

Aprendizaje Reforzado: Regulación cuadrática lineal

Lección 26

Dr. Pablo Alvarado Moya

CE5506 Introducción al reconocimiento de patrones
Área de Ingeniería en Computadores
Tecnológico de Costa Rica

II Semestre, 2019

Contenido

- 1 Repaso
- 2 Recompensas por estado-acción
- 3 MDP de horizonte finito
- 4 Sistemas dinámicos lineales
 - Regulación cuadrática lineal (LQR)
 - Modelos
 - Ecuación de Riccati en tiempo discreto

Repaso

(1)

- Definimos un MDP como la tupla $\langle \mathcal{S}, \mathcal{A}, \{P_{sa}\}, \gamma, R \rangle$ con estados.
- El factor de degradación $\gamma \in [0, 1)$
- Función de recompensa $R : \mathcal{S} \rightarrow \mathbb{R}$.
- En la **iteración de valor**, repetimos hasta convergencia

$$V(s) \leftarrow R(s) + \max_{a \in \mathcal{A}} \gamma \sum_{s'} P_{sa}(s') V(s')$$

que al final será equivalente a $V^*(s)$

- Con $V^*(s)$ encontramos π^*

$$\pi(s) = \arg \max_{a \in \mathcal{A}} \sum_{s'} P_{sa}(s') V(s')$$

Repaso

(2)

- En el caso de estados continuos, lo que hicimos fue aproximar $V(s)$.
- En esos casos usamos regresión para la aproximación del siguiente estado $\underline{s}_{(t+1)}$ y un simulador para poder estimar valores esperados.

Plan de acción

- Lo que haremos ahora es proponer algunas modificaciones comunes a MDP que
 - usan otra definición de la función de recompensa R
 - usan otra estrategia para degradar la recompensa
 - nos permiten calcular $V(s)$ exactamente a pesar de continuidad de s

Recompensas por estado acción

- Después de usar una función de recompensa $R : \mathcal{S} \rightarrow \mathbb{R}$, regresaremos a la definición que considera también la acción a tomar.
- De nuevo usaremos $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Con eso, redefinimos el saldo total como

$$R(s_{(0)}, a_{(0)}) + \gamma R(s_{(1)}, a_{(1)}) + \gamma^2 R(s_{(2)}, a_{(2)}) + \dots$$

- Queremos encontrar la política que maximiza el valor esperado del saldo total:

$$\mathbb{E} [R(s_{(0)}, a_{(0)}) + \gamma R(s_{(1)}, a_{(1)}) + \gamma^2 R(s_{(2)}, a_{(2)}) + \dots]$$

- Esto permite incorporar diferencias entre el costo que conlleva tomar cierta acción.

Por ejemplo, “**deténgase**” es distinto a “**acelere**” en términos energéticos.

Ecuaciones de Bellman

- Las ecuaciones de Bellman serían ahora

$$V^*(s) = \max_{a \in \mathcal{A}} \left(R(s,a) + \gamma \sum_{s'} P_{sa}(s') V^*(s') \right)$$

- En la **iteración de valor**, repetimos hasta convergencia

$$V(s) \leftarrow \max_{a \in \mathcal{A}} \left(R(s,a) + \gamma \sum_{s'} P_{sa}(s') V(s') \right)$$

que al final será equivalente a $V^*(s)$

- La política óptima estará dada por:

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} \left(R(s,a) + \gamma \sum_{s'} P_{sa}(s') V^*(s') \right)$$

MDP de horizonte finito

- Un MDP de **horizonte finito** es una tupla $\langle \mathcal{S}, \mathcal{A}, \{P_{sa}\}, T, R \rangle$, donde T es el tiempo de horizonte.
- Redefinimos el saldo total como

$$R(s_{(0)}, a_{(0)}) + R(s_{(1)}, a_{(1)}) + \cdots + R(s_{(T)}, a_{(T)})$$

- Queremos encontrar el valor esperado del saldo total, que ahora percibe solo una ventana de tiempo.
- La política óptima puede ser **no-estacionaria**, es decir, que depende del tiempo.
- ¡Las acciones óptimas dependerán de cuándo en el tiempo se deban tomar! pues existirán recompensas no visibles en la ventana de tiempo.

Probabilidades de transición no estacionarias

- Hasta ahora hemos supuesto que

$$\underline{s}_{(t+1)} \sim P_{\underline{s}_t, \underline{a}_t}$$

con $P_{\underline{s}_t, \underline{a}_t}$ estacionario (no depende del tiempo).

- Estas probabilidades pueden depender del tiempo (no-estacionarias):

$$\underline{s}_{(t+1)} \sim P_{\underline{s}_t, \underline{a}_t}^{(t)}$$

- Ejemplos:
 - 1 probabilidades de transición de un avión dependen de cuánto combustible tiene (pues la masa se reduce con el tiempo) o las condiciones climáticas.
 - 2 modelo de tránsito depende de la hora del día

Recompensas no-estacionarias

- Para generalizar el modelo anterior, también podemos suponer que la función de recompensa es no-estacionaria.
- Redefinimos el saldo total como

$$R_{(0)}(s_{(0)}, a_{(0)}) + R_{(1)}(s_{(1)}, a_{(1)}) + \dots + R_{(T)}(s_{(T)}, a_{(T)})$$

Estimación de función de valor

(1)

- Queremos un algoritmo que estime la política óptima
- Modifiquemos la definición de función de valor óptima:

$$V_{(t)}^*(\underline{s}) = E \left[R_{(t)}(\underline{s}_t, \underline{a}_t) + \dots + R_{(T)}(\underline{s}_{(T)}, \underline{a}_{(T)}) \middle| \pi^*, \underline{s}_{(t)} = \underline{s} \right]$$

- La **iteración de valor** para este caso es

$$V_{(t)}^*(\underline{s}) = \max_{\underline{a} \in \mathcal{A}} \left(R_{(t)}(\underline{s}, \underline{a}) + \sum_{\underline{s}' \in \mathcal{S}} P_{\underline{s}\underline{a}}^{(t)}(\underline{s}') V_{(t+1)}^*(\underline{s}') \right)$$

que es un problema de programación dinámica solucionable con recursión.

Estimación de función de valor

(2)

- Para iniciar la recursión usamos

$$V_{(T)}^*(\underline{s}) = \max_{\underline{a} \in \mathcal{A}} R_{(T)}(\underline{s}, \underline{a})$$

pues no consideramos nada luego de T , y retrocedemos en el tiempo.

- Luego encontramos la política óptica con

$$\pi_{(t)}^*(\underline{s}) = \arg \max_{\underline{a} \in \mathcal{A}} \left(R_{(t)}(\underline{s}, \underline{a}) + \sum_{\underline{s}' \in \mathcal{S}} P_{\underline{s}\underline{a}}(\underline{s}') V_{(t+1)}^*(\underline{s}') \right)$$

- Una diferencia de esta versión con tiempo de horizonte es que el resultado es directo, y no requiere iteraciones para convergencia asintótica.

Estimación de función de valor

(3)

- En la implementación, calculamos

$$V_{(T)}^*, V_{(T-1)}^*, V_{(T-2)}^*, \dots V_{(t)}^*$$

y con ellos calculamos

$$\pi_{(T)}^*, \pi_{(T-1)}^*, \pi_{(T-2)}^*, \dots \pi_{(t)}^*$$

Sistemas dinámicos lineales

- Vamos a usar los dos cambios anteriores: recompensas por estado-acción, y MDP con horizonte temporal finito, para resolver un tipo muy particular de sistemas, que sin embargo se encuentran a menudo.
- **Regulación cuadrática lineal** nos permitirá aplicar los conceptos anteriores a casos de espacios de estado y de acción continuos.

Planteando el nuevo problema LQR

- Un MDP se especifica con quintupla: $\langle \mathcal{S}, \mathcal{A}, \{P_{\underline{s}\underline{a}}\}, T, R \rangle$
- Asumiremos el espacio de estados $\mathcal{S} = \mathbb{R}^n$
- El espacio de acciones será $\mathcal{A} = \mathbb{R}^d$.
- Para las probabilidades de transición usaremos:

$$P_{\underline{s}\underline{a}}^{(t)} : \underline{s}_{(t+1)} = \mathbf{A}_{(t)}\underline{s}_{(t)} + \mathbf{B}_{(t)}\underline{a}_{(t)} + \underline{\mathbf{w}}_{(t)}$$

con el ruido gaussiano $\underline{\mathbf{w}}_t \sim \mathcal{N}(\underline{\mathbf{0}}, \underline{\Sigma}_{\underline{\mathbf{w}}})$, y las matrices conocidas *a-priori* $\mathbf{A}_{(t)} \in \mathbb{R}^{n \times n}$, $\mathbf{B}_{(t)} \in \mathbb{R}^{n \times d}$

- Posteriormente veremos que $\underline{\mathbf{w}}_t$ no es muy importante y lo podremos ignorar.

Recompensa cuadrática en LQR

- Supondremos que la recompensa está dada por

$$R_{(t)} \left(\underline{\mathbf{s}}_{(t)}, \underline{\mathbf{a}}_{(t)} \right) = - \left(\underline{\mathbf{s}}_{(t)}^T \mathbf{U}_{(t)} \underline{\mathbf{s}}_{(t)} + \underline{\mathbf{a}}_{(t)}^T \mathbf{V}_{(t)} \underline{\mathbf{a}}_{(t)} \right)$$

con $\mathbf{U}_{(t)} \in \mathbb{R}^{n \times n}$ y $\mathbf{V}_{(t)} \in \mathbb{R}^{d \times d}$, ambas positivas semidefinidas, por lo que $R_{(t)} \leq 0$.

- Por ejemplo:
 - Supongamos un sistema diseñado tal que queremos $\underline{\mathbf{s}}_{(t)} \approx 0$.
 - Podríamos elegir $\mathbf{U}_{(t)} = \mathbf{I}$, $\mathbf{V}_{(t)} = \mathbf{I}$ de modo que

$$R_{(t)} \left(\underline{\mathbf{s}}_{(t)}, \underline{\mathbf{a}}_{(t)} \right) = - \left(\|\underline{\mathbf{s}}_{(t)}\|^2 + \|\underline{\mathbf{a}}_{(t)}\|^2 \right)$$

donde la intención del término $\|\underline{\mathbf{a}}_{(t)}\|^2$ es desalentar un uso arbitrario y exagerado de acciones.

Dinámica no estacionaria

- Derivaremos ahora el caso general de dinámica no estacionaria.
- En este caso tendremos variación temporal en la recompensa y en la dinámica del sistema.
- Para simplificar la comprensión, las derivaciones pueden suponerse estacionarias, es decir, con $\mathbf{M} = \mathbf{M}_{(1)} = \dots = \mathbf{M}_{(T)}$ con todas las matrices involucradas \mathbf{A} , \mathbf{B} , \mathbf{U} y \mathbf{V} .

Modelo lineal

- La suposición principal en lo que sigue es que la **dinámica** del sistema es **lineal**:

$$\underline{s}_{(t+1)} = \mathbf{A}\underline{s}_{(t)} + \mathbf{B}\underline{a}_{(t)}$$

- Realizamos m experimentos con el sistema:

$$\underline{s}_{(0)}^{(1)} \xrightarrow{\underline{a}_{(0)}^{(1)}} \underline{s}_{(1)}^{(1)} \xrightarrow{\underline{a}_{(1)}^{(1)}} \underline{s}_{(2)}^{(1)} \xrightarrow{\underline{a}_{(2)}^{(1)}} \dots \underline{s}_{(T)}^{(1)} \dots$$

$$\underline{s}_{(0)}^{(2)} \xrightarrow{\underline{a}_{(0)}^{(2)}} \underline{s}_{(1)}^{(2)} \xrightarrow{\underline{a}_{(1)}^{(2)}} \underline{s}_{(2)}^{(2)} \xrightarrow{\underline{a}_{(2)}^{(2)}} \dots \underline{s}_{(T)}^{(2)}$$

⋮

$$\underline{s}_{(0)}^{(m)} \xrightarrow{\underline{a}_{(0)}^{(m)}} \underline{s}_{(1)}^{(m)} \xrightarrow{\underline{a}_{(1)}^{(m)}} \underline{s}_{(2)}^{(m)} \xrightarrow{\underline{a}_{(2)}^{(m)}} \dots \underline{s}_{(T)}^{(m)}$$

Estimando modelo dinámico

- Para encontrar el modelo dinámico entonces optimizamos

$$\arg \min_{\mathbf{A}, \mathbf{B}} \frac{1}{2} \sum_{i=1}^m \sum_{t=0}^{T-1} \left\| \underline{\mathbf{s}}_{(t+1)} - \left(\mathbf{A} \underline{\mathbf{s}}_{(t)} + \mathbf{B} \underline{\mathbf{a}}_{(t)} \right) \right\|^2$$

- Otra forma de llegar a un sistema lineal es **linealizando** un sistema no-lineal.

Linealización

- Supongamos que tenemos una función f no-lineal tal que

$$\underline{\mathbf{s}}_{(t+1)} = f\left(\underline{\mathbf{s}}_{(t)}, \underline{\mathbf{a}}_{(t)}\right)$$

- Usamos una serie de Taylor de primer orden alrededor de un punto $(\bar{\underline{\mathbf{s}}}_{(t)}, \bar{\underline{\mathbf{a}}}_{(t)})$ como aproximación lineal.
- Obviemos por un instante la dependencia de la acción:

$$\begin{aligned}\underline{\mathbf{s}}_{(t+1)} &= f(\underline{\mathbf{s}}_{(t)}) \\ &\approx f(\bar{\underline{\mathbf{s}}}_{(t)}) + \mathbf{J}_f(\bar{\underline{\mathbf{s}}}_{(t)})(\underline{\mathbf{s}}_{(t)} - \bar{\underline{\mathbf{s}}}_{(t)})\end{aligned}$$

con $\mathbf{J}_f(\bar{\underline{\mathbf{s}}}_{(t)})$ el jacobiano de f evaluado en $\bar{\underline{\mathbf{s}}}_{(t)}$.

- La linealización obviamente solo es válida en una pequeña vecindad alrededor de $(\bar{\underline{\mathbf{s}}}_{(t)}, \bar{\underline{\mathbf{a}}}_{(t)})$
- Por ello, empíricamente se debe linealizar alrededor de estado observado con frecuencia.

Ejemplo

- Por ejemplo, con el péndulo invertido

$$\underline{s}_{(t+1)} = f \left(\begin{bmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{bmatrix} \right)$$

si elegimos que el centro corresponde a $x = 0$, y queremos que la base esté estática ($\dot{x} = 0$), que el poste esté vertical ($\theta = 0$) y que no se esté cayendo ($\dot{\theta} = 0$), entonces queremos linealizar alrededor de $\underline{\bar{s}}_{(t)} = \underline{0}$.

Linealización general

- Considerando las acciones, la extensión con el jacobiano es directa:

$$\begin{aligned}\underline{\mathbf{s}}_{(t+1)} &= f(\underline{\mathbf{s}}_{(t)}, \underline{\mathbf{a}}_{(t)}) \\ &\approx f(\bar{\underline{\mathbf{s}}}_{(t)}, \bar{\underline{\mathbf{a}}}_{(t)}) + \mathbf{J}_f(\bar{\underline{\mathbf{s}}}_{(t)}, \bar{\underline{\mathbf{a}}}_{(t)}) \begin{bmatrix} \underline{\mathbf{s}}_{(t)} - \bar{\underline{\mathbf{s}}}_{(t)} \\ \underline{\mathbf{a}}_{(t)} - \bar{\underline{\mathbf{a}}}_{(t)} \end{bmatrix}\end{aligned}$$

donde de nuevo $(\bar{s}_{(t)}, \bar{a}_{(t)})$ corresponden al centro de la serie, que es constante en la aproximación, seleccionado adecuadamente para la aplicación concreta.

- Con la linealización, ¿podemos reexpresar lo anterior como

$$\underline{\mathbf{s}}_{(t+1)} = \mathbf{A}\underline{\mathbf{s}}_{(t)} + \mathbf{B}\underline{\mathbf{a}}_{(t)}?$$

Lidiando con el intercepto

- Si expresamos la formula anterior como productos de matrices

$$\begin{aligned}\underline{s}_{(t+1)} &\approx \underbrace{f(\bar{\underline{s}}_{(t)}, \bar{\underline{a}}_{(t)})}_{\underline{c}} + \mathbf{J}_{f(\bar{\underline{s}}_{(t)}, \bar{\underline{a}}_{(t)})} \begin{bmatrix} \underline{s}_{(t)} - \bar{\underline{s}}_{(t)} \\ \underline{a}_{(t)} - \bar{\underline{a}}_{(t)} \end{bmatrix} \\ &= \mathbf{A}\underline{s}_{(t)} + \mathbf{B}\underline{a}_{(t)} + \underbrace{\underline{c} - \mathbf{A}\bar{\underline{s}}_{(t)} - \mathbf{B}\bar{\underline{a}}_{(t)}}_{\bar{\underline{c}}}\end{aligned}$$

- De la teoría de señales y sistemas sabemos que un sistema con intercepto **no** es lineal (sino afín)
- Si $\bar{\underline{c}} \neq \underline{0}$ podemos extender el estado $\underline{s}_{(t)}$ en una dimensión con un 1 y así considerar el intercepto implícitamente en la matriz \mathbf{A} .
- Por ejemplo: $\underline{s} = [\mathbf{1} \quad x \quad \dot{x} \quad \theta \quad \dot{\theta}]^T$ y el intercepto estará en la primera columna de \mathbf{A}

Resumiendo el problema LQR

(1)

- Un MDP de horizonte finito para LQR se especifica con la quintupla: $\langle \mathcal{S}, \mathcal{A}, \{P_{sa}\}, T, R \rangle$ con $\mathcal{S} = \mathbb{R}^n$, $\mathcal{A} = \mathbb{R}^d$.
- Suponemos que para las probabilidades de transición se cumple

$$P_{\underline{s}\underline{a}}^{(t)} : \underline{s}_{(t+1)} = \mathbf{A}_{(t)}\underline{s}_{(t)} + \mathbf{B}_{(t)}\underline{a}_{(t)} + \underline{\mathbf{w}}_{(t)}$$

- Para las recompensas usamos

$$R_{(t)}(\underline{s}_{(t)}, \underline{a}_{(t)}) = - \left(\underline{s}_{(t)}^T \mathbf{U}_{(t)} \underline{s}_{(t)} + \underline{a}_{(t)}^T \mathbf{V}_{(t)} \underline{a}_{(t)} \right)$$

con $\mathbf{U}_{(t)} \in \mathbb{R}^{n \times n}$ y $\mathbf{V}_{(t)} \in \mathbb{R}^{d \times d}$, ambas positivas semidefinidas.

Resumiendo el problema LQR

(2)

- **A**, **B**, **U** y **V** son dadas en la especificación del problema LQR
- La meta es encontrar una política que maximice la recompensa en el tiempo de horizonte finito:

$$\text{máx} : R_{(0)}(s_{(0)}, a_{(0)}) + R_{(1)}(s_{(1)}, a_{(1)}) + \dots + R_{(T)}(s_{(T)}, a_{(T)})$$

Estrategia de solución

- La estrategia será primero encontrar $V_{(T)}^*(\underline{s}_{(T)})$ y luego con programación dinámica (recursivamente) encontrar los valores en tiempos anteriores a T .
- Luego encontraremos la política usando $V_{(t)}^*(\underline{s}_{(t)})$.

Encontrando la función de valor y la política

- El inicio de la recursión se encuentra con:

$$\begin{aligned} V_{(T)}^*(\underline{\mathbf{s}}_{(T)}) &= \max_{\underline{\mathbf{a}}_{(T)} \in \mathcal{A}} R_{(T)}(\underline{\mathbf{s}}_{(T)}, \underline{\mathbf{a}}_{(T)}) \\ &= \max_{\underline{\mathbf{a}}_{(T)} \in \mathcal{A}} \left(- \left(\underline{\mathbf{s}}_{(T)}^T \mathbf{U}_{(T)} \underline{\mathbf{s}}_{(T)} + \underline{\mathbf{a}}_{(T)}^T \mathbf{V}_{(T)} \underline{\mathbf{a}}_{(T)} \right) \right) \\ &= -\underline{\mathbf{s}}_{(T)}^T \mathbf{U}_{(T)} \underline{\mathbf{s}}_{(T)} \end{aligned}$$

donde se ha considerado que con $\underline{\mathbf{a}}_{(T)} = \underline{\mathbf{0}}$ se obtiene el máximo.

- Entonces, la acción óptima de la política en T es

$$\pi_{(T)}^* \left(\underline{\mathbf{s}}_{(T)} \right) = \arg \max_{\underline{\mathbf{a}}_{(T)}} R_{(T)} \left(\underline{\mathbf{s}}_{(T)}, \underline{\mathbf{a}}_{(T)} \right) = \underline{\mathbf{0}}$$

Paso de programación dinámica

- El paso de programación dinámica que buscamos produce $V_{(t)}^*$ a partir de $V_{(t+1)}^*$, como lo hicimos antes.
- Para el caso de estados finitos teníamos

$$V_{(t)}^*(\mathbf{s}_{(t)}) = \max_{\mathbf{a}_{(t)}} \left(R_{(t)}(\mathbf{s}_{(t)}, \mathbf{a}_{(t)}) + \sum_{\mathbf{s}'_{(t+1)}} P_{\mathbf{s}_{(t)}\mathbf{a}_{(t)}}^{(t)}(\mathbf{s}'_{(t+1)}) V_{(t+1)}^*(\mathbf{s}'_{(t+1)}) \right)$$

- Para el caso continuo, la suma se convierte en integral, pero podemos expresar directamente el resultado como el valor esperado:

$$V_{(t)}^*(\mathbf{s}_{(t)}) = \max_{\mathbf{a}_{(t)}} \left(R_{(t)}(\mathbf{s}_{(t)}, \mathbf{a}_{(t)}) + \mathbb{E}_{\mathbf{s}_{(t+1)} \sim P_{\mathbf{s}_{(t)}\mathbf{a}_{(t)}}} \left[V_{(t+1)}^*(\mathbf{s}_{(t+1)}) \right] \right)$$

- Por usar LQR, resulta que la expresión anterior se puede expresar de forma cuadrática.

Expresión cuadrática para función de valor

- Supongamos que

$$V_{(t+1)}^*(\mathbf{s}_{(t+1)}) = \mathbf{s}_{(t+1)}^T \mathbf{\Phi}_{(t+1)} \mathbf{s}_{(t+1)} + \psi_{(t+1)}$$

con $\mathbf{\Phi}_{(t+1)} \in \mathbb{R}^{n \times n}$, y $\psi_{(t+1)} \in \mathbb{R}$

- Se puede demostrar que en un paso de la programación dinámica anterior, entonces

$$V_{(t)}^*(\mathbf{s}_{(t)}) = \mathbf{s}_{(t)}^T \mathbf{\Phi}_{(t)} \mathbf{s}_{(t)} + \psi_{(t)}$$

para nueva matriz $\mathbf{\Phi}_{(t)}$ y valor $\psi_{(t)}$.

- $$V_{(T)}^*(\underline{\mathbf{s}}_{(T)}) = -\underline{\mathbf{s}}_{(T)}^T \mathbf{U}_{(T)} \underline{\mathbf{s}}_{(T)}$$

de modo que $\Phi_{(T)} = -\mathbf{U}_{(T)}$, $\psi_{(T)} = 0$ para entonces

$$V_{(T)}^*(\underline{\mathbf{s}}_{(T)}) = \underline{\mathbf{s}}_{(T)}^T \boldsymbol{\Phi}_{(T)} \underline{\mathbf{s}}_{(T)} + \psi_{(T)}$$

- Luego de varias derivaciones algebraicas se obtiene que

$$V_{(t)}^*(\underline{\mathbf{s}}_{(t)}) = \max_{\underline{\mathbf{a}}_{(t)}} \left(-\underline{\mathbf{s}}_{(t)}^T \mathbf{U}_{(t)} \underline{\mathbf{s}}_{(t)} - \underline{\mathbf{a}}_{(t)}^T \mathbf{V}_{(t)} \underline{\mathbf{a}}_{(t)} + \mathbb{E}_{\underline{\mathbf{s}}_{(t+1)} \sim P_{\underline{\mathbf{s}}_{(t)} \underline{\mathbf{a}}_{(t)}}} \underbrace{\left[\underline{\mathbf{s}}_{(t+1)}^T \boldsymbol{\Phi}_{(t+1)} \underline{\mathbf{s}}_{(t+1)} + \psi_{(t+1)} \right]}_{V_{(t+1)}^*(\underline{\mathbf{s}}_{(t+1)})} \right)$$

$$\text{con } P_{\underline{s}(t)\underline{a}(t)} = \mathcal{N}(\mathbf{A}_{(t)}\underline{s}(t) + \mathbf{B}_{(t)}\underline{a}(t), \boldsymbol{\Sigma}_{\underline{w}}).$$

- La expresión anterior es a su vez cuadrática.

Encontrando las acciones óptimas

- Queremos maximizar lo anterior, que es una función cuadrática.
- Derivamos e igualamos a cero, con lo que se obtiene

$$\begin{aligned}\underline{\mathbf{a}}_{(t)} &= \underbrace{\left(\mathbf{B}_{(t)}^T \boldsymbol{\Phi}_{(t+1)} \mathbf{B}_{(t)} - \mathbf{V}_{(t)} \right)^{-1} \mathbf{B}_{(t)}^T \boldsymbol{\Phi}_{(t+1)} \mathbf{A}_{(t)} \underline{\mathbf{s}}_{(t)}}_{\mathbf{L}_{(t)}} \\ &= \mathbf{L}_{(t)} \underline{\mathbf{s}}_{(t)}\end{aligned}$$

- Lo anterior implica que la acción óptima es una función lineal del estado $\underline{\mathbf{s}}_{(t)}$
- Para calcular la política óptima tenemos entonces:

$$\begin{aligned}\pi_{(t)}^*(\underline{\mathbf{s}}_{(t)}) &= \arg \max_{\underline{\mathbf{a}}_{(t)}} R_{(t)}(\underline{\mathbf{s}}_{(t)}, \underline{\mathbf{a}}_{(t)}) + \mathbb{E}_{\underline{\mathbf{s}}_{(t+1)}} \left[V_{(t+1)}^*(\underline{\mathbf{s}}_{(t+1)}) \right] \\ &= \mathbf{L}_{(t)} \underline{\mathbf{s}}_{(t)}\end{aligned}$$

Ecuación de Riccati de tiempo discreto

- Si tomamos lo anterior y lo introducimos en la recursión de la programación dinámica, encontramos:

$$V_{(t)}^*(\underline{s}_{(t)}) = \underline{s}_{(t)}^T \Phi_{(t)} \underline{s}_{(t)} + \psi_{(t)}$$

con la ecuación de Riccati de tiempo discreto:

$$\begin{aligned} \Phi_{(t)} &= \mathbf{A}_{(t)}^T \left(\Phi_{(t+1)} - \Phi_{(t+1)} \mathbf{B}_{(t)} (\mathbf{B}_{(t)}^T \Phi_{(t+1)} \mathbf{B}_{(t)} - \mathbf{V}_{(t)})^{-1} \mathbf{B}_{(t)}^T \Phi_{(t+1)} \right) \mathbf{A}_{(t)} \\ &\quad - \mathbf{U}_{(t)} \\ \psi_{(t)} &= -\text{tr} [\underline{\Sigma}_{\mathbf{w}} \Phi_{(t+1)}] + \psi_{(t+1)} \end{aligned}$$

- Observe la dependencia recursiva de los cálculos, que nos permite calcular la función de valor óptima a través de expresar $\Phi_{(t)}$ en términos de $\Phi_{(t+1)}$

Resumiendo...

- El algoritmo para encontrar la solución exacta de un LQR de horizonte finito se resume como sigue:
 - 1 Inicialice $\Phi_T = -\mathbf{U}_T$, $\psi_T = 0$.
 - 2 Calcule recursivamente Φ_t , ψ_t usando Φ_{t+1} , ψ_{t+1} con las ecuaciones de Riccati para $t = T-1, T-2, \dots, 0$.
 - 3 Calcule la matrix \mathbf{L}_t usando Φ_{t+1} , ψ_{t+1} .
 - 4 Calcule la política $\pi^*(\mathbf{s}_t) = \mathbf{L}_t \mathbf{s}_t$.
- Observe que considerando la linealidad del sistema, y las suposición de que la recompensa es cuadrática, podemos obtener de forma exacta la función de valor óptima, a pesar de que los espacios de acción y de estado son continuos.
- Este algoritmo escala con $\mathcal{O}(n^3)$ con n la dimensión del espacio de estados, que es mucho mejor que el crecimiento exponencial del método de discretización.

Efectos del ruido

- Nótese que la política óptima **no** depende de $\psi_{(t)}$ ni $\psi_{(t+1)}$.
- Como no depende de $\psi_{(t)}$, ¡el ruido no tiene efecto en la política óptima! pues $\Sigma_{\underline{w}}$ solo aparece en los $\psi_{(t)}$.
- Esta es una propiedad única para los sistemas LQR.
- Si cambiamos la suposición de linealidad del sistema o la forma cuadrática de la recompensa, ya no se cumple esto.
- Note que la función de valor **sí** depende de $\psi_{(t)}$ y por tanto sí es afectada por el ruido.

Resumen

- 1 Repaso
- 2 Recompensas por estado-acción
- 3 MDP de horizonte finito
- 4 Sistemas dinámicos lineales
 - Regulación cuadrática lineal (LQR)
 - Modelos
 - Ecuación de Riccati en tiempo discreto

Este documento ha sido elaborado con software libre incluyendo \LaTeX , Beamer, GNUPlot, GNU/Octave, XFig, Inkscape, GNU-Make y Subversion en GNU/Linux



Este trabajo se encuentra bajo una Licencia Creative Commons Atribución-NoComercial-LicenciarIgual 3.0 Unported. Para ver una copia de esta Licencia, visite <http://creativecommons.org/licenses/by-nc-sa/3.0/> o envíe una carta a Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

© 2017–2019 Pablo Alvarado-Moya Área de Ingeniería en Computadores Instituto Tecnológico de Costa Rica