

Učenje potkrepljivanje (eng. Reinforcement learning)
 Teksto za kurs veštice inteligencije na Matematičkom fakultetu

NEMANJA MIČIĆ

Markovov proces odlučivanja: (S, A, R, T, P, γ)

S - skup svih stanja

A - skup svih akcija

R - skup svih nagrada ($R \subseteq \mathbb{R}$)

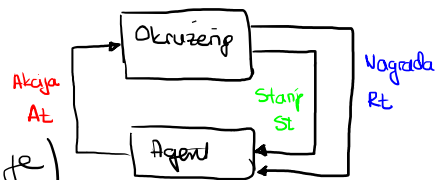
T - tranzicija

P - funkcija prelaska u okruženju

$$P: S \times R \times S \times A \rightarrow \mathbb{R}$$

$$P(s', r | s, a) = P(S_{t+1} = s', R_t = r | S_t = s, A_t = a)$$

γ - faktor umnogavanja $\in [0, 1]$



$S_0, A_0, R_0, S_1, A_1, R_1, S_2, A_2, R_2, \dots$

P je "velika" funkcija, a u praksi je zgodno da je svelimo na verovatnoću

	0	1
0		A
1	B	

Igrač je na $(0, 0)$

Nagrada uvek ista, +1

dostupne akcije: $a_0 = \downarrow$ $a_1 = \rightarrow$

$$P(\boxed{(0, 1)}, 1 | (0, 0), \rightarrow) = 1 = P(S_{t+1} = (0, 0), R_t = 1 | S_t = (0, 0), A_t = \rightarrow)$$

$$P(\boxed{(0, 1)}, 1 | (0, 0), \downarrow) = 0 = P(S_{t+1} = (0, 0), R_t = 1 | S_t = (0, 0), A_t = \downarrow)$$

Markovovo svojstvo

prethodna stanja

Nije važno kako smo došli

$$P(S_t, R_{t-1} | \cancel{S_0, A_0, R_0, \dots, S_{t-1}, A_{t-1}}) = P(S_t, R_{t-1} | \boxed{S_{t-1}, A_{t-1}})$$

Ako je poznato trenutno stanje procesa, za zaključivanje o budućnosti
poznavanje prošlosti nje neophodno

Ako napustim partiju šaha u stanju S , onda drugi igrač može da se završi

Kada agent interaguje sa okolinom generišu se epizode

Politika $\pi \rightarrow$ funkcija pomoću koje agent donosi odluke

\rightarrow raspodela verovatnoća

$\pi(a|s) \rightarrow$ verovatnoća da će agent preduzeti akciju a ako se nađe u stanju S

Mogu biti determinističke i nedeterminističke

Dobitak od koraka t

Tipično nas zanima G_t

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k} \quad (\gamma = 1)$$

Definicija MDP-a (faktor umnogavanja)

Zavis od politike
 Često više stanja nagrade nego daleke nagrade

$$G_t = \gamma^0 R_t + \gamma^1 R_{t+1} + \gamma^2 R_{t+2} + \dots$$

$$\boxed{G_t = R_t + \gamma G_{t+1}}$$

γ - faktor umnogavanja vagu značaj kratkoročnih i dugoročnih nagrada

$\gamma = 0 \rightarrow$ maksimalan fokus na kratkoročnu nagradu \rightarrow politična politika

U praksi na primer $\gamma = 0.9$; slično

$\gamma = 1$

Šta znači rešiti Markovov proces odlučivanja?

Naći politiku koja maksimizuje očekivanu dobit agenta

\rightarrow Teško biti egzaktan usled stohastičnosti

$$J = S_0, A_0, R_0, S_1, A_1, R_1, S_2, A_2, R_2$$

$$P(J) = P(S_0) \prod_{t=0}^{\infty} \pi(A_t | S_t) P(S_{t+1}, R_t | S_t, A_t)$$

Raspodela po svim mogućim trajektorijama

Raspodela po svim mogućim stanjima

Politika koju prati agent
 Verovatnoća da će agent da preduzme A_t kada dođe u S_t

Funkcija prelaska okoline
 Verovatnoća da će da se desi ono što je agent uradio

Očekivana nagrada koju agent dobije
prateći politiku π iz stanja s

$$U^\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$$

↑
Funkcija vrednosti stanja

$$q^\pi(s, a) = \mathbb{E}_\pi [G_t | S_t = s, A_t = a]$$

↑
Funkcija vrednosti
akcije u stanju

→ U^π i q^π se mogu povezati

kao i obrnuto

$$G_t = R_t + \gamma G_{t+1}$$

$$U^\pi(s) = \mathbb{E} [G_t | S_t = s]$$

$$q^\pi(s, a) = \mathbb{E} [G_t | S_t = s, A_t = a]$$

$U^\pi \rightarrow q^\pi$:

$$= \sum_a \pi(a|s) \mathbb{E} [G_t | S_t = s, A_t = a]$$

$$= \sum_a \pi(a|s) q^\pi(s, a)$$

Otežavamo verovatnoću
da ćemo izabrati
akciju a u stanju s

$q^\pi \rightarrow U^\pi$:

$$= \mathbb{E} [R_t + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$= \sum_{s', r} P(s', r | s, a) [r + \gamma \mathbb{E} [G_t | S_t = s']]$$

$$= \sum_{s', r} P(s', r | s, a) [r + \gamma U^\pi(s')]$$

Više nam
odgovara $U^\pi \rightarrow q^\pi$

Ovo je u praksi teže dobiti jer
zahteva da nam je poznato
kako funkcioniše okruženje

Kako da poređamo 2 politike po kvalitetu?

Komentar: Postoje pristupi koji modeluju i okruženje

Ali važn za sve stanja $s \in S$

$$U_{\pi_1}(s) \geq U_{\pi_2}(s), \text{ jasno je da } \pi_1 \text{ bolja od } \pi_2$$

(ili barem jednaka...)

Za MDP postoje optimalna politika π^* koja je bolja (ili jednaka) svim ostalim politikama.

$$U_*(s) = \max_\pi U_\pi(s) \quad q_*(s, a) = \max_\pi q_\pi(s, a)$$

→ Želimo (je) naći! → One dale U_* → koje ne bi bile optimalne
sili

Važno: $U_*(s) = \max_a q_*(s, a)$

Najbolja akcija a
u stanju s

Ali poznajući funkciju q_* , možemo izračunati bar jednu optimalnu politiku.

$$\pi_*(a|s) = \begin{cases} 1, & \text{ako nam odgovara } q_*(a|s) \\ 0, & \text{inače} \end{cases}$$

→ Deterministička politika

⇒ Želimo neku način da izračunavamo q_*

$$U_*(s) = \max_a q_*(s, a)$$

$$= \max_a \mathbb{E}_{\pi_*} [G_t | S_t = s, A_t = a]$$

$$= \max_a \mathbb{E}_{\pi_*} [R_t + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$= \max_a \left(\sum_{s', r} P(s', r | s, a) [r + \gamma \mathbb{E} [G_{t+1} | S_t = s', A_t = a]] \right)$$

$$= \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma U_*(s')]$$

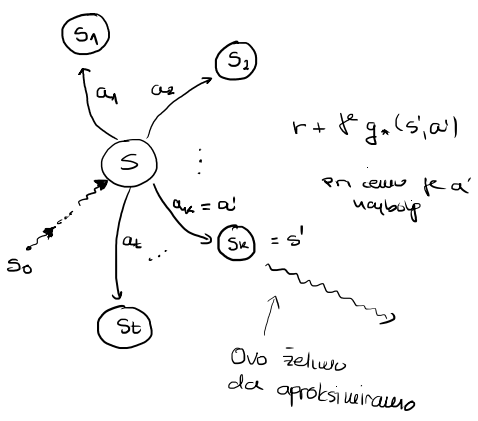
Belmanova (rekurentna)
jednačina za U_*

Belmanova jednačina za q je

$$q_*(s, a) = \sum_{s', r} P(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right]$$

Najbolje što možemo iz stanja
gde ćemo dospeti.

Otežavamo nagradu i buduću dobitku
verovatnoćama mogućih prelaza



Ponašanje pokazanog može doći do iterativnog postupka

$$U_{i+1}(s) \leftarrow \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma U_i(s')]$$

$$q_{i+1}(s, a) \leftarrow \sum_{s', r} P(s', r | s, a) [r + \gamma \max_{a'} q_i(s', a')]$$

13.2 Učenje u nepoznatom okruženju

Praktični iterativni postupak nam ne odgovara jer zahteva poznavanje funkcije okruženja P .

⇒ Agenci treba da istraži okruženje

→ novi element u našem razmišljanju

Poslednjom, Bellmanova jednačina za g glasi

$$g_*(s, a) = \sum_{s', r} P(s', r | s, a) \left[r + \gamma \max_{a'} g_*(s', a') \right]$$

Ekstremni slučaj kada
aproximujemo samo
pomoću prvog opitajanja



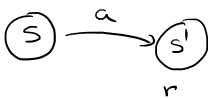
g nam nije poznata → pokušavamo da je ocenimo

Q learning

linearna aproksimacija pomoću
koje varijabilne važnosti

$$g(s, a) \leftarrow (1 - \alpha t) g(s, a) + \alpha t (r + \gamma \max_{a'} g(s', a'))$$

Popravljanje aproksimacije kroz 'vreme' (iteracije)



αt - korak učenja (eng. learning rate)

na primer: $\sum_{t=0}^{\infty} \alpha t = \infty$ $\sum_{t=0}^{\infty} \alpha t^2 < \infty$ (Robbins - Monro uslov)

Koliko varijabilno stari, toliko aproksimativni?
Modelujemo parametrom αt .

Na primer $\frac{1}{t}$ (t broj iteracija)

ϵ pohlepna politika

Politika koja je pohlepna "do na ϵ ". Postoji neka verovatnoća da politika ne bude pohlepna.
Preciznije, politika

$$\pi_{g, \epsilon}(a|s) = \begin{cases} 1 - \epsilon, & \text{ako } a = \arg \max_{a'} g(s, a') \\ \frac{\epsilon}{|A|-1}, & \text{inače} \end{cases}$$

Primer: Neka je S neko stanje
 $A = \{ \downarrow, \uparrow, \leftarrow, \rightarrow \}$ stanje

Alto je samo pohlepna
na primer

$$\pi_g(\downarrow, s) = 1$$

$$\pi_g(\uparrow, s) = 0$$

$$\pi_g(\rightarrow, s) = 0$$

$$\pi_g(\leftarrow, s) = 0$$

Alto je $\epsilon = 0.3$ pohlepna

$$\pi_g(\downarrow, s) = 1 - \epsilon = 1 - 0.3 = 0.7$$

$$\pi_g(\uparrow, s) = \frac{\epsilon}{|A|-1} = \frac{0.3}{4-1} = \frac{0.3}{3} = 0.1$$

$$\pi_g(\rightarrow, s) = 0.1$$

$$\pi_g(\leftarrow, s) = 0.1$$

Verovatnoća koja se ravnomerno raspoređuje na ostale
akcije u odnosu na (dosadašnje) najbolju.

Često ϵ smanjujemo tokom algoritma tako da imaš istraživački. Na primer $\epsilon_t = 1/t$, gde je t broj trenutne iteracije.

Algoritam g -učenja

Ulaz: Broj iteracija H

Izlaz: Aproksimativna funkcija g_*

1. Inicijalizaciju za sve parove s, a

2. Inicijalizaciju s na polazno stanje

3. $t \leftarrow 1$

4. ponavljanje

→ ϵ pohlepna politika sa $\epsilon = 1/t$

5. predumi akciju $a \sim \pi_{g, t}(a|s)$ i opazi nagradu r i novo stanje s'

6. $\alpha t \leftarrow 1/t$

// određuje se korak učenja

7. $g(s, a) \leftarrow (1 - \alpha t) g(s, a) + \alpha t (r + \gamma \max_{a'} g(s', a'))$

// onde se dešava "učenje"

8. Ako je s završno stanje onda

9. inicijalizaciju s na polazno stanje

10. inače

11. $s \leftarrow s'$

12. $t \leftarrow t + 1$

13. dok up ispunjen uslov $t = H$

14. Vrađi g kao rešenje

+ Vrlo jednostavno

- Memorijski (potencijalno) zahtevna tabela Q

- Nema uoč generalizacije sa stanje na stanje

Može li ovo primeniti na "sah"?

teorijski: da!

Praktično: NE! Može li mi čuvati g tabele u memoriji?

To je matrica dimenzija $|S| \times |A|$

X/astancé se! :)