Lecture 2

① Classical Planning. state Model S(P)
  ① state space S
  ② initial state s0 ∈ S
  ③ ~~set~~ set SG ⊆ S, goal states
  ④ actions A(s) ⊆ A
  ⑤ deterministic transition function $s' = f(a,s)$
  ⑥ action costs $c(a,s)$

Solutions: actions that map s0 into SG

② Blind search
  DFS. BFS. UCS ID.
  Heuristic Search
    A*. WA*. IDA*. Hill climbing.
    Best first search.

③ completeness: ~~goo~~ guaranteed to find a solution
  optimality: solution optimal?

  time complexity    Branching factor b 有几个分叉
  space complexity   goal depth d 搜几次

④ Breadth-first-search: 找到就了. 层层. FIFO queue
  Uniform cost search: 每次找 (g(s)) 最小后 ~~先扩~~ bfs pray queue
  depth-first search: 不撞南墙不回头 LIFO stack.
  Iterative deepening: combine bfs and dfs (stack.)
    limit = 0. BFS
    limit = 1  BFS
    limit = 2  BFS.

⑤. Systematic (Heuristic Search:)  h* perfect heuristic
                                    remaining cost.
    Greedy best first search ← 根据h排 priority queue
    WA*  g(s)+W×h(s).  h(s) 排序了好
    A*: g(s)+h(s) ←  priority queue
    IDA*            g: 移动成本.
                    h: 估算成本
  Local Heuristic Search:
    Hill-climbing
    Enforced-hill climbing

⑥ safe: $h^*(s) = \infty \Rightarrow h(s) = \infty$ for $s \in S$
  goal-aware: $h(s) = 0$ for goal states
  admissible: 表 $h(s) \le h^*(s)$ @ for all s
  consistent: $h(s) \le h(s') + c(a)$ for a, $s \to s'$
    父点的h ≤ 子点的 h + cost.
  据此 heuristic function 分台湾
  consistent + goal-aware ⇒ admissible
  admissible ⇒ goal-aware
  admissible ⇒ safe

⑦ ~~Greedy best~~ hill-climbing
  每次找局部最大 smaller h 右则停其
  Enforced hill-climbing   only if h(s)>0
    + bfs + smaller h-value.

⑧

| | DFS | BFS | ID | A* | HC | IDA* |
|---|---|---|---|---|---|---|
| Complete | No | Yes | Yes | Yes | No | Yes |
| Optimal | No | Yes* | Yes | Yes | No | Yes |
| Time | ∞ | $b^d$ | $b^d$ | $b^d$ | ∞ | $b^d$ |
| Space | b·d | $b^d$ | b·d | $b^d$ | b | b·d |

A*. IDA* optimal | if h is admissible
  $h \le h^*$

BFS optimal | if costs are same

① **Classical Planning**
state space $S$
initial state $s_0$
$S_G$
actions $A(s)$
deterministic $s' = f(a,s)$
action cost $c(a,s)$

**Conformant planning**
state space $S$
set of possible $S_0$
$S_G$
actions $A(s)$
non-deterministic $F(a,s)$
actions with $c(a,s)$

③ **Example** ⊓⊓⊓

**Pre** : on$(x,y)$, onTable$(x)$, clear$(x)$
holding$(x)$, armEmpty$()$

**Actions** : stack$(x,y)$, unstack$(x,y)$
Putdown$(x)$, pickup$(x)$

stack$(x,y)$: pre $\{$holding$(x)$, clear$(y)\}$
add $($on$(x,y)$, clear$(x)$, armEmpty$()$$)$
del $\{$~~holding~~ clear$(y)$, holding$(x)$$\}$

---

**MDPS**
state Space $S$
$s_0$
$S_G$
$A(s)$
transition probabilities $P_a(s'|s)$
$c(a,s)$

**POMDPS**
States $S$
actions $A(s)$
transition prob $P_a(s/s)$
initial belief state $b_0$
final belief state $b_f$
sensor model $P_a(o|s)$

① actions - deterministic + initial location known
$\Rightarrow$ classical

② action - stochastic + location observable
$\Rightarrow$ MDP

③ action stochastic + location partially observable
$\Rightarrow$ POMDP.

② **STRIPS** ⭐ $P = \langle F, O, I, G \rangle$

$F \Rightarrow$ atoms (boolean)
$O \Rightarrow$ operators (actions) $\longleftarrow$ add list Add$(o)$
$I \Rightarrow$ initial situation delete list del$(o)$
$G \Rightarrow$ goal situation. precondition list pre$(o)$

states $(s \in S)$ are collections of atoms $F \cdot S = 2^F$
initial state $s_0$ is $I$.
goal states. $G \subseteq S$
actions $a$ in $O$ st. prec$(a) \subseteq s$
next state $s' = S -$ Del$(a) +$ Add$(a)$
action costs $C(a,s) = 1$

---

**Lecture 12** Reward shaping
add additional rewards

$Q(s,a) = Q(s,a) + \alpha[r - F(s,s') + \gamma \max Q(s',a')$
$- Q(s,a)]$

$$\boxed{F(s,s') = \gamma \, \overline{\Phi}(s') - \overline{\Phi}(s)}$$
potential function

$\overline{\Phi}(s)$ $\dfrac{1}{|x(g) - x(s)| + |y(g) - y(s)|}$
(goal) Manhattan distance

**n-step temporal difference learning**
discounted future reward

$G(t) = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^n V(s_{t+n})$
n step

$G_t^n = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + v$

$Q(s,a) = Q(s,a) + \alpha[G_t + \gamma^n Q(s',a')$
$- Q(s,a)]$

$$G_t = \sum_{i=t+1}^{t+n} \gamma^{i-t-1} r_i$$

① $P$, $h^*$

define $P'$, $h'^*$ used to estimate $h^*$

[transformation $r$] $\Rightarrow$ $P \Rightarrow P'$

given $\Pi \in P$ abbr $h^*(\Pi)$ by $\boxed{h'^*(r(\Pi))}$

② [Relaxation] ☆

relaxation of $h^*$ is $R = (P', r, h'^*)$

$r : P \to P'$ and $h'^* : P' \to R_0^+ \cup \{\infty\}$

$h^R(\Pi) = h'^*(r(\Pi)) \Leftarrow \boxed{h^R(\Pi) \leq h^*(\Pi)}$ ☆

① [native] if $P' \subseteq P$ and $h'^* = h^*$

② [efficiently] constructible : $r(\Pi)$ 可构造

③ [efficiently] computable : $h'^*(\Pi')$ 可计算

③ $R = (P', r, h'^*) :$ remove preconditions and deletes, $h^*$

Native ? Yes.

eff. constructible : Yes.

eff. computable : NO, NP-hard

$\Rightarrow$ Approximate $h^*$

[lecture 5] delete Relaxation Heuristic

① delete Relaxation : $\boxed{A^+ : \text{set of relaxed actions}}$

STRIPS. planning task $\Pi = (F, A, c, I, G)$

$\Pi^+ = (F, A^+, c, I, G)$

② dominance $\boxed{s'^+ \text{ dominates } s^+}$ (delete)

③ $h^+(s)$ the cost of $\boxed{\text{optimal relaxed plan for } s}$

$h^*$ : optimal plan. $\boxed{h^* = 20}$ $\boxed{h^+ = 10}$

用 $h^{add}$ and $h^{max}$ to estimate $h^+$.

$h^{add}(s,g) = \begin{cases} 0 & g \subseteq s \\ \min_{a \in A, g \in add_a} c(a) + h^{add}(s, pre_a) & |g| = 1 \\ \sum_{g' \in g} h^{add}(s, \{g'\}) & |g| > 1 \end{cases}$

$h^{max}(s,\bar{s}) = \begin{cases} 0 \\ \min_{a \in A, g \in add.} \underline{c(a)} + \boxed{h^{max}}(s, pre_a) & |g| = 1 \\ \boxed{\max}_{g \in g} h^{max}(s, \{g'\}) & |g| > 1 \end{cases}$

$\boxed{h^{max} \leq h^+} \Rightarrow h^{max} \leq h^* \quad |\text{optimist}$

$h^{add} \geq h^+ \Rightarrow \text{exist } h^{add}_{(s)} > h^*_{(s)}$

$\boxed{\text{pessimistic}}$

④ $\boxed{\text{Best supporter}}$ from $h^{add}$, $h^{max}$.

$bs^{max}(P) = \boxed{arg} \min c(a) + h^{max}(s, pre}$

$bs^{add}(P) = \boxed{arg} \min c(a) + h^{add}(s, pre}$

$\boxed{\text{得出的 action}}$ 画走. 得出是小

cost 的 action

$\Rightarrow \boxed{h^{ff}(\Pi) \text{ 所有 action}}$

$h^{ff} > h^+ \Rightarrow \text{exist } h^{FF}(s) > h^*(s)$

$\boxed{\text{may be inadmissible}}$

[lecture 6] [Iterated width]

① $IW(k)$ BFS that prunes newly generated states whose $\boxed{\text{novelty}(s)} >$



没出现过的最多少

Novelty table

| | |
|---|---|
| P | 1 |
| a | 1 |
| r | 1 |
| Pa | 1 |
| PR | 1 |
| ar | 1 |

$\boxed{\text{Lecture } 8}$ MDP

Bellman equations

$$V(s) = \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s)[r(s,a,s') + \gamma V(s')]$$

next state

$$Q(s,a) = \sum_{s \in S} P_a(s'|s)[r(s,a,s') + \gamma V(s')]$$

action

$$V(s) = \max_{a \in A(v)} Q(s,a).$$

最大 Q值的 action.

$\boxed{\text{Lecture } 9}$ Value policy iteration

Value Iteration

$$V(s) = \max_{a \in A} \sum_{s \in S} P_a(s'|s)[r(s,a,s') + \gamma V(s')]$$

$$V^{\pi}(s) = \arg\max_{a \in A} Q^{\pi}(a,s)$$

action

$$Q^{\pi}(a,s) = \sum_{s \in S} P_a(s'|s)[r(s,a,s') + \gamma V^{\pi}(s')]$$

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

$\boxed{\text{Lecture } 10}$ Monte Carlo Tree Search

① Select ② expand ③ simulate ④ Backpropagation

$$V(s) = \max \sum P_a(s'|s)[r(s,a,s') + \gamma V(s')]$$

$\boxed{\text{UCB}}$ $\pi(s) = \arg\max Q(s,a) + \sqrt{\dfrac{2\ln N(s)}{N(s,a)}}$

$\boxed{\text{UCT}}$ $\pi(s) = \arg\max Q(s,a) + 2 \times C_P \cdot \sqrt{\dfrac{2\ln N(s)}{N(s,a)}}$

$\boxed{\text{Lecture } 11}$

Reinforcement learning
$\boxed{Q - \text{Learning}}$
$\boxed{\text{SARSA}}$

Q Learning : off policy

$$Q(s,a) = Q(s,a) + \alpha[r + \gamma \cdot \max Q(s',a') - Q(s,a)]$$

$\boxed{\text{SARSA}}$ on-policy

$$Q(s,a) = Q(s,a) + \alpha[r + \gamma \cdot Q(s',a') - Q(s,a)]$$

下一动作

$\boxed{\text{Lecture } 12}$

$\boxed{\text{Approximate Q-function}}$

$$f(s,a) = \begin{cases} f_1(s,a) \\ f_2(s,a) \\ \vdots \\ f_n(s,a) \end{cases}$$

feature vector

weight vector

$$W_i^a \leftarrow W_i^a + \alpha[r + \gamma \max Q(s',a') - Q(s,a)] \cdot f_i(s,a)$$

Q-learn

$$Q(s,a) = f_1(s,a) \cdot W_1^a + f_2(s,a) \cdot W_2^a + \cdots = \sum_{i=0}^{n} f_i(s,a) \cdot W_i^a$$

h_add

[A] [B] | [A/B/C]  →

| | cl(A) | cl(B) | cl(c) | oT(A) | oT(B) | oT(c) | on(A,c) | on(A,B) | on(B,c) | h(A) | h(B) | h(c) | onTree |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ∞ | ∞ | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | ∞ | 1 | ∞ | ∞ | 1 | 0 |
| 2 | 0 | 0 | 0 | 2 | 8 | 8 | 0 | 2 | 3/2 | 8 | 1 | 2 | 0 |

add / max

goal: { on(A,c), oT(c), on(B,c), oT(B) }

Initial: on(A,c), oT(c), on(B,c), oT(B), cl(A), cl(B)

cl(c) → unstack(A,c)  cost = 1
  └ pre: { on(A,c), cl(A), uT } ⇒ 1.

oT(A) → putdown(A)  cost = 1
  └ pre: h(A) ⇒ ∞ ⇒ 2
          ∞

on(B,c) → stack(B,c)
  └ pre: { cl(c), h(B) } =(∞) ⇒ 2
          0   ∞

on(A,B) → stack(A,B)  cost = 1
  └ pre: { cl(B), h(A) } ⇒ ∞ ⇒ 2
          0   ∞

h(A) → unstack(A,c)  cost = 1
  └ pre: { uT, cl(B), oT(B), on(A,c) } ⇒ 1
          0   0   0   0

h(B) → pickup(B)
  └ pre: { cl(B), cl(B) } ⇒ 1.
          0    0

h(c) → pickup(c)
  └ pre: { cl(c), uT } h(c)
          0    ∞   ⇒ ∞

bs(A) → { unstack(A,c) }
  └ pre: { on(A,c), cl(A), uT }

bscon(B,c) = { stack(B,c) }
  └ pre: { cl(B), h(A), h(B) }

bs(cl(c)) = { unstack(A,c) }
  └ pre: { on(A,c) }  ← support

bscon(A,B) = { stack(A,B) }
  └ pre: { cl(B), h(A) }

h^add = 5      bscon(A,B) = { stack(A,B) }
h_max = 2        └ pre: { cl(B), h(A) }

{ 3 : add
  2 : max  ⇒ 2

Relaxed plan → { unstack(A,c) ; pickup(B) ; stack(A,B) ; oT(B) ; stack(A,B) }

  └ pre: { uT, cl(B) }
  └ pickup(B)
    └ pre: { uf, cl(B), oT(B) } cost=1
             0    0    0    ⇒ 1.
    └ stack(A,B)  cost=1
      └ pre: { uf, cl(B) } oT(B)
                0    0   ⇒ 1.

h_del = 4

h_max hold. does it change last supported bs
  pickup. └ pre { cl(c), uT } h(c)
                 0    ∞   ⇒ ∞

| i | t(A) | t(B) | t(C) | t(D) | P(T) | P(A) | P(B) | P(C) | P(D) |
|---|------|------|------|------|------|------|------|------|------|
| 0 | 0 | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ |
| 1 | 0 | 2 | 2 | ∞ | ∞ | ∞ | 0 | 0 | ∞ |
| 2 | 0 | 2 | 3 | ∞ | 3 | ∞ | 0 | 0 | ∞ |
| 3 | 0 | 3 | 3 | 4 | ∞ | 3/4 | 0 | 0 | ∞ |
| 4 | 0 | 3 | 4 | ∞ | 7/4 | 0 | 0 | ∞ | ∞ |

Goal

Goal: P(D) t

1: t(A) P(C)

Goal: P(D) t
1: t(A) P(C)

$h(\text{add}) =$

Goal

$h(\text{max}) = 4$

load(C), t(a), t(b), t(d)

P(C D) 1+ ③(3 or ③)
└─ pre: t(C D), P(T)

at B.    min(∞,  d(A,B)   d(C,B).
        1+0 ,  1+∞ )

(C,B)   P(a)   min ∮ ∞,
t(B)                1+ h(at A)
t(B) 1+0
these   dir(A,B)  pre: {at A}.
P(T)         0
    └─ load C 1+2+0
        └─ pre: (at C C), P(C)

at PC T).       [P(B)]  unload(B)
min [∞,                1+ ⑤
    1+ h(at cc).P(c)]     t
    1+2+0              └─ pre: t(B), P(T))
                      L→ pre: t(B), P(T))

min ∮ ∞,   len load(A)    (a+t+∞, a+t+∞, a+n+∞
1+ 0 ⑬    1+ 0 ⑬         (D, t), t
       └─ pre: t t(A), t(a) )   └─ pre: t(B), P(C T)

1+   3+0

[P(B)]

P(C D)

1<10    ,     (t t, tt)
1=10 , (pre: s  t old (s, pre: s  s t)
   min ∮ A, g old C(C) + ∮ old (s, pre)
      Σ s; g old       t

at(s) at(A) at(Pe) at(Pe) at(Br) at(Da) v(sy) v(Da) v(Ba) v(Da) v(sy) v(Ba) v(De) v(Da)

| i | at(sy) | at(A) | at(Pe) | at(Pe) | at(Br) | at(Da) | v(sy) | v(Da) | v(Ba) | v(Da) | v(sy) | v(Ba) | v(De) | v(Da) |
|---|--------|-------|--------|--------|--------|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| 0 | 0 | ∞ | ∞ | ∞ | ∞ | ∞ | 0 | ∞ | ∞ | ∞ | 0 | ∞ | ∞ | ∞ |
| 1 | 0 | 1.5 | 1 | 1 | 1 | 1 | 0 | 1 | ∞ | ∞ | 1 | 1 | 1 | 1.5 |
| 2 | 0 | 1.5 | 1 | 1 | ∞ | 1 | 0 | 1.5 | 1.5 | 1.5 | 1 | 1.5 | 1 | 1.5 |
| 3 | 0 | 1.5 | 1.5 | 1 | 1 | 1.5 | 0 | 1.5 | 1.5 | 1.5 | 1 | 1.5 | 1.5 | 1.5 |

∮ old (1.5, 1.5) =
1.5, 1.5, 1.5

min (0)  old(Br, Br
   pre: at Br)
(s, s)

min (∞
pre: 1.5, A, t

2: at(sy), v(sy)
g: at(sy), v(×), v×p×al×

Ba
Da

1
-1
y  3.5  -1.5
Pe   Pe  3.5