

WSI – zadanie 6

- Wstęp

Celem zadania było zaimplementowanie algorytmu Q-learning, oraz użycie go do rozwiązania problemu FrozenLake8x8. Problem ten polega na tym, że mamy agenta, który ze startu musi dotrzeć do mety, omijając dziury w zamrożonym jeziorze. Domyślna punktacja jest następująca – za dotarcie do mety otrzymuje 100 punktów, za wpadnięcie do dziury -100, w przeciwnym wypadku 0.

- Trening

Trening polega na stworzeniu tzw. QTable, czyli tabeli stanów, w jakich może się znaleźć postać w danym momencie (w naszym przypadku stanów jest 64, każdy odpowiada jednemu polu na planszy), oraz określenie wartości dla każdej możliwej akcji (w naszym przypadku są 4 akcje – pójście w jeden z 4 kierunków). Następnie, algorytm może albo wykonać najlepszy ruch w danej pozycji (ruch o najwyższym wyniku), albo wykonać ruch losowy. Szansę na losowość sterujemy parametrem epsilon. Dodatkowym utrudnieniem jest poślizg, czyli dodatkowy losowy ruch, o ile pierwotny w danej epoce nie zakończył gry. Szansą na poślizg sterujemy parametrem `slippery_rate`. Po każdym ruchu, specjalna funkcja wylicza wartość danego ruchu.

Warte odnotowania jest, że poślizg bardzo zaburza ocenę ruchu. Przeanalizuję to na przykładzie: znajdujemy się na poniższej mapie na polu o 2 wyższym od mety. Wykonujemy ruch w dół. Niestety, pech chciał, że się poślizgnęliśmy w lewo. Wpadliśmy do dziury. Algorytm w takiej sytuacji potraktuje nasz pierwotny krok w dół jako bezpośrednie wpadnięcie do dziury. Dlatego ważne jest danie chociaż odrobiny losowości w ruchach, bo inaczej taki poślizg mógłby zablokować drogę do mety.

Na koniec treningu otrzymujemy przykładowy QTable:

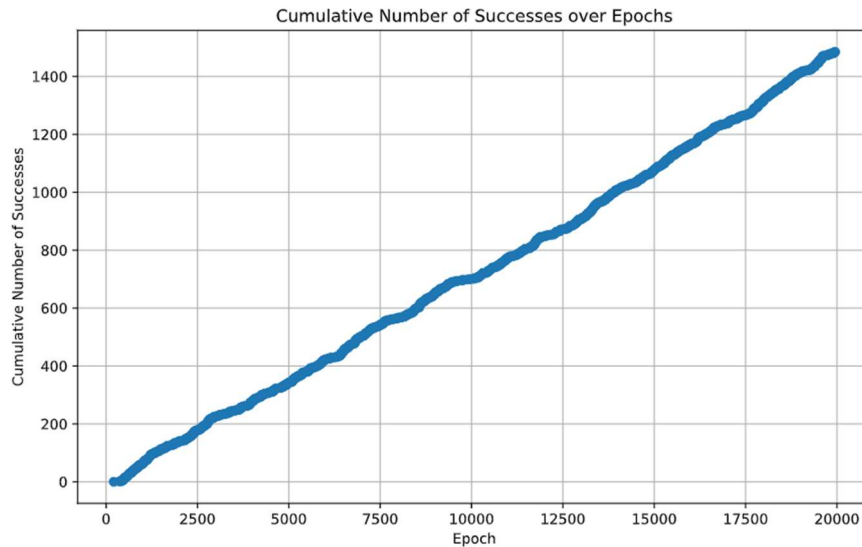
| | | | | | | | |
|------------------|-----------------|-------------------|-----------------|-------------------|-------------------|-------------------|-------------------|
| ↑(36.2) | ↑(38.3) | ↑(39.4) | ↑(39.6) | ↑(40.5) | ↑(42.9) | ↑(43.0) | ↑(43.6) |
| ↔(37.7) (39.3)→ | ↔(39.7) (47.8)→ | ↔(40.9) (53.9)→ | ↔(42.4) (56.0)→ | ↔(43.4) (56.2)→ | ↔(46.9) (57.0)→ | ↔(45.4) (45.4)→ | ↔(43.8) (43.9)→ |
| ↓(42.0) | ↓(39.6) | ↓(40.0) | ↓(16.9) | ↓(42.8) | ↓(44.5) | ↓(58.3) | ↓(59.5) |
| ↑(38.8) | ↑(40.3) | ↑(47.1) | ↑(41.9) | ↑(40.6) | ↑(54.4) | ↑(45.0) | ↑(47.3) |
| ↔(37.6) (43.8)→ | ↔(39.5) (45.8)→ | ↔(39.1) (25.5)→ | ↔(41.3) (45.1)→ | ↔(31.3) (48.5)→ | ↔(42.5) (43.7)→ | ↔(43.1) (59.0)→ | ↔(43.3) (47.7)→ |
| ↓(38.9) | ↓(39.9) | ↓(33.6) | ↓(-100.0) | ↓(35.0) | ↓(27.2) | ↓(46.2) | ↓(56.8) |
| ↑(40.1) | ↑(39.0) | ↑(40.7) | ↑(0.0) | ↑(42.2) | ↑(44.0) | ↑(47.2) | ↑(52.0) |
| ↔(28.3) (34.1)→ | ↔(27.5) (31.5)→ | ↔(31.9) (-100.0)→ | ↔(0.0) (0.0)→ | ↔(-100.0) (32.1)→ | ↔(7.8) (47.9)→ | ↔(36.6) (57.5)→ | ↔(51.7) (50.5)→ |
| ↓(24.2) | ↓(25.1) | ↓(24.9) | ↓(0.0) | ↓(31.2) | ↓(-100.0) | ↓(33.1) | ↓(55.9) |
| ↑(32.8) | ↑(41.4) | ↑(12.5) | ↑(-46.9) | ↑(45.3) | ↑(0.0) | ↑(49.5) | ↑(54.5) |
| ↔(21.6) (23.9)→ | ↔(21.9) (21.8)→ | ↔(32.1) (-0.6)→ | ↔(23.5) (-0.2)→ | ↔(-8.1) (-74.6)→ | ↔(0.0) (0.0)→ | ↔(-100.0) (58.1)→ | ↔(35.1) (51.3)→ |
| ↓(22.3) | ↓(13.2) | ↓(-3.1) | ↓(-19.0) | ↓(-23.6) | ↓(0.0) | ↓(31.0) | ↓(62.2) |
| ↑(25.7) | ↑(22.9) | ↑(20.2) | ↑(0.0) | ↑(6.0) | ↑(97.5) | ↑(39.4) | ↑(54.2) |
| ↔(3.2) (2.3)→ | ↔(9.5) (-15.7)→ | ↔(-8.9) (-27.1)→ | ↔(0.0) (0.0)→ | ↔(-46.9) (19.8)→ | ↔(12.7) (51.3)→ | ↔(26.9) (57.4)→ | ↔(39.2) (53.8)→ |
| ↓(-5.1) | ↓(-61.3) | ↓(-27.1) | ↓(0.0) | ↓(-8.7) | ↓(3.8) | ↓(-100.0) | ↓(67.6) |
| ↑(10.9) | ↑(0.0) | ↑(0.0) | ↑(-10.0) | ↑(-5.3) | ↑(12.1) | ↑(0.0) | ↑(50.3) |
| ↔(-5.4) (-19.0)→ | ↔(0.0) (0.0)→ | ↔(0.0) (0.0)→ | ↔(0.0) (0.0)→ | ↔(-10.1) (7.3)→ | ↔(-15.8) (-41.0)→ | ↔(0.0) (0.0)→ | ↔(-100.0) (49.4)→ |
| ↓(-9.7) | ↓(0.0) | ↓(0.0) | ↓(0.1) | ↓(-27.1) | ↓(-22.5) | ↓(0.0) | ↓(73.6) |
| ↑(0.7) | ↑(0.0) | ↑(-10.0) | ↑(-0.0) | ↑(0.0) | ↑(-10.0) | ↑(0.0) | ↑(68.2) |
| ↔(-1.0) (-10.0)→ | ↔(0.0) (0.0)→ | ↔(0.0) (0.0)→ | ↔(-0.1) (0.0)→ | ↔(0.0) (0.0)→ | ↔(-10.0) (-10.0)→ | ↔(0.0) (0.0)→ | ↔(-100.0) (65.7)→ |
| ↓(-0.0) | ↓(0.0) | ↓(-0.1) | ↓(-10.0) | ↓(0.0) | ↓(-0.2) | ↓(0.0) | ↓(100.0) |
| ↑(-10.0) | ↑(-10.0) | ↑(0.0) | ↑(0.0) | ↑(-10.0) | ↑(-10.0) | ↑(-10.0) | ↑(0.0) |
| ↔(-1.0) (-0.6)→ | ↔(-0.1) (0.3)→ | ↔(-0.1) (-10.0)→ | ↔(0.0) (0.0)→ | ↔(-10.0) (-0.4)→ | ↔(-10.1) (-7.9)→ | ↔(-0.0) (34.4)→ | ↔(0.0) (0.0)→ |
| ↓(-1.0) | ↓(-1.0) | ↓(-1.0) | ↓(0.0) | ↓(-1.0) | ↓(-1.8) | ↓(-0.0) | ↓(0.0) |

Kolorem zielonym oznaczony został start, czerwonym meta, a dziury są jasnoniebieskie.

- Testy

Testowanie polegało na wykonanie treningu z odpowiednimi parametrami (parametry dla poszczególnych treningów można znaleźć w plikach json). Następnie, wykonywałem 1000 prób, w których agent brał zawsze najlepszy ruch w pozycji, jednak dalej z występującym poślizgiem. Ustawiłem limit kroków w pojedynczym teście na 100, żeby uniknąć pętli.

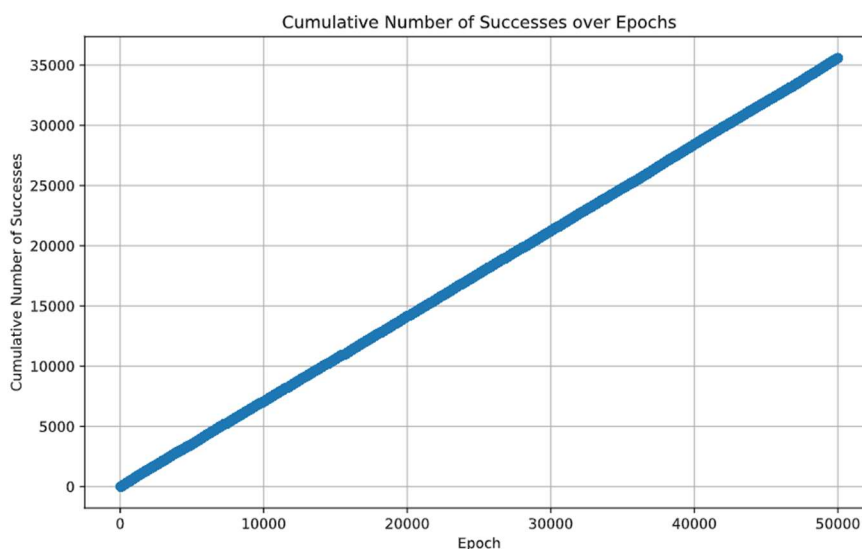
Moje rozwiązanie posiada 2 systemy przyznawania punktów. Pierwszy jest domyślny, czyli punkty za dotarcie do mety, ujemny punkty za wpadnięcie do dziury. Ten algorytm sprawdziłem na 2 mapach, jednej bardziej otwartej (map.txt), a drugiej bardziej zamkniętej (map3.txt). Są to treningi zawierające w nazwie pliku no_bonus. Dla każdej mapy dałem 2 treningi – z dużą losowością i z małą. Okazuje się, że tylko w jednym przypadku algorytm odnalazł drogę do mety – było to na otwartej mapie i z dużą losowością (pliki zaczynające się od 1_slip_no_bonus_high_random). Podczas treningu tylko ok. 7% epok kończyło się sukcesem. Podczas testów wypadł dużo lepiej, ponieważ aż 86,9% testów kończyło się sukcesem, co przy szansie na poślizg równej 30% jest naprawdę niezłym wynikiem. Natomiast wykres ilości sukcesów podczas treningu jest bardzo nieregularny, co wraz z niewielką ich ilością wskazuje na czyste szczęście algorytmu:



- Ulepszony algorytm

Ulepszając algorytm, dodałem mu nową nagrodę i nową karę. Algorytm otrzymuje dodatkową nagrodę za każdy ruch, który go przybliży do mety (i karę za oddalanie się), oraz otrzymuje karę za każdy ruch który by go zderzył z krawędzią mapy (w mojej implementacji ruch na ścianę liczy się jako wykonany ruch, ale sam agent nie zmienia pozycji). Ten algorytm również przetestowałem na tych samych dwóch mapach, znowu z dużą i małą losowością. Od razu mogłem odrzucić testy z dużą losowością, ponieważ nie dawały one się wykazać algorytmowi. Ciekawie zaczęło się dzieć przy niskiej losowości:

- Na otwartej mapie algorytm podczas treningu uzyskał 71% sukcesów, oraz podczas testów, 100%:



Nie dość, że dużo więcej sukcesów, to zdarzały się one praktycznie co ten sam interwał (dlatego wykres jest zbliżony do prostej).

- Zamknięta mapa okazała się jeszcze ciekawsza – ze wszystkich testów na 3 mapie, jako jednemu udało się znaleźć prawidłową drogę, osiągając 9% sukcesów na

treningu, i 53% na testach. I udało mu się tego dokonać mimo tego, że często musiał wykonywać ruchy sprzeczne z założeniami (musiał się oddalać od mety).

