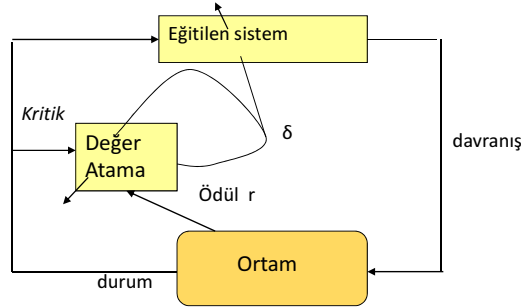


Eğitici-siz Öğrenme

- Pekiştirmeli Öğrenme (reinforcement learning)



Öğrenme işleminin her adımında istenilen yanıtı sağlayan bir eğitici yok

Eğitilen sistem, sonuçta elde edilecek yanıtı erişmek için gerekli davranışı eleştiriye gözönünde tutarak bulmak bulmak zorunda

1

Psikoloji açısından Pekiştirmeli öğrenme

- Biz kararlarımızı nasıl veriyoruz?
- Verdiğimiz kararlar daha sonraki davranışlarımızı nasıl etkiliyor?
- Verdiğimiz kararların sonuçları öğrenmemizi sağlıyor mu?

2

Şartlanma-Pekiştirmeli öğrenme

İlişkilendirme (association):

$$\begin{array}{lcl} O_1 & \longrightarrow & T_1 \\ O_2 & \longrightarrow & T_2 \\ O_1 & \longrightarrow & T_2 \end{array}$$

Klasik Şartlanma

Thordike'nin Yasası: $U_1 \longrightarrow Te_1$ 😊 $U_1 \longrightarrow Te_1$
 $U_2 \longrightarrow Te_2$ 😞 $U_2 \longrightarrow Te_2$

Etkin Şartlanma

δ

3

Psikolojide pekiştirmeli öğrenme

- Of several responses made to the same situation, those which are accompanied or closely followed by satisfaction to the animal will, other things being equal, be more firmly connected with the situation, so that, when it recurs, they will be more likely to recur; those which are accompanied or closely followed by discomfort to the animal will, other things being equal, have their connections with that situation weakened, so that, when it recurs, they will be less likely to occur. The greater the satisfaction or discomfort, the greater the strengthening or weakening of the bond. (Thorndike, 1911, p. 244)

4

Psikolojide pekiştirmeli öğrenme

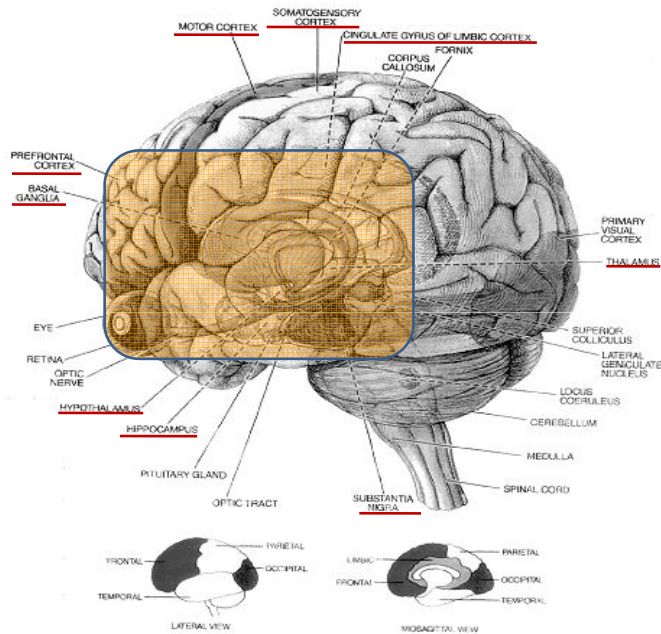
- Thondike (1898): uyarın-yanıt ilişkilendirmesi
(stimulus-response association)
- Skinner (1938): davranışsal düzenleme
(behavioral regulation)

5

Nörobilim açısından Pekiştirmeli öğrenme

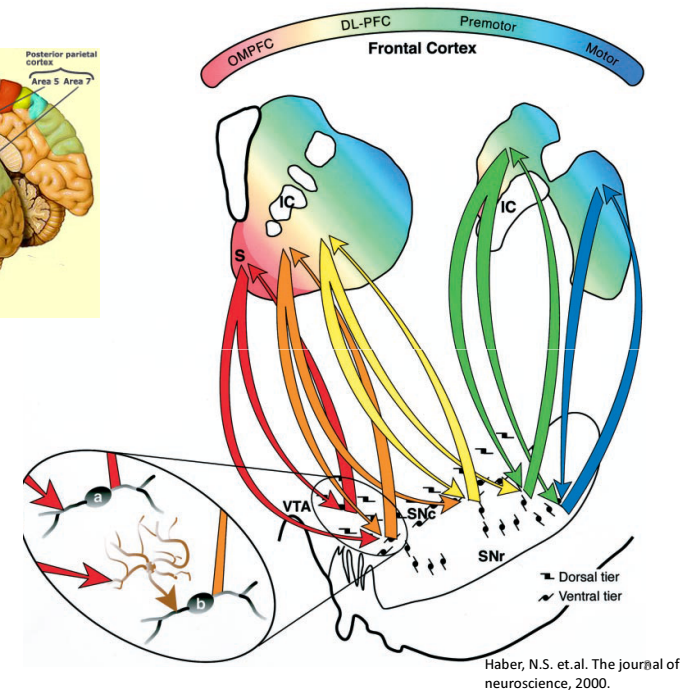
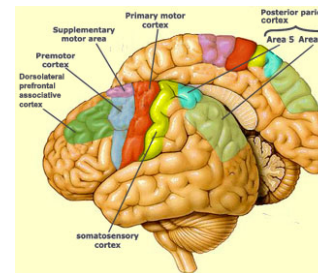
- Beyindeki hangi bölgeler yer alıyor?
- Bu bölgelerin birbirleriyle bağlantıları neler?
- Bağlantıları etkileyen mekanizmalar neler?

6



7

<http://thebrain.mcgill.ca/>



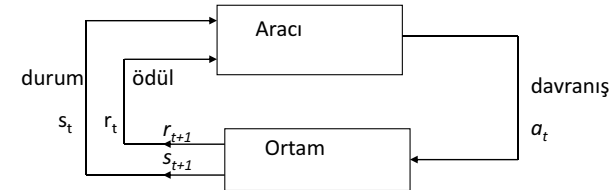
Makina öğrenmesinde pekiştirmeli öğrenme (Machine learning)

- Ortamdaki **belirsizliğe** rağmen bir **amaca erişmek** için aktif karar veren bir aracının ortamla ilişkisi inceleniyor.
- Aracı davranışlarını seçerken **yararlanma-arama** ikilemi ile yüzleşir. (exploit-explore)
- Pekiştirmeli öğrenme sistemi:
 - π yaklaşım (policy)
 - r ödül fonksiyonu (reward function)
 - Q^π, V^π değer fonksiyonu (value function)
 - s ortam modeli

9

Makina öğrenmesinde pekiştirmeli öğrenme

- öğrenen, karar veren aracı
- etkileşim içinde olduğu ortam



- amaç: π^* optimal yaklaşımı bulmak

$$V^*(s) = \max_{\pi} V^{\pi}(s) \quad Q^*(s,a) = \max_{\pi} Q^{\pi}(s,a)$$

10

Makina öğrenmesinde pekiştirmeli öğrenme

$$V^{\pi}(s) = E_{\pi} \{R_t | s_t = s\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s \right\}$$

$$= \sum_a \pi(s,a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^{\pi}(s')]$$

11

Makina öğrenmesinde pekiştirmeli öğrenme

$$Q^{\pi}(s,a) = E_{\pi} \{R_t | s_t = s, a_t = a\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right\}$$

$$= \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^{\pi}(s')]$$

12

Bir Pekiştirmeli öğrenme metodu: Dinamik programlama

- ortam modeli : Markov karar işlevi
(Markov Decision Process (MDP))

- yaklaşım belirleme : ardışıl

$$\begin{aligned} V_{k+1}(s) &= E_{\pi} \{ r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s \} \\ &= \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V_k(s')] \end{aligned}$$

13

Bir Pekiştirmeli öğrenme metodu: Dinamik programlama

- yaklaşım iyileştirme :

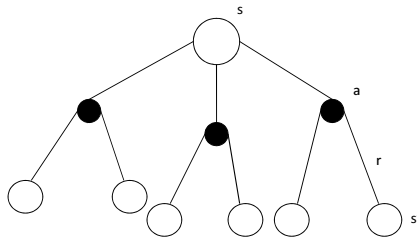
$$\begin{aligned} Q^{\pi}(s, a) &= E_{\pi} \{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s, a_t = a \} \\ &= \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^{\pi}(s')] \end{aligned}$$

- yaklaşım:

$$\pi'(s) = \arg \max_a Q^{\pi}(s, a) = \arg \max_a E \{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s, a_t = a \}$$

14

Bir Pekiştirmeli öğrenme metodu: Dinamik programlama



15

Bir pekiştirmeli öğrenme metodu: Monte Carlo

- ortam modeli: deneyim

gerçek deneyim

(on-line)

benzeşim deneyim

(simulated)

- yaklaşımla ve yaklaşım ötesinde

(on-policy)

(off-policy)

16

Bir pekiştirmeli öğrenme metodu: Monte Carlo

$$V^{\pi}(s) = (1 - \varepsilon) \max_a \sum_{s'} P_{ss'}^a \left[R_{ss'}^a + \gamma V^{\pi}(s') \right] \\ + \frac{\varepsilon}{|A(s)|} \sum_a \sum_{s'} P_{ss'}^a \left[R_{ss'}^a + \gamma V^{\pi}(s') \right]$$

17

Bir pekiştirmeli öğrenme metodu: Monte Carlo

$$Q^{\pi}(s, \pi'(s)) = \sum_a \pi'(s, a) Q^{\pi}(s, a) \\ = \frac{\varepsilon}{|A(s)|} \sum_a Q^{\pi}(s, a) + (1 - \varepsilon) \max_a Q^{\pi}(s, a)$$

18

Bir pekiştirmeli öğrenme metodu: Monte Carlo



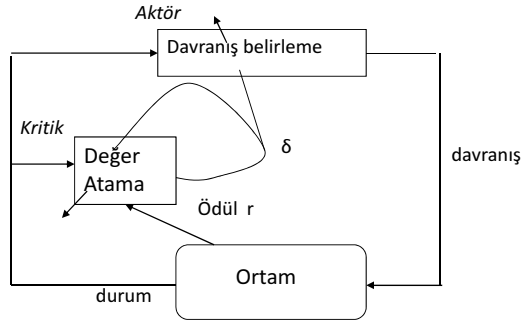
19

Bir pekiştirmeli öğrenme metodu: Zamansal fark

- Monte Carlo metoduna benziyor: ortamın tam modeline gereksinimi yok
- Dinamik programlamaya benziyor: en son çıktıyı beklemeden güncelleme yapabiliyor
- Yaklaşımla: Sarsa
yaklaşım ötesinde: Q-öğrenme (Q-learning)
- Aktör-Kritik

20

Bir zamansal fark metodu: Aktör-kritik



21

Bir Zamansal fark metodu: Aktör-kritik

$$V(t) = w^T(t)s(t)$$

$$\delta(t) = r(t+1) + \gamma V(t+1) - V(t)$$

$$w(k+1) = w(k) + \eta \delta(k)x(k)$$

$$\hat{w}(t+1) = \hat{w}(t) + \alpha \delta(t)e(t)$$

$$e(t+1) = \lambda e(t) + (1 - \lambda)a(t)s(t)$$

$$a(t) = f[\hat{w}^T(t)s(t) + n(t)]$$

22

Pekiştirmeli öğrenmeye ilişkin biliş bilimde bir uygulama

Biliş bilim ne ile ilgileniyor?

- **Davranışsal:** girişe karşılık gelen çıkış ne?
- **Fonksiyonel:** çıkış nasıl oluşuyor?
- **Fiziksel:** çıkışı ne üretiyor?

23

Pekiştirmeli öğrenme için geliştirilecek bir hesaplamalı modelde nelere dikkat edilmeli?

Davranışsal: uyarıcı → yanıt
yanıt → ödül/ceza
ödül → yararlan (exploit)
ceza → ara (explore)

Fonksiyonel: geçmiş değeri değerlendir → beklenti oluştur

Fiziksel: nöral yapıların/bağlantıların özellikleri

24

Pekiştirmeli öğrenme için önerilen bazı hesaplamalı modeller

- Barto & Sutton & Anderson (1983)
makina öğrenmesi
TD (temporal difference)
- Schultz & Dayan & Montague (1997)
Kritik, TD
Kritik: **VTA**
- Suri & Scultz (1998)
Aktör-Kritik, TD
Kritik: **nigrostriatal dopamin nöronları**
Aktör: **Striatum**

25

Bir pekiştirmeli öğrenme metodu: Zamansal fark (Temporal Difference(TD))

Barto, A.G.
IEEE, Syst.
Man&Cyber.1983

Gelecekteki ödülü öngörme

t anındaki
öngörü

$$P(t) = E\{r(t) + \gamma r(t+1) + \gamma^2 r(t+2) + \dots\}$$

$$P(t+1) = E\{r(t+1) + \gamma r(t+2) + \gamma^2 r(t+3) + \dots\}$$

t+1 anındaki
öngörü

$$P(t) = E\{r(t) + \gamma \overbrace{P(t+1)}^{P(t+1)}\}$$

$$P(t) = E\{r(t)\} + \gamma P(t+1)$$

$$\delta(t) \triangleq r(t) + \gamma P(t+1) - P(t) \leftarrow \text{Hata}$$

26

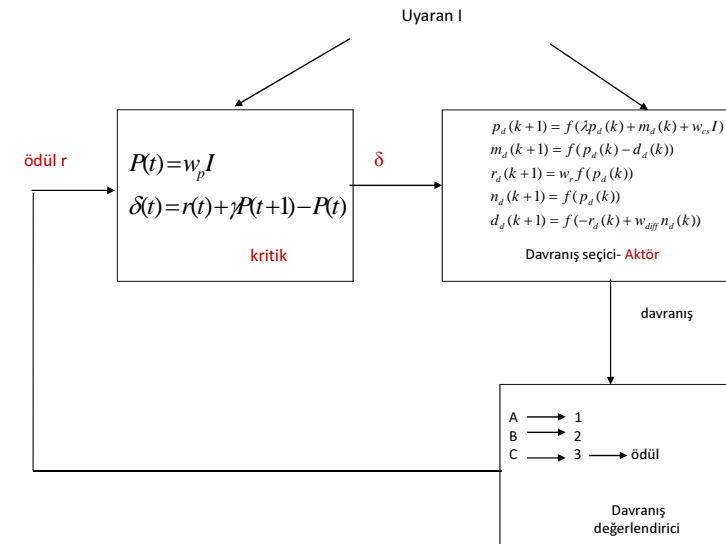
Ardışıl eşleştirme ödevi

- Amaç: Bir dizi öğrenmek

A	→	1
B	→	2
C	→	3
- Yöntem:

1) $U_1 = C$	$Te_1 = 3$	$U_1 \rightarrow Te_1$	ödül
2) $U_2 = B$	$Te_2 = 2$	$U_2 \rightarrow Te_2$	
		$U_1 \rightarrow Te_1$	ödül
3) $U_3 = A$	$Te_3 = 1$	$U_3 \rightarrow Te_3$	
		$U_2 \rightarrow Te_2$	
		$U_1 \rightarrow Te_1$	ödül

27



28

Davranış seçici sistem

$$p_d(k+1) = f(\lambda p_d(k) + m_d(k) + w_{cs} I)$$

$$m_d(k+1) = f(p_d(k) - d_d(k))$$

$$r_d(k+1) = w_r f(p_d(k))$$

$$n_d(k+1) = f(p_d(k))$$

$$d_d(k+1) = f(-r_d(k) + w_{diff} n_d(k))$$

$$f(x) = 0.5 \tanh(ax - \beta)$$

w_r ve w_{cs} öğrenme ile değiştirilecek

29

Güncelleme terimleri

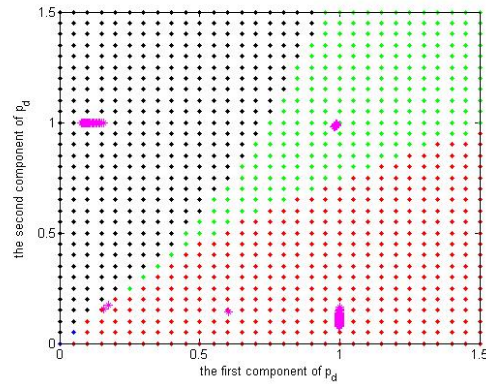
$$w_r(n+1) = w_r(n) + \eta_{w_r} \cdot \delta(n) \cdot w_r \cdot f(p_d(k))$$

$$w_v(n+1) = w_v(n) + \eta_{w_v} \cdot \delta(n) I^T(n-1)$$

$$w_{cs}(n+1) = w_{cs}(n) + \eta_{w_{cs}} \cdot \delta(n) d(n-1) I^T(n-1)$$

30

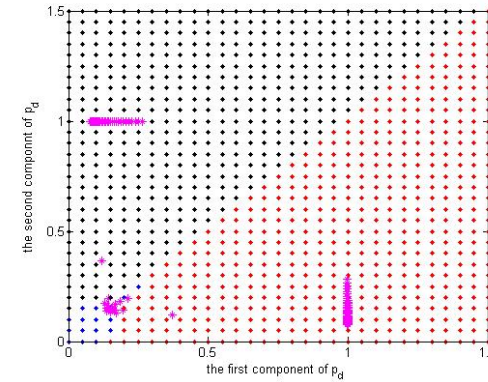
Dinamik sistemin davranışı



$$w_r = \begin{bmatrix} 1.25 & 0 \\ 0 & 1.25 \end{bmatrix}$$

31

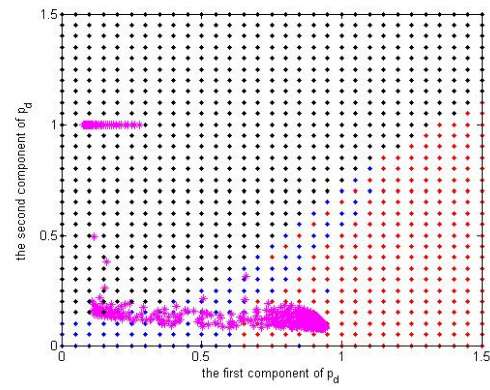
Dinamik sistemin davranışı



$$w_r = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

32

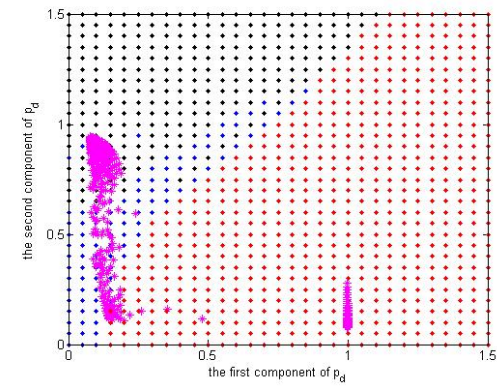
Dinamik sistemin davranışı



$$w_r = \begin{bmatrix} 0.5 & 0 \\ 0 & 1 \end{bmatrix}$$

33

Dinamik sistemin davranışı



$$w_r = \begin{bmatrix} 1 & 0 \\ 0 & 0.5 \end{bmatrix}$$

34