

Markov Decision Process

Why do we care about MDPs

Markov Decision Processes are the defining problem of RL and any way to solve them is a method of RL

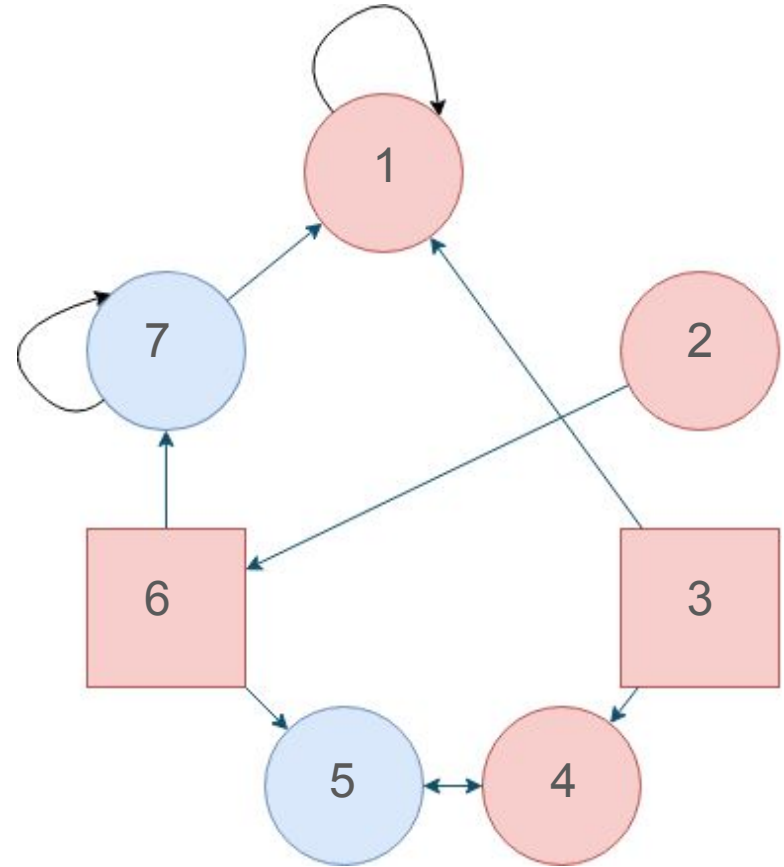
A Simple Zero-Sum Game

Players: Adam (Square) and Eve (Circle)

Objective for Eve:

Move a token to a blue vertex

Question - Given a starting vertex, does Eve have a winning strategy?



In General, we define a Game on a Graph as:

- $G = (V, E)$
- $V = V_{\text{Adam}} \uplus V_{\text{Eve}}$
- A Colouring Function: $E \rightarrow C$
- A Win Condition (Qualitative): $V_{\text{win}} \subset V$

In this case, our 2 players are Adam and Eve. Eve wins if there exists a single strategy such that no matter what Adam does, the token reaches $v \in V_{\text{win}}$

- Objective: Reachability, Safety, Buchi, co-Buchi, etc.

In General, we define a Game on a Graph as:

- A Win Condition (Quantitative): $f: \text{Paths} \rightarrow \mathbb{R} \cup \{\pm\infty\}$

In this case, our 2 players are Max and Min. If the cumulative total is below a certain threshold ϵ Min wins, else Max wins

- Alternative Formulation: Optimization!

We can think of these as rewards

MDP as Game of Agent vs Environment

The previous game is an example of a 2-player game: Adam and Eve.

Reinforcement Learning → agent vs the environment.

Games on Graph → 1.5 player game

Example of such a game: Blackjack

Simplified Blackjack Rules

Deck & Card Values:

- **Cards:** $\{1, 2, 3, \dots, 10, 11, 12, 13\}$
- **Draw:** With replacement ($\Pr(c)=1/13$)
- Both player and dealer get their cards face up. The dealer starts with 1 card.

Player's Turn:

- **Hit:** Draw one card
- **Stand:** End turn; player's score is the sum total of the cards that they have drawn

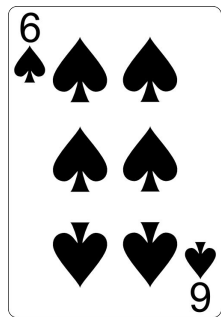
Dealer's Turn:

- Hits until $d \geq 17$
- $d > 21$: dealer busts → **player wins**.

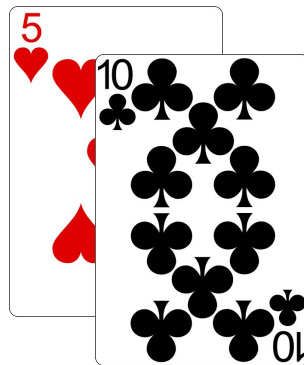
Outcome (if neither busts):

- $p > d$: **win**
- $p \leq d$: **loss**

Dealer

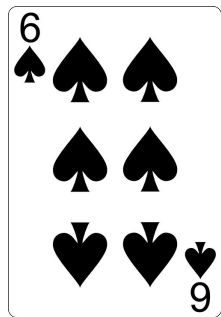


Player

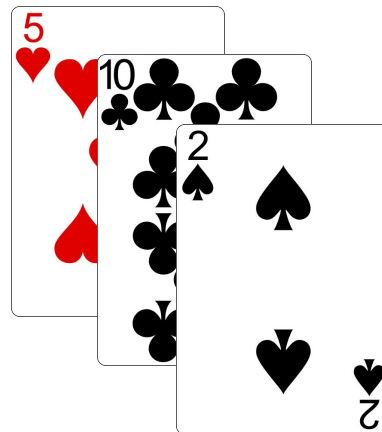


Hit or Stand?

Dealer

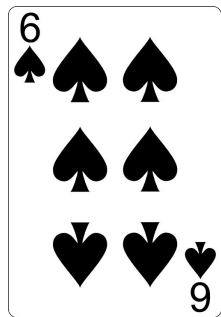


Player

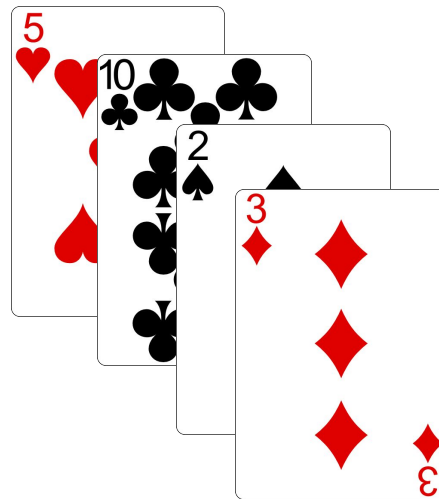


Hit or Stand?

Dealer



Player

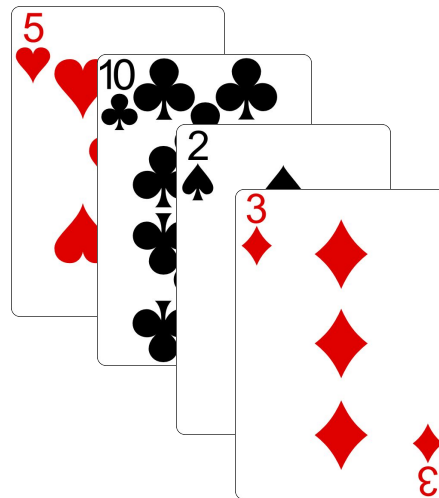


Players Total: 20

Dealer

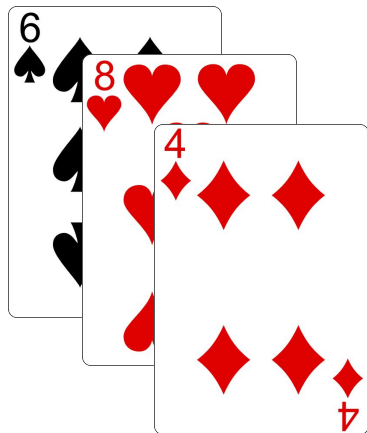


Player

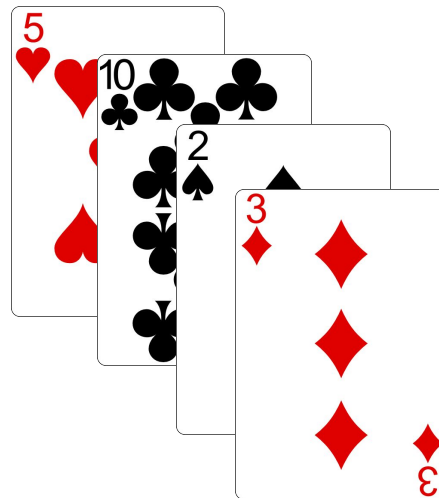


Players Total: 20

Dealer

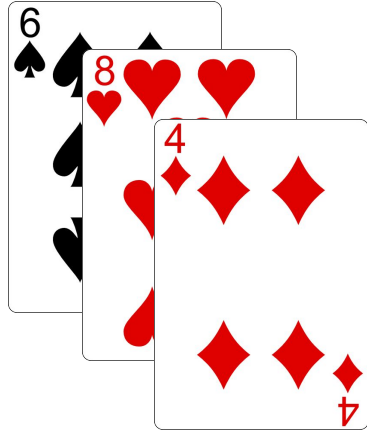


Player



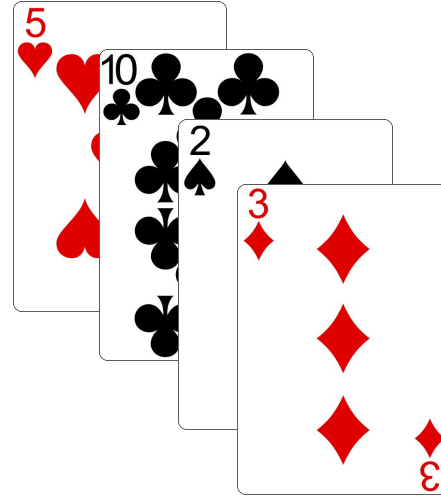
Players Total: 20

Dealer



Dealer's Total: 18

Player



Players Total: 20

Player Wins!

Modelling Blackjack as a MDP

$S = \{(p,d): p \text{ is the player's score, } d \text{ is the dealer's score}\} \cup \{\text{Win, Lose}\}$

Action Space = {Hit, Stand}

Goal

Win: 1-sink,

Lose: 0-sink

Transition Function

Hit:

$(p+x,d)$ where $P(x = \{1,2,\dots,13\}) = 1/13$

Lose if $p+x > 21$

Stand:

Player Stops

Dealer draws until $d \geq 17$

Game Formulation of an MDP - Markov Chain

Graph - $\{V_{\text{avg}}, E\}$

$V_{\text{avg}} - \{1, 2, \dots, n-1, n\}$

All vertices have two outgoing edges with probability $\frac{1}{2}$

v_{n-1} : 0-sink, v_n : 1-sink | Both edges are self loops

Transition matrix: $N \times N$ with all entries as 0, $\frac{1}{2}$, or 1

Said to be stopping if from every vertex you can reach a sink

Reachability of Markov Chains

We wish to reach the 1-sink and avoid the 0-sink

$$v_n = 1. \quad | \quad v_{n-1} = 0. \quad | \quad v_i = \frac{1}{2} v_j + \frac{1}{2} v_k$$

Let

$v = \{1, 2, \dots, n-2\}$: v is a column vector

$Q = (n-2) \times (n-2)$ transition matrix without 0-sink or 1-sink

$b = (n-2) \times 1$ vector of transitioning to 1-sink from the current vertex

$$v = Qv + b$$

Properties of Markov Chains

Properties that are useful to solving Markov Chains

- $\lim_{n \rightarrow \infty} Q^n = 0$ for a stopping matrix
- $(I - Q)^{-1} = I + Q + Q^2 + \dots$
- Unique solution given by $\nu = (I - Q)^{-1}b$

Game Formulation of Markov Decision Processes

$$G = \{V_{\text{avg}} \cup V_{\text{max}}, E\}$$

$$V_{\text{avg}} \cup V_{\text{max}} = \{1, 2, \dots, n-1, n\}$$

All vertices have two outgoing edges

v_{n-1} : 0-sink (lose), v_n : 1-sink (win) (both edges are self loops)

V_{avg} : Both vertices have probability $\frac{1}{2}$ (simplification)

V_{max} : Certain edges to $V_{\text{avg}} \cup V_{\text{max}}$, but no probabilities

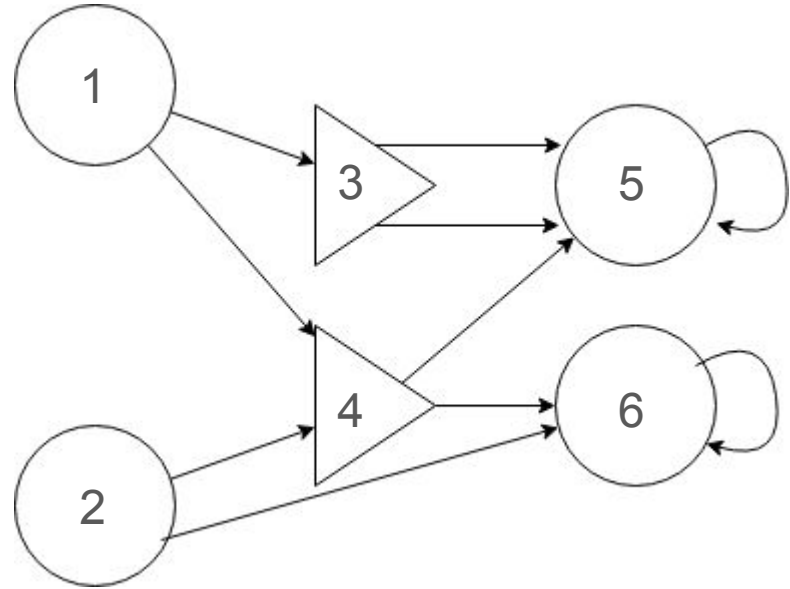
What is a strategy

Vertices: V_{\max} (Circles) and V_{\min} (Triangles)

Objective for Eve:

Avoid v_5 , Reach v_6

Question - What should be the agent's strategy from v_1 and v_2 .



Games on Graphs Terminology

Sequential: Players and the environment alternate turns while making decisions.

Fully Observable: The current state of the system is always known with certainty.

Stochastic: Outcomes of actions involve randomness, meaning the same action taken in the same state can lead to different results.

Markovian (RL Term): The decision-making process depends only on the current state and not on the sequence of past states.

Positionally Determined (Game Theory Term): The optimal strategy from a given state does not depend on the history of previous states.

A Comparison of Terminology

Games	Reinforcement Learning
Players	Agents
Moves	Actions
Strategy	Policy
Quantitative (Winning) Objective	Reward Function
Fog of War	Partial Observability (POMDP)
Game Balancing	Reward Shaping
Cheating	Adversarial Attack

Markov Decision Process to Markov Chain

Each strategy σ for G gives a Markov Chain G_σ . An MDP is said to be stopping if for every σ , the Markov Chain G_σ is stopping. Corresponding to G_σ is the vector denoting reachability probabilities defined in the previous slide

$$\bar{v} := \begin{bmatrix} \max_{\sigma} v_{\sigma}(1) \\ \max_{\sigma} v_{\sigma}(2) \\ \vdots \\ \max_{\sigma} v_{\sigma}(n) \end{bmatrix}$$

$$0 \leq w_i \leq 1 \quad \text{for every } i \tag{1.1}$$

$$w_n = 1$$

$$w_{n-1} = 0$$

$$\text{for each } i \leq n-2 \text{ in } V_{max} \quad w_i = \max(w_j, w_k) \quad \text{where } j \text{ and } k \text{ are children of } i$$

$$\text{for each } i \leq n-2 \text{ in } V_{avg} \quad w_i = \frac{1}{2}w_j + \frac{1}{2}w_k \quad \text{where } j \text{ and } k \text{ are children of } i$$

Policy Iteration Algorithm

Algorithm 1.1: Strategy improvement algorithm for MDPs, also known as policy iteration

```
1 algorithm strategy-improvement( $G$ )  
2    $\sigma \leftarrow$  an arbitrary positional strategy  
3    $v_\sigma \leftarrow$  probabilities to reach 1-sink in  $G_\sigma$   
4   repeat  
5     pick a node  $i \in V_{max}$  s.t.  $v_\sigma(i) < \max(v_\sigma(j), v_\sigma(k))$  where  $j, k$  are its children  
6      $\sigma'(i) \leftarrow \arg\max\{v_\sigma(j), v_\sigma(k)\}$   
7      $\sigma \leftarrow \sigma'$   
8      $v_\sigma \leftarrow$  probabilities to reach 1-sink in  $G_\sigma$   
9   until  $\sigma$  is optimal
```


Successive Approximation Algorithm

Algorithm 1.2: Value iteration algorithm, also known as Successive approximation algorithm

```
1 algorithm value-iteration( $G$ )
2   let  $\bar{u}$  be the vector assigning 1 to 1-sink and 0 to everything else
3   repeat
4     define  $\bar{u}'$  as follows:
5        $u'(i) = \max(u(j), u(k))$ , if  $i$  is a max node and  $j, k$  are its children
6        $u'(i) = \frac{1}{2}u(j) + \frac{1}{2}u(k)$ , if  $i$  is an average node and  $j, k$  are its children
7        $u'(n-1) = 0$ 
8        $u'(n) = 1$ 
9     let  $u \leftarrow u'$ 
10  until  $\bar{u}$  satisfies (1.1)
```