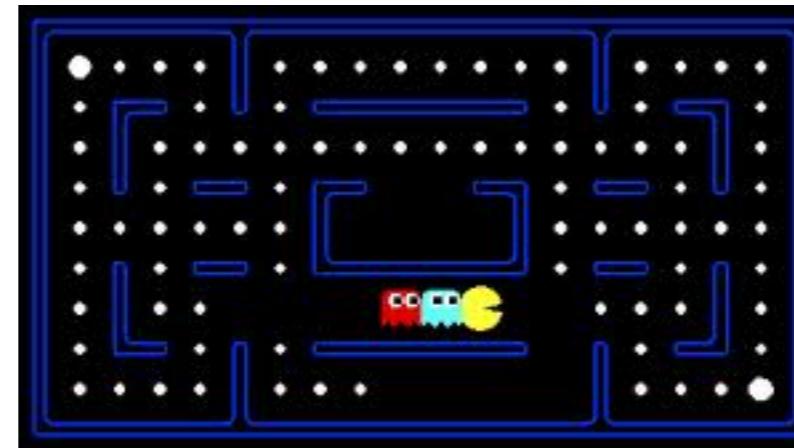
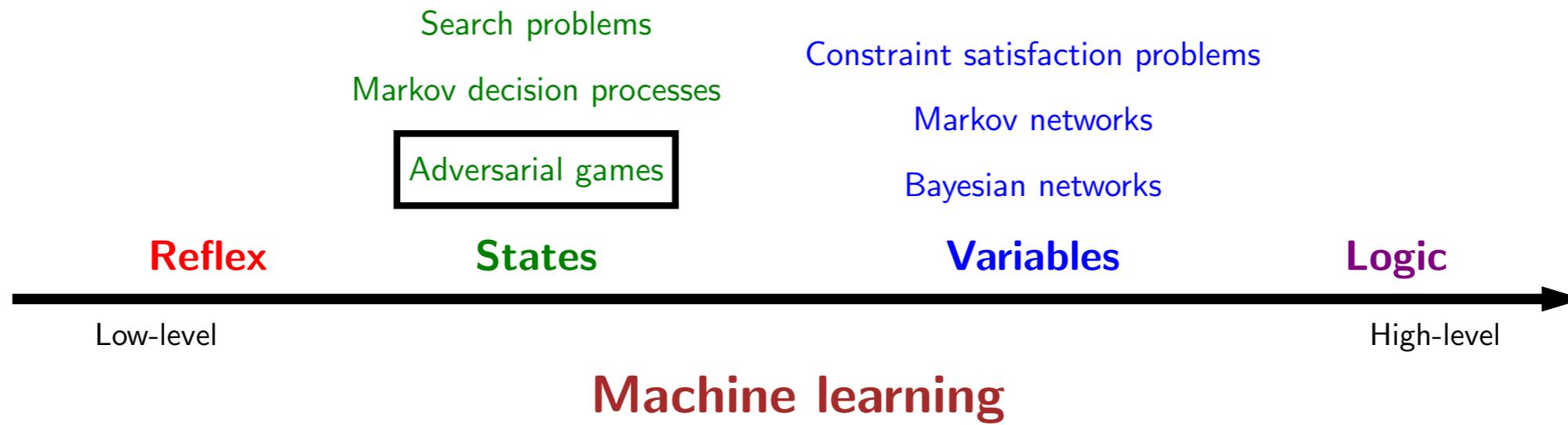




Games: overview



Course plan



- This lecture will be about games, which have been one of the main testbeds for developing AI programs since the early days of AI. Games are distinguished from the other tasks that we've considered so far in this class in that they make explicit the presence of other agents, whose utility is not generally aligned with ours. Thus, the optimal strategy (policy) for us will depend on the strategies of these agents. Moreover, their strategies are often unknown and adversarial. How do we reason about this?

A simple game

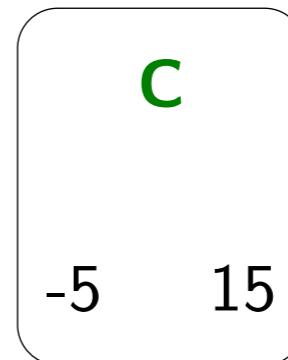
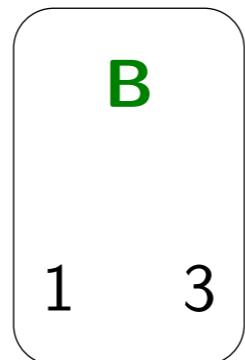
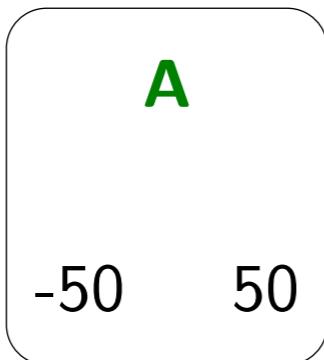


Example: game 1

You choose one of the three bins.

I choose a number from that bin.

Your goal is to maximize the chosen number.



working with

$$\max(A) = 50$$

working against

$$\min(B) = 1$$

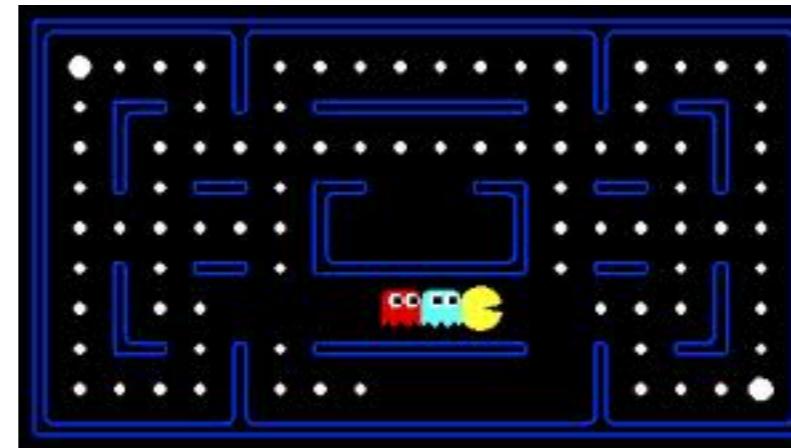
random

$$E[C] = 5$$

- Which bin should you pick? Depends on your mental model of the other player (me).
- If you think I'm working with you (unlikely), then you should pick A in hopes of getting 50. If you think I'm against you (likely), then you should pick B as to guard against the worst case (get 1). If you think I'm just acting uniformly at random, then you should pick C so that on average things are reasonable (get 5 in expectation).



Games: modeling



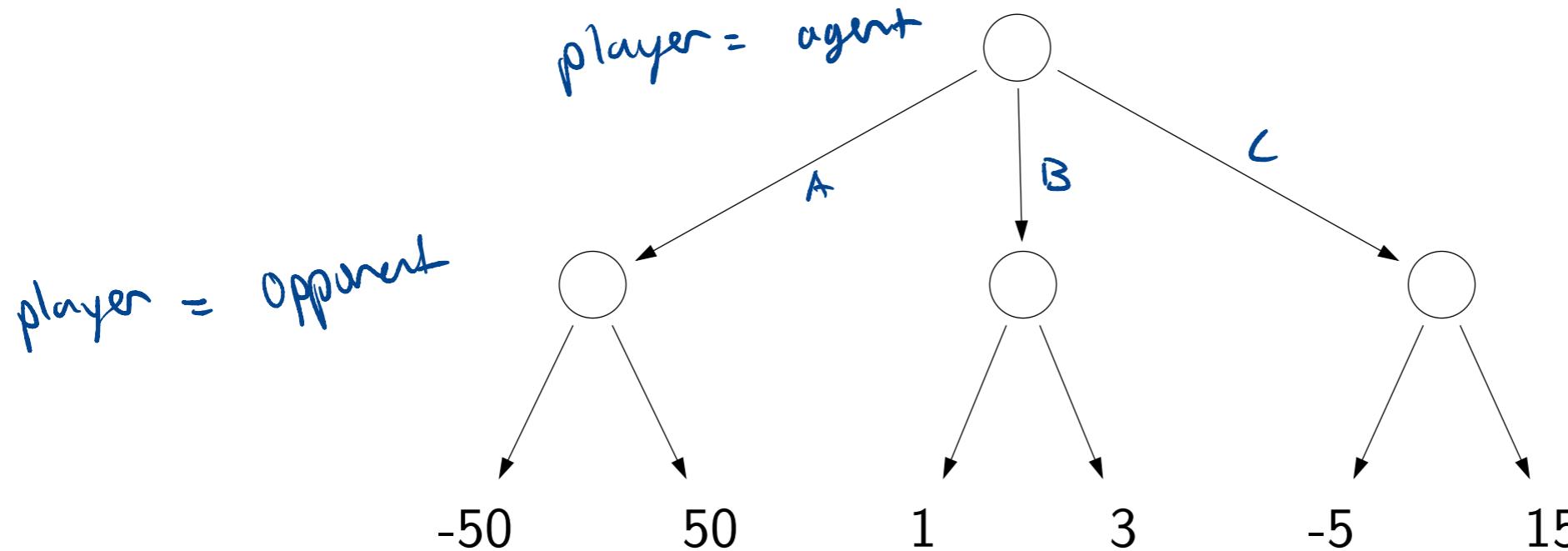
Game tree



Key idea: game tree

Each node is a decision point for a player.

Each root-to-leaf path is a possible outcome of the game.



Two-player zero-sum games

Players = {agent, opp}



Definition: two-player zero-sum game

s_{start} : starting state

$\text{Actions}(s)$: possible actions from state s

$\text{Succ}(s, a)$: resulting state if choose action a in state s

$\text{IsEnd}(s)$: whether s is an end state (**game over**)

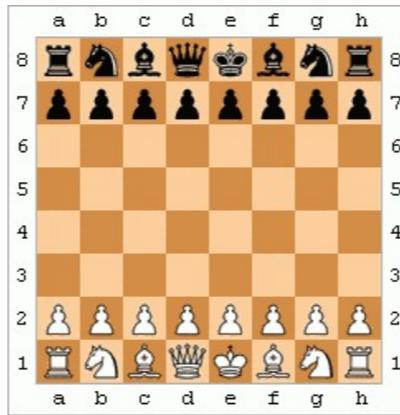
$\text{Utility}(s)$: agent's utility for end state s

$\text{Player}(s) \in \text{Players}$: player who controls state s

only get reward
at end state

- In this lecture, we will specialize to **two-player zero-sum** games, such as chess. To be more precise, we will consider games in which people take turns (unlike rock-paper-scissors) and where the state of the game is fully-observed (unlike poker, where you don't know the other players' hands). By default, we will use the term **game** to refer to this restricted form.
- We will assume the two players are named agent (this is your program) and opp (the opponent). Zero-sum means that the utility of the agent is negative the utility of the opponent (equivalently, the sum of the two utilities is zero).
- Following our approach to search problems and MDPs, we start by formalizing a game. Since games are a type of state-based model, much of the skeleton is the same: we have a start state, actions from each state, a deterministic successor state for each state-action pair, and a test on whether a state is at the end.
- The main difference is that each state has a designated Player(s), which specifies whose turn it is. A player p only gets to choose the action for the states s such that $\text{Player}(s) = p$.
- Another difference is that instead of having edge costs in search problems or rewards in MDPs, we will instead have a utility function Utility(s) defined only at the end states. We could have used edge costs and rewards for games (in fact, that's strictly more general), but having all the utility at the end states emphasizes the all-or-nothing aspect of most games. You don't get utility for capturing pieces in chess; you only get utility if you win the game. This ultra-delayed utility makes games hard.

Example: chess



Players = {white, black}

State s : (position of all pieces, whose turn it is)

Actions(s): legal chess moves that Player(s) can make

IsEnd(s): whether s is checkmate or draw

Utility(s): $+\infty$ if white wins, 0 if draw, $-\infty$ if black wins

↑
White is
agent

↑ opponent

Characteristics of games

- All the utility is at the end state



- Different players in control at different states



The halving game



Problem: halving game

Start with a number N .

Players take turns either decrementing N or replacing it with $\lfloor \frac{N}{2} \rfloor$.

The player that is left with 0 wins.

[semi-live solution: HalvingGame]

```
class HalvingGame(object):
    def __init__(self, N):
        self.N = N

    # state = (player, number)

    def startState(self):
        return (+1, self.N)

    def isEnd(self, state):
        player, number = state
        return number == 0

    def utility(self, state):
        player, number = state
        assert number == 0
        return player * float('inf')

    def actions(self, state):
        return ['-1', '/']

    def player(self, state):
        player, number = state
        return player

    def succ(self, state, action):
        player, number = state
        if action == '-1':
            return (-player, number-1)
        elif action == '/':
            return (-player, number//2)
```

players are +1 and -1
oppnent

current number

agent

{ if agent wins, utility is $+\infty$
if agent loses, utility is $-\infty$

{
-1 or $\frac{1}{2}$

player changes

policy & player p

Policies

Deterministic policies: $\pi_p(s) \in \text{Actions}(s)$

action that player p takes in state s

Stochastic policies $\pi_p(s, a) \in [0, 1]$:

probability of player p taking action a in state s

randomized

[semi-live solution: `humanPolicy`]

```
# general code
def humanPolicy(game, state):
    while True:
        action = input('Input action:')
        if action in game.actions(state):
            return action

# controller
policies = {+1: humanPolicy, -1: humanPolicy}
game = HalvingGame(N=15)
state = game.startState()

while not game.isEnd(state):
    print('*'*10, state)
    player = game.player(state)
    policy = policies[player]
    action = policy(game, state)
    state = game.succ(state, action)

print('utility = {}'.format(game.utility(state)))
```

user input
policy

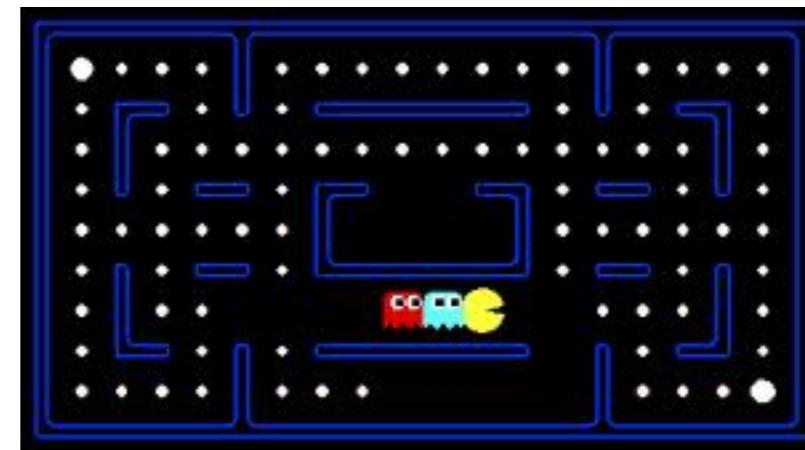
↳ policy printed only
at the end

- Following our presentation of MDPs, we revisit the notion of a **policy**. Instead of having a single policy π , we have a policy π_p for each player $p \in \text{Players}$. We require that π_p only be defined when it's p 's turn; that is, for states s such that $\text{Player}(s) = p$.
- It will be convenient to allow policies to be stochastic. In this case, we will use $\pi_p(s, a)$ to denote the probability of player p choosing action a in state s .
- We can think of an MDP as a game between the agent and nature. The states of the game are all MDP states s and all chance nodes (s, a) . It's the agent's turn on the MDP states s , and the agent acts according to π_{agent} . It's nature's turn on the chance nodes. Here, the actions are successor states s' , and nature chooses s' with probability given by the transition probabilities of the MDP: $\pi_{\text{nature}}((s, a), s') = T(s, a, s')$.



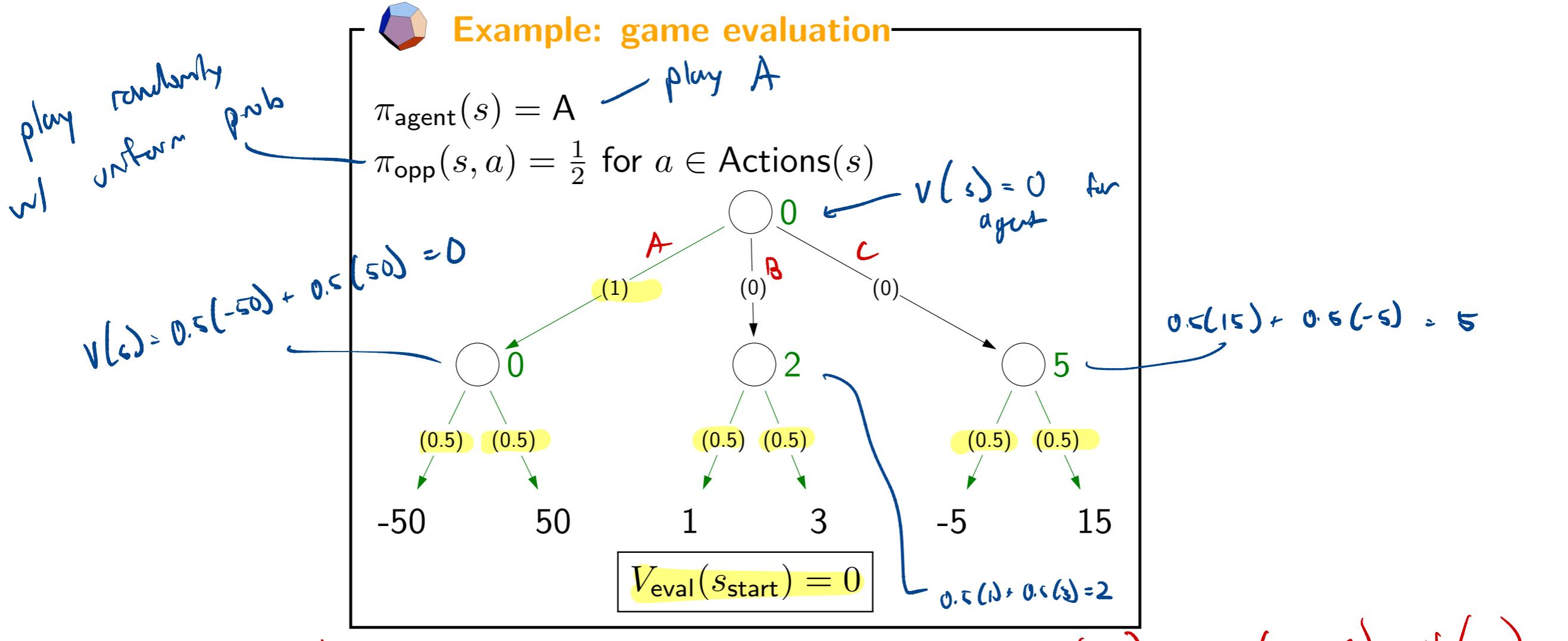
Games: game evaluation

→ policy evaluation



* know π_{agent} , π_{opp}

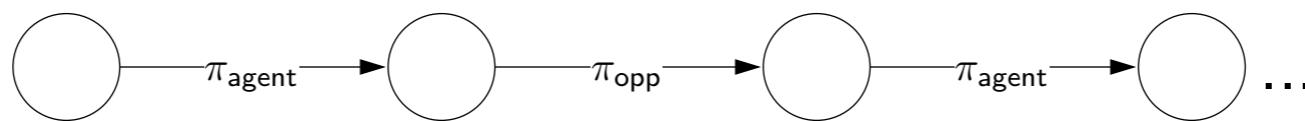
Game evaluation example



$$V_{\text{eva}}(s_{\text{start}}) = \underbrace{\pi(s, A)}_i \cdot V(s') + \underbrace{\pi(s, B)}_j \cdot V(s') + \underbrace{\pi(s, C)}_0 \cdot V(s')$$

Game evaluation recurrence

Analogy: recurrence for policy evaluation in MDPs



Value of the game:

* Utilities of
the agent
or from perspective

recursive
formula

$$V_{\text{eval}}(s) = \begin{cases} \text{Utility}(s) & \text{IsEnd}(s) \\ \sum_{a \in \text{Actions}(s)} \pi_{\text{agent}}(s, a) V_{\text{eval}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{agent} \\ \sum_{a \in \text{Actions}(s)} \pi_{\text{opp}}(s, a) V_{\text{eval}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{opp} \end{cases}$$

value
|
prob of action
|
value of successor state

- More generally, we can write down a recurrence for $V_{\text{eval}}(s)$, which is the **value** (expected utility) of the game at state s .
- There are three cases: If the game is over ($\text{IsEnd}(s)$), then the value is just the utility $\text{Utility}(s)$. If it's the agent's turn, then we compute the expectation over the value of the successor resulting from the agent choosing an action according to $\pi_{\text{agent}}(s, a)$. If it's the opponent's turn, we compute the expectation with respect to π_{opp} instead.

$V_{\text{eval}}(s) \rightarrow \begin{cases} \text{know } \pi_{\text{agent}}, & \text{know } \pi_{\text{opp}} \\ \end{cases}$

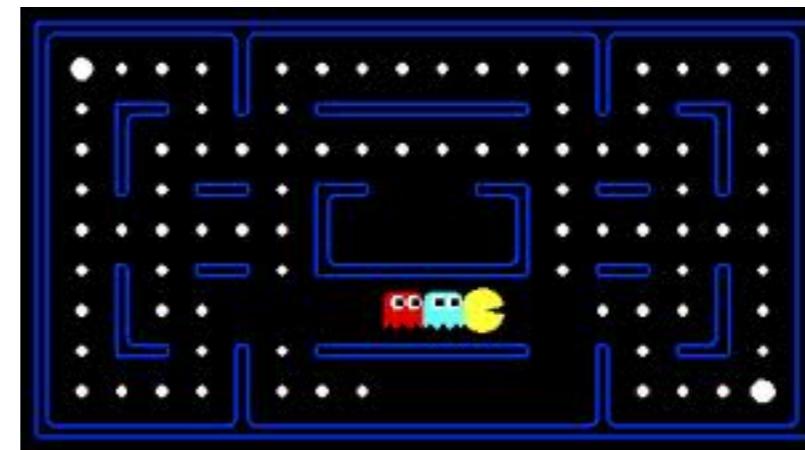
$V_{\text{expectimax}}(s) \rightarrow \text{know } \pi_{\text{opp}}$

$V_{\text{minimax}}(s) \rightarrow \text{do it } + \begin{cases} \text{know } \pi_{\text{opp}} \\ \text{either} \end{cases}$



Games: expectimax

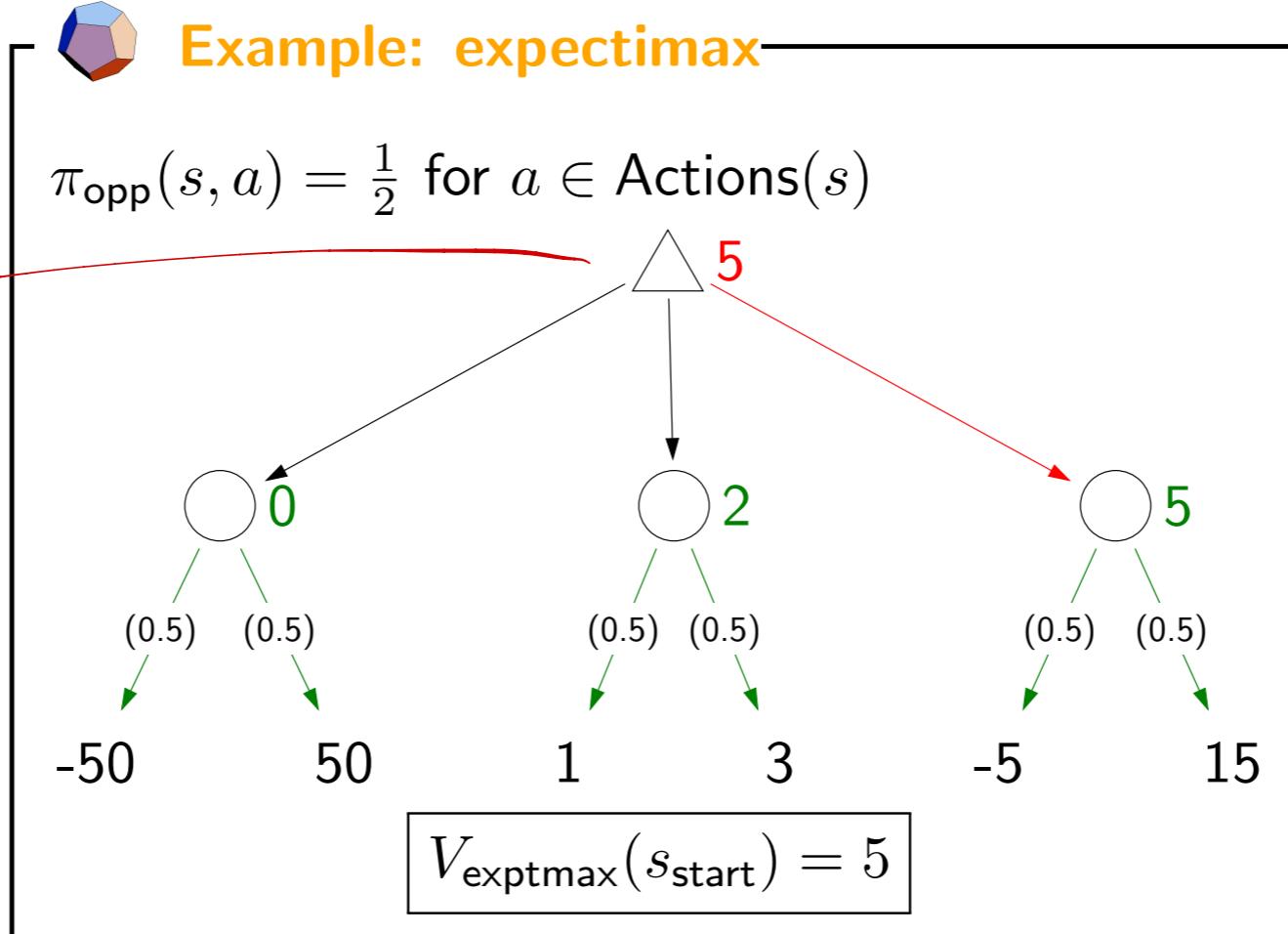
⇒ value iteration



Joint know π_{agent} , know R_{opp}

Expectimax example

take max over states
successor states

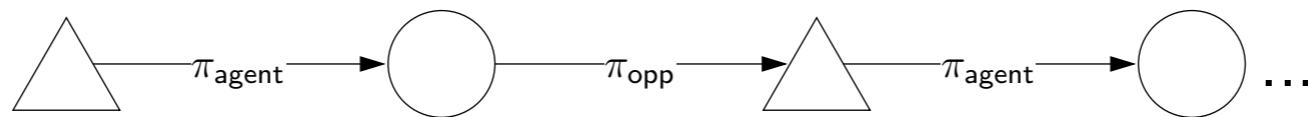


△ is a max node, where you maximize over actions

Expectimax recurrence

Analogy: recurrence for value iteration in MDPs

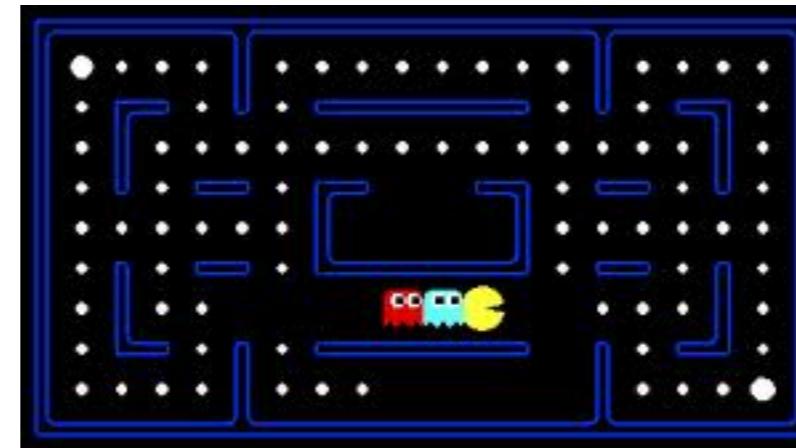
* agent is a maximising agent



$$V_{\text{exptmax}}(s) = \begin{cases} \text{Utility}(s) & \text{IsEnd}(s) \\ \max_{a \in \text{Actions}(s)} V_{\text{exptmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{agent} \\ \sum_{a \in \text{Actions}(s)} \pi_{\text{opp}}(s, a) V_{\text{exptmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{opp} \end{cases}$$



Games: minimax



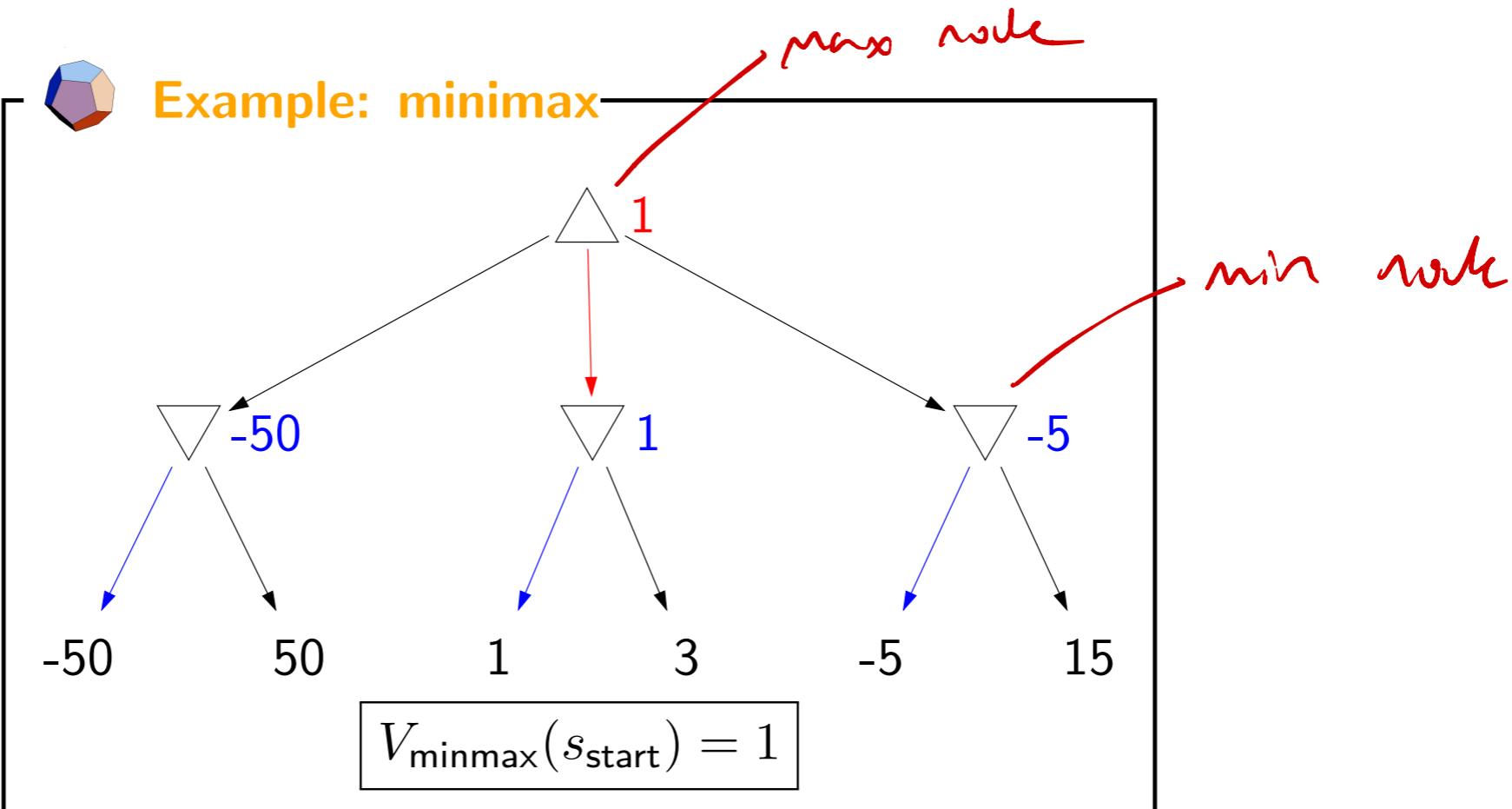
don't know π_{agent} or π_{opp}

Problem: don't know opponent's policy

Approach: assume the worst case



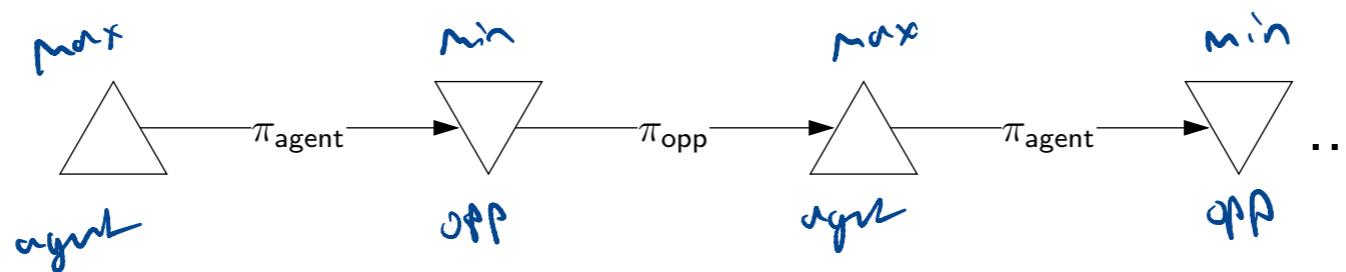
Minimax example



* assume worst case \rightarrow opponent is trying to minimize the agent's utility

Minimax recurrence

No analogy in MDPs:

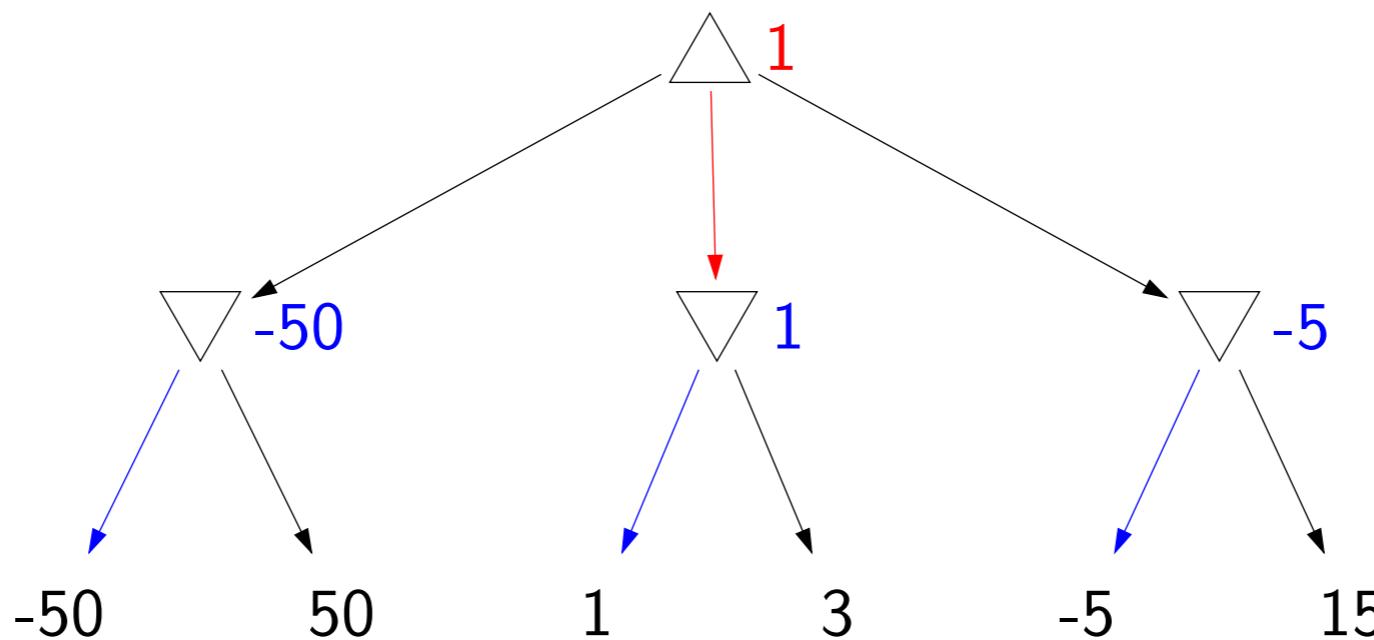


$$V_{\text{minmax}}(s) = \begin{cases} \text{Utility}(s) & \text{IsEnd}(s) \\ \max_{a \in \text{Actions}(s)} V_{\text{minmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{agent} \\ \min_{a \in \text{Actions}(s)} V_{\text{minmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{opp} \end{cases}$$

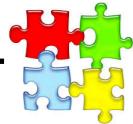
Extracting minimax policies

$$\pi_{\max}(s) = \arg \max_{a \in \text{Actions}(s)} V_{\min\max}(\text{Succ}(s, a))$$

$$\pi_{\min}(s) = \arg \min_{a \in \text{Actions}(s)} V_{\min\max}(\text{Succ}(s, a))$$



The halving game



Problem: halving game

Start with a number N .

Players take turns either decrementing N or replacing it with $\lfloor \frac{N}{2} \rfloor$.

The player that is left with 0 wins.

[semi-live solution: `minimaxPolicy`]

```

def minimaxPolicy(game, state):
    def recurse(state):
        # Return (utility of that state, action that achieves that
        # utility)
        if game.isEnd(state):
            return (game.utility(state), None)
        # List of (utility of succ, action leading to that succ)
        candidates = [
            (recurse(game.succ(state, action))[0], action)
            for action in game.actions(state)
        ]
        player = game.player(state)
        if player == +1:
            return max(candidates)
        elif player == -1:
            return min(candidates)
        assert False

    utility, action = recurse(state)
    print('minimaxPolicy: state {} => action {} with utility {}'.format(
        state, action, utility))
    return action

```

$Utility(s)$

\leftarrow gets the utility

$\max V_{\min}$

$(succ(s, a))$

agent is human

```

policies = {+1: humanPolicy, -1: minimaxPolicy}

state = game.startState()
while not game.isEnd(state):
    # Who controls this state?
    player = game.player(state)
    policy = policies[player]
    # Ask policy to make a move
    action = policy(game, state)
    # Advance state
    state = game.succ(state, action)

```

\leftarrow opponent is playing
optimally w) minimax

Face off

Recurrences produces policies:

*agent knows opponent
is playing π_7*

$$V_{\text{exptmax}} \Rightarrow \pi_{\text{exptmax}(7)}, \pi_7 \text{ (some opponent)}$$

$$V_{\min\max} \Rightarrow \pi_{\max}, \pi_{\min} \quad \text{--- opponent plays this}$$

Play policies against each other:

Agent plays this

<i>agent</i> <i>opponent</i>	π_{\min}	π_7
π_{\max}	$V(\pi_{\max}, \pi_{\min})$	$V(\pi_{\max}, \pi_7)$
$\pi_{\text{exptmax}(7)}$	$V(\pi_{\text{exptmax}(7)}, \pi_{\min})$	$V(\pi_{\text{exptmax}(7)}, \pi_7)$

What's the relationship between these values?

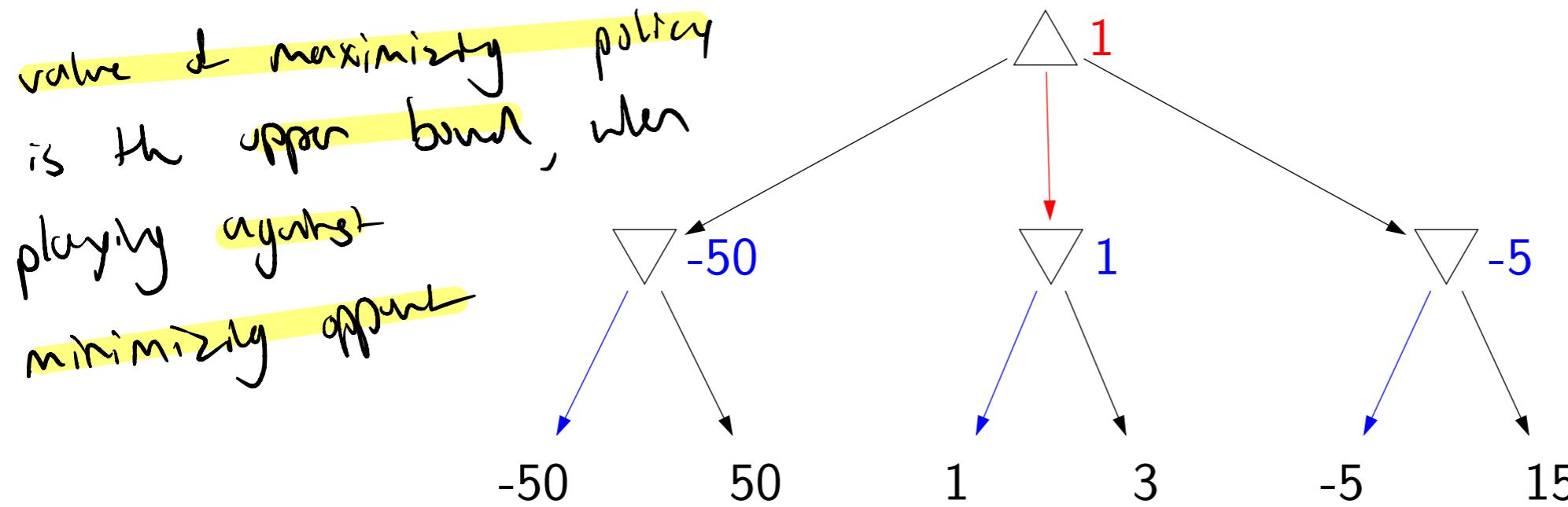
$$V(\pi_{\text{exptmax}(7)}, \pi_{\min}) \leq V(\pi_{\max}, \pi_{\min}) \leq V(\pi_{\max}, \pi_7) \leq V(\pi_{\text{exptmax}}, \pi_7)$$

Minimax property 1



Proposition: best against minimax opponent

$$V(\pi_{\max}, \pi_{\min}) \geq V(\pi_{\text{agent}}, \pi_{\min}) \text{ for all } \pi_{\text{agent}}$$



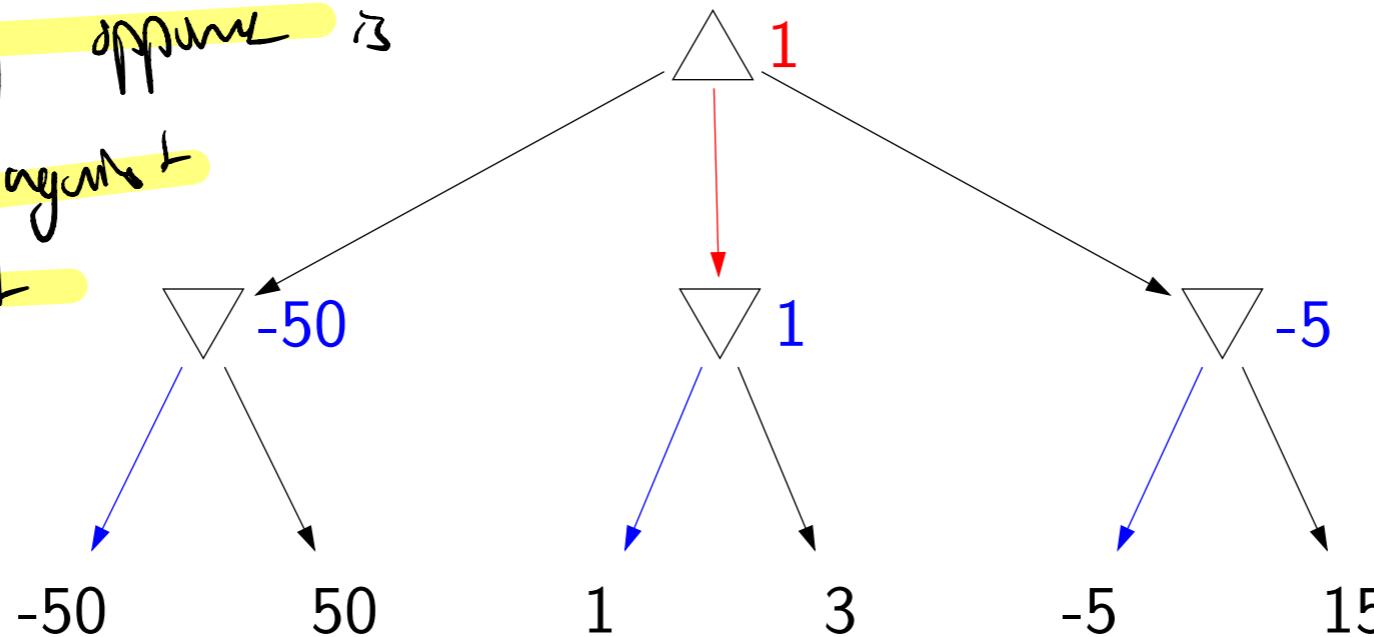
Minimax property 2



Proposition: lower bound against any opponent

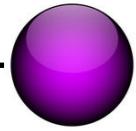
$$V(\pi_{\max}, \pi_{\min}) \leq V(\pi_{\max}, \pi_{\text{opp}}) \text{ for all } \pi_{\text{opp}}$$

value of mining opponent is
lower bound when agent is
a maximizing agent



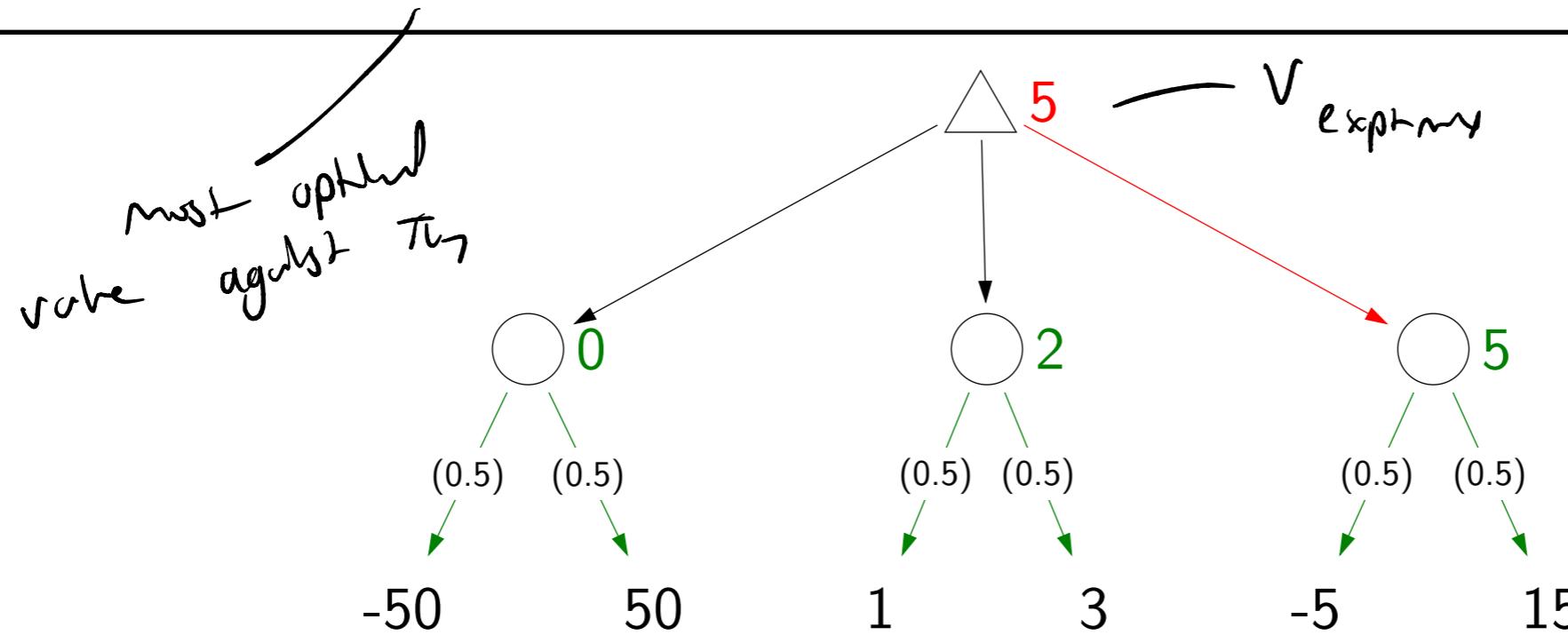
- The second property is the analogous statement for the opponent: if the opponent changes his policy from π_{\min} to π_{opp} , then he will be no better off (the value of the game can only increase).
- From the point of view of the agent, this can be interpreted as guarding against the worst case. In other words, if we get a minimax value of 1, that means no matter what the opponent does, the **agent is guaranteed at least a value of 1**. As a simple example, if the minimax value is $+\infty$, then the agent is guaranteed to win, provided it follows the minimax policy.

Minimax property 3



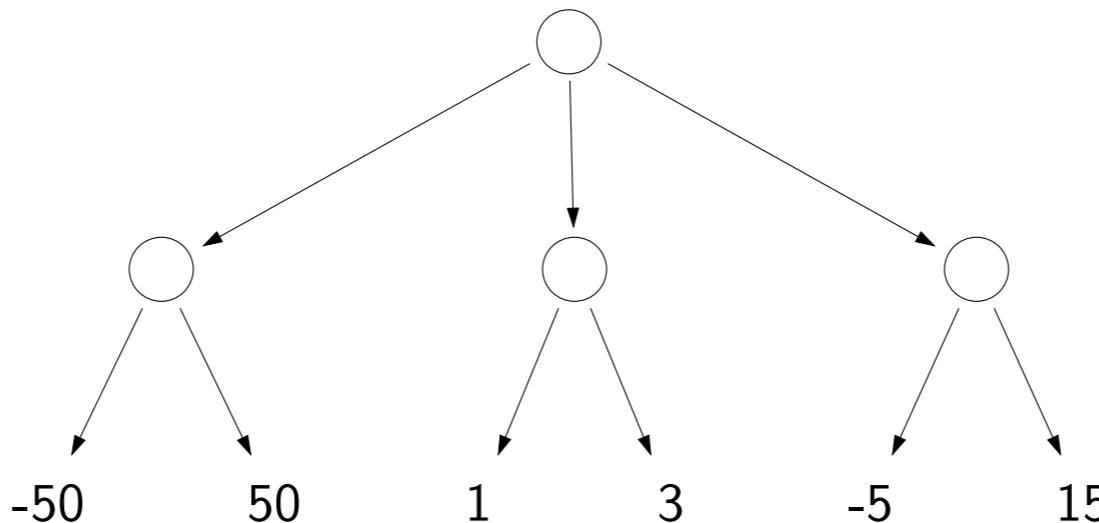
Proposition: not optimal if opponent is known

$$V(\pi_{\max}, \pi_7) \leq V(\pi_{\text{exptmax}(7)}, \pi_7) \text{ for opponent } \pi_7$$



- However, following the minimax policy might not be optimal for the agent if the opponent is known to be not playing the adversarial (minimax) policy.
- Consider the running example where the agent chooses A, B, or C and the opponent chooses a bin. Suppose the agent is playing π_{\max} , but the opponent is playing a stochastic policy π_7 corresponding to choosing an action uniformly at random.
- Then the game value here would be 2 (which is larger than the minimax value 1, as guaranteed by property 2). However, if we followed the expectimax $\pi_{\text{exptmax}}(7)$, then we would have gotten a value of 5, which is even higher.

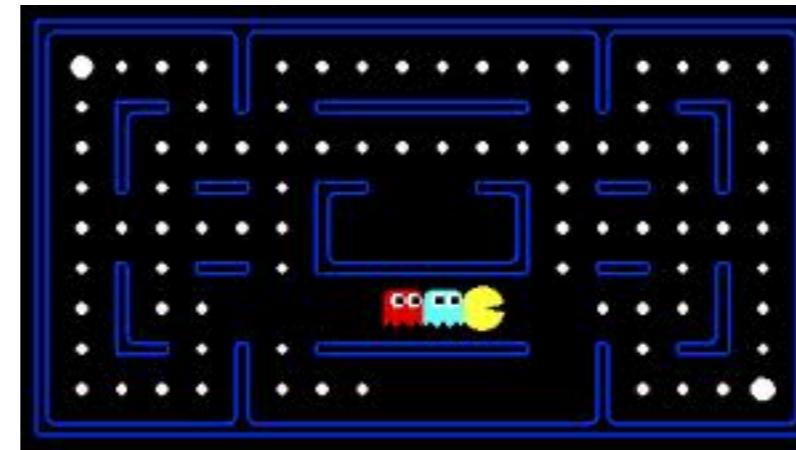
Relationship between game values



π_{\max}	π_{\min}	π_7
	$V(\pi_{\max}, \pi_{\min})$	$\leq V(\pi_{\max}, \pi_7)$
	1	2
	\vee	\wedge
$\pi_{\text{exptmax}(7)}$	$V(\pi_{\text{exptmax}(7)}, \pi_{\min})$	$V(\pi_{\text{exptmax}(7)}, \pi_7)$
	-5	5

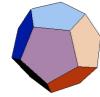


Games: expectiminimax



What if there is uncertainty in nature?

A modified game



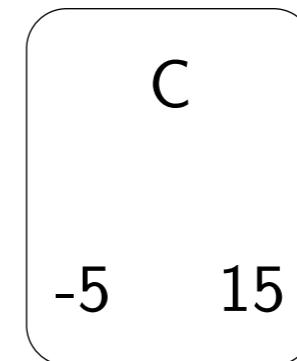
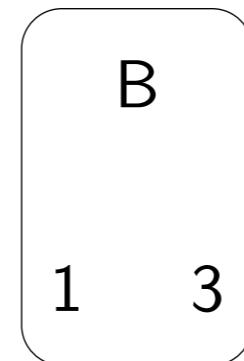
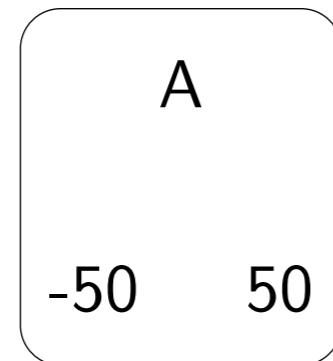
Example: game 2

You choose one of the three bins.

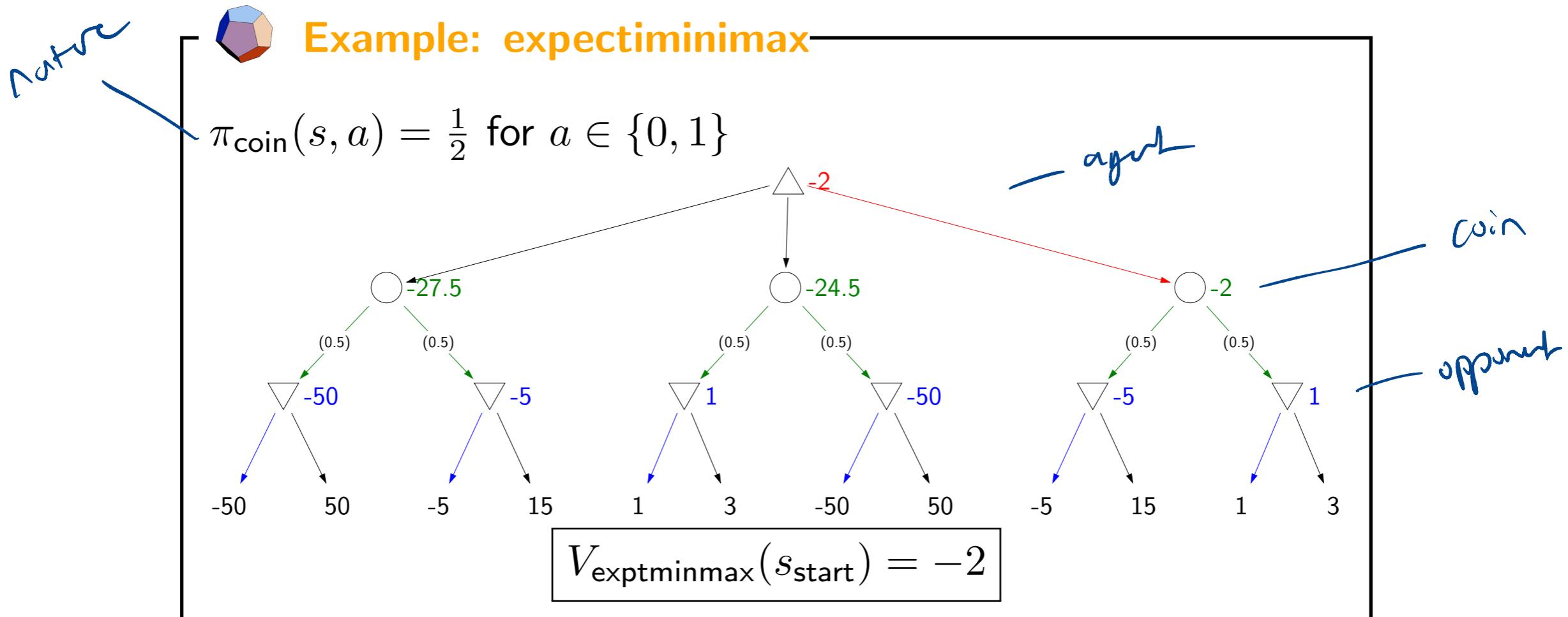
Flip a coin; if heads, then move one bin to the left (with wrap around).

I choose a number from that bin.

Your goal is to maximize the chosen number.



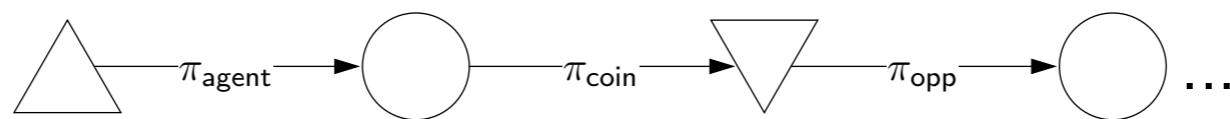
Expectiminimax example



now tries to avoid -50, picks
rightmost branch

Expectiminimax recurrence

Players = {agent, opp, coin}



$$V_{\text{exptminmax}}(s) = \begin{cases} \text{Utility}(s) & \text{IsEnd}(s) \\ \max_{a \in \text{Actions}(s)} V_{\text{exptminmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{agent} \\ \min_{a \in \text{Actions}(s)} V_{\text{exptminmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{opp} \\ \sum_{a \in \text{Actions}(s)} \pi_{\text{coin}}(s, a) V_{\text{exptminmax}}(\text{Succ}(s, a)) & \text{Player}(s) = \text{coin} \end{cases}$$

same as before

↑ we know policy of coin,
just evaluate it



Summary so far

Primitives: **max** nodes, **chance** nodes, **min** nodes

Composition: alternate nodes according to model of game

Value function $V_{\dots}(s)$: recurrence for expected utility

Scenarios to think about:

- What if you are playing against multiple opponents?
- What if you and your partner have to take turns (table tennis)?
- Some actions allow you to take an extra turn?