
INTRO TO AI

Rui Wu

CS Department, Rutgers University
rw761@scarletmail.rutgers.edu

1 Introduction and AI agent

In Artificial Intelligence, the central problem at hand is that of the creation of a rational agent. "rational" means that maximally achieving **predefined goals(Rutgers), or expected utility(Berkeley)**. And the "agent" could be a human, a robot, a software, etc, it can choose its actions based on only its previous actions and observations.

A reflex agent is one that does not think about the consequences of its actions, but rather selects an action based on the current state of the world. These agents are typically outperformed by planning agents, which make decisions based on (hypothesized) consequences of actions and have a model of how the world evolves in response to actions. A re-planning agent is similar to a planning agent but can adapt its plan when the environment changes or when new information becomes available.

2 Search

2.1 State Spaces and Search Problems

In order to create a rational planning agent, we need a way to mathematically express the given environment in which the agent will exist. To do this, we must formally express a search problem - given our agent's current state (its configuration within its environment), how can we arrive at a new state that satisfies its goals in the best possible way? A search problem consists of the following elements:

- A state space: The set of all possible states that are possible in your given world.
- A set of actions available in each state.
- A transition model: Outputs the next state when a specific action is taken at current state.
- An action cost: Incurred when moving from one state to another after applying an action.
- A start state: The state in which an agent exists initially.
- A goal test: A function that takes a state as input, and determines whether it is a goal state.

The difference between a **world state** and a **search state** is that a world state contains all information about a given state, whereas a search state contains only the information about the world that is necessary for planning.

2.2 State Space Graphs and Search Trees

A state space graph is constructed with states representing nodes, with directed edges existing from a state to its children. These edges represent actions, and any associated weights represent the cost of performing the corresponding action.

Unlike state space graphs, search trees have no such restriction on the number of times a state can appear. Search trees are also a class of graph with states as nodes and actions as edges between states, each state/node encodes not just the state itself, but the entire path (or plan) from the start state to the given state in the state space graph.

2.3 Uninformed Search

Uninformed search is like trying to find your way through a maze without a map. Here's how it works:

1. Starting point: You begin at the entrance of the maze. This is your "start state."
2. Making a plan: As you explore, you keep a list of possible paths to try. This list is called the "frontier."
3. Exploring: You pick a path from your list and follow it a bit further into the maze.
4. New choices: When you reach a point where the path branches, you add these new options to your list of paths to try later.
5. Keeping track: For each new path you discover, you remember: Where it leads? How you got there? How far you've traveled?
6. Repeat: You keep doing this - picking a path, exploring it, and adding new options to your list.
7. Finding the exit: You continue this process until you pick a path that leads to the exit (your "goal state").
8. Success!: Once you find the exit, you can trace back the path you took to get there.

The key thing about uninformed search is that you don't have any special knowledge about which paths might be better. You're just systematically trying different options until you find the goal.

When we have no knowledge of the location of goal states in our search tree, we are forced to select our strategy for tree search from one of the techniques that falls under the umbrella of uninformed search. We'll now cover three such strategies in succession: **depth-first search**, **breadth-first search**, and **uniform cost search**.

2.4 Kinds of Uninformed Search

Depth-first search(DFS) is a strategy for exploration that always selects the **deepest** frontier node from the start node for expansion. Removing the deepest node and replacing it on the frontier with its children necessarily means the children are now the new deepest nodes - their depth is one greater than the depth of the previous deepest node. This implies that to implement DFS, we require a structure that always gives the most recently added objects highest priority. DFS is not complete. If there exist cycles in the state space graph, this inevitably means that the corresponding search tree will be infinite in depth. DFS simply finds the "leftmost" solution in the search tree without regard for path costs, and so is not optimal.

Breadth-first search(BFS) is a strategy for exploration that always selects the **shallowest** frontier node from the start node for expansion. If we want to visit shallower nodes before deeper nodes, we must visit nodes in their order of insertion. If a solution exists, then the depth of the shallowest node s must be finite, so BFS must eventually search this depth. Hence, it's complete. BFS is generally not optimal because it simply does not take costs into consideration when determining which node to replace on the frontier.

Bidirectional Search is running two simultaneous BFS: one forward from the initial state, and one backward from the goal. They stop when the two meet in the middle.

Uniform cost search(UCS), is a strategy for exploration that always selects the lowest cost frontier node from the start node for expansion. To represent the frontier for UCS, the choice is usually a heap-based priority queue, where the priority for a given enqueued node v is the path cost from the start node to v , or the backward cost of v . UCS is complete. If a goal state exists, it must have some finite length shortest path. UCS is also optimal if we assume all edge costs are non negative.

2.5 Informed Search

Uniform cost search is good because it's both complete and optimal, but it can be fairly slow because it expands in every direction from the start state while searching for a goal. If we have some notion of the direction in which we should focus our search, we can significantly improve performance and "hone in" on a goal much more quickly. This is exactly the focus of **informed search**.

Heuristics are the driving force(functions) that allow estimation of distance to goal states. A common heuristic that is used to solve Pacman problem is **Manhattan distance**, which for two points (x_1, y_1) and (x_2, y_2) is defined as follows:

$$Manhattan(x_1, y_1, x_2, y_2) = |x_1 - x_2| + |y_1 - y_2|$$

Greedy search is a strategy for exploration that always selects the frontier node with the lowest heuristic value for expansion, which corresponds to the state it believes is nearest to a goal.

A* Search combining UCS and Greedy. It is a strategy for exploration that always selects the frontier node with the lowest estimated total cost for expansion (total backward cost + heuristic value), where total cost is entire cost from the start node to the goal node.

- $g(n)$ -cost computed by UCS
- $h(n)$ -cost computed by Greedy
- $f(n) = g(n) + h(n)$

A heuristic h is **admissible** if: $0 \leq h(n) \leq h^*(n)$
Where: $h^*(n)$ is the true cost to a nearest goal.

Graph Search is the tree search with the optimization that ensuring the each node isn't already in the set before expansion and add it to the set after expansion if it's not. To maintain optimality under A^* graph search, we need to an even stronger property than admissibility, **consistency**. Mathematically, the consistency constrain can be expressed as follow:

$$\text{for any } A, C, h(A) - h(C) \leq \text{cost}(A, C)$$

Theorem: For a given search problem, if the consistency constraint is satisfied by a heuristic function h , using A^* graph search with h on that search problem will yield an optimal solution.

3 Markov Decision Processes

A Markov Decision Process (MDP) is defined by the following components:

- **States** (S): The set of all possible states the agent can be in.
- **Actions** (A): The set of all possible actions the agent can take.
- **Transition Function** ($T(s, a, s')$): The probability that taking action a in state s leads to state s' , i.e., $T(s, a, s') = P(s'|s, a)$.
- **Reward Function** ($R(s, a, s')$): The immediate reward received after transitioning from state s to state s' via action a .
- **Start State:** The initial state where the agent begins.
- **Discount Factor** (γ): A factor $0 \leq \gamma \leq 1$ that discounts future rewards.

3.1 Markov Property

The Markov property asserts that the future state depends only on the current state and the action taken, not on the sequence of previous states. Formally, the transition from state s to s' depends only on s and a , not on the history of past states.

3.2 Objective of MDP

The goal in an MDP is to find an **optimal policy** $\pi^* : S \rightarrow A$, which maps each state s to an action a , such that the agent maximizes its expected utility (total reward over time). The utility is computed as the sum of (possibly discounted) rewards, taking into account the uncertainty of action outcomes.

3.3 Policy

A policy $\pi(s)$ defines the action a that the agent should take when in state s . The optimal policy π^* maximizes the expected utility over all states.

3.4 Utility in Sequential Decision Making

To define the agent's objective, we use the concept of utility, which is the sum of future rewards. Given a discount factor γ , the utility of a reward sequence $[r_1, r_2, r_3, \dots]$ is given by:

$$U([r_1, r_2, r_3, \dots]) = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$$

The discount factor γ ensures that rewards received sooner have higher value than those received later.

3.5 Bellman Equation

The value of a state $V^*(s)$ under an optimal policy π^* is given by the **Bellman equation**:

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

This recursive equation defines the value of a state as the maximum expected utility achievable by taking the best action a and following the optimal policy thereafter.

3.6 Solving MDPs

MDPs can be solved using different techniques, such as:

- **Value Iteration:** An iterative algorithm that computes the optimal value function by repeatedly applying the Bellman equation until the values converge.
- **Policy Iteration:** Alternates between policy evaluation (computing the value of a given policy) and policy improvement (updating the policy based on the computed values).

4 Reinforcement Learning

Reinforcement Learning (RL) is an area of machine learning where an agent learns to interact with an environment to maximize cumulative rewards over time. Unlike Markov Decision Processes (MDPs), where the transition model and rewards are known, RL deals with environments where the agent does not have prior knowledge of these factors and must learn from experience.

4.1 Components of Reinforcement Learning

The basic elements of RL are:

- **Agent:** The learner or decision-maker.
- **Environment:** Everything the agent interacts with.
- **State** ($s \in S$): A representation of the current situation of the environment.
- **Action** ($a \in A$): Choices available to the agent to interact with the environment.
- **Reward** ($r \in \mathbb{R}$): Feedback from the environment based on the agent's action.
- **Policy** (π): A mapping from states to actions. The agent's objective is to learn an optimal policy π^* that maximizes cumulative reward.
- **Value Function:** Measures the expected future rewards starting from a particular state or state-action pair. The value function helps guide the agent to learn better policies.

4.2 Exploration vs. Exploitation

A fundamental challenge in reinforcement learning is balancing exploration (trying new actions to discover their effects) and exploitation (choosing known actions that yield high rewards). The agent must explore sufficiently to learn about the environment while exploiting what it knows to maximize its rewards.

4.3 Model-Based vs. Model-Free Learning

- **Model-Based Learning:** The agent learns an approximate model of the environment (i.e., the transition probabilities and rewards) and then uses this model to solve the MDP using value or policy iteration.
- **Model-Free Learning:** The agent directly learns the value of actions and states from experience without explicitly learning the transition probabilities and rewards.

4.4 Q-Learning

Q-learning is a model-free reinforcement learning algorithm that seeks to learn the optimal action-value function, $Q^*(s, a)$. The update rule for Q-learning is:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

Where:

- $Q(s, a)$ is the current estimate of the action-value function.
- r is the reward received after taking action a in state s .
- s' is the next state.
- α is the learning rate.
- γ is the discount factor.

Q-learning is considered an off-policy algorithm because it updates its Q-values based on the action that maximizes future rewards, regardless of the current policy.

4.5 Temporal Difference Learning

Temporal difference (TD) learning is another approach to reinforcement learning where the agent updates its value estimates after each action, based on the difference between predicted and actual rewards. The TD update rule is:

$$V(s) \leftarrow V(s) + \alpha [r + \gamma V(s') - V(s)]$$

Here, the agent updates its estimate of the value of state s based on the reward r and the value of the next state s' .

4.6 Exploration Strategies

- **ϵ -Greedy:** The agent explores by selecting a random action with probability ϵ and exploits the current best-known action with probability $1 - \epsilon$. This balances exploration and exploitation.
- **Exploration Functions:** The agent uses additional strategies to encourage exploring less visited states by providing bonuses to underexplored actions or states.

4.7 Approximate Q-Learning

In environments with large or continuous state spaces, storing Q-values for every state-action pair becomes impractical. Approximate Q-learning addresses this by using a feature-based representation. The Q-function is represented as a linear combination of features:

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

Here, f_1, f_2, \dots, f_n are features representing the state-action pair, and w_1, w_2, \dots, w_n are weights that are learned through experience.