# Homework 3

## Group 30

a) States
$$X = \{0,1,2,3,4,5,6,7,8,9,10,11\}$$

Actions
$$A = \{p, s\}; \quad p = \text{"play"}, \quad s = \text{"stop"}$$

Observations
$$Z = \{g, w, t, e\}$$

Transition probabilities:

$P_p =$

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0.2 | 0.4 | 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.2 | 0 | 0.4 | 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.2 | 0 | 0 | 0.4 | 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.2 | 0 | 0 | 0 | 0.4 | 0.4 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.2 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0.2 | 0.4 | 0.4 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0.2 | 0 | 0.4 | 0.4 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.2 | 0.4 | 0.4 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.2 | 0 | 0.4 | 0.4 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.2 | 0 | 0 | 0.8 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

$P_s =$

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

Observation probabilities

$O_p = O_s =$

| g | w | t | e |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |

Immediate cost function

$C =$

| p | s |
|---|---|
| 1 | 1 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 0 | 0 |
| 1 | 1 |
| 1 | 0 |
| 1 | 1 |
| 1 | 0 |
| 1 | 0 |
| 0 | 0 |

.b)

.a)

$$\alpha_0 = \mu_0 = \begin{bmatrix} \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} \end{bmatrix}$$

(column indices: 0 1 2 3 4 5 6 7 8 9 10 11)

$$\hat{\alpha}_1 = \alpha_0 \, \rho_p = \begin{bmatrix} \frac{1}{12} & \frac{1}{30} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} & \frac{11}{60} & \frac{1}{12} & \frac{1}{30} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} & \frac{11}{60} \end{bmatrix}$$

$$\alpha_1^T = \hat{\alpha}_1 \, \text{diag}\left(\delta_p\left(w|\cdot\right)\right)$$

$$= \hat{\alpha}_1 \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{15} & 0 & 0 & 0 \end{bmatrix}$$

normalized $\Rightarrow = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$

**.b)**

$$\alpha_0 = M_0 = \begin{bmatrix} \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} & \frac{1}{12} \end{bmatrix}$$

(columns labeled: 0 1 2 3 4 5 6 7 8 9 10 11)

(b) After feeling the web with its legs and playing two turns, assuming that the spider made no observation after each step (i.e., it makes two empty observations).

$$\hat{\alpha}_1 = \alpha_0 \, P_p = \begin{bmatrix} \frac{1}{12} & \frac{1}{20} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} & \frac{11}{60} & \frac{1}{12} & \frac{1}{20} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} & \frac{11}{60} \end{bmatrix}$$

$$\hat{\alpha}_2 = \alpha_1 \, P_p$$

$$\alpha_1^T = \hat{\alpha}_1 \, \text{diag}(O_p(w | \cdot))$$

$$\alpha_2^T = \hat{\alpha}_2 \, \text{diag}(O_p(e | \cdot))$$

$$= \hat{\alpha}_1 \begin{bmatrix} 0 & & & & & & & & & & & \\ & 0 & & & & & & & & & & \\ & & 0 & & & & & & & & & \\ & & & 0 & & & & & & & & \\ & & & & 0 & & & & & & & \\ & & & & & 0 & & & & & & \\ & & & & & & 0 & & & & & \\ & & & & & & & 1 & & & & \\ & & & & & & & & 0 & & & \\ & & & & & & & & & 0 & & \\ & & & & & & & & & & 0 & \\ & & & & & & & & & & & 0 \end{bmatrix}$$

$$\hat{\alpha}_3 = \alpha_2 \, P_p$$

$$\alpha_3^T = \hat{\alpha}_3 \, \text{diag}(O_p(e | \cdot))$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{15} & 0 & 0 & 0 \end{bmatrix}$$

---

$$\hat{\alpha}_2 = \alpha_1 \, P_p$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{15} & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} .2 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & 0 & .8 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & 0 & .8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{75} & \frac{2}{75} & \frac{2}{75} & 0 \end{bmatrix}$$

$$\alpha_2^T = \hat{\alpha}_2 \, \text{diag}(O_p(e | \cdot))$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{75} & \frac{2}{75} & \frac{2}{75} & 0 \end{bmatrix} \begin{bmatrix} 0 & & & & & & & & & & & \\ & 1 & & & & & & & & & & \\ & & 1 & & & & & & & & & \\ & & & 1 & & & & & & & & \\ & & & & 1 & & & & & & & \\ & & & & & 0 & & & & & & \\ & & & & & & 1 & & & & & \\ & & & & & & & 1 & & & & \\ & & & & & & & & 0 & & & \\ & & & & & & & & & 1 & & \\ & & & & & & & & & & 1 & \\ & & & & & & & & & & & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{2}{75} & \frac{2}{75} & 0 \end{bmatrix}$$

$$\hat{\alpha}_3 = \alpha_2 P_P$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{2}{75} & \frac{2}{75} & 0 \end{bmatrix} \begin{bmatrix} .2 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & 0 & .8 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & 0 & .8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{4}{375} & 0 & \frac{4}{375} & \frac{4}{125} \end{bmatrix} \overset{\text{normalized}}{\implies} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{5} & 0 & \frac{1}{5} & \frac{3}{5} \end{bmatrix}$$

$$\alpha_3^T = \hat{\alpha}_3 \; diag(\partial p(e|.))$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{5} & 0 & \frac{1}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{5} & 0 \end{bmatrix}$$

(c) After starting and playing 3 times, assuming that the spider made no observation after each step (i.e., it makes three empty observations).

$$\hat{\alpha}_1 = \alpha_0 P_p = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} .2 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & 0 & .8 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & 0 & .8 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{5} & \frac{2}{5} & \frac{2}{5} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\alpha_1^T = \hat{\alpha}_1 \, \text{diag}\left(O_p \mid e(\cdot)\right)$$

$$= \begin{bmatrix} \frac{1}{5} & \frac{2}{5} & \frac{2}{5} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & \frac{2}{5} & \frac{2}{5} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

---

$$\hat{\alpha}_2 = \alpha_1 P_p$$

$$= \begin{bmatrix} 0 & \frac{2}{5} & \frac{2}{5} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} .2 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & 0 & .8 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & 0 & .8 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{4}{25} & 0 & \frac{4}{25} & \frac{8}{25} & \frac{4}{25} & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

normalized $\Rightarrow$ = $\begin{bmatrix} \frac{1}{5} & 0 & \frac{1}{5} & \frac{2}{5} & \frac{1}{5} & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$

$$\alpha_2^T = \hat{\alpha}_2 \, diag\left(0_p \,|\, e(\cdot)\right)$$

$$= \left[\tfrac{1}{5} \; 0 \; \tfrac{1}{5} \; \tfrac{3}{5} \; \tfrac{1}{5} \; 0\,0\,0\,0\,0\,0\,0\right] \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \left[0 \; 0 \; \tfrac{1}{5} \; \tfrac{2}{5} \; \tfrac{1}{5} \; 0\,0\,0\,0\;0\,0\,0\right]$$

$$\hat{\alpha}_3 = \alpha_2 \, P_p$$

$$= \left[0 \; 0 \; \tfrac{1}{5} \; \tfrac{2}{5} \; \tfrac{1}{5} \; 0 \; 0\,0\,0 \; 0\,0\,0\right] \begin{bmatrix} .2 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & .4 & .4 & 0 & 0 & 0 & 0 & 0 & 0 \\ .2 & 0 & 0 & 0 & 0 & .8 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & .4 & .4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & .4 & .4 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & .2 & 0 & 0 & .8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$= \left[\tfrac{4}{25} \; 0 \; 0 \; \tfrac{2}{25} \; \tfrac{6}{25} \; \tfrac{8}{25} \; 0\,0\,0\,0\,0\,0\right]$$

$$\overset{normalized}{\Longrightarrow} \;=\; \left[\tfrac{1}{5} \; 0 \; 0 \; \tfrac{1}{10} \; \tfrac{3}{10} \; \tfrac{4}{10} \; 0\,0\,0\,0\,0\,0\right]$$

$$\alpha_3^T = \hat{\alpha}_3 \, diag\left(0_p \,|\, e(\cdot)\right)$$

$$= \left[\tfrac{1}{5} \; 0 \; 0 \; \tfrac{1}{10} \; \tfrac{3}{10} \; \tfrac{4}{10} \; 0\,0\,0\,0\,0\right] \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \left[0 \; 0 \; 0 \; \tfrac{1}{10} \; \tfrac{3}{10} \; 0 \; 0\,0\,0 \; 0\,0\,0\right] \overset{normalize}{\Longrightarrow} \left[0\,0\,0 \; \tfrac{1}{4} \; \tfrac{3}{4} \; 0\,0\,0\,0\,0\,0\right]$$

c)

Belief: [0.2    0.08    0.24    0.32    0.16    0   0   0   0   0   0]

## MLS

- For MLS we have to identify the most likely state according to the belief
- The highest probability is $b(3) = 0.32$ so state 3 is the most likely state
- To know the best action to take we have to calculate which action has the lowest cost between playing and stopping, since with this heuristic we will never do the "identify" action.

→ The cost-to-go for this policy would be:

$$J^\pi(3) = c_\pi(3) + \gamma \sum_{y \in \chi} P_\pi(y|3) \times J^\pi(y)$$

↓
state we currently are

$$J^\pi(3) = 1 + 0.9\left(0.2\, J^\pi(0) + 0.4\, J^\pi(4) + 0.4\, J^\pi(5)\right)$$

$J^\pi(5) = 0$, because according to the homework any action after reaching level 5 will have a cost of 0

$$J^\pi(4) = 1 + 0.9\left(0.2\, J^\pi(0) + 0.8\, \overbrace{J^\pi(5)}^{0}\right) =$$

$$= 1 + 0.9 \times 0.2\, J^\pi(0) = 1 + 0.18\, J^\pi(0)$$

$$J^\pi(3) = 1 + 0.9\left(0.2\, J^\pi(0) + 0.4\, J^\pi(4) + 0.4\, J^\pi(5)\right) =$$

$$= 1 + 0.9\left(0.2\, J^\pi(0) + 0.4\,(1 + 0.18\, J^\pi(0))\right) = 1.36 + 0.2448\, J^\pi(0)$$

$$J^\pi(2) = 1 + 0.9\left(0.2\, J^\pi(0) + 0.4\,(1.36 + 0.2448\, J^\pi(0)) + 0.4\,(1 + 0.18\, J^\pi(0))\right) =$$

$$= 1.8496 + 0.3329\, J^\pi(0)$$

$$J^\pi(1) = 1 + 0.9\left(0.2\, J^\pi(0) + 0.4\, J^\pi(2) + 0.4\, J^\pi(3)\right) = 2.1554 + 0.3880\, J^\pi(0)$$

$$J^\pi(0) = 1 + 0.9\left(0.2\, J^\pi(0) + 0.4\, J^\pi(1) + 0.4\, J^\pi(2)\right) = 2.4418 + 0.4396\, J^\pi(0)$$

$$J^\pi(0) = 2.4418 + 0.4396\, J^\pi(0) \iff J^\pi(0) = 4.3572$$

$$J^\pi(3) = 1.36 + 0.2448\, J^\pi(0) = 2.4266 \quad \longrightarrow \text{cost for policy play}$$

Now for the action "stop"

$$J_i^\pi(3) = 0 + 0.9\,(1 \times J^\pi(0))$$

$\underbrace{\phantom{J}}_{stop}$

We will assume that in all other states the policy is optimal, and will do the "play" action, for state 0

$$J^\pi(3) = 0.9 \times 4.3572 = 3.921$$

→ According to MLS heuristic and with the given belief the best action to make is "Play" since $J^{play}(3) < J^{stop}(3)$

                                                ↓         ↓

                                               2.4266      3.921

## QMDP:

We have to calculate the action that minimizes $\sum_{x \in \mathcal{X}} b(x)\, Q^*_{MDP}(x,a)$

                              ↳ calculated the same way as above

$Q^*(0, Play) = 1 + 0.9(0.2\,J^*(0) + 0.4\,J^*(1) + 0.4\,J^*(2)) =$

          $= 1 + 0.9\,(0.2 \times 4.3572 + 0.4 \times 3.8460 + 0.4 \times 2.4953 = 4.0672$

$Q^*(0, Stop) = 1 + 0.9 \times (1 \times J^*(0)) = 4.92148$ → for state 2 the optimal policy is choosing the stop action

$Q^*(1, Play) = 1 + 0.9(0.2\,J^*(0) + 0.4\,J^*(2) + 0.4\,J^*(3)) =$

          $= 1 + 0.9\,(0.2 \times 4.3572 + 0.4 \times 2.4953 + 0.4 \times 2.4266 = 3.5562$

$Q^*(1, Stop) = 0 + 0.9 \times (1 \times J^*(0)) = 3.92148$

$Q^*(2, Play) = 1 + 0.9(0.2\,J^*(0) + 0.4\,J^*(3) + 0.4\,J^*(4)) =$

          $= 1 + 0.9\,(0.2 \times 4.3572 + 0.4 \times 2.4266 + 0.4 \times 1.7843 = 3.3002$

$Q^*(2, Stop) = 0 + 0.9 \times (1 \times J^*(8)) = 2.4953$

$Q^*(3, Play) = 1 + 0.9(0.2\,J^*(0) + 0.4\,J^*(4) + 0.4\,J^*(5)) =$

          $= 1 + 0.9\,(0.2 \times 4.3572 + 0.4 \times 1.7843 + 0.4 \times 0 = 2.4266$

$Q^*(3, Stop) = 0 + 0.9 \times (1 \times J^*(0)) = 3.92148$

$Q^*(4, Play) = 1 + 0.9(0.2\,J^*(0) + 0.8\,J^*(5))$

          $= 1 + 0.9\,(0.2 \times 4.3572 + 0) = 1.7843$

$Q^*(4, Stop) = 0 + 0.9 \times (1 \times J^*(0)) = 3.92148$

$\sum_{x \in \mathcal{X}} b(x)\, Q^*_{MDP}(x, Play) = 0.2 \times 4.0672 + 0.08 \times 3.5562 + 0.24 \times 3.3002 +$
                        $+ 0.32 \times 2.4266 + 0.16 \times 1.7843 = 2.9520$

$\sum_{x \in \mathcal{X}} b(x)\, Q^*_{MDP}(x, Stop) = 0.2 \times 4.9215 + 0.08 \times 3.9215 + 0.24 \times 2.4953 +$
                        $+ 0.32 \times 3.9215 + 0.16 \times 3.9215 = 3.7792$

- Since our objective is to minimize that function, the QMDP heuristic will choose the action "play" for the given belief