

Visszacsatolós tanulási módszerek a Mario játékra alkalmazva

Schnebli Zoltán

Babeş–Bolyai Tudományegyetem, Kolozsvár

2018 május 26

A gépi tanulás

- ▶ Felügyelt tanulás

A gépi tanulás

- ▶ Felügyelt tanulás
 - ▶ regresszió, osztályozás

A gépi tanulás

- ▶ Felügyelt tanulás
 - ▶ regresszió, osztályozás
- ▶ Felügyeletlen tanulás

A gépi tanulás

- ▶ Felügyelt tanulás
 - ▶ regresszió, osztályozás
- ▶ Felügyeletlen tanulás
 - ▶ klaszterezés

A gépi tanulás

- ▶ Felügyelt tanulás
 - ▶ regresszió, osztályozás
- ▶ Felügyeletlen tanulás
 - ▶ klaszterezés
- ▶ **Félig felügyelt tanulás**

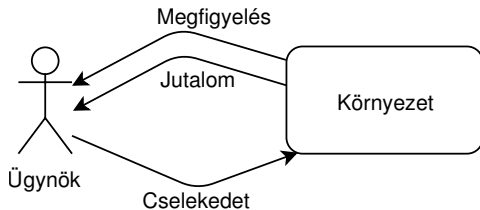
A gépi tanulás

- ▶ Felügyelt tanulás
 - ▶ regresszió, osztályozás
- ▶ Felügyeletlen tanulás
 - ▶ klaszterezés
- ▶ **Félig felügyelt tanulás**
 - ▶ robotika

A visszacsatolós tanulás

Működési elv:

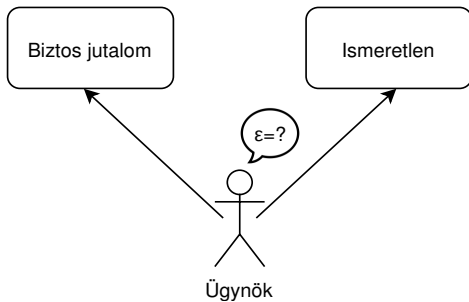
- ▶ Ügynök
- ▶ Környezet
- ▶ Állapot - Cselekedet
- ▶ Jutalom



A visszacsatolós tanulás

Felfedezés és kizsákmányolás

- ▶ ϵ - mohó stratégia
- ▶ csökkenő ϵ - paraméter



A visszacsatolós tanulás

Elemei:

- ▶ Irányelv
- ▶ Jutalomfüggvény
- ▶ Értékfüggvény
- ▶ Környezet modellje

A visszacsatolós tanulás

Markov döntési folyamat

- ▶ $(S, A, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$

A visszacsatolós tanulás

Markov döntési folyamat

- ▶ $(S, A, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$
 - ▶ S - állapottér

A visszacsatolós tanulás

Markov döntési folyamat

- ▶ $(S, A, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$
 - ▶ S - állapottér
 - ▶ A - cselekvéstér

A visszacsatolós tanulás

Markov döntési folyamat

- ▶ $(S, A, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$
 - ▶ S - állapottér
 - ▶ A - cselekvéstér
 - ▶ $\mathcal{P}_a(s, s')$ - s' valószínűségét határozza meg

A visszacsatolós tanulás

Markov döntési folyamat

- ▶ $(S, A, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$
 - ▶ S - állapottér
 - ▶ A - cselekvéstér
 - ▶ $\mathcal{P}_a(s, s')$ - s' valószínűségét határozza meg
 - ▶ $\mathcal{R}_a(s, s')$ - s' -be vezető cselekedet jutalmát határozza meg

A visszacsatolós tanulás

Markov döntési folyamat

- ▶ $(S, A, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$
 - ▶ S - állapottér
 - ▶ A - cselekvéstér
 - ▶ $\mathcal{P}_a(s, s')$ - s' valószínűségét határozza meg
 - ▶ $\mathcal{R}_a(s, s')$ - s' -be vezető cselekedet jutalmát határozza meg
 - ▶ γ - engedményfaktor

A visszacsatolós tanulás

Megoldási módszerek:

- ▶ Dinamikus programozás
 - ▶ kipróbál minden lehetőséget mielőtt lép
 - ▶ sok memória
- ▶ Monte Carlo módszerek
 - ▶ tapasztalat
 - ▶ pontos környezeti modell
- ▶ **Időbeli-differencia tanulás** (ID tanulás)
 - ▶ előző két elv egyesítése

Időbeli-differencia tanulás

A legegyszerűbb ID tanulás

- ▶ TD(0)
- ▶ $V(s_t) \leftarrow V(s_t) + \alpha[r_t + \gamma V(s_{t+1}) - V(s_t)]$
- ▶ α - tanulási ráta

Időbeli-differencia tanulás

Q - tanulás

- ▶ $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$
- ▶ nem igényel környezeti modellt
- ▶ képes optimális irányelvet találni

Időbeli-differencia tanulás

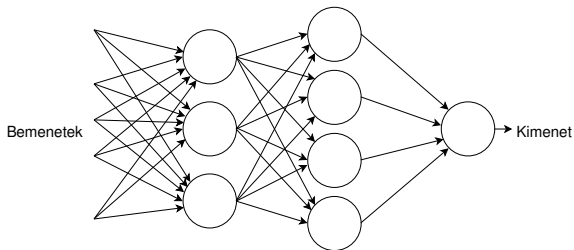
SARSA - tanulás

- ▶ Állapot-Cselekvés-Jutalom-Állapot-Cselekvés
- ▶ Q - tanulásból származik
- ▶ $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$

Időbeli-differencia tanulás

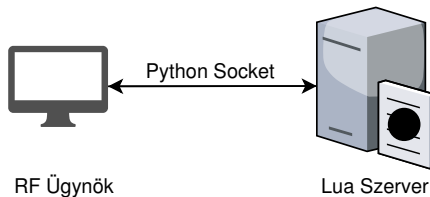
Deep Q - tanulás

- ▶ Neurális háló
- ▶ Tapasztalat
- ▶ 3 réteg



Alkalmazás

- ▶ Visszacsatolósos ügynök
- ▶ Lua szerver
- ▶ Környezeti modell



Alkalmazás

- ▶ JSON üzenetek
- ▶ a játék reprezentációja

[illegible]

Eredmények

- ▶ 55+ százalékos sikeresség
- ▶ maximálisan elért távolság

