

Path planning using reinforcement learning

Luka Cerovć, Lazar Lukić

Matematički fakultet

1. Postavka problema

Problem koji rešavamo porazumeva agenta koja treba da se kreće kroz dati prostor za skladištenje, da pokupi kutiju i da je odnese na zadato mesto. Prostor je predstavljen matricom velicine $N \times M$ na kojoj postoje polja po kojima agent može da se kreće, dok ostala polja predstavljaju zidove preko kojih agent ne može da prelazi.

2. Učenje potkrepljivanjem

Za rešavanje ovog problema koristićemo algoritam *Q-learning* koji je iz familije algoritama učenja potkrepljivanjem. Algoritam se zasniva na principu nagrade i kazne koje agent dobija da osnovu trenutnog stanja i akcije koju izvrši iz datog stanja. Cilj nam je da agent posle dovoljno iteracija nauči koja akcija je najbolja u nekom proizvoljnom stanju tako da reši zadati problem, sa ciljem da maksimizuje ukupnu nagradu (u našem slučaju da to uradi najkraćim mogućim putem).

2.1 Stanja i akcije

Jedno stanje agenta sadrži:

- 3x3 pogled u okolini agenta
- X koordinata agenta
- Y koordinata agenta
- X koordinata ciljne lokacije
- Y koordinata ciljne lokacije
- Prethodna akcija

Agent iz nekog stanja prelazi u sledeće tako što se pomeri gore, dole, levo ,desno ili dijagonalno za jedno polje u matrici.

Postoje i takozvana terminalna stanja, koja prekidaju iteraciju ukoliko se agent nadje u jednom od njih. U našem slučaju imamo dva takva stanja:

- Ako je agent udario u zid
- Ako je agent došao na cilj

2.2 Nagrade i kazne

Kada agent izvrši neku akciju on dobija odredjenu nagradu (kaznu) koja se akumulira tokom jedne iteracije.

Takođe, cilj nam je da agent uvek ide najkraćim putem i zato ćemo mu davati malu kaznu svaki put kad se pomeri za jedno polje.

Moramo biti oprezni kod politike davanja nagrada (kazni) jer neka stanja su važnija od drugih. Npr: Ne želimo da ga kažnjavamo previše za pomeraj jer neće nikad stići da istraži svoje okruženje. Tako će kazna za zid biti mnogo veća od ostalih.

3. Treniranje agenta

Agenta ćemo trenirati kroz više iteracija. Jedna iteracija traje dok ne dođe u terminalno stanje ili predje dozvoljeni maksimalni broj poteza u jednoj iteraciji. Na početku iteracije agent se postavlja u nasumično stanje tako da se njegova pozicija ne poklopi sa nekim zidom. To podrazumeva poziciju agenta i poziciju cilja.

Tokom jedne iteracije, agent u svakom potezu radi sledeće:

1. Agent bira koju akciju da izvrši u trenutnom stanju
2. Predje u sledeće stanje I dobije odgovarajuću nagradu
3. Ažurira vrednosti *Q-matrice*

3.1 *Q-matrica*

Već smo pomenuli da imamo 7 dimenzija jednog stanja. U svakom od tih stanja agent treba da može da odluči koja akcija je najbolja. Takva matrica koja ima vrednost za svaku akciju u svakom stanju je *Q-matrica*.

3.2 Izbor akcije (ϵ -greedy)

Sada znamo da će agent koristiti Q -matricu za izbor akcije iz nekog stanja. Međutim, šta ako agent još uvek ništa nije naučio? U tom slučaju možemo da biramo nasumičnu akciju. Pitanje je: kako da znamo kada je agent spreman da primeni naučeno? Zato uvodimo hiperparametar ϵ koji će se smanjivati tokom iteracija. ϵ predstavlja verovatnoću da agent izabere nasumičnu akciju u trenutnom stanju. Na ovaj način podstičemo agenta da istražuje u ranim iteracijama, a kasnije da koristi to što je naučio.

3.3 Belmanova jednačina

Za ažuriranje vrednosti Q -matrice koristimo Belmanovu jednačinu.

$$Q(S_t) = Q(S_{t,a}) + \alpha * [R_{t+1} + \gamma * \max Q(S_{t+1}) - Q(S_{t,a})]$$

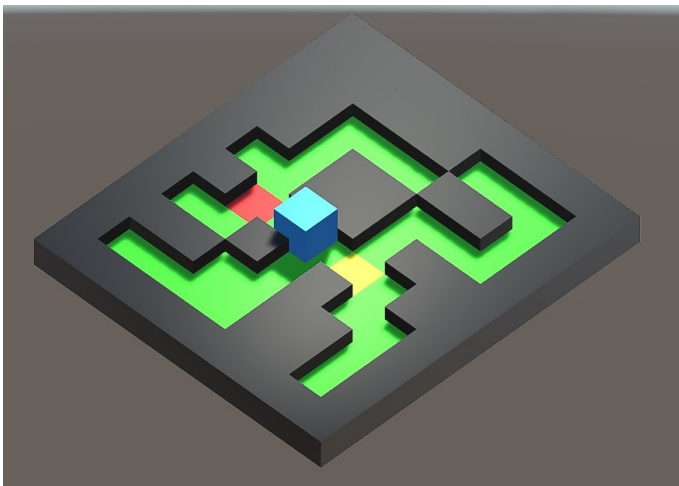
Gde su:

- $Q(S_{t,a})$ – vrednost Q -matrice za stanje t i akciju a
- $\max Q(S_{t+1})$ – vrednost Q -matrice za najbolju akciju u stanju u koje prelazimo
- R_{t+1} – nagrada za prelaz u stanje $t+1$
- α – koeficijent učenja (koliko brzo agent uči)
- γ – discount factor (koliko uzimamo od već naučenog)

4. Evaluacija algoritma

Sada ćemo primeniti algoritam na naš problem i koristeći različite vrednosti parametara. Na svim testovima prikazaćemo rezultate na mapama sa različitim brojem zidova (prepreka). Problem donošenja kutije na cilj se svodi na traženje dva puta (pokrenemo algoritam još jednom). Imamo malo ϵ i u testu da bi probali da izbegnemo što više petlji u kojima agent može da se nađe (to je najveći problem ovog modela).

Mapa okruženja (primer jednog stanja):



- zidovi su crne boje
- pozicije za kretanje su zelene boje
- agent je kocka
- kutija je crveni kvadrat
- cilj je žuti kvadrat

TEST 1 – Mapa veličine 10x10,maksimalan broj prepreka 10 u 100000 iteracija.

```
Map with 0 obstacles. eps: 0.0
Success rate: 99.18%

Map with 1 obstacles. eps: 0.011
Success rate: 95.15%

Map with 2 obstacles. eps: 0.022
Success rate: 91.5%

Map with 3 obstacles. eps: 0.033
Success rate: 87.74%

Map with 4 obstacles. eps: 0.044
Success rate: 83.86%

Map with 5 obstacles. eps: 0.05499999999999999
Success rate: 79.39%

Map with 6 obstacles. eps: 0.06599999999999999
Success rate: 75.48%

Map with 7 obstacles. eps: 0.07699999999999999
Success rate: 71.16%

Map with 8 obstacles. eps: 0.08799999999999998
Success rate: 66.52%

Map with 9 obstacles. eps: 0.09899999999999998
Success rate: 62.89%

Map with 10 obstacles. eps: 0.10999999999999997
Success rate: 59.34%
```

TEST 2 – Mapa 10x10,maksimalni broj prepreka 10 u 50000 iteracija.

```
Map with 0 obstacles. eps: 0.0
Success rate: 72.85%

Map with 1 obstacles. eps: 0.011
Success rate: 68.76%

Map with 2 obstacles. eps: 0.022
Success rate: 65.41%

Map with 3 obstacles. eps: 0.033
Success rate: 63.0%

Map with 4 obstacles. eps: 0.044
Success rate: 59.41%

Map with 5 obstacles. eps: 0.05499999999999999
Success rate: 55.38%

Map with 6 obstacles. eps: 0.06599999999999999
Success rate: 52.94%

Map with 7 obstacles. eps: 0.07699999999999999
Success rate: 50.11%

Map with 8 obstacles. eps: 0.08799999999999998
Success rate: 47.23%

Map with 9 obstacles. eps: 0.09899999999999998
Success rate: 44.02%

Map with 10 obstacles. eps: 0.10999999999999997
Success rate: 42.95%
```

TEST 3 - Mapa 10x10, maksimalni broj prepreka 10 u 1000000 iteracija.

```
Map with 0 obstacles. eps: 0.0
Success rate: 100.0%

Map with 1 obstacles. eps: 0.011
Success rate: 98.39%

Map with 2 obstacles. eps: 0.022
Success rate: 97.08%

Map with 3 obstacles. eps: 0.033
Success rate: 94.8%

Map with 4 obstacles. eps: 0.044
Success rate: 92.46%

Map with 5 obstacles. eps: 0.05499999999999999
Success rate: 90.31%

Map with 6 obstacles. eps: 0.06599999999999999
Success rate: 87.22%

Map with 7 obstacles. eps: 0.07699999999999999
Success rate: 83.63%

Map with 8 obstacles. eps: 0.08799999999999998
Success rate: 80.67%

Map with 9 obstacles. eps: 0.09899999999999998
Success rate: 76.79%

Map with 10 obstacles. eps: 0.10999999999999997
Success rate: 73.21%
```

5. Literatura:

- <http://incompleteideas.net/book/RLbook2018.pdf>