \sum_a \pi(a|s) q^*(s, a)$
74. The analysis of Value Iteration serves as a bridge to understanding more advanced, model-free algorithms like:
- a. Policy Iteration.
 - b. Q-learning and SARSA.
 - c. Linear Regression.
 - d. Support Vector Machines.
75. The principles of contraction and fixed points have been extended to handle:
- a. Only finite, discrete MDPs.
 - b. Problems with no rewards.
 - c. Continuous state and action spaces.
 - d. Non-Markovian environments.

Part II

Explanatory Questions (30 Questions)

1. Explain the three fundamental theoretical questions (existence, uniqueness, computability) concerning the optimal value function and why they are important for establishing the reliability of value-based RL.
2. Describe all five components of a finite Markov Decision Process (MDP) tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ and the role each plays.
3. What is the Markov Property and why is it a critical assumption for the entire MDP framework?
4. Explain the dual role of the discount factor γ . How does it function both as a mathematical necessity and as a parameter that shapes the agent's behavior?
5. Differentiate between a state-value function (V^π) and an action-value function (q^π). Why is the action-value function often more useful for control?
6. Derive the Bellman expectation equation for $V^\pi(s)$ starting from its definition, $V^\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$. Explain each step.
7. What is the "Principle of Optimality" and how is it mathematically embodied in the Bellman optimality equation for $q^*(s, a)$?
8. Compare and contrast the Bellman expectation equation for q^π with the Bellman optimality equation for q^* . What is the single most important difference and what does it signify?
9. What is a fixed point of an operator? Provide a geometric interpretation for a real-valued function and explain why fixed-point iteration ($x_{k+1} = f(x_k)$) does not always converge.
10. Define a metric space and list its four axiomatic properties. Why is this formal structure necessary for defining a contraction mapping?
11. What is a contraction mapping? Explain the role and critical importance of the contraction constant α being strictly less than 1.
12. What does it mean for a metric space to be "complete"? Provide an example of a complete space and an incomplete space, and explain why completeness is a necessary condition for the Banach Fixed-Point Theorem.
13. State the Banach Fixed-Point Theorem. What two crucial properties does it guarantee for a contraction mapping on a non-empty complete metric space?

14. Outline the proof of *existence* in the Banach Fixed-Point Theorem. How does this proof constructively provide an algorithm for finding the fixed point?
15. Outline the proof of *uniqueness* in the Banach Fixed-Point Theorem.
16. How is the problem of finding the optimal value function q^* reframed as a fixed-point problem? Define the specific operator used, \mathcal{T}^* .
17. What metric is used to define the distance between two value functions in the proof of convergence, and why is this specific metric chosen? What powerful guarantee does it provide?
18. Provide a step-by-step proof that the Bellman Optimality Operator \mathcal{T}^* is a contraction mapping in the infinity norm.
19. In the proof of contraction for \mathcal{T}^* , a key inequality regarding the max function is used: $|\max_x f(x) - \max_x g(x)| \leq \max_x |f(x) - g(x)|$. Explain intuitively why this inequality holds.
20. Describe the Value Iteration algorithm. How does its update rule directly correspond to the repeated application of the Bellman Optimality Operator?
21. Synthesize the entire argument of the report into a logical syllogism that proves Value Iteration converges to the unique optimal value function. State the premises and the conclusion clearly.
22. Why does the uniqueness of the optimal value function q^* provide an "unambiguous target for learning"? How is the optimal policy π^* extracted once q^* is known?
23. The report states that the theoretical analysis serves as a "bridge to understanding more advanced algorithms" like Q-learning. Explain this connection. How do model-free methods relate to the Bellman operator?
24. What is meant by "uniform convergence" in the context of Value Iteration, and why is this a much stronger guarantee than, for example, convergence of the average error?
25. If the discount factor γ were equal to 1, the Bellman operator would not be a contraction. What would be the mathematical and conceptual problems in an infinite-horizon MDP?
26. Explain the relationship between a deterministic policy and a stochastic policy. How can a deterministic policy be seen as a special case of a stochastic one?
27. The report transforms an "intractable infinite problem into a solvable one-step lookahead problem." Explain which concept achieves this transformation and how.
28. Why is the Bellman optimality equation a system of *non-linear* equations, while the Bellman expectation equation is a system of *linear* equations?

29. The proof of the Banach Fixed-Point Theorem relies on showing a sequence is a "Cauchy sequence." Explain in simple terms what a Cauchy sequence is and why this concept is important for proving convergence.
30. What are the broader implications of the theoretical guarantees for deploying RL systems in critical real-world applications?

Answers to Multiple-Choice Questions

1. c (To learn a policy that maximizes a cumulative reward signal.)
2. b (Quantify the long-term desirability of states or actions.)
3. c (Stability)
4. c (The Banach Fixed-Point Theorem)
5. b (The set of all possible states.)
6. c (Independent of the past given the present state.)
7. b (The probability of transitioning to state s' and receiving reward r , given the current state s and action a .)
8. c ($\gamma \in [0, 1)$)
9. b ("Myopic," focusing on immediate rewards.)
10. b (The probability that the agent "survives" to the next time step.)
11. b (Each state to a single, specific action.)
12. c (The sum of discounted future rewards.)
13. c (The expected return when starting in state s and following policy π .)
14. d ("How good is it to take action a in state s , and then commit to policy π ?")
15. b ($V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q^\pi(s, a)$)
16. c ($(\mathcal{S}, \mathcal{A}, P, R, \gamma)$)
17. c (Allows future predictions based only on the current state, not the full history.)
18. b (1)
19. a (The value functions could diverge to infinity.)
20. a (Formulate an effective policy without a model of the environment.)
21. b (A recursive structure.)
22. b (Policy evaluation.)
23. c (A weighted average of the values of its potential successor states.)
24. c ($|\mathcal{S}|$ linear equations in $|\mathcal{S}|$ variables.)
25. d (Temporal-Difference (TD) learning and SARSA.)
26. b (Bellman optimality equation.)

- 27. a ($V^{\pi^*}(s) \geq V^{\pi}(s)$ for all $s \in \mathcal{S}$ and all policies π .)
- 28. c ($\max_{\pi} V^{\pi}(s)$.)
- 29. b (Constitute an optimal policy with regard to the state resulting from the first decision.)
- 30. b (Maximization operator.)
- 31. c (The optimal policy itself.)
- 32. b (Solving a system of non-linear equations.)
- 33. b (It describes the value of a fixed, given strategy.)
- 34. a (It defines the value of the best possible strategy.)
- 35. c (Choose the action a that maximizes $q^*(s, a)$.)
- 36. b ($\mathcal{T}(x) = x$)
- 37. c (An equilibrium state.)
- 38. c (Fixed-point iteration.)
- 39. c (Linearity)
- 40. a (Identity of Indiscernibles)
- 41. c ($0 \leq \alpha < 1$)
- 42. b (Non-expansive)
- 43. b (It guarantees the existence and uniqueness of a fixed point.)
- 44. c (Every Cauchy sequence converges to a limit within the space.)
- 45. b (The set of rational numbers, \mathbb{Q} .)
- 46. b (Exactly one fixed point.)
- 47. b (Provides the algorithm for finding the fixed point.)
- 48. c (A Cauchy sequence.)
- 49. c (Proof by contradiction.)
- 50. a (The sequence is Cauchy.)
- 51. b (A complete metric space.)
- 52. b (Non-expansive with no fixed point.)
- 53. c (The line $y = x$.)
- 54. a (Contraction Mapping Theorem.)

- 55. b ($d(x^*, \mathcal{T}(x^*)) = 0$.)
- 56. b (A search for a fixed point of an operator.)
- 57. b (Action-value function)
- 58. c ($\mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$)
- 59. c (Infinity norm (L_∞ -norm).)
- 60. c (Largest absolute difference between q_1 and q_2 across all state-action pairs.)
- 61. b (γ)
- 62. b (The expected immediate reward terms.)
- 63. a ($|\max f(x) - \max g(x)| \leq \max |f(x) - g(x)|$)
- 64. c (It guarantees uniform convergence across the entire state-action space.)
- 65. c (The repeated application of the Bellman Optimality Operator, $q_{k+1} \leftarrow \mathcal{T}^* q_k$.)
- 66. b (The Banach Fixed-Point Theorem applied to the Bellman Optimality Operator.)
- 67. c (The unique optimal action-value function, q^* .)
- 68. b (The initial value function estimate q_0 .)
- 69. c (The value estimates improve across all states and actions simultaneously.)
- 70. a (The space is complete and the operator is a contraction.)
- 71. c (Trust)
- 72. a (Truly optimal.)
- 73. c ($\pi^*(s) = \arg \max_a q^*(s, a)$)
- 74. b (Q-learning and SARSA.)
- 75. c (Continuous state and action spaces.)

Answers to Explanatory Questions

1. **Answer:** The three fundamental questions are:

- **Existence:** Does an optimal value function (q^*) even exist? This is important because if no such function exists, the entire goal of finding it is ill-defined.
- **Uniqueness:** If q^* exists, is it the only one? This is crucial for reliability. If multiple, different optimal value functions existed, it would lead to ambiguity about which one to target and whether a policy derived from one is truly optimal.
- **Computability:** If q^* exists and is unique, can we design an algorithm that is guaranteed to find it? This connects theory to practice. Without a guaranteed method of computation, the existence and uniqueness of q^* would be purely academic.

Answering these affirmatively establishes that the problem is well-posed, has a single correct solution, and that we have a reliable method (Value Iteration) to find that solution.

2. **Answer:** The 5-tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ defines an MDP:

- \mathcal{S} : The set of **States**. This is the set of all possible configurations of the environment. It defines "where" the agent can be.
- \mathcal{A} : The set of **Actions**. This is the set of all choices the agent can make. It defines "what" the agent can do.
- P : The **Transition Probability Function**, $p(s', r|s, a)$. This is the dynamics model of the environment. It defines the probability of ending up in state s' and receiving reward r after taking action a in state s .
- R : The **Reward Function**. This function defines the goal of the agent. It provides a scalar feedback signal, r , for each transition. The agent's objective is to maximize the cumulative sum of these rewards.
- γ : The **Discount Factor**. A scalar in $[0, 1)$ that determines the present value of future rewards. It ensures long-term returns are finite and models the agent's preference for immediate versus delayed gratification.

3. **Answer:** The **Markov Property** posits that the future is independent of the past, given the present. Formally, the probability of the next state S_{t+1} depends only on the current state S_t and the current action A_t , not on any previous states or actions. This assumption is critical because it dramatically simplifies the problem. Without it, an agent would need to consider the entire history of interactions to make an optimal

decision, which is computationally intractable. The Markov property allows the agent to make decisions based solely on its current observation, as the state s_t is assumed to "encapsulate" all relevant information from the history.

4. **Answer:** The discount factor γ has a dual role:

- (a) **Mathematical Necessity:** In infinite-horizon problems, an agent could accumulate rewards forever. Without discounting ($\gamma < 1$), the total return ($G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$) could be infinite, making the value functions diverge. This would render the optimization problem ill-defined. The discount factor ensures the sum converges to a finite value.
- (b) **Behavioral Parameter:** Conceptually, γ models the agent's "patience." A value close to 1 makes the agent **farsighted**, giving significant weight to long-term consequences. A value close to 0 makes the agent **myopic**, focusing almost exclusively on maximizing its immediate reward. It can also be interpreted as a survival probability, where $1 - \gamma$ is the chance the process terminates at any step.

5. **Answer:**

- The **state-value function**, $V^\pi(s)$, gives the expected return from *starting* in state s and then following policy π . It tells you how good a *state* is under a given policy.
- The **action-value function**, $q^\pi(s, a)$, gives the expected return from *starting* in state s , taking a specific action a , and *then* following policy π . It tells you how good a specific *action* is in a state.

The action-value function q^π is often more useful for control, especially in model-free settings. To improve a policy using V^π , you need a model of the environment ($p(s'|s, a)$) to see which action leads to the best next state. With q^π , you don't need a model. To find the best action in a state s , you can simply compare the $q^\pi(s, a)$ values for all actions a and choose the one with the highest value. This allows for direct policy improvement.

6. **Answer:** The derivation of the Bellman expectation equation for $V^\pi(s)$ is as follows:

(a) **Start with the definition:**

$$V^\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

(b) **Expand the return G_t into its first term and the rest:** The return is $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = R_{t+1} + \gamma G_{t+1}$.

$$V^\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

(c) **Split the expectation:** Using the linearity of expectation.

$$V^\pi(s) = \mathbb{E}_\pi[R_{t+1}|S_t = s] + \gamma \mathbb{E}_\pi[G_{t+1}|S_t = s]$$

(d) **Expand over actions and successor states:** We use the law of total expectation, averaging over all possible actions a (according to policy π) and all possible outcomes (s', r) (according to the environment dynamics p).

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1}|S_{t+1} = s']]$$

(e) **Recognize the recursive structure:** The term $\mathbb{E}_\pi[G_{t+1}|S_{t+1} = s']$ is, by definition, the value of the successor state s' , which is $V^\pi(s')$.

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V^\pi(s')]$$

This final equation expresses the value of a state recursively in terms of the expected values of its successor states.

7. **Answer:** The **Principle of Optimality**, in the context of RL, states that 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 principle is embodied in the Bellman optimality equation for $q^*(s, a)$ through the maximization operator:

$$q^*(s, a) = \sum_{s', r} p(s', r|s, a) \left[r + \gamma \max_{a' \in \mathcal{A}} q^*(s', a') \right]$$

The term $\max_{a' \in \mathcal{A}} q^*(s', a')$ asserts that after the initial action a and transition to state s' , an optimal policy must follow by choosing the action a' that maximizes the value from that new state. It mathematically enforces the idea that the "remaining decisions" must be optimal.

8. **Answer:**

- **Bellman Expectation for q^π :**

$$q^\pi(s, a) = \sum_{s', r} p(s', r|s, a) \left[r + \gamma \sum_{a' \in \mathcal{A}} \pi(a'|s') q^\pi(s', a') \right]$$

- **Bellman Optimality for q^* :**

$$q^*(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a' \in \mathcal{A}} q^*(s', a') \right]$$

The single most important difference is the replacement of the policy-weighted average $\sum_{a' \in \mathcal{A}} \pi(a' | s')$ with the direct maximization $\max_{a' \in \mathcal{A}}$.

This signifies the conceptual leap from **policy evaluation** to **optimal control**. The expectation equation is passive; it calculates the value of a *given, fixed* policy π . The optimality equation is active; it defines the value of the *best possible* strategy by implicitly selecting the best action at the next step.

9. **Answer:** A **fixed point** of an operator \mathcal{T} is a point x in the domain that is mapped to itself, i.e., $\mathcal{T}(x) = x$.

The **geometric interpretation** for a real-valued function $f(x)$ is the point where the graph of the function $y = f(x)$ intersects the line $y = x$.

Fixed-point iteration, $x_{k+1} = f(x_k)$, does not always converge because the behavior depends on the properties of the function near the fixed point. If the function is "too steep" (the absolute value of its slope is greater than 1), the iterations will spiral away from the fixed point (divergence). If the slope is gentle (absolute value less than 1), it will converge. This is precisely the condition captured by contraction mappings.

10. **Answer:** A **metric space** (X, d) is a set X paired with a distance function (metric) d that provides a notion of distance between elements of the set. The four axioms are:
 - (a) **Non-negativity:** $d(x, y) \geq 0$. (Distance is never negative).
 - (b) **Identity of Indiscernibles:** $d(x, y) = 0 \iff x = y$. (Distance is zero only if the points are the same).
 - (c) **Symmetry:** $d(x, y) = d(y, x)$. (The distance from x to y is the same as from y to x).
 - (d) **Triangle Inequality:** $d(x, z) \leq d(x, y) + d(y, z)$. (The direct path is the shortest).

This formal structure is necessary because the definition of a contraction mapping, $d(\mathcal{T}(x), \mathcal{T}(y)) \leq \alpha d(x, y)$, is fundamentally about how an operator affects *distances*. Without a well-defined metric, the concept of "shrinking distances" would be meaningless.

11. **Answer:** A **contraction mapping** is an operator \mathcal{T} on a metric space that uniformly shrinks the distance between any two points by at least a constant factor α . The role of the contraction constant α is to quantify this "shrinking factor."

The condition that α must be **strictly less than 1** is absolutely critical. If $\alpha = 1$ (a non-expansive map), the operator only guarantees not to increase distances. This is not enough to ensure convergence to a unique point. For example, a simple translation $f(x) = x + 1$ is non-expansive but has no fixed point. The strict inequality $\alpha < 1$ ensures that the operator actively pulls all points in the space closer together, which forces any iterative sequence to converge towards a single, unique equilibrium point.

12. **Answer:** A metric space is **complete** if every Cauchy sequence in the space converges to a limit that is also an element of the space. Intuitively, a complete space has no "holes" or "missing points."

- **Complete Example:** The set of real numbers \mathbb{R} . Any sequence of real numbers that gets arbitrarily close to itself (Cauchy) will converge to another real number.
- **Incomplete Example:** The set of rational numbers \mathbb{Q} . The sequence of rational approximations to $\sqrt{2}$ (e.g., 1, 1.4, 1.41, 1.414, ...) is a Cauchy sequence, but it converges to $\sqrt{2}$, which is irrational and thus not in the space \mathbb{Q} .

Completeness is necessary for the Banach Fixed-Point Theorem because the proof of existence first shows that the iterative sequence $x_{k+1} = \mathcal{T}(x_k)$ is a Cauchy sequence. Completeness is the property that then guarantees this sequence has a limit *within the space*, which can then be proven to be the fixed point.

13. **Answer:** The **Banach Fixed-Point Theorem** states: Let (X, d) be a non-empty complete metric space and let $\mathcal{T} : X \rightarrow X$ be a contraction mapping on X . Then, \mathcal{T} has **exactly one** fixed point $x^* \in X$.

The two crucial properties it guarantees are:

- (a) **Existence:** A fixed point is guaranteed to exist.
- (b) **Uniqueness:** This fixed point is guaranteed to be the only one.

14. **Answer:** The proof of existence is constructive:

- (a) **Construct a Sequence:** Start with an arbitrary point x_0 and generate a sequence using fixed-point iteration: $x_{k+1} = \mathcal{T}(x_k)$.
- (b) **Show it is a Cauchy Sequence:** By repeatedly applying the contraction property, it's shown that the distance between successive terms shrinks geometrically: $d(x_{k+1}, x_k) \leq \alpha^k d(x_1, x_0)$. Using the triangle inequality, this is extended to show that the distance between any two terms $d(x_m, x_n)$ can be made arbitrarily small by choosing large enough m, n . This satisfies the definition of a Cauchy sequence.
- (c) **Show the Limit is a Fixed Point:** Because the space is complete, the Cauchy sequence must converge to a limit, x^* , within the space. By showing that $d(x^*, \mathcal{T}(x^*)) \rightarrow$

0, we conclude that $x^* = \mathcal{T}(x^*)$, meaning the limit is a fixed point.

This proof provides an algorithm because the sequence construction in step 1 is precisely the algorithm of fixed-point iteration (which is the basis for Value Iteration). The proof guarantees that this very algorithm will converge to the desired solution.

15. **Answer:** The proof of uniqueness is a classic proof by contradiction:

- (a) **Assume two distinct fixed points exist:** Let's call them x^* and y^* , where $x^* \neq y^*$. This means $d(x^*, y^*) > 0$.
- (b) **Apply the definition of a fixed point:** Since they are fixed points, $\mathcal{T}(x^*) = x^*$ and $\mathcal{T}(y^*) = y^*$.
- (c) **Apply the contraction property:**

$$d(x^*, y^*) = d(\mathcal{T}(x^*), \mathcal{T}(y^*)) \leq \alpha d(x^*, y^*)$$

- (d) **Derive the contradiction:** The inequality implies $(1 - \alpha)d(x^*, y^*) \leq 0$. Since $\alpha < 1$, the term $(1 - \alpha)$ is positive. Since we assumed $d(x^*, y^*) > 0$, their product must be positive. This contradicts the inequality.
- (e) **Conclude:** The only way to resolve the contradiction is if the initial assumption was false. Therefore, $d(x^*, y^*) = 0$, which implies $x^* = y^*$. The fixed point is unique.

16. **Answer:** The problem is reframed by defining an operator whose fixed point is the solution to the Bellman optimality equation.

- **The Space:** The set of all bounded action-value functions, $q(s, a)$.
- **The Operator:** The **Bellman Optimality Operator**, \mathcal{T}^* , is defined by the right-hand side of the Bellman optimality equation. It takes an action-value function q as input and outputs a new action-value function \mathcal{T}^*q :

$$(\mathcal{T}^*q)(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a' \in \mathcal{A}} q(s', a') \right]$$

A function q^* is the optimal value function if and only if it satisfies $q^* = \mathcal{T}^*q^*$, which is the definition of a fixed point.

17. **Answer:** The metric used is the **infinity norm** (L_∞ -norm), defined as:

$$\|q_1 - q_2\|_\infty = \max_{s \in \mathcal{S}, a \in \mathcal{A}} |q_1(s, a) - q_2(s, a)|$$

This metric is chosen because it measures the **worst-case error** between two value

functions across the entire state-action space. Proving that the Bellman operator is a contraction in this specific norm provides a guarantee of **uniform convergence**. This means that with each iteration, the maximum possible error between the current estimate and the true optimal value function is guaranteed to shrink, ensuring that the estimates improve everywhere simultaneously.

18. **Answer:** We want to show $\|\mathcal{T}^*q_1 - \mathcal{T}^*q_2\|_\infty \leq \gamma\|q_1 - q_2\|_\infty$.

(a) **Start with the definition of the operator difference for a single (s,a):**

$$|(\mathcal{T}^*q_1)(s, a) - (\mathcal{T}^*q_2)(s, a)| = \left| \sum p(\dots)[r + \gamma \max q_1] - \sum p(\dots)[r + \gamma \max q_2] \right|$$

(b) **Simplify:** The expected reward terms cancel, and we factor out γ and the probability.

$$= \left| \gamma \sum p(s', r|s, a) (\max_{a'} q_1(s', a') - \max_{a'} q_2(s', a')) \right|$$

(c) **Apply triangle inequality and max property:** Move the absolute value inside the sum and use the property $|\max f - \max g| \leq \max |f - g|$.

$$\leq \gamma \sum p(s', r|s, a) |\max_{a'} q_1(s', a') - \max_{a'} q_2(s', a')| \leq \gamma \sum p(s', r|s, a) \max_{a'} |q_1(s', a') - q_2(s', a')|$$

(d) **Bound the local difference by the global maximum (the infinity norm):**

The term $\max_{a'} |q_1(s', a') - q_2(s', a')|$ is less than or equal to the max difference over ALL states, which is $\|q_1 - q_2\|_\infty$.

$$\leq \gamma \sum p(s', r|s, a) \|q_1 - q_2\|_\infty$$

(e) **Use the property of probability distributions:** The sum of probabilities is 1.

$$= \gamma \|q_1 - q_2\|_\infty$$

(f) **Take the maximum over all (s,a):** Since this holds for any (s,a), it holds for the maximum.

$$\|\mathcal{T}^*q_1 - \mathcal{T}^*q_2\|_\infty = \max_{s,a} |(\mathcal{T}^*q_1)(s, a) - (\mathcal{T}^*q_2)(s, a)| \leq \gamma \|q_1 - q_2\|_\infty$$

This completes the proof that \mathcal{T}^* is a γ -contraction.

19. **Answer:** The inequality is $|\max_x f(x) - \max_x g(x)| \leq \max_x |f(x) - g(x)|$. Let x_f be the point where $f(x)$ is maximized, so $\max_x f(x) = f(x_f)$. The difference is $\max f - \max g = f(x_f) - \max g$. Since $\max g \geq g(x_f)$, we have $f(x_f) - \max g \leq f(x_f) - g(x_f)$.

So, $\max f - \max g \leq f(x_f) - g(x_f) \leq |f(x_f) - g(x_f)|$. Since $|f(x_f) - g(x_f)|$ is the difference at a specific point x_f , it must be less than or equal to the maximum difference over all points, $\max_x |f(x) - g(x)|$. Therefore, $\max f - \max g \leq \max_x |f(x) - g(x)|$. By symmetry, $\max g - \max f \leq \max_x |f(x) - g(x)|$. Combining these gives the final inequality. Intuitively, the difference between the peaks of two functions cannot be greater than the greatest vertical gap between the functions anywhere.

20. **Answer:** The **Value Iteration** algorithm is an iterative method to find the optimal value function.

- (a) **Initialization:** Start with an arbitrary action-value function estimate, q_0 .
- (b) **Iteration:** For each step k , compute the next estimate q_{k+1} from the previous one q_k for all state-action pairs using the update rule:

$$q_{k+1}(s, a) \leftarrow \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a' \in \mathcal{A}} q_k(s', a') \right]$$

This update rule is precisely the definition of the Bellman Optimality Operator, \mathcal{T}^* , applied to the function q_k . Therefore, the entire algorithm can be written concisely as the sequence $q_{k+1} \leftarrow \mathcal{T}^* q_k$. It is the concrete implementation of the abstract fixed-point iteration procedure.

21. **Answer:** The logical syllogism is as follows:

- **Premise 1:** The set of all bounded action-value functions on a finite MDP, equipped with the infinity norm metric, forms a **complete metric space**.
- **Premise 2:** The Bellman Optimality Operator, \mathcal{T}^* , is a **contraction mapping** on this space with a contraction constant $\gamma < 1$.
- **Conclusion from Banach's Theorem:** Therefore, according to the Banach Fixed-Point Theorem, the operator \mathcal{T}^* is guaranteed to have **one and only one fixed point**, which is the optimal action-value function, q^* .
- **Final Implication:** The theorem further guarantees that the sequence generated by repeatedly applying the operator ($q_{k+1} = \mathcal{T}^* q_k$), which is precisely the Value Iteration algorithm, will converge to this unique fixed point q^* .

22. **Answer:** The uniqueness of q^* provides an "unambiguous target for learning" because it confirms there is a single, correct answer that all convergent, value-based algorithms should find. There is no ambiguity about what the "optimal value" is.

Once this unique q^* is found, the optimal policy π^* is extracted **greedily**. For any given state s , the agent simply chooses the action a that has the highest learned optimal

action-value. This is because $q^*(s, a)$ represents the maximum possible expected return if you start by taking action a . A rational agent will naturally choose the action that promises the best outcome. Mathematically:

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} q^*(s, a)$$

23. **Answer:** The analysis of Value Iteration, a model-based method, provides the theoretical core for understanding model-free methods like Q-learning. The Bellman Optimality Operator, \mathcal{T}^* , represents an ideal "full backup" that requires a model ($p(s', r|s, a)$) to average over all possible next states.

Model-free algorithms like Q-learning work without this model. Instead of a full backup, they perform a **sample backup**. After taking action a_t in state s_t and observing a single reward r_{t+1} and next state s_{t+1} , Q-learning updates its estimate of $q(s_t, a_t)$ towards a target: $r_{t+1} + \gamma \max_{a'} q(s_{t+1}, a')$. This target is a noisy, sampled estimate of the true Bellman backup $(\mathcal{T}^*q)(s_t, a_t)$. Q-learning can thus be viewed as a stochastic approximation method that uses samples from the environment to move its value function estimate, on average, in the direction prescribed by the Bellman Optimality Operator. The fixed-point analysis provides the "ground truth" that these stochastic updates are trying to reach.

24. **Answer: Uniform convergence**, guaranteed by using the infinity norm, means that the maximum error between the current value estimate q_k and the optimal value function q^* shrinks with every iteration across all state-action pairs. That is, $\|q_k - q^*\|_\infty \rightarrow 0$.

This is a much stronger guarantee than convergence of the average error. Average error could be low even if the error for a few critical states is very large. Uniform convergence prevents this; it ensures that our value estimates improve everywhere, preventing a situation where the policy might be nearly optimal for most of the space but dangerously poor for a few critical states where the value estimate has failed to converge.

25. **Answer:** If $\gamma = 1$, the Bellman operator is only non-expansive, not a contraction, so the Banach Fixed-Point Theorem does not apply.

- **Mathematical Problem:** The proof of convergence for Value Iteration fails. More fundamentally, in an infinite-horizon MDP with ongoing positive rewards, the total return $\sum R_{t+k+1}$ would diverge to infinity. The value functions would be infinite, and the optimization problem becomes meaningless.
- **Conceptual Problem:** An agent with $\gamma = 1$ is infinitely patient. It would be

willing to accept any amount of short-term pain for an infinitesimally larger reward in the distant future. This makes it difficult to compare policies. For example, a policy that gets a reward of +1 at every step is indistinguishable from one that gets +1 at every step except the first, where it gets -1,000,000, as both have infinite total reward.

26. **Answer:**

- A **stochastic policy**, $\pi(a|s)$, defines a probability distribution over the available actions for each state. It specifies the probability of taking action a while in state s .
- A **deterministic policy**, $a = \pi(s)$, maps each state directly to a single action.

A deterministic policy is a special case of a stochastic policy where, for the chosen action $\hat{a} = \pi(s)$, the probability is 1, and for all other actions $a \neq \hat{a}$, the probability is 0. That is, $\pi(\hat{a}|s) = 1$ and $\pi(a|s) = 0$ for $a \neq \hat{a}$.

27. **Answer:** The concept that achieves this transformation is the **Bellman equation**.

The original definition of a value function, $V^\pi(s) = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]$, involves an infinite sum over all future time steps, which is intractable to compute directly.

The Bellman equation provides a recursive relationship that decomposes this infinite problem into a solvable one-step lookahead. For example, $V^\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma V^\pi(S_{t+1}) | S_t = s]$. This equation expresses the value of a state in terms of the expected immediate reward and the discounted value of the *next* state. This transforms the problem from summing over an infinite future to solving a system of equations where each equation only looks one step ahead.

28. **Answer:**

- The **Bellman expectation equation** for a given policy π is a system of **linear** equations. For V^π , the unknowns are the values $V^\pi(s)$ for each state s . The equation $V^\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(\dots)[r + \gamma V^\pi(s')]$ is linear with respect to the unknown variables $V^\pi(s')$.
- The **Bellman optimality equation** is a system of **non-linear** equations because of the max operator. The equation $V^*(s) = \max_a \sum_{s',r} p(\dots)[r + \gamma V^*(s')]$ involves a maximization over the unknown variables $V^*(s')$, which is a non-linear operation. This non-linearity is why direct solution methods (like matrix inversion) cannot be used, necessitating iterative approaches like Value Iteration.

29. **Answer:** A **Cauchy sequence** is a sequence where the terms become arbitrarily close to each other as the sequence progresses. Formally, for any small distance $\epsilon > 0$, there

exists a point in the sequence after which the distance between any two terms is less than ϵ .

This concept is important because it is a condition for convergence without knowing the limit itself. In the proof of the Banach theorem, we first show that the iterative sequence $x_{k+1} = \mathcal{T}(x_k)$ is Cauchy. This tells us the sequence is "settling down" and trying to converge. Then, the property of *completeness* of the space guarantees that there is a point *within the space* for it to converge to.

30. **Answer:** The theoretical guarantees are crucial for deploying RL in the real world because they provide a **foundation of trust**. For critical applications (e.g., autonomous driving, medical treatment, robotics), we cannot rely on algorithms that are just heuristics that "seem to work." The proofs of existence, uniqueness, and computability assure us that:

- The problem has a single, correct optimal solution.
- The algorithm we are using (or a derivative of it) is mathematically guaranteed to find that solution under the given assumptions.
- The resulting policy is truly optimal and not just a locally optimal or dangerously suboptimal one.

This confidence is essential for safety, reliability, and certification of AI systems in high-stakes environments.