

Q-Learning con grafos de coordinación

Requiere grafo de coordinación del sistema $G = (\mathcal{V}, \mathcal{E})$; orden de eliminación de los agentes $\mathcal{O} \subseteq \mathcal{V}$; conjunto de vecinos para cada agente $\Gamma(i)$

1: **Para** cada arco $(i, j) \in \mathcal{E}$ **hacer**

2: Inicialice $Q_{ij}(s_{ij}, a_i, a_j)$ de manera optimista

3: **Fin Para**

4: **Para** cada episodio **hacer**

5: Inicialice el estado observado s_i para todos los agentes $i \in \mathcal{V}$

6: **Para** cada periodo de decisión k **hacer**

7: Obtenga la acción conjunta \mathbf{a} con ϵ -greedy como método de selección

8: Para cada agente $i \in \mathcal{V}$: aplique $a_i \in \mathbf{a}$, observe r_i y s_i^k

9: Para todos los arcos $(i, j) \in \mathcal{E}$: forme $s_{ij}^k = s_i^k \cup s_j^k$

10: Obtener $\mathbf{a}^* = \max_{\mathbf{a}} Q(\mathbf{s}^k, \mathbf{a}) \rightarrow \text{ELIMINACIÓN VARIABLE}(s_{ij}^k \forall (i, j) \in \mathcal{E})$

11: **Para** cada arco $(i, j) \in \mathcal{E}$ **hacer**

12: $Q_{ij}(s_{ij}^{k-1}, a_i^{k-1}, a_j^{k-1}) := (1 - \alpha)Q_{ij}(s_{ij}^{k-1}, a_i^{k-1}, a_j^{k-1}) + \alpha \left[\frac{r_i^k}{|\Gamma(i)|} + \frac{r_j^k}{|\Gamma(j)|} + \gamma Q_{ij}(s_{ij}^k, a_i^*, a_j^*) \right]$

13: **Fin Para**

14: Para cada agente $i \in \mathcal{V}$ actualice $s_i^{k-1} \leftarrow s_i^k$

15: **Fin Para**

16: **Fin Para**

Algoritmo de eliminación de variable

```
1: Función ELIMINACIÓNVARIABLE( $s_{ij}^k$ )
2:   Para cada agente  $i \in \mathcal{O}$  hacer
3:     Calcular función de pago condicional  $\phi$ 
4:     si el agente  $i$  es el primero a eliminar entonces
5:        $k \equiv$  siguiente agente a eliminar,  $k \in \Gamma(i)$ 
6:        $\phi_{kj} = \max_{a_j} \left[ \sum_{j \in \Gamma(i)} Q_{ij}(s_{ij}^k, a_i, a_j) \right]$ 
7:     sino si el agente  $i$  es el último a eliminar entonces
8:        $\phi_i = \max_{a_i} [\phi_i(a_i)]$ 
9:     sino
10:       $k \equiv$  siguiente agente a eliminar,  $k \in \Gamma(i)$ 
11:       $\phi_{kj} = \max_{a_j} \left[ \sum_{j \in \Gamma(i)} Q_{ij}(s_{ij}^k, a_i, a_j) + \phi_{ik}(a_i, a_k) \right]$ 
12:     Fin si
13:     Transmitir  $\phi$  al siguiente agente a eliminar
14:   Fin Para

15: Para cada agente  $i$  en el orden inverso de  $\mathcal{O}$  hacer
16:   si el agente  $i$  fue el último eliminado entonces
17:     Obtener  $a_i^* = \arg\max_{a_i} [\phi_i(a_i)]$ 
18:     Transmitir al agente anterior  $a_i^*$ 
19:     Transmitir  $\max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}) = \max_{a_i} \phi_i(a_i)$ 
20:   sino
21:     Esperar de los agentes vecinos  $a_k^*, a_j^*$  y  $\max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$ 
22:     Obtener  $a_i^* = \arg\max_{a_i} [\phi_{kj}(a_k^*, a_j^*)]$ 
23:   Fin si
24: Fin Para
25: Retornar  $a_i^* \quad \forall \quad i \in \mathcal{V}$ 
26: Fin Función
```

Q-Learning con *best response*

Requiere conjunto de agentes N , conjunto de vecinos para cada agente NB_i

1: **Para** cada agente $i \in \{1, 2, \dots, |N|\}$ **hacer**

2: **Para** cada agente vecino $j \in \{1, 2, \dots, |NB_i|\}$ **hacer**

3: Inicializar $Q_{i,NB_i[j]}(s_{i,NB_i[j]}, a_i, a_{NB_i[j]})$ de manera optimista

4: Inicializar modelo para la estimación de política $\theta_{i,NB_i[j]}(s_{i,NB_i[j]}, a_{NB_i[j]}) = 1/|\mathcal{A}_{NB_i[j]}|$

5: **Fin Para**

6: **Fin Para**

7: **Para** cada episodio **hacer**

8: Inicializar el estado observado s_i para todos los agentes $i \in N$

9: **Para** cada periodo de decisión k **hacer**

10: Obtener la acción conjunta \mathbf{a} con ϵ -greedy como método de selección

11: Aplicar $a_i \in \mathbf{a}$ para cada agente $i \in \mathcal{V}$

12: **Para** cada agente $i \in \{1, 2, \dots, |N|\}$ **hacer**

13: **Para** cada agente vecino $j \in \{1, 2, \dots, |NB_i|\}$ **hacer**

14: Observar $s_i^k, r_i^k, s_{NB_i[j]}^k$ y $a_{NB_i[j]}^{k-1}$

15: Formar estado conjunto $s_{i,NB_i[j]}^k = s_i^k \cup s_{NB_i[j]}^k$ y acción conjunta $a_{i,NB_i[j]}^{k-1} = a_i^{k-1} \cup a_{NB_i[j]}^{k-1}$

16: Actualizar modelo de estimación de la política para el vecino $NB_i[j]$:

$$\theta_{i,NB_i[j]}(s_{i,NB_i[j]}^{k-1}, a_{NB_i[j]}^{k-1}) = \frac{v(s_{i,NB_i[j]}^{k-1}, a_{NB_i[j]}^{k-1})}{\sum_{a_{NB_i[j]} \in \mathcal{A}_{NB_i[j]}} v(s_{i,NB_i[j]}^{k-1}, a_{NB_i[j]})}$$

17: Encontrar la *best response* respecto al vecino $NB_i[j]$:

$$br_i^k = \max_{a_i \in \mathcal{A}_i} \left[\sum_{a_{NB_i[j]} \in \mathcal{A}_{NB_i[j]}} Q_{i,NB_i[j]}(s_{i,NB_i[j]}^k, a_i, a_{NB_i[j]}) \times \theta_{i,NB_i[j]}(s_{i,NB_i[j]}^k, a_{NB_i[j]}) \right]$$

18:

Actualizar factores Q :

$$Q_{i, \text{NB}_i[j]}^k(s_{i, \text{NB}_i[j]}^{k-1}, a_{i, \text{NB}_i[j]}^{k-1}) = (1 - \alpha)Q_{i, \text{NB}_i[j]}^{k-1}(s_{i, \text{NB}_i[j]}^{k-1}, a_{i, \text{NB}_i[j]}^{k-1}) + \alpha [r_i^k + \gamma b_{r_i}^k]$$

19:

Actualizar $s_{i, \text{NB}_i[j]}^{k-1} \leftarrow s_{i, \text{NB}_i[j]}^k$

20:

Fin Para

21:

Acción que corresponde a *best response* con todos los vecinos:

$$a_i^k = \underset{a_i \in \mathcal{A}_i}{\operatorname{argmax}} \left[\sum_{j \in \{1, 2, \dots, |\text{NB}_i|\}} \sum_{a_{\text{NB}_i[j]} \in \mathcal{A}_{\text{NB}_i[j]}} Q_{i, \text{NB}_i[j]}(s_{i, \text{NB}_i[j]}^k, [a_i \cup a_{\text{NB}_i[j]}]) \times \theta_{i, \text{NB}_i[j]}(s_{i, \text{NB}_i[j]}^k, a_{\text{NB}_i[j]}) \right]$$

22:

Fin Para

23:

Fin Para

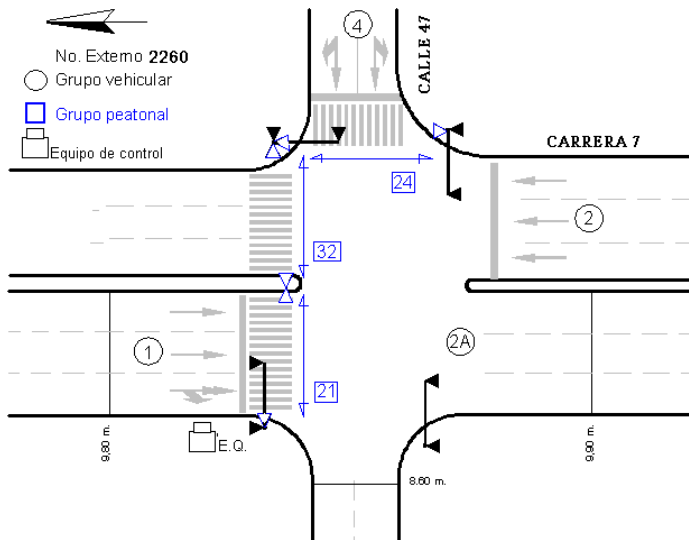
24:

Fin Para

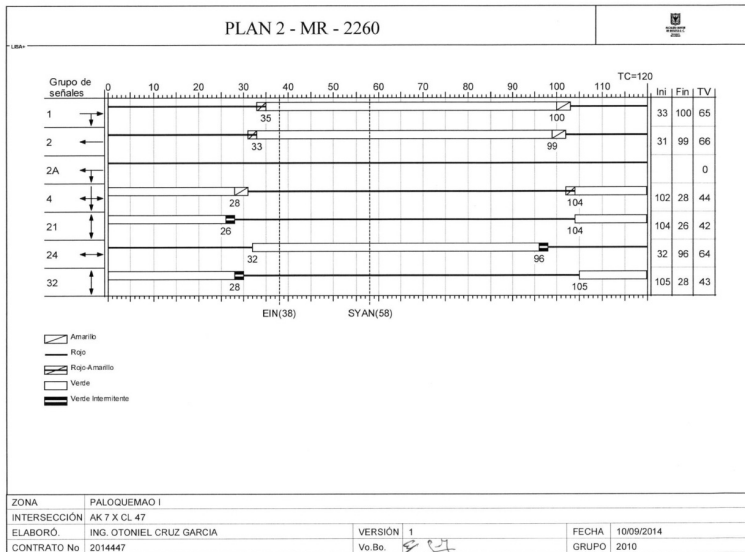
Discretización del espacio de estados

- 1: **Procedimiento** VECTOR QUANTIZATION
- 2: **Para** cada hora $h \in 0, 1, \dots, 23$ **hacer**
- 3: **Para** cada fase k **hacer**
- 4: Observar pasivamente el sistema durante n días
- 5: Guardar todos los vectores estado $s \in \mathbb{R}^{2(2+2i)-1}$ y concatenar las n observaciones en una matriz S ,
 donde cada columna j represente una dimensión del estado
- 6: Calcular la desviación estándar de cada columna $\text{Stdev}(S)_j$
- 7: Calcular el número de centroides $C_{h,k} = \sum_j \text{Stdev}(S)_j$
- 8: Aplicar K-Means sobre S usando el número de centroides calculado $\rightarrow VQ_{h,k}(S)$
- 9: **Fin Para**
- 10: **Fin Para**
- 11: **Fin Procedimiento**
- 12: **Procedimiento** DISCRETIZACIÓN(estado continuo s_{ij})
- 13: **Retornar** $\arg\min_{y \in VQ_{h,k}(s_{ij})} \{d_E(s_{ij}, y)\}$ ▷ euclidiana
- 14: **Fin Procedimiento**

Ejemplo de un plan de fases



Ejemplo de un plan de fases



De control óptimo a aprendizaje por refuerzo

La política π corresponde a un conjunto de funciones $\pi = \{\mu_0, \mu_1, \dots\}$ en donde cada μ_k mapea el estado s_k en una acción de control $a_k = \mu_k(s_k)$, tal que $\mu_k \in A(s_k)$ para todo s_k .

Costo de seguir la política π :

$$J_\pi(s) = \lim_{N \rightarrow \infty} \mathbf{E} \left[\sum_{k=0}^N \gamma^k r(s_k, \mu_k(s_k), s_{k+1}) | s = s_0 \right]$$

La política óptima $\pi^* \in \Pi$ corresponde a aquella que maximiza el funcional:

$$J_\pi^*(s_0) = \max_{\pi \in \Pi} J_\pi(s_0)$$

De control óptimo a aprendizaje por refuerzo

Solucionar la ecuación de Bellman para encontrar el valor óptimo del funcional para una iteración del problema:

$$J_{\pi}^*(s_i) = \max_{a \in A(s_i)} \sum_{j=1}^n p_{s_i, s_j}(a) (r(s_i, a, s_j) + \gamma J(s_j)), \quad i = 1 \cdots n$$

D. Bertsekas muestra que se puede obtener el funcional a partir de factores Q óptimos:

$$Q^*(s_i, a) = \sum_{j=1}^n p_{s_i, s_j}(a) (r(s_i, a, s_j) + \gamma J^*(s_j)), \quad \forall (s_i, a)$$

De control óptimo a aprendizaje por refuerzo

Reemplazando $J^*(s_i) = \max_{a \in A(s_i)} Q^*(s_i, a)$, se obtiene una ecuación de Bellman para sistemas de estados aumentados:

$$Q^*(s_i, a) = \sum_{j=1}^n p_{s_i, s_j}(a) \left(r(s_i, a, s_j) + \gamma \max_{a' \in A(s_j)} Q^*(s_j, a') \right), \quad \forall (s_i, a)$$

Una vez calculados los $Q^*(s_i, a) \forall (s_i, a)$, es directo obtener la política óptima π^* :

$$\mu^*(s_i) = \operatorname{argmax}_{a \in A(s_i)} Q^*(s_i, a) \quad \forall s_i$$