

Agents: concept

right function: performance measure for environment

rationality: maximum "expected utility"

rational agents: no omniscient, no clairvoyant

PEAS: performance measure, environment
actuator (执行), sensor (感应器)

environment type: 判定方法:

deterministic (决定的), episodic (不连贯的)

reflex agent can be rational

NP-hard solutions can be verified in polynomial time

PSPACE can be solved in polynomial amount of space

EXPTIME can be solved in exponential time $O(2^{n^m})$

Search:

search problem consists of:

state space action set Actions(s)

transition model Succ(s,a) start state

step cost function c(s,a,s') goal test

state space graph: search problem 教学表示

Nodes: world states Arcs: successors

General tree search

problem, strategy \Rightarrow solution / failure

initialization

loop do

no candidates \Rightarrow failure

choose a leaf node to expand (strategy)

node contains goal state \Rightarrow return

expand the node and adding -

-resulting nodes in the tree

DFS: $O(b^m)$ space: $O(bm)$ no optimal

BFS: $O(b^s)$ space: $O(b^s)$ optimal iff cost $\equiv 1$

Uniform Search: (find the least-cost)

永远探索最近的点

time $O(b^{\frac{s}{\epsilon}})$ C*: solution cost, ϵ : arc cost at least

C^*/ϵ tiers: space $O(b^{\frac{C}{\epsilon}})$ optimal!

Tree: explore options in every direction

no information about goal location

Uniform Search

每次都像BFS一样同层扫描, explore cost ≤ 1 的点

same property for search algorithm:

1. all fringe are priority queues

2. DFS and BFS can avoid $\log(n)$

3. one implementation \Rightarrow variable queueing object

Heuristic (manhattan, Euclidean)

admissible heuristic: $h^*(n)$: true cost for

$0 \leq h(n) \leq h^*(n)$ nearest goal

Dominance: $h_a \geq h_b$ if $\forall n: h_a(n) \geq h_b(n)$

max of admissible heuristic is admissible

consistency of heuristic:

estimated cost \leq actual costs

heuristic design: often use relaxed problems

A* Search

$g(n)$: path cost $h(n)$: goal proximity
 $\Rightarrow f(n) = g(n) + h(n)$

A* search is optimal iff $h(n)$ is admissible

Graph Search

像 tree search - 样, 但是不再访问 visited nodes
consistent \Rightarrow A* expanded nodes increase \Rightarrow A* graph
heuristic in its suboptimal 要先到达 \Rightarrow search is optimal

Search Games:

Initial state: $s_0 \in S$ Players: Players(s) {1, 2, ..., N}

Actions: $a \in \text{Actions}(s)$ Transition: $\text{Succ}(s, a) \in S$

Terminal test: Terminal-Test(s)

Terminal value: Utility(s, p) for player P

Minimax Search: DFS

① minimax-decision(s):

action = argmax(min-value((succ(s, a))))

return action (in Actions(s))

② min-value(s):

if terminal-Test(s) \Rightarrow return (Utility(s))

$V = +\infty$

for a in Actions(s):

$V = \min(V, \text{max-value}(\text{succ}(s, a)))$

return (V) 在所有值中取 min (对手)

③ max-value(s):

if terminal-Test(s) \Rightarrow return (Utility(s))

$V = +\infty$

for a in Actions(s):

$V = \max(V, \text{min-value}(\text{succ}(s, a)))$

return (V) (对手)

Another way:

① minimax-decision(s):

return argmax_a(value(succ(s, a)))

② value(s):

if terminal-test(s) \Rightarrow Utility(s)

if Player(s) = MAX \Rightarrow max(Value(Succ(s, a)))

if Player(s) = MIN \Rightarrow min()

minimax properties:

adversary assumed to be rational

minimax is optimal against perfect opponent

time $O(b^m)$ Space $O(bm)$

处理 limits: 层 level 用评估函数来赋值

use iteration deepening for anytime algorithm ??
deeper is better, evaluation function is confusing

2. Game tree pruning: 发现没用就剪掉

α : MAX's best choice β : MIN's best choice

max-value(state, α, β)

$v = -\infty$

for Successor in $\text{succ}(s, a)$:

$v = \max(v, \text{min-value}(\text{Successor}, \alpha, \beta))$

if $v \geq \beta \Rightarrow \text{return } v$

$\alpha = \max(\alpha, v), \text{return } v$

min 也一样, 把 $v \geq \beta$ 改成 $v \leq \alpha$, max \Rightarrow min, $\alpha \leq \beta$

time ($O(b^2)$) double solvable depth!

Search Chance

1. Independent Decision-making 每人分别行动

2. Joint State/Action Spaces: 结合所有 agents

3. Coordinated Decision Making: confirm joint plans

Each Agent proposes their actions and computer

4. Alternate Searching One Agent a Time

Generalized minimax: 每人都有 Utility, 不同的维度

Expectimax

① decision(s):

return $\operatorname{argmax}_a \text{Value}(\text{succ}(s, a))$

value(s):

If Terminal-test(s) \Rightarrow return Utility(s)

If Player(s)=MAX \Rightarrow return max among actions

If Player(s)=CHANCE \Rightarrow sum_a Pr(a) * Value(succ(s, a))

Expectiminimax

前代码加一行: If = MIN \Rightarrow return min

Property: minimax are invariant with respect to monotonic

expectiminimax are with positive affine transformation

should take real value

Maximum Expected Utility:

Outcome: 結果 lotteries: distribution Preference: \approx

completeness: $L \approx L'$ or $L' \approx L$ Transitivity: 可传递

Independence: $(L \approx L') \Rightarrow [p, L; 1-p, L'] \approx [p, L; 1-p, L']$

Continuity: $(L \approx L' \approx L'') \Rightarrow [p, L; 1-p, L''] \approx [p, L; 1-p, L'']$

\approx satisfies 4 axioms iff $L \approx L' \Leftrightarrow EV(L) \geq EV(L')$

但大部分的EV不会满足 $EV(L) = \sum_{o \in O} p(o) V(o) \Rightarrow$ convex

Decision Theory: pick a strategy to maximize utility

given world outcomes

Game Theory: given other player's strategy

Game Theory

除3下以外策略合集

dominant: $\forall k \neq i, u_i(\pi_{i,k}, \bar{\pi}_{-i}) > u_i(\pi_{i,k'}, \bar{\pi}_{-i}) \Rightarrow$ strictly dominated

对于玩家来说,无论其他人出什么,选择策略都 max

Nash Equilibrium: $\forall i, u_i(\bar{\pi}) \geq u_i(\pi'_i, \bar{\pi}_{-i})$

当其它人不改变策略时,没有人会改变

strict nash equilibrium: $\forall i, u_i(\bar{\pi}) > u_i(\pi'_i, \bar{\pi}_{-i})$

寻找 pure nash 方法: 找出 dominant \Rightarrow 剔掉被 dominant 的

下划线法: 划出每一行 and 每一列最大值,有重合即为纳什均衡

Theorem: finite number of players+actions $\Rightarrow \exists$ nash

混合 nash: 行序法: 设 A 为 p_1, \dots, p_n , B 为 q_1, \dots, q_n

$\Rightarrow p_1 + \dots + p_n = 1, q_1 + \dots + q_n = 1$ 且 A 的配置需要让 B 觉得每件事期望相同
correlated Equilibrium: 即站在上帝视角

$\forall \pi'_i \in \Pi_i, \sum_{\pi_{-i}} P_{\pi_i}(\pi_i, \pi_{-i}) u_i(\pi_i, \pi_{-i}) \geq \sum_{\pi'_i} P_{\pi_i}(\pi_i, \pi_{-i}) u_i(\pi'_i, \pi_{-i})$

Pareto 最优: 每个人都达到了最大 outcome

Markov Decision Process

s: set of states a: set of actions

T(s, a, s'): transition probability (如果采取 a, 到 s' 的可能性是多少)

R(s, a, s'): reward

MDP: non-deterministic problems, present 和 future 是独立的

π^* : π give an action for each state. π^* 达到最优

discounting: r 每走一步就乘以 r

如何解决 infinite Utilities?

1. 限制 horizon (finite step or finite \Rightarrow policy for π)

2. Discounting

$V^*(s)$: value of utility starting in s

$Q^*(s, a)$: value of q-state(s, a)

$V_k(s)$: the optimal value of s in k steps

$O(s^2A)$ and will converge to optimal

公式大全

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

$$Q^*(s, a) = \sum_s T(s, a, s')[R(s, a, s') + r V^*(s')]$$

$$\text{Bellman: } V^*(s) = \max_a \sum_s T(s, a, s')[R(s, a, s') + r V^*(s')]$$

$$\text{Value: } V_{k+1}(s) \leftarrow \max_a \sum_s T(s, a, s')[R(s, a, s') + r V_k(s')]$$

$$Q\text{-iteration: } Q_{k+1}(s, a) = \sum_s T(s, a, s)[R(s, a, s') + r \max_a Q_k(s', a')]$$

$$\text{Policy extraction: } \pi_{V_{k+1}} = \arg \max_a \sum_s T(s, a, s)[R(s, a, s') + r V_{k+1}(s')]$$

$$\text{Policy evaluation: } V_{k+1}^{\pi}(s) = \sum_s P(s'|s, \pi(s)) [R(s, \pi(s), s') + r V_k^{\pi}(s')]$$

$$\text{Policy improve: } \pi_{\text{new}}(s) = \arg \max_a \sum_s P(s'|s, a) [R(s, a, s') + r V_{k+1}^{\pi}(s')]$$

Policy evaluation: $O(S^2)$ per iteration

value iteration: $O(AS^2)$ slow / policy converges before max

Reinforcement Learning

receive feedback in the form of rewards, max expected difference from MDP: $X \rightarrow R$

Direct evaluation: average all the sampled values under waste information about state connections, long-time

Temporal Difference Learning

实时更新 $V(s)$, policy fix

$$\text{sample} = R(s, \pi(s), s') + r V^\pi(s')$$

$$V^\pi(s) = (1-\alpha)V^\pi(s) + \alpha \text{sample} = V^\pi(s) + \alpha(\text{sample} - V^\pi(s))$$

exponential moving average: $\alpha \downarrow \Rightarrow$ give converging

TD value Learning problems: 也可用改成 Q learning

$$\text{sample} = R(s, a, s') + r \max_a Q(s, a)$$

$$Q(s, a) = (1-\alpha)Q(s, a) + \alpha \text{sample} \quad \text{suboptimally}$$

Q -learning converges to optimal policy even if acting

前提: explore enough; learning rate small enough

Explore Scheme

random actions: ϵ : act randomly $1-\epsilon$: follow policies

\Rightarrow lower ϵ overtime

Exploration Function: u -value estimate n -visit count

$$f(u, n) = u + \frac{1}{n} \text{ modified Q-Update}$$

$$Q(s, a) = \alpha R(s, a, s') + r \max_{a'} f(Q(s, a'), N(s, a'))$$

Constraint Satisfaction Problems

constraint Graphs:

① each constraint relates 2 variables, arcs are constraints

Search Method:

DFS-BFS - Backtracking Search

Backtracking search is the basic uninformed algorithm
每次检查一个 node, 如果不符合 constraint 就后退

Backtracking Search 伪代码

backtracking-search(csp):

return recursive-backtracking(f, csp)

recursive-backtracking(assignment, csp)

if assignment is complete \Rightarrow return assignment

var = select-unassigned-variable(variable[csp], csp)
 $\xrightarrow{\text{assignment}}$

for value in order-domain(var, assignment, csp):

if value is consistent $\xleftarrow{\text{constraints}} \text{constraints}[csp]$ \Rightarrow

add [var = value] to assignment

result = recursive-backtracking(assignment, csp)

if result \neq failure \Rightarrow return result

remove {var = value} from assignment

return failure

改进算法方法:

ordering, filter, structure

① filter: Forward checking: enforcing consistency of arcs pointing to each new assignment

遍历每一个点, 把不符合的踢出去

Consistency

Arc: $X \rightarrow Y$ is consistent iff for every $x \in X$ there is $y \in Y$ such that

Arc consistency of an Entire CSP: all arcs are consistent

伪代码:

AC-3(csp): inputs: csp, a binary CSP

local variables: queue, a queue of arcs

while queue is not empty:

(A, B) = remove-first(queue)

if remove-inconsistent(A, B) \Rightarrow

for node in neighbours[A] then

add (node, A) to queue

remove-first(A, B) returns true iff succeeds

remove = false

for x in Domain[A]

if no value y in Domain[B] allows(x, y) \Rightarrow consistent
then delete x from Domain[A]

removed = true

return removed

$O(n^2d^3) \rightarrow$ reduce $O(n^2d^2)$

Ordering: minimum reducing value

倾向于去掉可以留下最少的点,

Least constraining value

倾向于去掉 constraints 最小的点,

K-consistency: for each $k-1$ nodes, could be extended
strong k-consistency: $k=1, k=2, \dots, n$ consistent

Tree-structured CSPs

if the constraint graph has no loop, CSP $O(nd^2)$

① order: order variables from parents to children

② remove backward:

for $i=n:2$, RemoveInconsistent(Parent(x_i), x_i)

for $i=1:n$, assign x_i consistently with Parent(x_i)

Cutset Conditioning

选一个点, 使得去掉这个点, 就可以变成 tree

再当此点为所处情况时遍历考虑.

Tree Decomposition

分割成一个个小问题, 用共同边连接

Iterative Algorithm

variable selection: 随机分配

choose heuristic with min constraints

Climb hill: local and global and flat

simulated annealing: 给一定几率往 next 随机

simulated-annealing(problem, schedule)
inputs: problem, a mapping from time to "temperature"
for t from 1 to ∞ :

$T = \text{schedule}(t)$

if $T=0 \Rightarrow$ return current

next = randomly selected successor of current

$\Delta E = \text{value}[\text{next}] - \text{value}[\text{current}]$

$\Delta > 0 \Rightarrow \text{current} = \text{next}$

else: $\text{current} = \text{next}$ only if probability $e^{\frac{\Delta E}{T}}$