

学号: 20202202823

姓名: 朱相颐

班级: 电信 2001

1) 搜索

| Search Algorithm     | A-B-D-G | A-C-D-G | AB-A-B-C-D-G | A-B-C-D-F-G |
|----------------------|---------|---------|--------------|-------------|
| Depth first search   | ✓       | ✓       | ✓            | ✓           |
| Breadth first search | ✓       | ✓       | ✗            | ✗           |
| UCS                  | ✗       | ✗       | ✗            | ✗           |
| $A^*$ with $h_1$     | ✗       | ✗       | ✗            | ✗           |
| $A^*$ with $h_2$     | ✗       | ✗       | ✗            | ✓           |

2) ① [0, 12]

② [9, 10]

③ (~~12, 13~~) (12.5, 13)

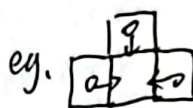
12)

1)  $(s_1, s_2, s_3, \dots, s_n) \quad S \in \{1, 2, \dots, M\}$

$M^n$

2)  $5^{\frac{MN}{2}}$

3) 是不可接受的



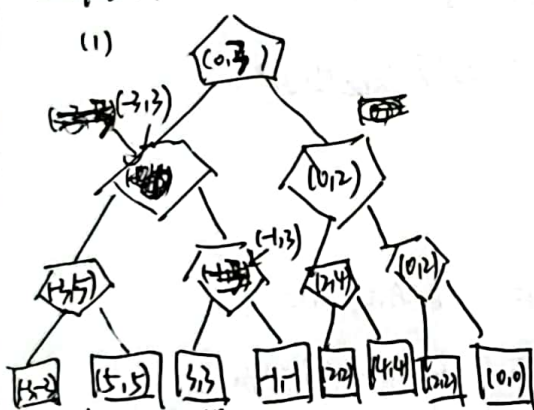
在左图中, 有 3 对 pacman 没有成对, 但只需要走一步便可成对, 预估值大于实际值, 所以是不可接受的

二、树搜索

(1)

12)  $u \leftarrow \max(uList)$

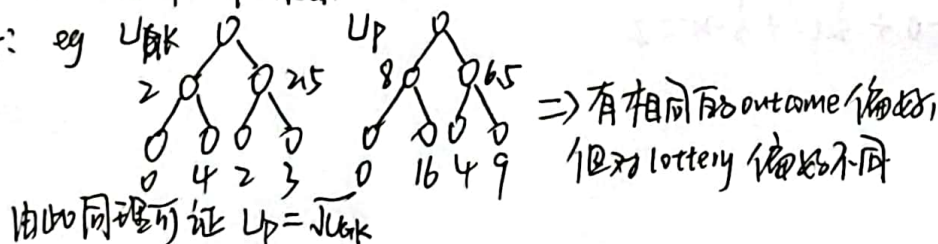
$v \leftarrow \max(\max(uList), \min(vList))$



三、效用函数

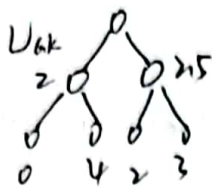
(1) 选择  $U_p = U_{GK}^2$  和  $U_p = \sqrt{U_{GK}}$

解释:

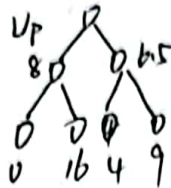


(2) 选择  $U_P = 0 U_{GK}$

解释:



此时 GK 会选择 outcome 肯定有收益的右支路



此时, Pacman 会选择能得到最大收益的右支路, 但也同时承担得到 0 收益的风险。  
1. 更冒险

同理可证  $U_P = 2 U_{GK}$  时, GK 更冒险。

四、马尔科夫决策过程。

| (1) State | $s_1$ | $s_2$ | $s_3$ | $s_4$ | $s_5$ | $s_6$ |
|-----------|-------|-------|-------|-------|-------|-------|
| $\pi(s)$  | Roll  | Roll  | Stop  | Stop  | Stop  | Stop  |
| $V(s)$    | 3     | 3     | 3     | 4     | 5     | 6     |

由题可知 1, 2 点时两状态会继续 Roll  
因此最终状态下  $V(s_1) = V(s_2)$

$$V(s_1) = -1 + \frac{1}{6}(V(s_1) + V(s_2) + 4 + 5 + 6)$$

$$V(s_2) = -1 + \frac{1}{6}(V(s_1) + V(s_2) + 3 + 4 + 5 + 6)$$

$$\text{可得 } V(s_1) = V(s_2) = 3$$

| (2) State | $s_1$ | $s_2$ | $s_3$     | $s_4$ | $s_5$ | $s_6$ |
|-----------|-------|-------|-----------|-------|-------|-------|
| $\pi(s)$  | Roll  | Roll  | Stop      | Stop  | Stop  | Stop  |
| $\pi'(s)$ | Roll  | Roll  | Roll/stop | Stop  | Stop  | Stop  |

4 点时 stop  $R(s_4) = 4 > 3$   
同理可证 5 点

$$\text{old policy average reward: } -1 + \frac{4+4+4+5+6}{6} = 3$$

new policy: 若 1 点 stop, 则  $R(s_1) = 1 < 3$   
若 2 点 stop, 则  $R(s_2) = 2 < 3$   
若 3 点 stop, 则  $R(s_3) = 3 = 3$

综上, 1, 2 点 Roll, 3 点可 Roll 可 stop  
4 点 stop, 5 点 stop

(3) 是最优的, 当前决策下平均收益为 3, 1 点, 2 点 stop 的收益为 0 或 1  
而, 4, 5, 6 点 stop 的直接收益要大于 3。  
1. 是最优的

五、强化学习

$$(1) Q(A, Down) = 2 \quad Q(B, Up) = 2$$

$$Q(A, Down) = \frac{1}{2} \cdot 0 + \frac{1}{2} (4 + \frac{1}{2} \cdot 0) = 2$$

$$Q(B, Down) = \frac{1}{2} \cdot 0 + \frac{1}{2} (-4 + \frac{1}{2} \cdot 0) = -2$$

$$Q(B, Up) = \frac{1}{2} \cdot 0 + \frac{1}{2} (0 + \frac{1}{2} \cdot 0) = 0$$

$$Q(B, Up) = \frac{1}{2} \cdot 0 + \frac{1}{2} (3 + \frac{1}{2} \cdot 2) = 2$$

$$(2) \hat{r}(A, Up, A) = 1 \quad \hat{r}(A, Up, A) = -1$$

$$\hat{r}(A, Up, B) = 0 \quad \hat{r}(A, Up, B) = \eta/a$$

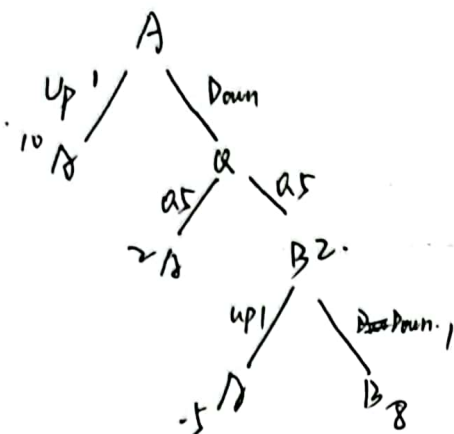
$$\hat{r}(B, Up, A) = \frac{1}{2} \quad \hat{r}(B, Up, A) = 3$$

$$\hat{r}(B, Up, B) = \frac{1}{2} \quad \hat{r}(B, Up, B) = 0$$



(3)

(i)  $\pi^*(A) = \text{Up}$   $\pi^*(B) = \text{Down}$



$V^*(A) = 20$

$V^*(B) = 16$

$A - \text{Down} - B : 2 + V(B)$

$A - \text{Down} - B : 0 + \frac{1}{2}(2 + \frac{1}{2}(0)) + \frac{1}{2}(2 + V(B)) \times$

$A - \text{Up} - B : 10 + V(B)$

$A - \text{Down} - B - \text{Up} - B : 1 + \frac{1}{2}(2 + \frac{1}{2}(0)) + \frac{1}{2}(2 + V(B)) \times$

$A - \text{Up} - B - \text{Up} - B : 10 + 10 + V(B) \checkmark$

$A - \text{Down} - B - \text{Down} - B : 2 + 8 + V(B) \checkmark$

$B - \text{Down} - B : 8 \checkmark$

$B - \text{Up} - B : 5 \times$

$\therefore \pi^*(B) = \text{Down} \cdot \pi^*(A) = \text{Up}$

② 迭代下去可得

$V^*(A) = 10 + 10 \times \frac{1}{2} + 10 \times \frac{1}{2}^2 + \dots = 20$

$V^*(B) = 8 + 8 \times \frac{1}{2} + 8 \times \frac{1}{2}^2 + \dots = 16$

(ii) 最终会收敛到上面所求得的  $V^*$

将其代入  $Q$ -learning 算法, 可以找到局部最优解  $V^*$  但不能找到全局最优解  $V^*$ . 要想找到全局最优解, 需要更多的关于这个世界的信息

