

1. Policy & Value Function

① Value Function: probability transition matrix

$$v = l + \lambda P(v)$$

↑
value vector
= $l + \lambda P + \lambda P^2 + \dots$

↑
expected under this iteration
↔ expected under j -th step

$$v_{cj} = l_{cj} + \sum_{i=1}^n v_{ci} p_{ci} \underbrace{s_{ci}=i | S_t=j, v}_{\substack{\text{value at } j\text{-th step} \\ \text{next step value}}} \underbrace{p_{ci}}_{\text{probability for } j \rightarrow i}$$

② Policy

- In terms of value:

$$p(v)_{ij} = \frac{p(v)_{ij} \exp(\lambda v_{ci})}{\sum_k p(v)_{kj} \exp(\lambda v_{ck})}$$

uncontrolled transition matrix (restriction)
↓ sensitivity

- In terms of attraction: $A_{ij} = \frac{1}{\lambda} v_{ci} + \ln p(v)_{ij}$

$$p(A)_{ij} = \frac{\exp(\lambda A_{ij})}{\sum_k \exp(\lambda A_{kj})}$$

Intuition: $v_{ci} \uparrow$, more attractive, easier to be selected

③ The problem is:

$$\begin{aligned} v_{t+1} &= l + \lambda P(v_t) \\ \text{s.t. } v - P(v) &= l \end{aligned}$$

2. Multiple agents

① Joint-state: (two agents)

$$S = S_1 \times S_2$$

② Policy:

$$P_i(v_i)_{ij} = \frac{\pi_i(c_i) \downarrow j \exp(\vec{v}_i c_i)}{\sum_k \pi_i(c_i) \exp(\vec{v}_i c_k)}, \quad P_2(v_2)_{ij} = \frac{\pi_2(c_2) \downarrow j \exp(\vec{v}_2 c_i)}{\sum_k \pi_2(c_2) \exp(\vec{v}_2 c_k)}$$

$$P(v_1, v_2) = P_2(v_2) P_1(v_1)$$

③ Value:

$$v_1^{(i)} = \vec{l}_1 + \vec{\vartheta}_1 \uparrow p(v_1^{(i)}, o)$$

$$v_2^{(i)} = \vec{l}_2 + \vec{\vartheta}_2 \uparrow p(o, v_2^{(i)})$$

$$v_1^{(i+1)} = \vec{l}_1 + \vec{\vartheta}_1 \uparrow p(v_1^{(i)}, v_2^{(i-1)})$$

$$v_2^{(i+1)} = \vec{l}_2 + \vec{\vartheta}_2 \uparrow p(v_1^{(i-1)}, v_2^{(i)})$$

until $v_K^{(i)} \approx v_K^{(i+1)}$

3. Stag - Hunt Game

① Process: stag → subject → computer

$$\textcircled{2} \text{ Update: } \text{PCA}(B, C) = \frac{\text{PCA}(A, C) \text{PCA}(c)}{\text{PCB}(c)}$$

↑ prior ↓

$$p(k_{\text{com}}(T)|y, k_{\text{sub}}) \propto p(y|1, \dots, T) | k_{\text{sub}}(1, \dots, T), k_{\text{com}}) p(k_{\text{com}})$$

↑ object: posterior of T step likelihood: Assume k_{com} is the opponent's strategy,
 the probability of generating trajectory $y(1, \dots, T)$
 = $\prod_{t=1}^{T-1} K^{T-t} p(s_{t+1}|s_t, k_{\text{sub}}(t), k_{\text{com}})$
 ↑ forget discount

③ We should solve:

$$\hat{k}_{\text{com}}(t) = \underset{k_{\text{com}} \in \{1, \dots, k_{\text{sub}}\}}{\operatorname{argmax}} p(k_{\text{com}}(t)|y, k_{\text{sub}}),$$

$$k_{\text{sub}}(t+1) = \hat{k}_{\text{com}}(t) + 1$$