

## Alpha-Beta Implementation

$\alpha$ : MAX's best option on path to root  
 $\beta$ : MIN's best option on path to root

```
def max-value(state,  $\alpha$ ,  $\beta$ ):  
    initialize v = -∞  
    for each successor of state:  
        v = max(v, min-value(successor,  $\alpha$ ,  $\beta$ )) if v  
            ≥  $\beta$   
        return v  
     $\alpha$  = max( $\alpha$ , v)  
    return v
```

```
def min-value(state,  $\alpha$ ,  $\beta$ ):  
    initialize v = +∞  
    for each successor of state:  
        v = min(v, max-value(successor,  $\alpha$ ,  $\beta$ )) if v  
            ≤  $\alpha$   
        return v  
     $\beta$  = min( $\beta$ , v)  
    return v
```



Algorithm

Week 5.

→ Minimax + Utility with Depth Cutoff Pseudocode

```
def value(state, depth):  
    if the state is a terminal state or depth > cutoff: return the state's  
        utility value  
    if the next agent is MAX: return max-value(state)  
    if the next agent is MIN: return min-value(state)
```

```
def max-value(state):  
    initialize v = -∞  
    for each successor of state:  
        v = max(v, value(successor))  
    return v
```

  
$$V(s) = \max_{s' \in \text{successors}(s)} V(s')$$
  
  

```
def min-value(state):  
    initialize v = +∞  
    for each successor of state:  
        v = min(v, value(successor))  
    return v
```

  
$$V(s') = \min_{s \in \text{successors}(s')} V(s)$$

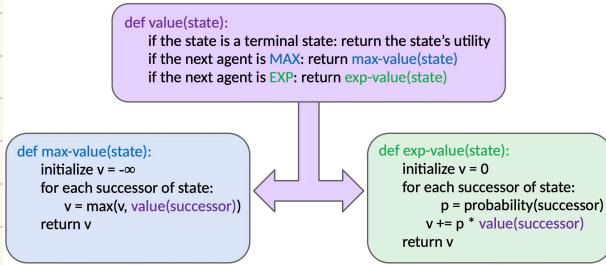
→ Uncertain Outcome

- Average - Case  $\longleftrightarrow$  Expectimax
- Worst - Case  $\longleftrightarrow$  Minimax

↳ Expectimax Search

- Calculate expected utilities

## Expectimax Pseudocode



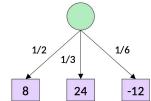
Probability

```

def exp-value(state):
    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v

```

$$v = (1/2)(8) + (1/3)(24) + (1/6)(-12) = 10$$



Tips: no pruning

Instead → Depth-limited Expectimax

深度 limit → 有限时间 而非 optimal

1. depth of ↗
  2. Max node ↗
  3. chance node ↗
  4. Depth-limit ↗
- ↗ Heuristic evaluation function

# → Markov Decision Processes

- **Dangers** { Dangerous Optimism  
Dangerous Pessimism
- Utility → functions from outcome to real number of agent's preference
- Maximum expected utility:

A rational agent chose action  $\rightarrow$  Maximizes its expected utility, given its knowledge

## ↳ Utilities: Uncertain Outcome

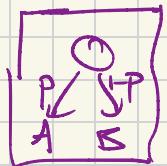
### ↳ Preferences

- Prizes, A, B . . .
- Lotteries,  $L = [P, A; (1-P), B]$

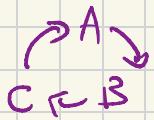
◦ Notion:  $A \succ B$  Preference  
 $A \sim B$  Indifference

- Rational Preferences

↳ Constraint: **Axiom of Transitivity**:



↓  
intransitive preference



## The Axioms of Rationality

### Orderability

$$(A \succ B) \vee (B \succ A) \vee (A \sim B)$$

### Transitivity

$$(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$$

### Continuity

$$A \succ B \succ C \Rightarrow \exists p [p, A; 1-p, C] \sim B$$

### Substitutability

$$A \sim B \Rightarrow [p, A; 1-p, C] \sim [p, B; 1-p, C]$$

### Monotonicity

$$A \succ B \Rightarrow$$

$$(p \geq q \Leftrightarrow [p, A; 1-p, B] \succeq [q, A; 1-q, B])$$

Maximum Expected Utility  
MEU Principle:

$$U(A) \geq U(B) \Leftrightarrow A \succ B$$

$$U([p_1, s_1, \dots; p_n, s_n]) =$$

$$\sum_i p_i U(s_i)$$

## Utility Scale

◦ Normalized Utilities:  $U_+ = 1.0$     $U_- = 0.0$

◦ behavior is invariant under positive linear transformation

$$U'(x) = k_1 U(x) + k_2 \quad \text{where } k_1 > 0$$

## Human Utility

◦ Standard lottery (e.g.  $p + (1-p) = 1$ )

## Money

◦ expected monetary value:  $EMV(L) = p \times X + (1-p) \times Y$

$$U(L) < V(EMV(L))$$

- Risk-averse & risk prone

## ⇒ Non-Deterministic Search

- Markov Decision Processes

↳ A set of states  $S \in S$

↳ A set of actions  $a \in A$

↳ Transition Function  $T(s, a, s') \sim P(s'|s, a)$   
 (also called the model or dynamics)

↳ Reward function  $R(s, a, s')$

↳ Start state / terminal state

- MPS

only care present state (action outcomes depends only on the current state)

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \dots, S_0 = s_0)$$

$$\overset{=} P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

$\overbrace{\quad\quad\quad}$  accumulate  
 ||  
 Present

## ◦ Policies

↳ plan, or sequence of actions, from start to goal

↳ For MDPs, want to optimal policy  $\pi^*: S \rightarrow A$



One that maximizes Expected Utility it follows

↳  $\pi^*$ : gives an action for each state

(Expectimax didn't compute entire policies)

## ◦ MDP Search Trees

Each MDP state  $\longrightarrow$  expectimax-like search tree

$\Delta s$

|

$\bigcirc s, a \quad (s, a, s')$  called a transition

$A s, a, s' \quad T(s, a, s') = P(s'|s, a) \quad R(s, a, s')$

## ◦ Utilities of Sequence

More  $>$  less Utilities / first Come U  $>$  latter Come U

## ◦ Discounting

Value of rewards decay exponentially

$1 \rightarrow r \rightarrow r^2 \dots$

$$1 \begin{cases} A \\ 0 \\ A \end{cases} \rightarrow \text{Sooner reward probability} > \text{later rewards}$$

$$\gamma \begin{cases} A \\ 0 \\ A \end{cases} \rightarrow \text{help algorithms converge}$$

↑  
e.g.  $V([1, 2, 3]) = 1x1 + 0.9x2 + 0.81x3$

$$V([1, 2, 3]) < V([3, 2, 1])$$

## o Stationary Preferences

$$[a_1, a_2, \dots] \succ [b_1, b_2, \dots]$$

↑  
[r, a, a, ...]  $\succ$  [r, b, b, ...]

↳ Additive utility:  $V([r_0, r_1, r_2, \dots]) = r_0 + r_1 + r_2 + \dots$

↳ Discounted Utility:  $V([r_0, r_1, r_2, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

## o Infinite Utilities

game last forever  $\rightarrow$  infinite reward ?

↳ Finite horizon (similar to depth limited search)

↳ Terminate episodes after a fixed T steps

↳ Gives nonstationary policies

↳ Discounting use  $0 < \gamma < 1$  (Smaller  $\gamma$  means smaller "horizon")

$$V(r_0, \dots, r_{\infty}) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max} / (1 - \gamma)$$

最大累积奖赏  $R_{\max}$

## o Optimal Quantities

↳ The Value (utility) of a state  $s$ :

$V^*(s)$  = expected utility starting in  $s$  and acting optimally

↳ The Value (utility) of a  $q$ -state  $(s, a)$

$Q^*(s, a)$  = expected Utility starting out having taken action  $a$  from state  $s$  and acting optimally

Optimal policy

$\pi^*(s)$  = optimal action from state  $s$

## o Values of States

Fundamental operation:

Expected utility under optimal action + Average sum of rewards (discount)

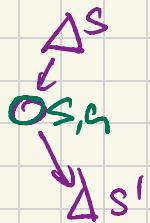
= Expectimax Calculation

## o Recursive definition of Value:

$$V^*(s) = \max_a Q^*(s, a) \quad \text{if } \max (-10, 10)$$

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

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



## ◦ Time-Limited Value

通过深度为  $k$  的 expectation 计算

$$V_k(s)$$

如果  $k$  步结束时，状态  $s$  的最优值

$\left. \begin{array}{l} k \text{ 步数到多少为止} \\ s \text{ 是当前 state} \end{array} \right\}$

## ◦ Value Iteration

↳ Bellman equation characterize the optimal value

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

↳ Value iteration Computes:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

↳ Start with  $V_0(s) = 0$ :  $\rightarrow$  expected reward sum

↳ 给定  $V_k(s)$  值，对每个状态执行一次 expectation

↳ 重复直到 Converges  $\rightarrow$  states 变化足够小，已经 converges optimal

↳ Complexity Iteration  $\mathcal{O}(S^2 A)$  (time)  $S \rightarrow A \rightarrow S$

Converge  $\Rightarrow$  unique optimal value

## o Convergence

- 基础：1. 最大深度  $M$  的树

2. 折扣因子小于 1

-  $V_k \neq V_{k+1} \rightarrow$  深度  $k+1$  的 expectimax 结果

$V_k$ :  $k$  层  $\overline{E}$ , reward  $\rightarrow 0$

$V_{k+1}$ :  $k+1$  层, reward  $\rightarrow$  实际奖励 ( $V_k$ )



## → Fixed Policies

### o Fixed Policy's Expectimax tree

$$S \rightarrow \pi(s) \rightarrow S'$$

固定于某个特定 policy, 可 not optimal, 但 follow policy  $\pi$

### o Utility for a Fixed Policy

Basic operation: Compute the Utility of a state  $S$  under a fix (non-optimal) policy

$$V^\pi(s) = \sum_{S'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

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

## Utilities for a fixed Policy



$$V^\pi(s) = \sum_{s', r} P(s', r | s, \pi(s)) [r + \gamma V^\pi(s')]$$

- $P(s', r | s, \pi(s))$  是在状态  $s$  下, 根据策略  $\pi$  采取行动  $\pi(s)$  后, 转移到状态  $s'$  并获得即时奖励  $r$  的概率。
- $\pi(s)$  是策略  $\pi$  在状态  $s$  下推荐的行动。
- $r$  是即时奖励。
- $\gamma$  是折扣因子, 用于减少未来奖励的当前价值。
- $V^\pi(s')$  是遵循策略  $\pi$  时状态  $s'$  的期望效用。

Utility

## Policy Evaluation

How to calculate  $V$ 's for fixed policy  $\pi$ ?

→ 1. recursive Bellman equations into updates (like iteration)

$$V_0^\pi(s) = 0$$

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

- $V_{k+1}^\pi(s)$ : 这是在策略  $\pi$  下, 状态  $s$  在  $k+1$  次迭代后的价值。它代表从状态  $s$  出发, 遵循策略  $\pi$  所期望获得的总折扣奖励的估计。
- $T(s, \pi(s), s')$ : 这是在状态  $s$  下, 按策略  $\pi$  执行行动  $\pi(s)$  后转移到状态  $s'$  的概率。
- $R(s, \pi(s), s')$ : 这是执行行动  $\pi(s)$  从状态  $s$  转移到状态  $s'$  时获得的即时奖励。
- $\gamma$ : 这是折扣因子, 它决定了未来奖励相对于即时奖励的重要性。
- $V_k^\pi(s')$ : 这是在策略  $\pi$  下, 状态  $s'$  在  $k$  次迭代后的价值。

公式中的求和操作将每个可能的下一个状态  $s'$  的概率乘以在该状态下可能获得的奖励, 然后将所有这些乘积累加起来。这样做是为了计算当前状态  $s$  在遵循固定策略  $\pi$  下的期望总折扣奖励。

→ 2. linear system

Matrix

## → Policy Extraction

### ◦ Computing Actions from Values

每个行动的期望回报:

动作  
回报

$$\underline{Q}(s, a) = \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma V^*(s')]$$

最佳行动:

$$\pi^*(s) = \arg \max_a Q(s, a)$$

最大      s state do a , 所期望得到的总回报

### ◦ Computing Actions from Q-Value

$$\pi^*(s) = \arg \max Q^*(s, a)$$

## → Policy Iteration

### ◦ Problem with Value Iteration

→ slow -  $O(S^2 A)$  per iteration

→ The "max" at each state rarely changes

→ Policy converges before the value

## 0 Approach for Optimal Value

o Policy Evaluation: Calculate utilities for fixed policy until Convergence

$$+ \quad V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

o Improvement:

update policy using one step look-ahead with resulting converged utilities as future value

Repeat steps until policy converges

$$\pi_{i+1}(s) = \arg \max \sum_{s'} T(s, a, s') [\hat{r}(s, a, s') + \gamma V^{\pi_i}(s')]$$

→ optimal, Converge faster

## → Comparison (Value & policy)

Value      {  
o Every Value Iteration updates both Value and (implicitly) policy  
o Don't track policy, but take max over actions

Policy      {  
o policy iteration { evaluation - do several time, Everytime do one action  
Improvement - chose new policy

