

Tema 2 - Învățare prin recompensa

Data publicare: 9.04.2020

Deadline: 26.04.2020

Se acceptă teme trimise până la data de 29.04.2020. Se aplica 0.5 puncte depunere în prima zi de întârziere, și câte 1 punct pentru fiecare din următoarele două. (din max 10 puncte)

Se consideră o lume de tip grid, de dimensiune $N \times M$. În interiorul gridului pot exista bucăți de brânză (plasate în celulele gridului). În această lume există un șoarece și o pisică. Scopul șoarecelui este să adune toate bucățile de brânză fără a fi mâncat de pisică. În cazul în care pisica întâlnește șoarecele jocul se termină. Șoarecele și pisica se pot deplasa în direcțiile N, E, S, V, fără a putea trece celulele-obstacol. Dacă pisica ajunge prea aproape de șoarece (maxim A pași între ei) aceasta îl va urmări. În caz contrar, pisica execută mișcări aleatoare. Pisica ignoră bucățile de brânză de pe hartă.

Modelați lumea descrisă în problemă (gridul și mișcarea pisicii) și ajutați șoarecele să învețe cum să adune toate bucățile de brânză fără a fi mâncat de pisică, folosind algoritmul de învățare prin recompensă Q-learning. Programul trebuie să poată fi rulat în două moduri: pas cu pas sau execuție continuă. Afișarea poate fi făcută grafic sau în mod text, dar să fie cât mai realistă: să fie vizibile celulele gridului, obstacolele, bucățile de brânză, șoarecele și pisica, precum și deplasările celor două animale.

Se vor experimenta 4 strategii de exploatare / explorare:

- Max First: exploatare pură, acțiunea cu utilitatea maximă va fi aleasă de fiecare dată
- Random: ignoră tabela de utilități și alege aleator o acțiune posibilă
- Exploatare: atâta timp cât sunt acțiuni nefolosite într-o stare, se va alege aleator dintre acestea
- Explorare / exploatare ponderată: permite echilibrarea explorării cu exploatarea folosind o probabilitate pentru fiecare acțiune bazată pe valoarea utilității acesteia

Descrierea gridului este preluată dintr-un fișier text cu următorul format:

$N \ M$

$N \times M$ valori 0, 1 sau 2 // 0 reprezintă celula liberă, 1 reprezintă obstacol,
// 2 reprezintă bucata de brânză

A // numărul de pași maxim pentru care putem spune că șoarecele este prea aproape

de pisică

xs ys // poziția inițială a soarecelui în grid

xp yp // poziția inițială a pisicii în grid

Se cer:

- [2.5p] Implementarea jocului (Pentru implementare puteți folosi [Gym](#))
 - Numărul de bucăți de brânză și dimensiunea habitatului să poată fi variat
 - Implementarea trebuie să permită rularea pas cu pas (după antrenare)
- [6.5p] Implementarea sistemului
 - [3.5p] Algoritm Q-Learning
 - [1p] Parametri variabili (rata de învățare, factor de discount)
 - [0.5p] MaxFirst
 - [0.5p] Random
 - [0.5p] Exploatare
 - [1p] Explorarea ponderată
 - [3p] Grafice (pentru toate strategiile de exploatare)
 - [1p] Evoluția scorului în funcție de numărul episodului de antrenament
 - [2p] Procentul de jocuri câștigate în funcție de valoarea (jocurile vor fi rulate în batchuri):
 - factorului de învățare
 - factorului de discount
- [1p] Grafice comparative:
 - [0.5p] Cum afectează numărul de episoade de antrenament valorile din tabela de utilități în cazul strategiei maxfirst?
 - [0.5p] Care sunt diferențele între tabela de utilități din cazul strategiei maxfirst și random?

Observație:

- pozițiile inițiale ale soarecelui, pisicii și bucăților de brânză vor fi aleatoare fără a exista coliziuni inițiale
- se considera că există cel puțin o bucată de brânză pe hartă

Bonus [2 puncte]:

Implementarea algoritmului de învățare prin recompensă [SARSA](#) și trasarea curbei de învățare Q-learning - SARSA (de exemplu număr episod - recompensă totală pentru acel episod - sau alte reprezentări pentru a putea pune în evidență eventualele

diferențe dintre cei doi algoritmi pentru acest caz). Curba de învățare poate fi reprezentată și sub forma tabelară.

Arhiva: Se va încărca o arhivă care conține 2 fișiere:

- codul python / notebook-ul aferent rezolvării temei
- un fișier PDF în care sunt trecute: un readme al implementării, graficele cerute și analiza acestora (text explicativ). Numele arhivei trebuie să fie de forma: Tema2_<nume>_<prenume>_<grupa>.