# Ex 10-1

## (a)

|          | roof | sleigh | terminal |
|----------|------|--------|----------|
| roof     | 0.5  | 0.9    | 0.1      |
| sleigh   | 0.4  | 0.5    | 0.1      |
| terminal | 0    | 0      | 0        |

$$P = \begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.4 & 0.5 & 0.1 \\ 0 & 0 & 0 \end{pmatrix}$$

## (b)

$$U = R + \gamma P U$$

$$\Rightarrow (I - \gamma P) U = R \quad \text{← expensive .}$$

$$\Rightarrow \underbrace{U = (I - \gamma P)^{-1} R}_{} \quad O(a^3)$$

$$\Rightarrow \left[ \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} - 1 \cdot \begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.4 & 0.5 & 0.1 \\ 0 & 0 & 0 \end{pmatrix} \right] \begin{bmatrix} u_r \\ u_s \\ u_t \end{bmatrix} = \begin{bmatrix} -2 \\ -1 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{pmatrix} 0.5 & -0.4 & -0.1 \\ -0.4 & 0.5 & -0.1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} u_r \\ u_s \\ u_t \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \\ 0 \end{pmatrix}$$
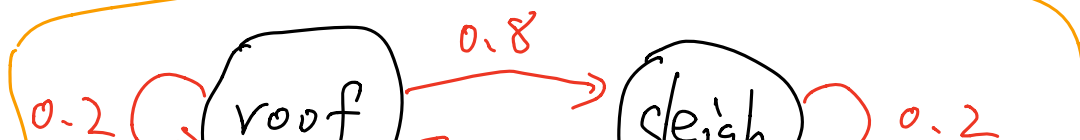
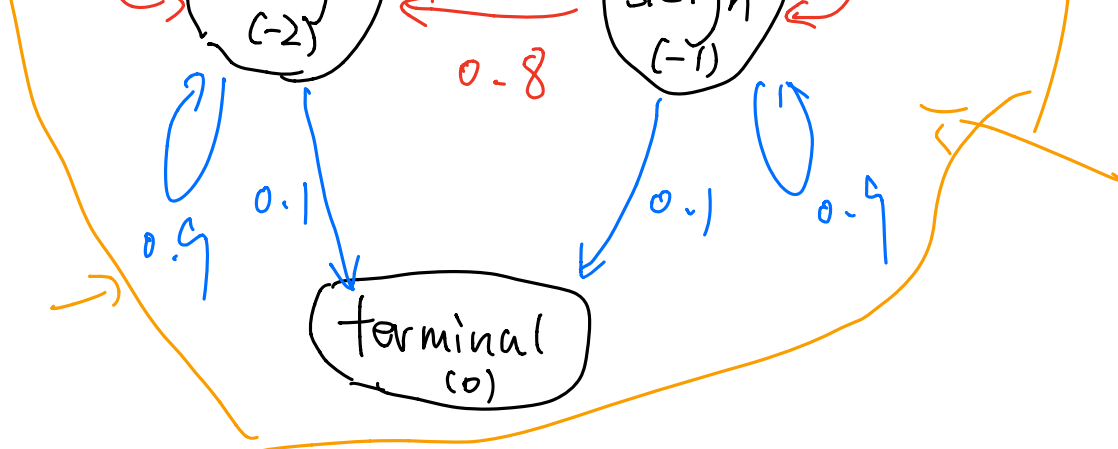$$\Rightarrow \begin{cases} u_s = -\dfrac{130}{9} \\ u_r = -\dfrac{140}{9} \\ u_t = 0 \end{cases}$$

---

# Ex 10-2

action: C (change)

action: t (throw)

## a)



$$0.2 \quad \bigcirc \quad \text{roof} \xrightarrow{0.8} \text{sleigh} \quad \bigcirc \quad 0.2$$

(−2)

0.8

(−1)

0.1 0.9

0.9

0.1

terminal
(0)

$$P(X_s \mid X_r, c) = 0.8 \qquad P(X_t \mid X_r, t) = 0.1$$

$$P(X_s \mid X_s, c) = 0.2 \qquad P(X_r \mid X_r, t) = 0.9$$

$$P(X_r \mid X_s, c) = 0.8 \qquad P(X_t \mid X_s, t) = 0.1$$

$$P(X_r \mid X_r, c) = 0.2 \qquad P(X_s \mid X_s, t) = 0.9$$

---

c)

$$\pi_0(X_s) = t \qquad \leftarrow c$$

$$\pi_0(X_r) = t \qquad (\text{initial})$$

$$\leftarrow c$$

① policy evaluation:

$$U_s = R_s + \gamma \sum_{s' \in S} P(s' \mid s, a) U(s')$$

$$= -1 + 0.9 \cdot U_s + 0.1 \cdot U_t$$

similarly:

$$U_r = -2 + \gamma (0.9 \, U_r + 0.1 \cdot U_t) \qquad \} \text{ unsolved.}$$

$$U_t = 0$$

$$\Rightarrow \begin{cases} u_s = -10 \\ u_r = -20 \\ u_t = 0 \end{cases} \qquad \arg\max_x f(x)$$

② policy improvement.

$$U = \boxed{R} + \arg\max_a \left( \boxed{r} \underbrace{\sum_{s'} p(s'|s,a) \; U}_{C_t} \right)$$

$$\Rightarrow \qquad T(x_i, a) = \sum_{s'} p(s'|s,a) \; U$$

$$T(x_s, c) = 0.8 \times (-20) + 0.2 \times (-10) = -18$$

$$T(x_s, t) = 0.9 \times (-10) + 0.1 \times (0) = -9$$

$$T(x_r, c) = 0.8 \times (-10) + 0.2 \times (-20) = -12$$

$$T(x_r, t) = 0.9 \times (-20) + 0.1 \times (0) = -18$$

$$T(x_s, c) \quad < \quad T(x_s, t) \Rightarrow \boxed{\pi_1(x_s) = t}$$

$$T(x_r, c) \quad > \quad T(x_r, t) \Rightarrow \boxed{\pi_1(x_r) = c}$$

③ policy evaluation

$$\begin{cases} u_s = -1 + 0.9 \, u_s + 0.1 \, u_t \\ u_\cdots = \cdots + 0.0 \, u + 0.2 \, u \end{cases} \Rightarrow \begin{cases} u_s = -10 \\ u_r = -12.5 \end{cases}$$

$u_r = -21 + 0.8\, u_s + 0.2\, u_r$

$u_t = 0$

④ policy inprovment.

$T(X_r, c) = 0.8 \times (-10) + 0.2 \times (-12.5) = -10.5$

$T(X_r, t) = 0.9 \times (-12.5) + 0.1 \times 0 = -11.25$

$T(X_s, c) = 0.8 \times (-12.5) + 0.2 \times (-10) = -12$

$T(X_s, t) = 0.9 \times (-10) + 0.1 \times 0 = -9$

$T(X_s, c) \quad < \quad T(X_s, t) \qquad \pi_2(X_s) = t$

$T(X_r, c) \quad > \quad T(X_r, t) \qquad \pi_2(X_r) = c$

terminate . $\pi_2$ is optimal