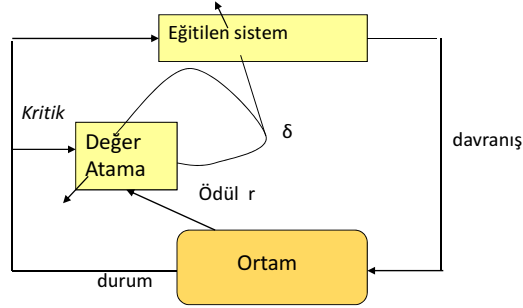


Eğiticiyiz Öğrenme

- Pekiştirmeli Öğrenme (reinforcement learning)



Öğrenme işleminin her adımında istenilen yanıtı sağlayan bir eğitici yok

Eğitilen sistem, sonuçta elde edilecek yanıtı erişmek için gerekli davranışı eleştiriyi 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): $O_1 \longrightarrow T_1$
 $O_2 \longrightarrow T_2$
 $O_1 \longrightarrow T_2$

Klasik Şartlanma

Thordike'nin Yasası: $U_1 \longrightarrow T_{e1}$ 😊 $U_1 \longrightarrow T_{e1}$
 $U_2 \longrightarrow T_{e2}$ 😞 $U_2 \longrightarrow T_{e2}$

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

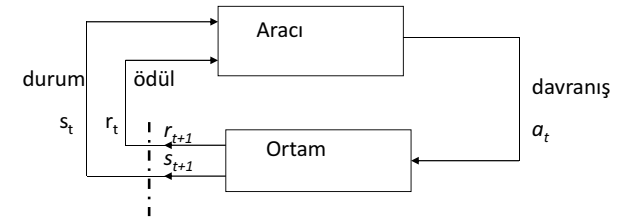
Makina öğrenmesinde pekiştirmeli öğrenme (Machine learning)

- 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



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

yaklaşım

$$\pi_t(s, a)$$

durum davranış

10

Makina öğrenmesinde pekiştirmeli öğrenme

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

Genel

$$\Pr\{s_{t+1} = s', r_{t+1} = r | s_t, a_t, r_t, s_{t-1}, a_{t-1}, r_{t-1}, \dots, s_0, a_0, r_0\}$$

Markov

$$\Pr\{s_{t+1} = s', r_{t+1} = r | s_t, a_t\}$$

$$P_{ss'}^a = \Pr\{s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a\} \quad \text{Durum geçiş}$$

$$R_{ss'}^a = \Pr\{r_{t+1} = r | s_t = s, a_t = a, s_{t+1} = s'\} \quad \text{Yanıt}$$

Daha öncede benzerini görmüştük, nerede? Durum denklemleri

11

Makina öğrenmesinde pekiştirmeli öğrenme

Markov karar işlevi (MDP) ele alındığında değer fonksiyonları ne oluyor?

$$V^{\pi}(s) = E_{\pi}\{R_t | s_t = s\}$$

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

$$Q^{\pi}(s, a) = E_{\pi}\{R_t | s_t = s, a_t = a\}$$

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

Bu değerler, deneyimlere dayalı olarak belirlenebilir. **Monte Carlo Metodu**

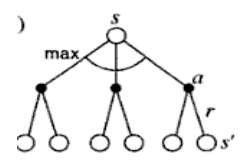
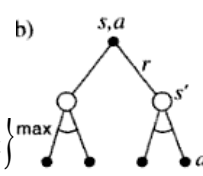
12

Bir pekiştirmeli öğrenme metodu: Monte Carlo

- ortam modeli: deneyim
 - gerçek deneyim (on-line)
 - benzeşim deneyim (simulated)
- yaklaşımla ve yaklaşım ötesinde
 - (on-policy) (off-policy)

13

Optimal değerleri belirleme:

$$\begin{aligned}
 V^*(s) &= \max_{a \in A(s)} Q^{\pi^*}(s, a) \\
 &= \max_a E_{\pi^*} \{R_t | s_t = s, a_t = a\} \\
 &= \max_a E_{\pi^*} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right\} \\
 &= \max_a E_{\pi^*} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s, a_t = a \right\} \\
 &= \max_a E_{\pi^*} \{r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a\} \\
 &= \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')] \\
 Q^*(s, a) &= E \{r_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') | s_t = s, a_t = a\} \\
 &= \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma \max_{a'} Q^*(s', a')]
 \end{aligned}$$



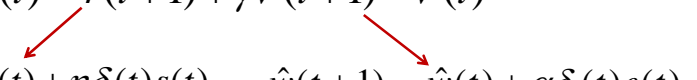
14

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

15

Aktör-kritik için bir uygulama

$$\begin{aligned}
 V(t) &= w^T(t)s(t) \quad \text{Değer fonksiyonu} \\
 a(t) &= f[\hat{w}^T(t)s(t) + n(t)] \quad \text{Davranış} \\
 \delta(t) &= r(t+1) + \gamma V(t+1) - V(t) \quad \text{Hata} \\
 w(t+1) &= w(t) + \eta \delta(t)s(t) \quad \hat{w}(t+1) = \hat{w}(t) + \alpha \delta(t)e(t) \\
 e(t+1) &= \lambda e(t) + (1 - \lambda)a(t)s(t)
 \end{aligned}$$


16

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?

17

Pekiştirmeli öğrenme için geliştirilecek bir hesaplamalı modelde nelere dikkat edilmeli?

Davranışsal: uyarın \longleftrightarrow yanıt
yanıt \longleftrightarrow ödül/ceza
ödül \longrightarrow yararlan (exploit)
ceza \longrightarrow ara (explore)

Fonksiyonel: geçmiş değeri değerlendir \longrightarrow beklenti oluştur

Fiziksel: nöral yapıların/bağlantıların özellikleri

18

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**

19

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

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

Gelecekteki ödülü öğrenme

t anındaki
öngörü

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

t+1 anındaki
öngörü

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

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

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

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

20

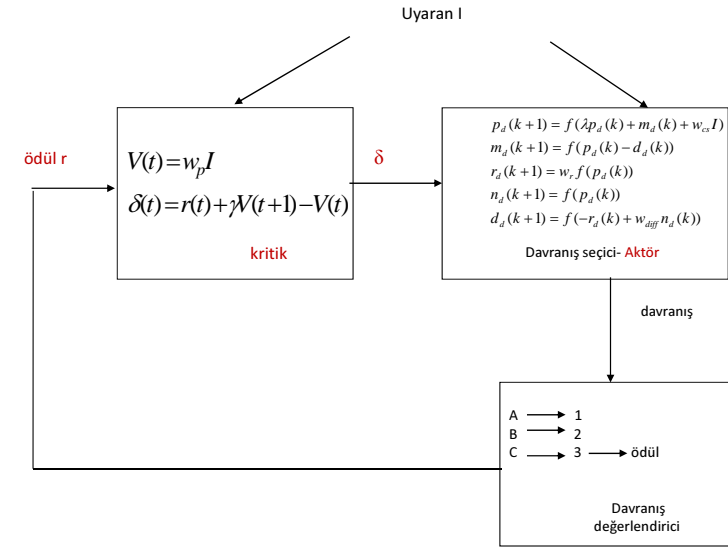
Ardışıl eşleştirme ödevi

- **Amaç:** Bir dizi öğrenmek

A	→	1
B	→	2
C	→	3
- **Yöntem:**

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

21



22

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

23

Güncelleme terimleri

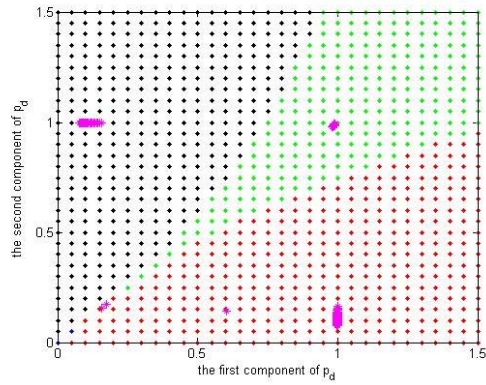
$$w_r(n+1) = w_r(n) + \eta_{w_r} \cdot \delta(n) \cdot w_r \cdot f(p_d(k))$$

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

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

24

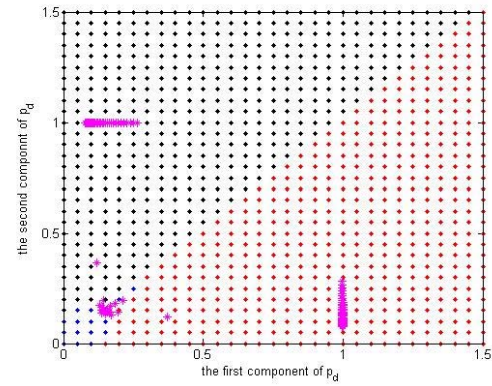
Dinamik sistemin davranışı



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

25

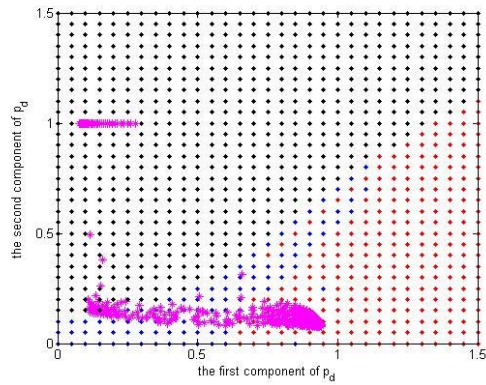
Dinamik sistemin davranışı



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

26

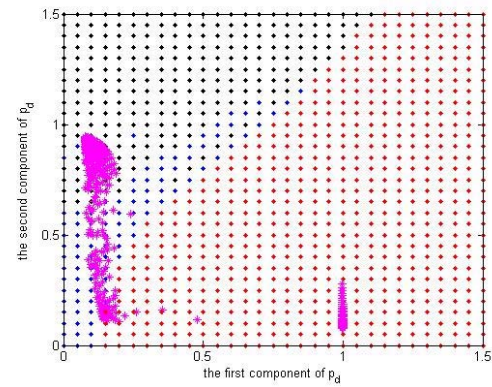
Dinamik sistemin davranışı



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

27

Dinamik sistemin davranışı



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

28