

2015-04-18

Problemy Decyzyjne Markowa

Problemy Decyzyjne Markowa
na podstawie AIMA ch17 i slajdów S. Russel'aWojciech Jaśkowski
Instytut Informatyki,
Politechnika Poznańska
18 kwietnia 2015

Problemy Decyzyjne Markowa

na podstawie AIMA ch17 i slajdów S. Russel'a

Wojciech Jaśkowski

Instytut Informatyki,
Politechnika Poznańska

18 kwietnia 2015

Sekwencyjne problemy decyzyjne

Sekwencyjny problem decyzyjny

ocena (użyteczność) agenta zależy od **sekwencji decyzji**, a nie od pojedynczej decyzji.

Problemy planowania i przeszukiwania — szczególny przypadek sekwencyjnych problemów decyzyjnych (bez niepewności)

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

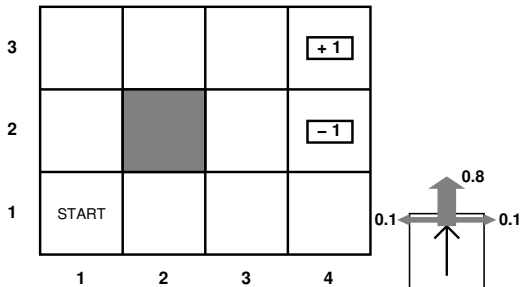
└ Sekwencyjne problemy decyzyjne

Sekwencyjny problem decyzyjny

ocena (użyteczność) agenta zależy od **sekwencji decyzji**, a nie od pojedynczej decyzji.

Problemy planowania i przeszukiwania — szczególny przypadek sekwencyjnych problemów decyzyjnych (bez niepewności)

Środowisko 4x3



- Agent rozpoczyna na polu „Start”.
- W każdym kroku wykonuje akcję Góra, Dół, Lewo, Prawo.
- Stany terminalne: pola +1 lub -1.
- „Nagrody” za wejście na pole wynoszą: +1, -1 lub -0.04 (pozostałe pola)
- Model ruchu agenta: (0.8, 0.1, 0.1)
- Agent zna swoją pozycję.

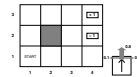
2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Środowisko 4x3

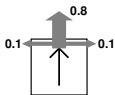
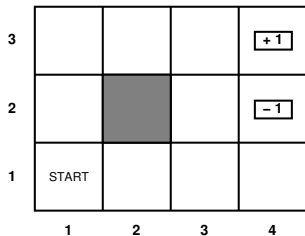
Środowisko 4x3



- Agent rozpoczyna na polu „Start”.
- W każdym kroku wykonuje akcję Góra, Dół, Lewo, Prawo.
- Stany terminalne: pola +1 lub -1.
- „Nagrody” za wejście na pole wynoszą: +1, -1 lub -0.04 (pozostałe pola)
- Model ruchu agenta: (0.8, 0.1, 0.1)
- Agent zna swoją pozycję.

1. Interakcja ze środowiskiem kończy się, gdy agent dotrze do stanów terminalnych
2. Czyli stan środowiska jest znany.

Cechy środowiska



Środowisko deterministyczne \implies proste rozwiązanie GGPPP.

- Ile wynosi prawd. dotarcia do +1 dla GGPPP? [zadanie 2]

Środowisko jest: [zadanie 1]

1. Całkowicie vs. częściowo obserwowalne?
2. Deterministyczne vs. stochastyczne?
3. Epizodyczne vs. sekwencyjne?
4. Statyczne vs. dynamiczne vs. semidynamiczne?
5. Dyskretne vs. ciągłe?
6. Znane vs. nieznanne?

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Cechy środowiska

1. Całkowicie
2. Stochastyczne
3. Sekwencyjne
4. Statyczne
5. Dyskretne
6. Znane
7. Prawd. dotarcia do +1 wynosi $0.8^5 + 0.1^4 \times 0.8 = 0.32776$ (dla strategii GGPPP).

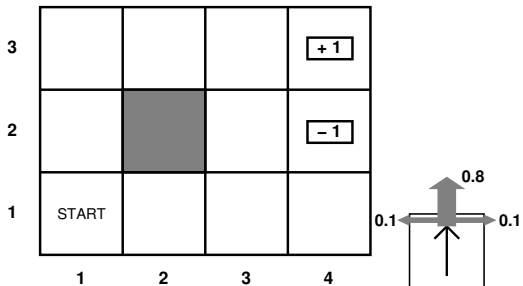


Środowisko jest: [zadanie 1]

- Całkowicie vs. częściowo obserwowalne?
- Deterministyczne vs. stochastyczne?
- Epizodyczne vs. sekwencyjne?
- Statyczne vs. dynamiczne vs. semidynamiczne?
- Dyskretne vs. ciągłe?
- Znane vs. nieznanne?

Środowisko deterministyczne \implies proste rozwiązanie GGPPP.
• Ile wynosi prawd. dotarcia do +1 dla GGPPP? [zadanie 2]

Proces decyzyjny Markowa (MDP)



MDP

to sekwencyjny proces decyzyjny dla środowiska:

- 1 stochastycznego
- 2 całkowicie obserwowalnego
- 3 z „markowskim” modelem przejść
- 4 z addytywną funkcją nagrody

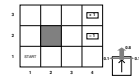
2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Proces decyzyjny Markowa (MDP)

Proces decyzyjny Markowa (MDP)



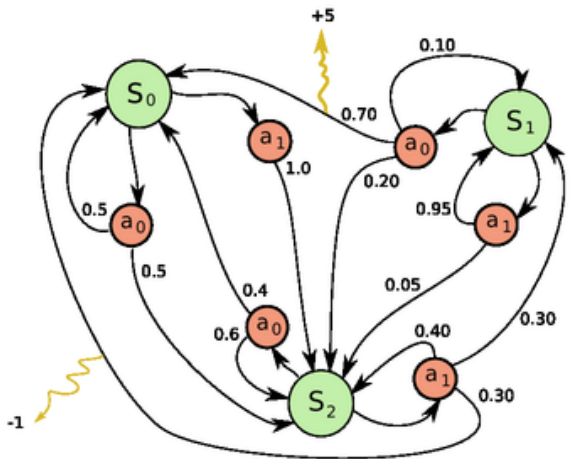
MDP

to sekwencyjny proces decyzyjny dla środowiska:

- 1 stochastycznego
- 2 całkowicie obserwowalnego
- 3 z „markowskim” modelem przejść
- 4 z addytywną funkcją nagrody

1. **MDP = Markov decision process.** W literaturze anglojęzycznej bardzo często używa się skrótu „MDP”, więc warto go pamiętać.

Proces decyzyjny Markowa (MDP) — definicja

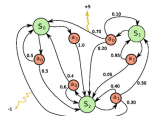


2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Proces decyzyjny Markowa (MDP) — definicja



- Stany $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- Akcje $a \in A$, zbiór akcji ze stanu s , to $A(s)$.

Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- **Stany** $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- **Akcje** $a \in A$, zbiór akcji ze stanu s , to $A(s)$.

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Proces decyzyjny Markowa (MDP) — definicja

1. Nazewnictwo: **model przejść** = **model tranzycji** = **model**
2. **Własność Markowa**: prawd. przejścia ze stanu s do stanu s' zależy tylko od stanu s a nie od historii poprzednich stanów. Przyszłe stany procesu są warunkowo niezależne od stanów przeszłych:

$$P(s^k | a, s^{k-1}) = P(s^k | a, s^{k-1}, s^{k-2}, \dots, s^0).$$

Dzięki własności Markowa, aby podejmować racjonalne (optymalne) decyzje wystarczy aktualny stan (nie trzeba znać całej historii stanów, *vide* definicja racjonalności agenta).

3. Uogólniona nagroda jest za parę (s, a) czyli wykonanie danej akcji w danym stanie, a nie za samą obecność w danym stanie. A jak interpretować nagrodę $R(s, a, s')$?
4. Świat jest niedeterministyczny, ale całkowicie obserwowalny (częściowo jedynie w przypadku POMDP)
5. (PO)MPD jest rozszerzeniem (ukrytych) łańcuchów Markowa o możliwość decyzji (poprzez akcje) i nagrody (motywacja).

Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- **Stany** $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- **Akcje** $a \in A$, zbiór akcji ze stanu s , to $A(s)$.
- **Model (przejść)** $T(s, a, s') \equiv P(s'|s, a)$ — prawd., że akcja a w stanie s prowadzi do stanu s' .
 - **własność Markowa**

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Proces decyzyjny Markowa (MDP) — definicja

Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- Stany $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- Akcje $a \in A$, zbiór akcji ze stanu s , to $A(s)$.
- Model (przejść) $T(s, a, s') \equiv P(s'|s, a)$ — prawd., że akcja a w stanie s prowadzi do stanu s' .
 - własność Markowa

1. Nazewnictwo: **model przejść** = **model tranzycji** = **model**
2. **Własność Markowa**: prawd. przejścia ze stanu s do stanu s' zależy tylko od stanu s a nie od historii poprzednich stanów. Przyszłe stany procesu są warunkowo niezależne od stanów przeszłych:

$$P(s^k | a, s^{k-1}) = P(s^k | a, s^{k-1}, s^{k-2}, \dots, s^0).$$

Dzięki własności Markowa, aby podejmować racjonalne (optymalne) decyzje wystarczy aktualny stan (nie trzeba znać całej historii stanów, *vide* definicja racjonalności agenta).

3. Uogólniona nagroda jest za parę (s, a) czyli wykonanie danej akcji w danym stanie, a nie za samą obecność w danym stanie. A jak interpretować nagrodę $R(s, a, s')$?
4. Świat jest niedeterministyczny, ale całkowicie obserwowalny (częściowo jedynie w przypadku POMDP)
5. (PO)MPD jest rozszerzeniem (ukrytych) łańcuchów Markowa o możliwość decyzji (poprzez akcje) i nagrody (motywacja).

Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- **Stany** $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- **Akcje** $a \in A$, zbiór akcji ze stanu s , to $A(s)$.
- **Model (przejść)** $T(s, a, s') \equiv P(s'|s, a)$ — prawd., że akcja a w stanie s prowadzi do stanu s' .
 - **własność Markowa**
- **Funkcja nagrody** (ang. reward) $R(s)$
 - np. $R(s) = \begin{cases} -0,04 & \text{dla stanów nieterminalnych (kara)} \\ \pm 1 & \text{dla stanów terminalnych} \end{cases}$
 - Można też uogólnić nagrodę do $R(s, a)$ lub $R(s, a, s')$, ale nie zmienia to podstawowych cech problemu.

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- Stany $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- Akcje $a \in A$, zbiór akcji ze stanu s , to $A(s)$.
- Model (przejść) $T(s, a, s') \equiv P(s'|s, a)$ — prawd., że akcja a w stanie s prowadzi do stanu s' .
 - własność Markowa
- Funkcja nagrody (ang. reward) $R(s)$
 - np. $R(s) = \begin{cases} -0,04 & \text{dla stanów nieterminalnych (kara)} \\ \pm 1 & \text{dla stanów terminalnych} \end{cases}$
 - Można też uogólnić nagrodę do $R(s, a)$ lub $R(s, a, s')$, ale nie zmienia to podstawowych cech problemu.

1. Nazewnictwo: **model przejść** = **model tranzycji** = **model**
2. **Własność Markowa**: prawd. przejścia ze stanu s do stanu s' zależy tylko od stanu s a nie od historii poprzednich stanów. Przyszłe stany procesu są warunkowo niezależne od stanów przeszłych:

$$P(s^k | a, s^{k-1}) = P(s^k | a, s^{k-1}, s^{k-2}, \dots, s^0).$$

Dzięki własności Markowa, aby podejmować racjonalne (optymalne) decyzje wystarczy aktualny stan (nie trzeba znać całej historii stanów, *vide* definicja racjonalności agenta).

3. Uogólniona nagroda jest za parę (s, a) czyli wykonanie danej akcji w danym stanie, a nie za samą obecność w danym stanie. A jak interpretować nagrodę $R(s, a, s')$?
4. Świat jest niedeterministyczny, ale całkowicie obserwowalny (częściowo jedynie w przypadku POMDP)
5. (PO)MPD jest rozszerzeniem (ukrytych) łańcuchów Markowa o możliwość decyzji (poprzez akcje) i nagrody (motywacja).

Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- **Stany** $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- **Akcje** $a \in A$, zbiór akcji ze stanu s , to $A(s)$.
- **Model (przejść)** $T(s, a, s') \equiv P(s'|s, a)$ — prawd., że akcja a w stanie s prowadzi do stanu s' .
 - **własność Markowa**
- **Funkcja nagrody** (ang. reward) $R(s)$
 - np. $R(s) = \begin{cases} -0,04 & \text{dla stanów nieterminalnych (kara)} \\ \pm 1 & \text{dla stanów terminalnych} \end{cases}$
 - Można też uogólnić nagrodę do $R(s, a)$ lub $R(s, a, s')$, ale nie zmienia to podstawowych cech problemu.
- **Użyteczność** agenta jest (addytywną) funkcją uzyskanych nagród (np. suma nagród w czasie życia).

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Proces decyzyjny Markowa (MDP) — definicja

Elementy MDP:

- Stany $s \in S$, stan początkowy $s_0 \in S$ i stany terminalne
- Akcje $a \in A$, zbiór akcji ze stanu s , to $A(s)$.
- Model (przejść) $T(s, a, s') \equiv P(s'|s, a)$ — prawd., że akcja a w stanie s prowadzi do stanu s' .
 - własność Markowa
- Funkcja nagrody (ang. reward) $R(s)$
 - np. $R(s) = \begin{cases} -0,04 & \text{dla stanów nieterminalnych (kara)} \\ \pm 1 & \text{dla stanów terminalnych} \end{cases}$
 - Można też uogólnić nagrodę do $R(s, a)$ lub $R(s, a, s')$, ale nie zmienia to podstawowych cech problemu.
- Użyteczność agenta jest (addytywną) funkcją uzyskanych nagród (np. suma nagród w czasie życia).

1. Nazewnictwo: **model przejść = model tranzycji = model**
2. **Własność Markowa**: prawd. przejścia ze stanu s do stanu s' zależy tylko od stanu s a nie od historii poprzednich stanów. Przyszłe stany procesu są warunkowo niezależne od stanów przeszłych:

$$P(s^k | a, s^{k-1}) = P(s^k | a, s^{k-1}, s^{k-2}, \dots, s^0).$$

Dzięki własności Markowa, aby podejmować racjonalne (optymalne) decyzje wystarczy aktualny stan (nie trzeba znać całej historii stanów, *vide* definicja racjonalności agenta).

3. Uogólniona nagroda jest za parę (s, a) czyli wykonanie danej akcji w danym stanie, a nie za samą obecność w danym stanie. A jak interpretować nagrodę $R(s, a, s')$?
4. Świat jest niedeterministyczny, ale całkowicie obserwowalny (częściowo jedynie w przypadku POMDP)
5. (PO)MPD jest rozszerzeniem (ukrytych) łańcuchów Markowa o możliwość decyzji (poprzez akcje) i nagrody (motywacja).

Rozwiązywanie MDP

Cel dla:

- problemów przeszukiwania: znalezienie optymalnej sekwencji akcji.
- MDP: znalezienie **optymalnej polityki** $\pi(s)$

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Rozwiązywanie MDP

Rozwiązywanie MDP

Cel dla:
 • problemów przeszukiwania: znalezienie optymalnej sekwencji akcji.
 • MDP: znalezienie **optymalnej polityki** $\pi(s)$

1. Nie może być niezależna, ponieważ środowisko jest niedeterministyczne.
2. **Polityka** (ang. policy) = **strategia**
3. **optymalna polityka** = **racjonalna polityka** (czyli agent ma być racjonalny!)
4. Polityka *explicite* reprezentuje funkcję agenta i w ten sposób opisuje prostego agenta odruchowego.
5. Na rysunku: Optymalna polityka, gdy kara $R(s)$ wynosi $-0,04$. Zwróćmy uwagę na konserwatywny wybór (3,1).

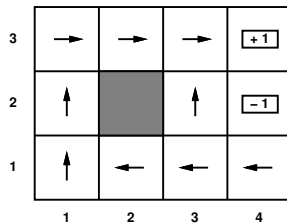
Rozwiązywanie MDP

Cel dla:

- problemów przeszukiwania: znalezienie optymalnej sekwencji akcji.
- MDP: znalezienie **optymalnej polityki** $\pi(s)$

Polityka: $\pi : S \rightarrow A$ (wybrana akcja dla każdego stanu s)

- dlaczego polityka nie może być niezależna od stanu? [zadanie 3]



2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Rozwiązywanie MDP

Rozwiązywanie MDP

- Cel dla:
- problemów przeszukiwania: znalezienie optymalnej sekwencji akcji.
 - MDP: znalezienie **optymalnej polityki** $\pi(s)$
- Polityka:** $\pi : S \rightarrow A$ (wybrana akcja dla każdego stanu s)
- dlaczego polityka nie może być niezależna od stanu? [zadanie 3]



1. Nie może być niezależna, ponieważ środowisko jest niedeterministyczne.
2. **Polityka** (ang. policy) = **strategia**
3. **optymalna polityka** = **racjonalna polityka** (czyli agent ma być racjonalny!)
4. Polityka *explicite* reprezentuje funkcję agenta i w ten sposób opisuje prostego agenta odruchowego.
5. Na rysunku: Optymalna polityka, gdy kara $R(s)$ wynosi $-0,04$. Zwróćmy uwagę na konserwatywny wybór (3,1).

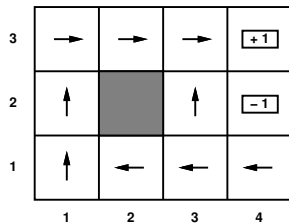
Rozwiązywanie MDP

Cel dla:

- problemów przeszukiwania: znalezienie optymalnej sekwencji akcji.
- MDP: znalezienie **optymalnej polityki** $\pi(s)$

Polityka: $\pi : S \rightarrow A$ (wybrana akcja dla każdego stanu s)

- dlaczego polityka nie może być niezależna od stanu? [zadanie 3]



Optymalna polityka π^* to polityka, która maksymalizuje oczekiwaną wartość funkcji użyteczności.

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Rozwiązywanie MDP

Rozwiązywanie MDP

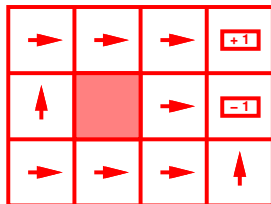
- Cel dla:
- problemów przeszukiwania: znalezienie optymalnej sekwencji akcji.
 - MDP: znalezienie **optymalnej polityki** $\pi(s)$
- Polityka:** $\pi : S \rightarrow A$ (wybrana akcja dla każdego stanu s)
- dlaczego polityka nie może być niezależna od stanu? [zadanie 3]



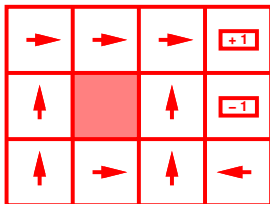
Optymalna polityka π^* to polityka, która maksymalizuje oczekiwaną wartość funkcji użyteczności.

1. Nie może być niezależna, ponieważ środowisko jest niedeterministyczne.
2. **Polityka** (ang. policy) = **strategia**
3. **optymalna polityka** = **racjonalna polityka** (czyli agent ma być racjonalny!)
4. Polityka *explicite* reprezentuje funkcję agenta i w ten sposób opisuje prostego agenta odruchowego.
5. Na rysunku: Optymalna polityka, gdy kara $R(s)$ wynosi $-0,04$. Zwróćmy uwagę na konserwatywny wybór (3,1).

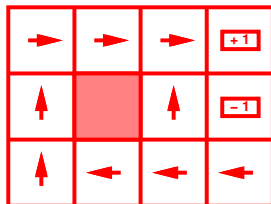
Ryzyko kary vs. nagroda



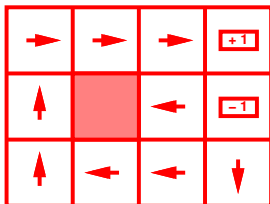
$$r = [-\infty : -1.6284]$$



$$r = [-0.4278 : -0.0850]$$



$$r = [-0.0480 : -0.0274]$$



$$r = [-0.0218 : 0.0000]$$

- r - kara dla stanów nieterminalnych

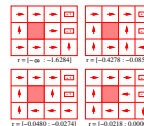
2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Ryzyko kary vs. nagroda

Ryzyko kary vs. nagroda

• r - kara dla stanów nieterminalnych

1. Polityka zmienia się w zależności od r . W tym przypadku manipulujemy tylko karą za ruch. Jeśli kara jest duża, należy jak najszybciej dostać się do stanu terminalnego niezależnie od tego, czy jest to pole $+1$ czy -1 (życie w tym świecie jest bardzo bolesne, więc **w tym świecie** samobójstwo jest jak najbardziej racjonalną opcją). Jeśli r jest małe (da się żyć), robimy wszystko, żeby tylko nie trafić na -1 .

Ryzyko kary vs. nagroda c.d.

- Jak wygląda optymalna polityka dla $r > 0$ [zadanie 5]

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Ryzyko kary vs. nagroda c.d.

1. Życie w tym świecie jest bardzo przyjemne, więc trzeba robić wszystko, byleby żyć. Ucieczka od pól terminalnych. Pozostałe pola, obojętnie. Interesujące: w polu [4,3] należy iść w dół (ryzyko wpadnięcia na -1).

Ryzyko kary vs. nagroda c.d.

- Jak wygląda optymalna polityka dla $r > 0$ [zadanie 5]
- Kompromis pomiędzy **ryzykiem kary** a **szansą na nagrodę** jest cechą charakterystyczną MDP.
 - Nie występuje w problemach deterministycznych.
 - Dlatego MDP rozważane są w wielu dziedzinach:
 - AI
 - badania operacyjne
 - ekonomia
 - teoria sterowania

2015-04-18

Problemy Decyzyjne Markowa

└ MDP

└ Ryzyko kary vs. nagroda c.d.

- Jak wygląda optymalna polityka dla $r > 0$ [zadanie 5]
- Kompromis pomiędzy **ryzykiem kary** a **szansą na nagrodę** jest cechą charakterystyczną MDP.
- Nie występuje w problemach deterministycznych.
- Dlatego MDP rozważane są w wielu dziedzinach:
 - AI
 - badania operacyjne
 - ekonomia
 - teoria sterowania

1. Życie w tym świecie jest bardzo przyjemne, więc trzeba robić wszystko, byleby żyć. Ucieczka od pól terminalnych. Pozostałe pola, obojętnie. Interesujące: w polu [4,3] należy iść w dół (ryzyko wpadnięcia na -1).

Stacjonarność preferencji sekwencji stanów

- W MDP **ocenie podlega sekwencja stanów** → musimy zrozumieć preferencje pomiędzy sekwencjami stanów.
- Naturalne założenie:
*preferencje względem sekwencji stanów są **stacjonarne**, czyli:*

$$[s_1, s_2, \dots] \succ [s'_1, s'_2, \dots] \implies [s_0, s_1, s_2, \dots] \succ [s_0, s'_1, s'_2, \dots]$$

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Stacjonarność preferencji sekwencji stanów

- W MDP ocenie podlega sekwencja stanów → musimy zrozumieć preferencje pomiędzy sekwencjami stanów.
- Naturalne założenie:
preferencje względem sekwencji stanów są **stacjonarne**, czyli:

$$[s_1, s_2, \dots] \succ [s'_1, s'_2, \dots] \implies [s_0, s_1, s_2, \dots] \succ [s_0, s'_1, s'_2, \dots]$$

1. Dotychczas zakładaliśmy, że użyteczność jest sumą nagród za przejścia do stanów. Teraz zajmiemy się tym głębiej tym zagadnieniem i rozważymy inne możliwe definicje funkcji użyteczności.
2. Sekwencja stanów = historia obserwacji z pierwszego wykładu (jest pełna obserwowalność).
3. Rozrysować.

Użyteczność sekwencji stanów

Twierdzenie

Przy założeniu **stacjonarności** sekwencji stanów istnieją tylko dwie możliwości, aby przyporządkować użyteczności do sekwencji stanów:

- **Addytywna funkcja użyteczności**

$$U([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$$

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczność sekwencji stanów

Przy założeniu **stacjonarności** sekwencji stanów istnieją tylko dwie możliwości, aby przyporządkować użyteczności do sekwencji stanów

- Addytywna funkcja użyteczności

$$U([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$$

1. ang. discounted utility function
2. współczynnik dyskontowy (ang. discount factor)
3. a , bo $U(a) = 2$, $U(b) = 0.9^4 \times 3 = 1.96$.
4. Zdyskontowana f. użyteczności ma uzasadnienie w życiu, w ekonomii, w ludzkich zachowaniach. Współczynnik dyskontowy decyduje, czy ważniejsze są nagrody w odległej przyszłości, czy dzisiaj. Czy lepiej dzisiaj zjeść lody, czy zbierać na maszynę do robienie lodów? Czy oszczędzać na telewizor, czy kupić go na raty?

Użyteczność sekwencji stanów

Twierdzenie

Przy założeniu **stacjonarności** sekwencji stanów istnieją tylko dwie możliwości, aby przyporządkować użyteczności do sekwencji stanów:

- **Addytywna funkcja użyteczności**

$$U([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$$

- **Zdyskontowana funkcja użyteczności:**

$$U([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots,$$

gdzie $\gamma \in (0, 1)$ jest **współczynnikiem dyskontowym**.

Która sekwencja nagród $a = [2, 0, 0, 0, 0]$ czy $b = [0, 0, 0, 0, 3]$ jest preferowana przy $\gamma = 0.9$? [\[zadanie 6\]](#)

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Użyteczność sekwencji stanów

1. ang. discounted utility function
2. współczynnik dyskontowy (ang. discount factor)
3. a , bo $U(a) = 2$, $U(b) = 0.9^4 \times 3 = 1.96$.
4. Zdyskontowana f. użyteczności ma uzasadnienie w życiu, w ekonomii, w ludzkich zachowaniach. Współczynnik dyskontowy decyduje, czy ważniejsze są nagrody w odległej przyszłości, czy dzisiaj. Czy lepiej dzisiaj zjeść lody, czy zbierać na maszynę do robienie lodów? Czy oszczędzać na telewizor, czy kupić go na raty?

Twierdzenie

Przy założeniu **stacjonarności** sekwencji stanów istnieją tylko dwie możliwości, aby przyporządkować użyteczności do sekwencji stanów

- **Addytywna funkcja użyteczności**

$$U([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$$

- **Zdyskontowana funkcja użyteczności**

$$U([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$$

gdzie $\gamma \in (0, 1)$ jest **współczynnikiem dyskontowym**

Która sekwencja nagród $a = [2, 0, 0, 0, 0]$ czy $b = [0, 0, 0, 0, 3]$ jest preferowana przy $\gamma = 0.9$? [\[zadanie 6\]](#)

Użyteczności — problem nieskończoności

Addytywna użyteczność może być nieskończona. **Nieskończona przyszłość** jest problemem, bo

jak porównać dwie sekwencje stanów, gdy użyteczności obu wynoszą $+\infty$?

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczności — problem nieskończoności

Addytywna użyteczność może być nieskończona. **Nieskończona przyszłość** jest problemem, bo

jak porównać dwie sekwencje stanów, gdy użyteczności obu wynoszą $+\infty$?

1. Agent może krążyć pomiędzy dwoma nieterminalnymi stanami w nieskończoność zdobywając dodatnie nagrody.

Użyteczności — problem nieskończoności

Możliwe rozwiązania:

- 1 **Skończony horyzont:** koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczności — problem nieskończoności

1. Współczynnik dyskontowy: za „horyzontem”, nagrody wynoszą już prawie 0.
2. Jeśli mamy stany absorbujące to oczekiwana użyteczność każdej polityki jest skończona.
3. Średnia nagroda: tak też można i to ma sens, ale tym się tutaj nie zajmujemy
4. Zwykle wybiera się opcję ze współczynnikiem dyskontowym lub też jeśli to możliwe 3), ponieważ 1) nie spełnia założenia o stacjonarności.

- **Skończony horyzont:** koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.

Użyteczności — problem nieskończoności

Możliwe rozwiązania:

- 1 **Skończony horyzont:** koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.

- 2 Używanie **współczynnika dyskontowego** $\gamma < 1$

- Jeśli $\forall s \in S |R(s)| \leq R_{max}$, to

$$U([s_0, \dots, s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{max}/(1 - \gamma)$$

- Mniejsze $\gamma \implies$ krótszy **horyzont**

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczności — problem nieskończoności

1. Współczynnik dyskontowy: za „horyzontem”, nagrody wynoszą już prawie 0.
2. Jeśli mamy stany absorbujące to oczekiwana użyteczność każdej polityki jest skończona.
3. Średnia nagroda: tak też można i to ma sens, ale tym się tutaj nie zajmujemy
4. Zwykle wybiera się opcję ze współczynnikiem dyskontowym lub też jeśli to możliwe 3), ponieważ 1) nie spełnia założenia o stacjonarności.

- **Skończony horyzont:** koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.
- Używanie **współczynnika dyskontowego** $\gamma < 1$
 - Jeśli $\forall s \in S |R(s)| \leq R_{max}$, to

$$U([s_0, \dots, s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{max}/(1 - \gamma)$$
 - Mniejsze $\gamma \implies$ krótszy **horyzont**

Użyteczności — problem nieskończoności

Możliwe rozwiązania:

1 Skończony horyzont: koniec po liczbie kroków T

- wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
- przykład: pole [3,1] i $T=3$ vs. $T=100$.

2 Używanie współczynnika dyskontowego $\gamma < 1$

- Jeśli $\forall s \in S |R(s)| \leq R_{max}$, to

$$U([s_0, \dots, s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{max}/(1 - \gamma)$$

- Mniejsze $\gamma \implies$ krótszy **horyzont**

3 Stan(y) absorbujące \implies z prawd. 1 agent w końcu zakończy działanie dla każdej polityki π

- polityka, która zawsze prowadzi do stanu terminalnego, to **polityka właściwa**
- możemy używać $\gamma = 1$

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczności — problem nieskończoności

- Współczynnik dyskontowy: za „horyzontem”, nagrody wynoszą już prawie 0.
- Jeśli mamy stany absorbujące to oczekiwana użyteczność każdej polityki jest skończona.
- Średnia nagroda: tak też można i to ma sens, ale tym się tutaj nie zajmujemy
- Zwykle wybiera się opcję ze współczynnikiem dyskontowym lub też jeśli to możliwe 3), ponieważ 1) nie spełnia założenia o stacjonarności.

Użyteczności — problem nieskończoności

Możliwe rozwiązania:

- Skończony horyzont**: koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.
- Używanie **współczynnika dyskontowego** $\gamma < 1$
 - Jeśli $\forall s \in S |R(s)| \leq R_{max}$, to

$$U([s_0, \dots, s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{max}/(1 - \gamma)$$
 - Mniejsze $\gamma \implies$ krótszy **horyzont**
- Stan(y) absorbujące** \implies z prawd. 1 agent w końcu zakończy działanie dla każdej polityki π
 - polityka, która zawsze prowadzi do stanu terminalnego, to **polityka właściwa**
 - możemy używać $\gamma = 1$

Użyteczności — problem nieskończoności

Możliwe rozwiązania:

- 1 **Skończony horyzont:** koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.
- 2 Używanie **współczynnika dyskontowego** $\gamma < 1$
 - Jeśli $\forall_{s \in S} |R(s)| \leq R_{max}$, to

$$U([s_0, \dots, s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{max}/(1 - \gamma)$$
 - Mniejsze $\gamma \implies$ krótszy **horyzont**
- 3 **Stan(y) absorbujące** \implies z prawd. 1 agent w końcu zakończy działanie dla każdej polityki π
 - polityka, która zawsze prowadzi do stanu terminalnego, to **polityka właściwa**
 - możemy używać $\gamma = 1$
- 4 **Średnia nagroda** — maksym. średniej wypłaty na krok

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczności — problem nieskończoności

1. Współczynnik dyskontowy: za „horyzontem”, nagrody wynoszą już prawie 0.
2. Jeśli mamy stany absorbujące to oczekiwana użyteczność każdej polityki jest skończona.
3. Średnia nagroda: tak też można i to ma sens, ale tym się tutaj nie zajmujemy
4. Zwykle wybiera się opcję ze współczynnikiem dyskontowym lub też jeśli to możliwe 3), ponieważ 1) nie spełnia założenia o stacjonarności.

Użyteczności — problem nieskończoności

Możliwe rozwiązania:

- **Skończony horyzont:** koniec po liczbie kroków T
 - wtedy polityka π jest niestacjonarna: tzn. zależy od pozostałego do końca czasu.
 - przykład: pole [3,1] i $T=3$ vs. $T=100$.
- Używanie **współczynnika dyskontowego** $\gamma < 1$
 - Jeśli $\forall_{s \in S} |R(s)| \leq R_{max}$, to

$$U([s_0, \dots, s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{max}/(1 - \gamma)$$
 - Mniejsze $\gamma \implies$ krótszy **horyzont**
- **Stan(y) absorbujące** \implies z prawd. 1 agent w końcu zakończy działanie dla każdej polityki π
 - polityka, która zawsze prowadzi do stanu terminalnego, to **polityka właściwa**
 - możemy używać $\gamma = 1$
- **Średnia nagroda** — maksym. średniej wypłaty na krok

Optymalna polityka i oczekiwana użyteczność

Jak porównywać polityki?

Porównywanie polityk → porównywanie **oczekiwanych wartości użyteczności** sekwencji stanów.

Niech agent realizuje **politykę** π zaczynając od **stanu** s . Wtedy:

- S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Optymalna polityka i oczekiwana użyteczność

1. Oczekiwana wartość użyteczności sekwencji stanów wynika z realizowania danej polityki.
2. Tak. Optymalna polityka dla MDP **nie zależy od stanu początkowego**. Dlatego po prostu piszemy π^* . Ale to jest tylko prawdą jeśli używamy zdyskontowanej użyteczności z nieskończonym horyzontem.

Intuicyjnie: Jeśli π_a^* jest optymalna przy rozpoczęciu ze stanu a , a π_b^* przy rozpoczęciu ze stanu b , to jeśli obie dotrą do stanu c , obie muszą mieć tę samą użyteczność ze stanu c .

Dlatego będziemy po prostu pisali $\pi^* = \pi_a^* = \pi_b^*$.

Optymalna polityka i oczekiwana użyteczność

Jak porównywać polityki?

Porównywanie polityk → porównywanie **oczekiwanych wartości użyteczności** sekwencji stanów.

Niech agent realizuje **politykę** π zaczynając od **stanu** s . Wtedy:

- S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)
- **Oczekiwana użyteczność** przy stosowanie polityki π od stanu s :

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Optymalna polityka i oczekiwana użyteczność

1. Oczekiwana wartość użyteczności sekwencji stanów wynika z realizowania danej polityki.
2. Tak. Optymalna polityka dla MDP **nie zależy od stanu początkowego**. Dlatego po prostu piszemy π^* . Ale to jest tylko prawdą jeśli używamy zdyskontowanej użyteczności z nieskończonym horyzontem.
Intuicyjnie: Jeśli π_a^* jest optymalna przy rozpoczęciu ze stanu a , a π_b^* przy rozpoczęciu ze stanu b , to jeśli obie dotrą do stanu c , obie muszą mieć tę samą użyteczność ze stanu c .
Dlatego będziemy po prostu pisali $\pi^* = \pi_a^* = \pi_b^*$.

- S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)
- **Oczekiwana użyteczność** przy stosowanie polityki π od stanu s :

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

Optymalna polityka i oczekiwana użyteczność

Jak porównywać polityki?

Porównywanie polityk → porównywanie **oczekiwanych wartości użyteczności** sekwencji stanów.

Niech agent realizuje **politykę** π zaczynając od **stanu** s . Wtedy:

- S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)
- **Oczekiwana użyteczność** przy stosowanie polityki π od stanu s :

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

- Optymalna polityka ze stanu s to

$$\pi_s^* = \operatorname{argmax}_\pi U^\pi(s)$$

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Optymalna polityka i oczekiwana użyteczność

Optymalna polityka i oczekiwana użyteczność

Jak porównywać polityki?

Porównywanie polityk → porównywanie **oczekiwanych wartości użyteczności** sekwencji stanów.Niech agent realizuje **politykę** π zaczynając od **stanu** s . Wtedy:• S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)• **Oczekiwana użyteczność** przy stosowanie polityki π od stanu s :

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

• Optymalna polityka ze stanu s to

$$\pi_s^* = \operatorname{argmax}_\pi U^\pi(s)$$

1. Oczekiwana wartość użyteczności sekwencji stanów wynika z realizowania danej polityki.
2. Tak. Optymalna polityka dla MDP **nie zależy od stanu początkowego**. Dlatego po prostu piszemy π^* . Ale to jest tylko prawdą jeśli używamy zdyskontowanej użyteczności z nieskończonym horyzontem.

Intuicyjnie: Jeśli π_a^* jest optymalna przy rozpoczęciu ze stanu a , a π_b^* przy rozpoczęciu ze stanu b , to jeśli obie dotrą do stanu c , obie muszą mieć tę samą użyteczność ze stanu c .

Dlatego będziemy po prostu pisali $\pi^* = \pi_a^* = \pi_b^*$.

Optymalna polityka i oczekiwana użyteczność

Jak porównywać polityki?

Porównywanie polityk → porównywanie **oczekiwanych wartości użyteczności** sekwencji stanów.

Niech agent realizuje **politykę** π zaczynając od **stanu** s . Wtedy:

- S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)
- **Oczekiwana użyteczność** przy stosowanie polityki π od stanu s :

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

- Optymalna polityka ze stanu s to

$$\pi_s^* = \operatorname{argmax}_\pi U^\pi(s)$$

- Zadanie: mamy dwa stany $a \neq b$. Czy $\pi_a^* = \pi_b^*$? [zadanie 7]

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Optymalna polityka i oczekiwana użyteczność

- S_t — zmienna losowa, oznaczająca stan osiągnięty w momencie t (czyli $S_0 = s$)
- **Oczekiwana użyteczność** przy stosowanie polityki π od stanu s :

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

- Optymalna polityka ze stanu s to $\pi_s^* = \operatorname{argmax}_\pi U^\pi(s)$
- Zadanie: mamy dwa stany $a \neq b$. Czy $\pi_a^* = \pi_b^*$? [zadanie 7]

1. Oczekiwana wartość użyteczności sekwencji stanów wynika z realizowania danej polityki.
2. Tak. Optymalna polityka dla MDP **nie zależy od stanu początkowego**. Dlatego po prostu piszemy π^* . Ale to jest tylko prawdą jeśli używamy zdyskontowanej użyteczności z nieskończonym horyzontem.

Intuicyjnie: Jeśli π_a^* jest optymalna przy rozpoczęciu ze stanu a , a π_b^* przy rozpoczęciu ze stanu b , to jeśli obie dotrą do stanu c , obie muszą mieć tę samą użyteczność ze stanu c .

Dlatego będziemy po prostu pisali $\pi^* = \pi_a^* = \pi_b^*$.

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczność stanu

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.
- Czyli możemy ją interpretować jako **użyteczność stanu** s .

2015-04-18

Problemy Decyzyjne Markowa

- └ Użyteczność

- └ Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.
- Czyli możemy ją interpretować jako **użyteczność stanu** s .

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.
- Czyli możemy ją interpretować jako **użyteczność stanu** s .

Użyteczność stanu

$U(s)$ jest oczekiwaną, zdyskontowaną sumą nagród uzyskanych przy realizowaniu optymalnej polityki od stanu s

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczność stanu

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.
- Czyli możemy ją interpretować jako **użyteczność stanu** s .

Użyteczność stanu

$U(s)$ jest oczekiwaną, zdyskontowaną sumą nagród uzyskanych przy realizowaniu optymalnej polityki od stanu s

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.
- Czyli możemy ją interpretować jako **użyteczność stanu** s .

Użyteczność stanu

$U(s)$ jest oczekiwaną, zdyskontowaną sumą nagród uzyskanych przy realizowaniu optymalnej polityki od stanu s

Porównanie R i U :

- $R(s)$ — nagroda krótkoterminowa
- $U(s)$ — nagroda długoterminowa

2015-04-18

Problemy Decyzyjne Markowa

└ Użyteczność

└ Użyteczność stanu

Użyteczność stanu

$U^{\pi^*}(s)$ jest (oczekiwaną, zdyskontowaną) **użytecznością optymalnej polityki** realizowanej od stanu s .

- Będziemy ją oznaczać $U(s)$.
- Czyli możemy ją interpretować jako **użyteczność stanu** s .

Użyteczność stanu

$U(s)$ jest oczekiwaną, zdyskontowaną sumą nagród uzyskanych przy realizowaniu optymalnej polityki od stanu s

Porównanie R i U :

- $R(s)$ — nagroda krótkoterminowa
- $U(s)$ — nagroda długoterminowa

Użyteczność stanów a optymalna polityka

3	0.812	0.868	0.912	+1	3
2	0.762		0.660	-1	2
1	0.705	0.655	0.611	0.388	1
	1	2	3	4	

→	→	→	+1
↑		↑	-1
↑	←	←	←
1	2	3	4

Założmy, że znamy użyteczności $U(s)$ każdego stanu. Czy znamy wtedy (optymalną) politykę? [\[zadanie 8\]](#)

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Użyteczność stanów a optymalna polityka

3	0.812	0.868	0.912	+1	→	→	→	+1
2	0.762		0.660	-1	↑		↑	-1
1	0.705	0.655	0.611	0.388	↑	←	←	←
	1	2	3	4				

Założmy, że znamy użyteczności $U(s)$ każdego stanu. Czy znamy wtedy (optymalną) politykę? [\[zadanie 8\]](#)

- Uwaga: nie wybieramy akcji, które prowadzą po prostu do stanów o najwyższej użyteczności! Musimy wziąć pod uwagę także model tranzycji, czyli to jest niepoprawne:

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} U(s')$$

- Ważne: ten sposób określania polityki możemy zastosować nawet, gdy użyteczności stanów nie są „optymalne”. Wtedy oczywiście „odczytana” w ten sposób polityka też nie będzie optymalna, w ogólności.

Użyteczność stanów a optymalna polityka

3	0.812	0.868	0.912	$\boxed{+1}$	3
2	0.762		0.660	$\boxed{-1}$	2
1	0.705	0.655	0.611	0.388	1
	1	2	3	4	

	→	→	→	$\boxed{+1}$
	↑		↑	$\boxed{-1}$
	↑	←	←	←
	1	2	3	4

Założmy, że znamy użyteczności $U(s)$ każdego stanu. Czy znamy wtedy (optymalną) politykę? [\[zadanie 8\]](#)

Użyteczności stanów jednoznacznie definiują politykę:

- wystarczy znaleźć akcję, która ma maksymalną oczekiwaną użyteczność (**polityka zachłanna ze względu na U**):

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Użyteczność stanów a optymalna polityka

3	0.812	0.868	0.912	$\boxed{+1}$	→	→	→	$\boxed{+1}$
2	0.762		0.660	$\boxed{-1}$	↑		↑	$\boxed{-1}$
1	0.705	0.655	0.611	0.388	↑	←	←	←
	1	2	3	4				

Założmy, że znamy użyteczności $U(s)$ każdego stanu. Czy znamy wtedy (optymalną) politykę? [\[zadanie 8\]](#)

Użyteczności stanów jednoznacznie definiują politykę:
 • wystarczy znaleźć akcję, która ma maksymalną oczekiwaną użyteczność (polityka zachłanna ze względu na U):

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

- Uwaga: nie wybieramy akcji, które prowadzą po prostu do stanów o najwyższej użyteczności! Musimy wziąć pod uwagę także model tranzycji, czyli to jest niepoprawne:

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} U(s')$$

- Ważne: ten sposób określania polityki możemy zastosować nawet, gdy użyteczności stanów nie są „optymalne”. Wtedy oczywiście „odczytana” w ten sposób polityka też nie będzie optymalna, w ogólności.

Optymalna akcja

3	0.812	0.868	0.912	+1	3
2	0.762		0.660	-1	2
1	0.705	0.655	0.611	0.388	1
	1	2	3	4	

→	→	→	+1
↑		↑	-1
↑	←	←	←
1	2	3	4

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

Dlaczego z $s_{3,1}$ idziemy w lewo a nie w górę? [\[zadanie 9\]](#)

2015-04-18

Problemy Decyzyjne Markowa

- Użyteczność

- Optymalna akcja

Optymalna akcja

3	0.812	0.868	0.912	+1
2	0.762		0.660	-1
1	0.705	0.655	0.611	0.388
	1	2	3	4

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

Dlaczego z $s_{3,1}$ idziemy w lewo a nie w górę? [\[zadanie 9\]](#)

1. Obserwacja szczegółowa: winne jest pole (4,1), bo $U(4,1)$ jest male (0.388)

Optymalna akcja

3	0.812	0.868	0.912	+1	3
2	0.762		0.660	-1	2
1	0.705	0.655	0.611	0.388	1
	1	2	3	4	

→	→	→	+1
↑		↑	-1
↑	←	←	←
1	2	3	4

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

Dlaczego z $s_{3,1}$ idziemy w lewo a nie w górę? [zadanie 9]

- z $s_{3,1}$ w górę = $\langle 0.8, 0.1, 0.1 \rangle \times \langle 0.660, 0.655, 0.388 \rangle = 0.6323$
- z $s_{3,1}$ w lewo = $\langle 0.8, 0.1, 0.1 \rangle \times \langle 0.655, 0.611, 0.660 \rangle = 0.6511$

Obserwacja: $\sum_{s'} P(s'|s, a) U(s')$ oznacza (sumaryczną, oczekiwaną) nagrodę jaką otrzymamy idąc do s' .

2015-04-18

Problemy Decyzyjne Markowa

↳ Użyteczność

↳ Optymalna akcja

Optymalna akcja

3	0.812	0.868	0.912	+1	→	→	→	+1
2	0.762		0.660	-1	↑		↑	-1
1	0.705	0.655	0.611	0.388	↑	←	←	←
	1	2	3	4	1	2	3	4

$$v^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

Dlaczego z $s_{3,1}$ idziemy w lewo a nie w górę? [zadanie 9]• z $s_{3,1}$ w górę = $\langle 0.8, 0.1, 0.1 \rangle \times \langle 0.660, 0.655, 0.388 \rangle = 0.6323$ • z $s_{3,1}$ w lewo = $\langle 0.8, 0.1, 0.1 \rangle \times \langle 0.655, 0.611, 0.660 \rangle = 0.6511$ **Obserwacja:** $\sum_{s'} P(s'|s, a) U(s')$ oznacza (sumaryczną, oczekiwaną) nagrodę jaką otrzymamy idąc do s' .

1. Obserwacja szczegółowa: winne jest pole (4,1), bo $U(4,1)$ jest male (0.388)

Programowanie dynamiczne: równanie Bellmana

Dane:

- Stanie s ma użyteczność $U(s)$.
- Ze stanem s związana jest nagroda $R(s)$.
- Model tranzycji to $P(s'|s, a)$.

Jaka jest zależność pomiędzy $U(s)$ a $U(s')$? [zadanie 10]

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Programowanie dynamiczne: równanie Bellmana

Dane:

- Stanie s ma użyteczność $U(s)$.
 - Ze stanem s związana jest nagroda $R(s)$.
 - Model tranzycji to $P(s'|s, a)$.
- Jaka jest zależność pomiędzy $U(s)$ a $U(s')$? [zadanie 10]

1. Podpowiedź: najpierw rozważ sytuację, w której szukamy U^π , dla danej polityki π (czyli agent wykonuje akcję $a = \pi(s)$).
2. Definicja użyteczności stanów prowadzi do prostej zależności pomiędzy użytecznościami sąsiadujących ze sobą stanów:
3. Najlepiej to sobie rozrysować
4. Nagroda = nagroda krótkoterminowa (teraz) + nagroda długoterminowa (w przyszłości)

Programowanie dynamiczne: równanie Bellmana

Dane:

- Stanie s ma użyteczność $U(s)$.
- Ze stanem s związana jest nagroda $R(s)$.
- Model tranzycji to $P(s'|s, a)$.

Jaka jest zależność pomiędzy $U(s)$ a $U(s')$? [zadanie 10]

Równanie Bellman'a (1957)

oczekiwana suma wypłat = aktualna wypłata
+ γ × oczekiwana suma wypłat po wybraniu najlepszej akcji:

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U(s') P(s'|s, a)$$

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Programowanie dynamiczne: równanie Bellmana

Dane:

- Stanie s ma użyteczność $U(s)$.
- Ze stanem s związana jest nagroda $R(s)$.
- Model tranzycji to $P(s'|s, a)$.

Jaka jest zależność pomiędzy $U(s)$ a $U(s')$? [zadanie 10]

Równanie Bellman'a (1957)

oczekiwana suma wypłat = aktualna wypłata
+ γ × oczekiwana suma wypłat po wybraniu najlepszej akcji:

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U(s') P(s'|s, a)$$

1. Podpowieź: najpierw rozważ sytuację, w której szukamy U^π , dla danej polityki π (czyli agent wykonuje akcję $a = \pi(s)$).
2. Definicja użyteczności stanów prowadzi do prostej zależności pomiędzy użytecznościami sąsiadujących ze sobą stanów:
3. Najlepiej to sobie rozrysować
4. Nagroda = nagroda krótkoterminowa (teraz) + nagroda długoterminowa (w przyszłości)

Równanie Bellmana: przykład

Równanie Bellman'a (1957)

oczekiwana suma wypłat = aktualna wypłata

+ $\gamma \times$ oczekiwana suma wypłat po wybraniu najlepszej akcji:

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U(s') P(s'|s, a)$$

Przykład:

$$U(1, 1) = -0.04 + \gamma \max\{$$

$$0.8 \times U(1, 2) + 0.1 \times U(2, 1) + 0.1 \times U(1, 1),$$

$$0.9 \times U(1, 1) + 0.1 \times U(1, 2),$$

$$0.9 \times U(1, 1) + 0.1 \times U(2, 1),$$

$$0.8 \times U(2, 1) + 0.1 \times U(1, 2) + 0.1 \times U(1, 1)\}$$

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Równanie Bellmana: przykład

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U(s') P(s'|s, a)$$

$$U(1, 1) = -0.04 + \gamma \max\{$$

$$0.8 \times U(1, 2) + 0.1 \times U(2, 1) + 0.1 \times U(1, 1),$$

$$0.9 \times U(1, 1) + 0.1 \times U(1, 2),$$

$$0.9 \times U(1, 1) + 0.1 \times U(2, 1),$$

$$0.8 \times U(2, 1) + 0.1 \times U(1, 2) + 0.1 \times U(1, 1)\}$$

Układ równań Bellmana

Układ równań:

$$\begin{cases} U(s_0) = R(s_0) + \gamma \max_{a \in A(s_0)} \sum_{s'} U(s') P(s' | s_0, a) \\ U(s_1) = R(s_1) + \gamma \max_{a \in A(s_1)} \sum_{s'} U(s') P(s' | s_1, a) \\ \vdots \\ U(s_n) = R(s_n) + \gamma \max_{a \in A(s_n)} \sum_{s'} U(s') P(s' | s_n, a) \end{cases}$$

- 1 równanie na stan + n stanów $\implies n$ równań z n niewiadomymi.
- Równania są liniowe czy nieliniowe? [\[zadanie 11\]](#)

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Układ równań Bellmana

Układ równań:

$$\begin{cases} U(s_0) = R(s_0) + \gamma \max_{a \in A(s_0)} \sum_{s'} U(s') P(s' | s_0, a) \\ U(s_1) = R(s_1) + \gamma \max_{a \in A(s_1)} \sum_{s'} U(s') P(s' | s_1, a) \\ \vdots \\ U(s_n) = R(s_n) + \gamma \max_{a \in A(s_n)} \sum_{s'} U(s') P(s' | s_n, a) \end{cases}$$

- 1 równanie na stan + n stanów $\implies n$ równań z n niewiadomymi.
- Równania są liniowe czy nieliniowe? [\[zadanie 11\]](#)

1. Nieliniowe (max). Nie tak łatwo jest rozwiązać \implies podejście iteracyjne

Algorytm iteracji wartości

Algorytm iteracji wartości

- 1 Rozpocznij z losowymi wartościami użyteczności U_0

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji wartości

Algorytm iteracji wartości

Algorytm iteracji wartości

- 1 Rozpocznij z losowymi wartościami użyteczności U_0

1. Ang. Value Iteration

Algorytm iteracji wartości

Algorytm iteracji wartości

- 1 Rozpocznij z losowymi wartościami użyteczności U_0
- 2 W kolejnych krokach i , uaktualniaj użyteczności zgodnie z układem równań Bellmana.

Dla wszystkich stanów s :

$$U_{i+1}(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U_i(s') P(s'|s, a)$$

- Uaktualnienie wykonujemy synchronicznie (kopia tablicy U).

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji wartości

1. Ang. Value Iteration

- Rozpocznij z losowymi wartościami użyteczności U_0
- W kolejnych krokach i , uaktualniaj użyteczności zgodnie z układem równań Bellmana.

Dla wszystkich stanów s :

$$U_{i+1}(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U_i(s') P(s'|s, a)$$

- Uaktualnienie wykonujemy synchronicznie (kopia tablicy U).

Algorytm iteracji wartości

Algorytm iteracji wartości

- 1 Rozpocznij z losowymi wartościami użyteczności U_0
- 2 W kolejnych krokach i , uaktualniaj użyteczności zgodnie z układem równań Bellmana.

Dla wszystkich stanów s :

$$U_{i+1}(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U_i(s') P(s'|s, a)$$

- Uaktualnienie wykonujemy synchronicznie (kopia tablicy U).
- 3 Jeśli osiągnęliśmy równowagę (brak zmian), to mamy **globalne** optimum.

Warto zauważyć: $\max_{a \in A(s)} \sum_{s'} U_i(s') P(s'|s, a)$ jest wybraniem **zachłannej akcji ze względu** na U_i .

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji wartości

Algorytm iteracji wartości

Algorytm iteracji wartości

- Rozpocznij z losowymi wartościami użyteczności U_0
- W kolejnych krokach i , uaktualniaj użyteczności zgodnie z układem równań Bellmana.

Dla wszystkich stanów s :

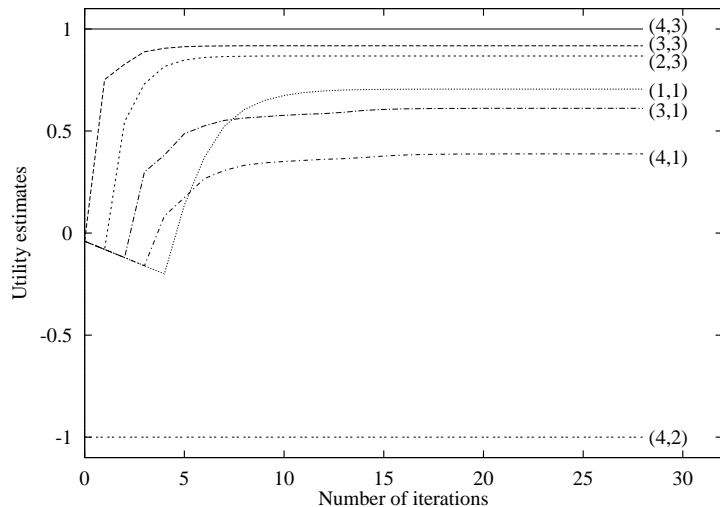
$$U_{i+1}(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} U_i(s') P(s'|s, a)$$

- Uaktualnienie wykonujemy synchronicznie (kopia tablicy U).
- Jeśli osiągnęliśmy równowagę (brak zmian), to mamy **globalne** optimum.

Warto zauważyć: $\max_{a \in A(s)} \sum_{s'} U_i(s') P(s'|s, a)$ jest wybraniem **zachłannej akcji ze względu** na U_i .

1. Ang. Value Iteration

Algorytm iteracji wartości — wykresy



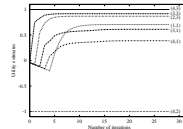
2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji wartości — wykresy

Algorytm iteracji wartości — wykresy



1. Można o tym algorytmie myśleć jak o propagowaniu się informacji poprzez przestrzeń stanów za pomocą lokalnych uaktualnień.

Algorytm iteracji wartości — uwagi i zbieżność

- Algorytm jest zbieżny do globalnego optimum
- Liczba wymaganych iteracji wynosi:

$$N = \lceil \log(2R_{\max}/\epsilon(1 - \gamma)) / \log(1/\gamma) \rceil,$$

gdzie:

- ϵ jest maksymalnym błędem pomiędzy obliczoną użytecznością stanu a użytecznością rzeczywistą
- R_{\max} jest maksymalną wartością nagrody
- Stosowane kryterium stopu:

$$\|U_{i+1} - U_i\| < \epsilon(1 - \gamma)/\gamma$$

- Optymalna polityka jest dostępna zanim wartości użyteczności zbiegną do idealnych.**
- Powyższe działa, gdy $\gamma < 1$. Jeśli $\gamma = 1$, używamy po prostu $< \epsilon$.

2015-04-18

Problemy Decyzyjne Markowa └ Rozwiązywanie MDP

└ Algorytm iteracji wartości — uwagi i zbieżność

- Algorytm jest zbieżny do globalnego optimum
- Liczba wymaganych iteracji wynosi:

$$N = \lceil \log(2R_{\max}/\epsilon(1 - \gamma)) / \log(1/\gamma) \rceil.$$

gdzie:

- ϵ jest maksymalnym błędem pomiędzy obliczoną użytecznością stanu a użytecznością rzeczywistą
- R_{\max} jest maksymalną wartością nagrody
- Stosowane kryterium stopu:

$$\|U_{i+1} - U_i\| < \epsilon(1 - \gamma)/\gamma$$

- Optymalna polityka jest dostępna zanim wartości użyteczności zbiegną do idealnych.
- Powyższe działa, gdy $\gamma < 1$. Jeśli $\gamma = 1$, używamy po prostu $< \epsilon$.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji polityki

1. (ang. Policy Iteration)
2. Czyli użyteczność polityki π .
3. Krok Ulepszanie polityki: To jest proste (i już to widzieliśmy).
4. Chociaż użyteczności wcale nie są poprawne.

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji polityki

1. (ang. Policy Iteration)
2. Czyli użyteczność polityki π .
3. Krok Ulepszanie polityki: To jest proste (i już to widzieliśmy).
4. Chociaż użyteczności wcale nie są poprawne.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji polityki

1. (ang. Policy Iteration)
2. Czyli użyteczność polityki π .
3. Krok Ulepszanie polityki: To jest proste (i już to widzieliśmy).
4. Chociaż użyteczności wcale nie są poprawne.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:
 - 1 (ewaluacja polityki) Oblicz użyteczność $U^\pi(s)$ dla wszystkich stanów $s \in S$.

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji polityki

1. (ang. Policy Iteration)
2. Czyli użyteczność polityki π .
3. Krok Ulepszanie polityki: To jest proste (i już to widzieliśmy).
4. Chociaż użyteczności wcale nie są poprawne.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:
 - 1 (ewaluacja polityki) Oblicz użyteczność $U^\pi(s)$ dla wszystkich stanów $s \in S$.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:
 - 1 (ewaluacja polityki) Oblicz użyteczność $U^\pi(s)$ dla wszystkich stanów $s \in S$.
 - 2 (ulepszenie polityki) Oblicz nową politykę π jako zachłanną względem $U(s)$:

$$\pi(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

2015-04-18

Problemy Decyzyjne Markowa

— Rozwiązywanie MDP

— Algorytm iteracji polityki

1. (ang. Policy Iteration)
2. Czyli użyteczność polityki π .
3. Krok Ulepszanie polityki: To jest proste (i już to widzieliśmy).
4. Chociaż użyteczności wcale nie są poprawne.

Algorytm iteracji polityki

- **Obserwacja:** można otrzymać optymalną politykę nawet gdy wartości funkcji użyteczności są niedokładne.
- **Pomysł:** Szukaj optymalnej polityki oraz wartości funkcji użyteczności równocześnie.

Algorytm iteracji polityki (Howard, 1960)

- 1 $\pi \leftarrow$ losowa polityka
 - 2 Powtarzaj dopóki π się zmienia:
 - 1 (ewaluacja polityki) Oblicz użyteczność $U^\pi(s)$ dla wszystkich stanów $s \in S$.
 - 2 (ulepszenie polityki) Oblicz nową politykę π jako zachłanną względem $U(s)$.
- $$\pi(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

Algorytm iteracji polityki — ocena polityki

Aby obliczyć użyteczności polityki π wystarczy dla wszystkich stanów policzyć:

$$U^\pi(s) = R(s) + \gamma \sum_{s'} U^\pi(s') P(s'|s, \pi(s))$$

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji polityki — ocena polityki

Aby obliczyć użyteczności polityki π wystarczy dla wszystkich stanów policzyć:

$$U^\pi(s) = R(s) + \gamma \sum_{s'} U^\pi(s') P(s'|s, \pi(s))$$

1. Nie ma już argmax'a, ponieważ liczymy użyteczność dla znanej polityki.
2. $O(n^3)$, a nawet w $O(n^{2.3})$

Algorytm iteracji polityki — ocena polityki

Aby obliczyć użyteczności polityki π wystarczy dla wszystkich stanów policzyć:

$$U^\pi(s) = R(s) + \gamma \sum_{s'} U^\pi(s') P(s'|s, \pi(s))$$

Tzn, mamy układ n liniowych równań i n niewiadomych. Można go rozwiązać w czasie: [\[zadanie 12\]](#)

- $O(n)$?
- $O(n^2)$?
- $O(n^3)$?

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Algorytm iteracji polityki — ocena polityki

Aby obliczyć użyteczności polityki π wystarczy dla wszystkich stanów policzyć:

$$U^\pi(s) = R(s) + \gamma \sum_{s'} U^\pi(s') P(s'|s, \pi(s))$$

Tzn, mamy układ n liniowych równań i n niewiadomych. Można go rozwiązać w czasie: [\[zadanie 12\]](#)

- $O(n)$?
- $O(n^2)$?
- $O(n^3)$?

1. Nie ma już argmax'a, ponieważ liczymy użyteczność dla znanej polityki.
2. $O(n^3)$, a nawet w $O(n^{2.3})$

Zmodyfikowany algorytm iteracji polityki

Algorytm iteracji polityki:

- bardzo szybko zbiega do optimum, ale
- każdy jego krok jest kosztowny obliczeniowo, gdy stanów jest dużo:
 - $O(n^3)$ zaczyna być bolesne

2015-04-18

Problemy Decyzyjne Markowa

└─ Rozwiązywanie MDP

└─ Zmodyfikowany algorytm iteracji polityki

Algorytm iteracji polityki:

- bardzo szybko zbiega do optimum, ale
- każdy jego krok jest kosztowny obliczeniowo, gdy stanów jest dużo:
 - $O(n^3)$ zaczyna być bolesne

Zmodyfikowany algorytm iteracji polityki

Algorytm iteracji polityki:

- bardzo szybko zbiega do optimum, ale
- każdy jego krok jest kosztowny obliczeniowo, gdy stanów jest dużo:
 - $O(n^3)$ zaczyna być bolesne

Pomysł: W kroku ewaluacji polityki, obliczmy przybliżoną $U(s)$

- (Przybliżoną) użyteczność $U(s)$ obliczamy wykonując k kroków algorytmu iteracji wartości (ze stałą polityką π) rozpoczynając od ostatnio znanego $U(s)$, czyli:

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

2015-04-18

Problemy Decyzyjne Markowa
└ Rozwiązywanie MDP

└ Zmodyfikowany algorytm iteracji polityki

Algorytm iteracji polityki:

- bardzo szybko zbiega do optimum, ale
- każdy jego krok jest kosztowny obliczeniowo, gdy stanów jest dużo:
 - $O(n^3)$ zaczyna być bolesne

Pomysł: W kroku ewaluacji polityki, obliczmy przybliżoną $U(s)$

- (Przybliżoną) użyteczność $U(s)$ obliczamy wykonując k kroków algorytmu iteracji wartości (ze stałą polityką π) rozpoczynając od ostatnio znanego $U(s)$, czyli:

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

Zmodyfikowany algorytm iteracji polityki

Algorytm iteracji polityki:

- bardzo szybko zbiega do optimum, ale
- każdy jego krok jest kosztowny obliczeniowo, gdy stanów jest dużo:
 - $O(n^3)$ zaczyna być bolesne

Pomysł: W kroku ewaluacji polityki, obliczmy przybliżoną $U(s)$

- (Przybliżoną) użyteczność $U(s)$ obliczamy wykonując k kroków algorytmu iteracji wartości (ze stałą polityką π) rozpoczynając od ostatnio znanego $U(s)$, czyli:

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

- Zwykle zbiega znacznie szybciej niż „czysty” algorytm iteracji wartości lub iteracji polityki.

2015-04-18

Problemy Decyzyjne Markowa

- └ Rozwiązywanie MDP

- └ Zmodyfikowany algorytm iteracji polityki

Algorytm iteracji polityki:

- bardzo szybko zbiega do optimum, ale
- każdy jego krok jest kosztowny obliczeniowo, gdy stanów jest dużo:
 - $O(n^3)$ zaczyna być bolesne

Pomysł: W kroku ewaluacji polityki, obliczmy przybliżoną $U(s)$

- (Przybliżoną) użyteczność $U(s)$ obliczamy wykonując k kroków algorytmu iteracji wartości (ze stałą polityką π) rozpoczynając od ostatnio znanego $U(s)$, czyli:

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

- Zwykle zbiega znacznie szybciej niż „czysty” algorytm iteracji wartości lub iteracji polityki.

Zmodyfikowany algorytm iteracji polityki

Zmodyfikowany algorytm iteracji polityki (van Nunen, 1976)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:
 - 1 (ewaluacja polityki) Wykonaj k kroków iteracji wartości z polityką π :

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

- 2 (ulepszenie polityki) Oblicz nową politykę π jako zachłanną względem U :

$$\pi(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

A co będzie jeśli $k = 1$? [\[zadanie 13\]](#)

2015-04-18

Problemy Decyzyjne Markowa

└ Rozwiązywanie MDP

└ Zmodyfikowany algorytm iteracji polityki

Zmodyfikowany algorytm iteracji polityki (van Nunen, 1976)

- 1 $\pi \leftarrow$ losowa polityka
- 2 Powtarzaj dopóki π się zmienia:
 - 1 (ewaluacja polityki) Wykonaj k kroków iteracji wartości z polityką π :

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$
 - 2 (ulepszenie polityki) Oblicz nową politykę π jako zachłanną względem U :

$$\pi(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

1. Jeśli $k = 1$, to mamy algorytm iteracji wartości!

Rozszerzenia

Algorytmy synchronicznie aktualizowały użyteczności poszczególnych stanów.

W praktyce nie jest to konieczne. Można:

- Wybrać jakikolwiek podzbiór stanów i
- zaaplikować do niego którąkolwiek aktualizację (**ulepszenie polityki** lub **iterację wartości**)
- \implies algorytm **Asynchronicznej Iteracji Polityki**
 - Pod pewnymi warunkami dot. początkowych użyteczności i początkowej polityki jest zbieżny
 - Umożliwia dobranie heurystyki wyboru stanów do aktualizacji, np. algorytm, który koncentruje się na ocenie użyteczności stanów, które z dużym prawd. mają szansę być odwiedzone (**Prioritized Sweeping**)

2015-04-18

Problemy Decyzyjne Markowa

- └ Rozwiązywanie MDP

- └ Rozszerzenia

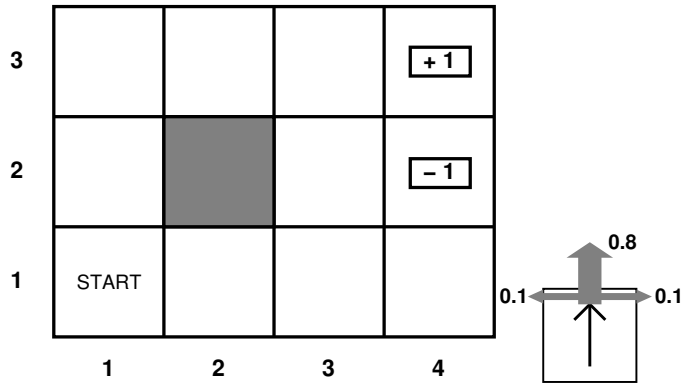
Rozszerzenia

Algorytmy synchronicznie aktualizowały użyteczności poszczególnych stanów.

W praktyce nie jest to konieczne. Można:

- Wybrać jakikolwiek podzbiór stanów i
- zaaplikować do niego którąkolwiek aktualizację (**ulepszenie polityki** lub **iterację wartości**)
- \implies algorytm **Asynchronicznej Iteracji Polityki**
 - Pod pewnymi warunkami dot. początkowych użyteczności i początkowej polityki jest zbieżny
 - Umożliwia dobranie heurystyki wyboru stanów do aktualizacji, np. algorytm, który koncentruje się na ocenie użyteczności stanów, które z dużym prawd. mają szansę być odwiedzone (**Prioritized Sweeping**)

Co zrobić, gdy agent nie wie gdzie jest?



[zadanie 14]

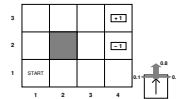
2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ Co zrobić, gdy agent nie wie gdzie jest?

Co zrobić, gdy agent nie wie gdzie jest?



[zadanie 14]

Rzecz o częściowej obserwowalności

- Agent nie zna aktualnego stanu środowiska („nie wie gdzie jest”) \implies nie ma sensu mówić o polityce $\pi(s)$!
- Musi zbierać informacje i wnioskować na temat możliwych stanów środowiska (**stan przekonań**, czyli rozkład prawd. nad stanami) \implies **estymacja stanu** (filtrowanie).

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ Rzecz o częściowej obserwowalności

- Agent nie zna aktualnego stanu środowiska („nie wie gdzie jest”) \implies nie ma sensu mówić o polityce $\pi(s)$!
- Musi zbierać informacje i wnioskować na temat możliwych stanów środowiska (**stan przekonań**, czyli rozkład prawd. nad stanami) \implies **estymacja stanu** (filtrowanie).

- zbiór stanów S ,
- zbiór akcji A ,
- model przejść $P(s'|s, a)$,
- funkcja nagrody $R(s)$.

POMDP:

- **zbiór stanów S ,**
- **zbiór akcji A ,**
- **model przejść $P(s'|s, a)$,**
- **funkcja nagrody $R(s)$,**

1. Model sensoryczny w ogólności może przybrać formę $P(e|s, a, s')$

- zbiór stanów S ,
- zbiór akcji A ,
- model przejść $P(s'|s, a)$,
- funkcja nagrody $R(s)$,
- model sensoryczny $P(e|s)$ czyli prawd., że agent notuje obserwację e będąc w stanie s .

POMDP:

- **zbiór stanów** S ,
- **zbiór akcji** A ,
- **model przejść** $P(s'|s, a)$,
- **funkcja nagrody** $R(s)$,
- **model sensoryczny** $P(e|s)$ czyli prawd., że agent notuje obserwację e będąc w stanie s ,

1. Model sensoryczny w ogólności może przybrać formę $P(e|s, a, s')$

- zbiór stanów S ,
- zbiór akcji A ,
- model przejść $P(s'|s, a)$,
- funkcja nagrody $R(s)$,
- model sensoryczny $P(e|s)$ czyli prawd., że agent notuje obserwację e będąc w stanie s ,
- początkowy stan przekonań b_0 (nieznany jest stan początkowy s_0).

POMDP:

- **zbiór stanów** S ,
- **zbiór akcji** A ,
- **model przejść** $P(s'|s, a)$,
- **funkcja nagrody** $R(s)$,
- **model sensoryczny** $P(e|s)$ czyli prawd., że agent notuje obserwację e będąc w stanie s ,
- **początkowy stan przekonań** b_0 (nieznany jest stan początkowy s_0).

1. Model sensoryczny w ogólności może przybrać formę $P(e|s, a, s')$

Stany przekonań

Aktualizacja stanu przekonań — **problem filtrowania** (estymacji stanu systemu):

$$b'(s') = \alpha P(e|s') \sum_s P(s'|s, a) b(s),$$

gdzie:

- $b(s)$ oznacza prawd., że agent jest w stanie s wg stanu przekonań b ,
- e jest poczynioną obserwacją.
- α jest współczynnikiem normalizującym ($\sum_{s \in S} b(s) = 1$),

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ Stany przekonań

Stany przekonań

Aktualizacja stanu przekonań — **problem filtrowania** (estymacji stanu systemu):

$$b'(s') = \alpha P(e|s') \sum_s P(s'|s, a) b(s),$$

gdzie:

- $b(s)$ oznacza prawd., że agent jest w stanie s wg stanu przekonań b ,
- e jest poczynioną obserwacją.
- α jest współczynnikiem normalizującym ($\sum_{s \in S} b(s) = 1$).

Stany przekonań

Twierdzenie (Astrom, 1965)

Optymalna polityka w POMDP jest funkcją $\pi : B \rightarrow A$, gdzie B jest zbiorem stanów przekonań). **Optymalna polityka nie zależy od aktualnego stanu, w którym jest agent.**

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ Stany przekonań

Stany przekonań

Twierdzenie (Astrom, 1965)

Optymalna polityka w POMDP jest funkcją $\pi : B \rightarrow A$, gdzie B jest zbiorem stanów przekonań). **Optymalna polityka nie zależy od aktualnego stanu, w którym jest agent.**

Jak zachowuje się agent realizujący politykę π ?

Powtarzaj:

- 1 Wykonaj akcję $a = \pi(b)$, gdzie b jest aktualnym stanem przekonań agenta
- 2 Otrzymaj obserwację środowiska e
- 3 Zaktualizuj swój stan przekonań obliczając b' (filtrowanie)

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ Stany przekonań

Stany przekonań

Twierdzenie (Astrom, 1965)

Optymalna polityka w POMDP jest funkcją $\pi : B \rightarrow A$, gdzie B jest zbiorem stanów przekonań). **Optymalna polityka nie zależy od aktualnego stanu, w którym jest agent.**

Jak zachowuje się agent realizujący politykę π ?

Powtarzaj:

- ◆ Wykonaj akcję $a = \pi(b)$, gdzie b jest aktualnym stanem przekonań agenta
- ◆ Otrzymaj obserwację środowiska e
- ◆ Zaktualizuj swój stan przekonań obliczając b' (filtrowanie)

POMDP → MDP

Wniosek

Można przekształcić POMDP w MDP operujący w przestrzeni stanów przekonań, gdzie

- $P(b'|a, b)$ jest prawd., że nowym stanem przekonań będzie b' pod warunkiem, że aktualny stan przekonań to b i agent wykonuje akcję a .

$P(b'|a, b)$ można łatwo wyprowadzić.

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ POMDP → MDP

Wniosek

Można przekształcić POMDP w MDP operujący w przestrzeni stanów przekonań, gdzie

- $P(b'|a, b)$ jest prawd., że nowym stanem przekonań będzie b' pod warunkiem, że aktualny stan przekonań to b i agent wykonuje akcję a .

$P(b'|a, b)$ można łatwo wyprowadzić.

POMDP → MDP

Ale, zauważmy, że przestrzeń stanów przekonań B :

- jest przestrzenią ciągłą (rozkłady prawd.),
- ma bardzo wiele wymiarów.
 - Jeśli mamy n stanów $\implies b$ jest n -wymiarowym wektorem liczb rzeczywistych

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ POMDP → MDP

Ale, zauważmy, że przestrzeń stanów przekonań B :

- jest przestrzenią ciągłą (rozkłady prawd.),
- ma bardzo wiele wymiarów.
 - Jeśli mamy n stanów $\implies b$ jest n -wymiarowym wektorem liczb rzeczywistych

POMDP → MDP

Ale, zauważmy, że przestrzeń stanów przekonań B :

- jest przestrzenią ciągłą (rozkłady prawd.),
- ma bardzo wiele wymiarów.
 - Jeśli mamy n stanów $\implies b$ jest n -wymiarowym wektorem liczb rzeczywistych

Istnieje **algorytm iteracji wartości** dla POMDP (1970), ale jest zbyt wolny nawet dla 4×3

- rozwiązywanie POMDP jest bardzo trudne obliczeniowo (PSPACE-trudne).

Istnieją algorytmy przybliżone oparte na **dynamicznych sieciach baysowskich**.

2015-04-18

Problemy Decyzyjne Markowa

└ POMDP

└ POMDP → MDP

Ale, zauważmy, że przestrzeń stanów przekonań B :

- jest przestrzenią ciągłą (rozkłady prawd.),
- ma bardzo wiele wymiarów.
 - Jeśli mamy n stanów $\implies b$ jest n -wymiarowym wektorem liczb rzeczywistych

Istnieje **algorytm iteracji wartości** dla POMDP (1970), ale jest zbyt wolny nawet dla 4×3

- rozwiązywanie POMDP jest bardzo trudne obliczeniowo (PSPACE-trudne).

Istnieją algorytmy przybliżone oparte na **dynamicznych sieciach baysowskich**.