

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 2857

Automatizirani razvoj agenata za različita okruženja

Paulo Sanković

Zagreb, lipanj 2022.

Zagreb, 11. ožujka 2022.

DIPLOMSKI ZADATAK br. 2857

Pristupnik: **Paulo Sanković (0036509900)**

Studij: Računarstvo

Profil: Računarska znanost

Mentor: doc. dr. sc. Marko Đurasević

Zadatak: **Automatizirani razvoj agenata za različita okruženja**

Opis zadatka:

Proučiti različita okruženja poput jednostavnih igara ili fizikalnih simulatora za koje je moguće razviti agente. Proučiti različite metode i modele poput umjetnih neuronskih mreža koje se mogu iskoristiti za upravljanje agentom u takvim okruženjima. Proučiti postojeća programska rješenja koja pružaju simulaciju pojedinih okruženja te istražiti mogućnosti povezivanja s takvim okruženjima. Razviti programski okvir koji omogućuje automatski razvoj agenata za odabrana okruženja te grafički prikazuje ponašanje razvijenih agenata. Ocijeniti efikasnost razvijenih agenata za odabrana okruženja. Predložiti moguća poboljšanja predloženog postupka s ciljem razvoja efikasnijih agenata. Radu priložiti izvorne programske kodove, dobivene rezultate uz potrebna objašnjenja i korištenu literaturu.

Rok za predaju rada: 27. lipnja 2022.

Zahvaljujem se mentoru doc. dr. sc. Marku Đuraseviću na pomoći, dobrim savjetima i razumijevanju tijekom izrade ovog diplomskog rada, kao i svim profesorima i kolegama s kojima sam proveo svoje osnovnoškolske, srednjoškolske i studentske dane. Posebno bih se zahvalio svojoj obitelji na pruženoj podršci i razumijevanju.

SADRŽAJ

1. Uvod	1
2. Podržano učenje	2
2.1. Ključni koncepti	3
2.2. Duboki modeli	4
2.2.1. Unaprijedni potpuno povezani modeli	5
2.2.2. Konvolucijski modeli	6
2.3. Algoritmi podržanog učenja	8
2.3.1. Deep Q Learning	8
2.3.2. Double Deep Q Learning	8
2.3.3. Actor Critic	9
3. OpenAI Gym	10
3.1. Struktura	10
3.1.1. Okolina	11
3.1.2. Interakcija s okolinom	11
3.1.3. Prostor akcija i prostor stanja	12
3.1.4. Omotači	13
3.1.5. Vektorizirana okruženja	13
3.2. Determinizam atari okruženja	14
3.3. Okruženja	15
3.3.1. Okruženje CartPole	15
3.3.2. Okruženje Breakout	17
4. Implementacija	19
4.1. Okolina CartPole	19
4.1.1. Implementacija Deep Q Learning algoritma	19
4.1.2. Implementacija Double Deep Q Learning algoritma	19

4.1.3.	Implementacija Actor Critic algoritma	19
4.1.4.	Usporedba	19
4.2.	Okolina Breakout	19
4.2.1.	Implementacija Deep Q Learning algoritma	19
4.2.2.	Implementacija Double Deep Q Learning algoritma	19
4.2.3.	Implementacija Actor Critic algoritma	20
4.2.4.	Usporedba	20
4.3.	Usporedba algoritama	20
5.	Moguća poboljšanja	21
6.	Zaključak	22
	Literatura	23
	Popis slika	26
	Popis tablica	27
	Popis programskih isječaka	28

1. Uvod

Inteligentni sustavi sposobni su analizirati, razumjeti i učiti iz dostupnih podataka putem posebno dizajniranih algoritama umjetne inteligencije. Zahvaljujući dostupnosti velikih skupova podataka, razvoja tehnologije i napretka u algoritmima došli smo do razine gdje nam razvijeni modeli mogu uvelike poslužiti u svakodnevnom životu.

Posebno je zanimljiva problematika u kojoj je bez znanja o pravilima i funkcioniranju specifične okoline, potrebno konstruirati specijaliziranog agenta koji se nalazi u određenom stanju okoline i ponavlja korake izvršavanja optimalne akcije i prijelaza u novo stanje okoline. Za svaku akciju agent prima određenu nagradu - mjeru koja označava koliko su akcije agenta ispravne za tu okolinu i koliko je napredak agenta kvalitetan. Agent izvršava akcije i prelazi u nova stanja sve dok se ne nađe u terminalnom (završnom) stanju. Dakle, cilj agenta u okolini jest pronaći optimalnu strategiju koja će maksimizirati očekivanu dobit (nagradu) u određenom vremenskom okviru.

Područje strojnog učenja koje se bavi prethodno navedenom problematikom naziva se podržano učenje (engl. *reinforcement learning*). Agenti koji se prilično dobro ponašaju u takvim okolinama možemo implementirati pomoću različitih algoritama podržanog učenja koji se temelje na umjetnim neuronskim mrežama (engl. *artificial neural networks*). Umjetne neuronske mreže su dobri aproksimatori funkcija i najbolje se ponašaju u okolini koja ima kompozitnu strukturu gdje vrlo kvalitetno duboki model predstave kao slijed naučenih nelinearnih transformacija.

U sklopu ovog rada bilo je potrebno proučiti i razumjeti metode i algoritme podržanog učenja i funkcioniranje umjetnih neuronskih mreža. Nadalje, bilo je potrebno istražiti, proučiti i naposljetku implementirati neke od algoritama podržanog učenja koji se zasnivaju na umjetnim neuronskim mrežama i koje je trebalo uklopiti u okoline koje su prikladne za simulaciju i testiranje ponašanja naučenih agenta.

Programsko rješenje implementirano je u programskom jeziku *Python*, primarno koristeći *PyTorch* biblioteku (engl. *library*), te ostale korisne biblioteke poput *numpy*, *tqdm*, *stable-baselines3*, *tensorboard*... Za simulaciju i testiranje ponašanja agenata u posebnom okruženju korištena je biblioteka *OpenAI Gym*.

2. Podržano učenje

Strojno učenje (engl. *Machine learning*) je grana umjetne inteligencije (engl. *artificial intelligence*) koje se može definirati kao skup metoda koje u podacima mogu automatski otkrivati obrasce, i potom te otkrivene obrasce iskoristavati pri budućem predviđanju podataka, ili obavljati druge zadatke odlučivanja u prisustvu nesigurnosti [11]. Drugim riječima, bez eksplicitnog programiranja moguće je napraviti sustave koji funkcioniraju kao ljudski mozak - imaju pristup podacima, koriste ih za učenje i samim time bolje razumiju entitete, domene i veze između podataka.

Strojno učenje dijeli se na 3 podvrste: nadzirano učenje, nenadzirano učenje i podržano (ojačano) učenje. Nadzirano učenje (engl. *supervised learning*) karakterizira učenje modela nad testnim podacima koji su označeni. Model točno zna da za određeni ulaz mora vratiti izlaz koji je istovjetan unaprijed pridruženoj oznaci. Algoritam mjeri točnost kroz funkciju gubitka, prilagođavajući se sve dok se izračunata razlika izlaza modela i stvarnog izlaza (pogreška) ne smanji do određene mjere. S druge strane, u nenadziranom učenju (engl. *unsupervised learning*) posjedujemo podatke bez zadatog izlaza - podatci su dani bez ciljane vrijednosti i u tim situacijama treba pronaći određenu pravilnost. Postupci poput grupiranja, smanjenja dimenzionalnosti, otkrivanja veza između primjeraka... pripadaju nenadziranom učenju.

Posebna i nama najzanimljivija podvrsta strojnog učenja jest podržano učenje (engl. *reinforcement learning*). Podržano učenje bavi se optimizacijom ponašanja agenta koji je u interakciji s okolinom (u kojoj se nalazi) i koji izvršava akcije na temelju informacija koje dobiva iz okoline. Agent pri svakom koraku od okoline dobiva povratnu informaciju u obliku nagrade ili kazne. Za razliku od prethodne dvije navedene podvrste koje mapiraju ulazne podatke na određeni format izlaza, u podržanom učenju je najizraženije učenje iz iskustva koje je čovjeku i drugim živim bićima veoma blisko.

2.1. Ključni koncepti

Za potpuno razumijevanje podržanog učenja, bitno je u navesti i pojasniti ključne koncepte i terminologiju. Okolina (engl. *environment*) označava svijet u kojem se agent nalazi i s kojim interaktira. Stanje s_t (engl. *state*) reprezentira kompletni opis okoline u određenom trenutku t . S druge strane, opservacija okoline o_t (engl. *observation*) predstavlja prilagođeni i ograničeni opis okoline koji agent dobije u nekom trenutku. Kada se agent nađe u nekom stanju s_t i kada od okoline dobije opservaciju o_t , tada agent poduzima akciju a_t (engl. *action*) i samim time u idućem vremenskom koraku inicira promjenu stanja s_{t+1} i opservacije okoline o_{t+1} .

Način na koji agent odabire akciju iz skupa svih dostupnih akcija naziva se politika (engl. *policy*). Politika ovisi o parametrima modela θ i može biti deterministička μ_θ ili stohastička π_θ . Veza između akcije, politike i stanja prikazana je izrazima 2.1. U našem slučaju koristiti ćemo politike koji su zapravo duboki modeli - aproksimacije funkcije odluke čiji se parametri uče optimizacijskim algoritmima.

$$\begin{aligned} a_t &= \mu_\theta(s_t) \\ a_t &\sim \pi_\theta(\cdot \mid s_t) \end{aligned} \tag{2.1}$$

Nadalje, putanja τ (engl. *trajectory*) je pojam koji označava niz stanja i pripadajućih akcija i predstavljen je izrazom 2.2. Interakcija s okolinom započinje u trenutku $t = 0$ kada pomoću funkcije za inicijalizaciju okoline ρ_0 poštujući pravila okoline, nasumično generiramo stanje s_0 .

$$\tau = (s_0, a_0, s_1, a_1, \dots) \tag{2.2}$$

Najvažnija i najkorisnija informacija koju agent dobiva od okoline jest nagrada r_t (engl. *reward*) koju generira funkcija R (engl. *reward function*) i koja u obzir uzima trenutno i iduće stanje te akciju koja je izazvala promjenu stanja. Povezanost između generiranja iznosa nagrade i same nagrade prikazana je izrazom 2.3.

$$r_t = R(s_t, a_t, s_{t+1}) \tag{2.3}$$

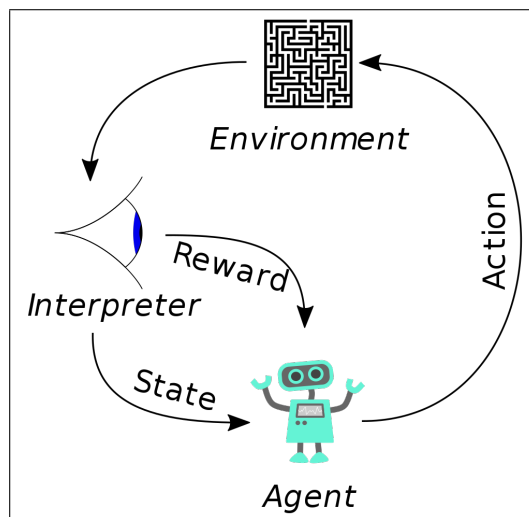
Želimo dobiti što bolji pregled koliko su bile dobre akcije koje je agent poduzeo. Tu informaciju možemo predstaviti na dva različita načina. Sumom nekorisniranih nagrada zbrajamo samo nagrade koje su dobivene u fiksnom vremenskom intervalu T (izraz 2.4). Na taj način dobivamo informaciju o tome kolika je bila prosječna nagrada u zadnjih T koraka. Ako pak želimo pokazati da su nam nagrade u trenutnom stanju

vrjednije nego nagrade koje ćemo dobiti u budućim stanjima, moramo uvesti korekcijski faktor $\gamma \in (0, 1)$ (engl. *discount factor*). Navedeni pristup sumiranja korigiranih nagrada prikazan je izrazom 2.5.

$$R(\tau) = r_0 + r_1 + \dots + r_T = \sum_{t=0}^T r_t \quad (2.4)$$

$$R(\tau) = r_0 + \gamma \cdot r_1 + \gamma^2 \cdot r_2 + \gamma^3 \cdot r_3 + \dots = \sum_{t=0}^{\infty} \gamma^t r_t \quad (2.5)$$

Dakle, cilj podržanog učenja jest naći optimalnu strategiju (niz optimalnih akcija) koje maksimiziraju ukupnu (kumulativnu) nagradu. U svakom koraku interakcije agenta s okolinom, agent prima opis stanja okoline u kojoj se nalazi. S obzirom na to stanje, izvršava akciju koja vrši neku promjenu nad okolinom i prebacuje ju u novo stanje. Agent prima povratnu informaciju od okoline koja reprezentira koliko je odabrana akcija dobra. Opisana interakcija agenta i okoline prikazana je na slici 2.1.



Slika 2.1: Ciklus interakcije agenta s okolinom [9]

2.2. Duboki modeli

Duboko učenje (engl. *Deep learning*) jest tip strojnog učenja (točnije, podskup strojnog učenja) koje nastoji oponašati način zaključivanja i obrasce koje ljudski mozak koristi za učenje i donošenje odluka. Veliku ulogu u cijeloj ideji dubokog učenja imaju duboke neuronske mreže (engl. *deep neural networks*, *DNN*) pomoću kojih se povezivanjem više slojeva procesnih elemenata (čvorova, neurona), dobivaju duboki modeli koji su sposobni učiti i baratati s podacima kompozitne strukture. Primjenom dubokih

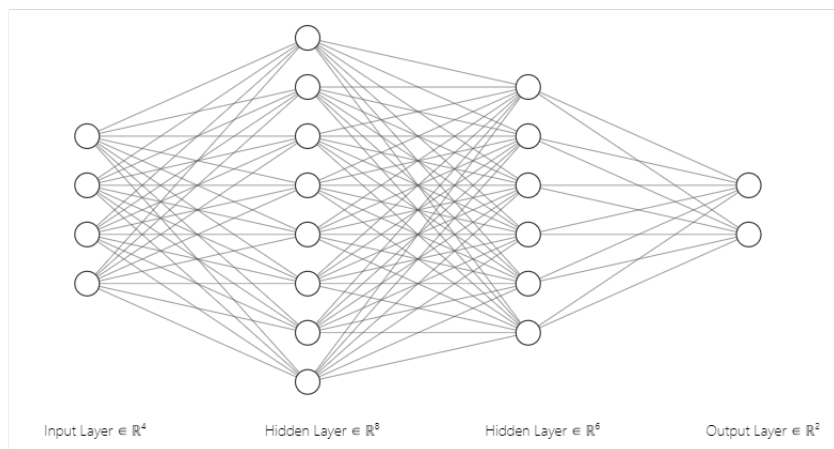
modela dolazimo do slijeda naučenih nelinearnih transformacija kojima aproksimiramo funkciju odluke, učimo mapiranje između ulaznih podataka i izlazih podataka, te nastojimo postići dobru generalizaciju nad stvarnim podacima [5].

2.2.1. Unaprijedni potpuno povezani modeli

Unaprijedni potpuno povezani modeli (engl. *Fully connected neural network*) (poznatiji i pod nazivom višeslojni perceptron (engl. *Multi-layer perceptron*)) sastoje se od lanaca potpuno povezanih slojeva. Svaki neuron iz prethodnog sloja povezan je s neuronom idućeg sloja.

Sastoje se od tri vrste slojeva - ulaznog sloja, izlaznog sloja i skrivenih slojeva. Ulaznom sloju dovode se podatci koje je potrebno obraditi. Izlaz neuronske mreže (u najosnovnijem obliku) predstavljen je logitima (engl. *logits*) - vektorima logaritama nenormaliziranih vrijednosti. Specifičnije, za slučaj da želimo provesti klasifikaciju podataka ili drugačije organizirati izlazne vrijednosti na izlaz dodajemo posebni sloj (npr. *Softmax* funkcija za klasifikaciju). Samo su ulaz i izlaz specificirani dimenzijama. Model ima slobodu da iskoristi skrivene slojeve na način koji osigurava najbolju aproksimaciju funkcije. Neuronskim mrežama želimo izgraditi modele koji nisu linearno odvojivi i zato koristimo nelinearnu aktivacijsku funkciju - najčešće ReLU (Rectified Linear Unit). Svaki od slojeva modelira jednu nelinearnu transformaciju.

Slika 2.4 [6] prikazuje arhitekturu potpuno povezanog modela koji je sastavljen od sveukupno 4 potpuno povezana sloja - ulaznog (dimenzije 4), izlaznog (dimenzije 2) i dva skrivena sloja (svaki dimenzije 8). Kodom 1 prikazana je implementacija navedenog modela u biblioteci *PyTorch*.



Slika 2.2: Arhitektura potpuno povezanog modela

```

1 import torch.nn as nn
2
3 model = nn.Sequential(
4     nn.Linear(in_features=4, out_features=8, bias=True),
5     nn.ReLU(),
6     nn.Linear(in_features=8, out_features=6, bias=True),
7     nn.ReLU(),
8     nn.Linear(in_features=6, out_features=2, bias=True),
9 )

```

Kôd 1: Implementacija potpuno povezanog modela na slici 2.4 koristeći biblioteku *PyTorch*

2.2.2. Konvolucijski modeli

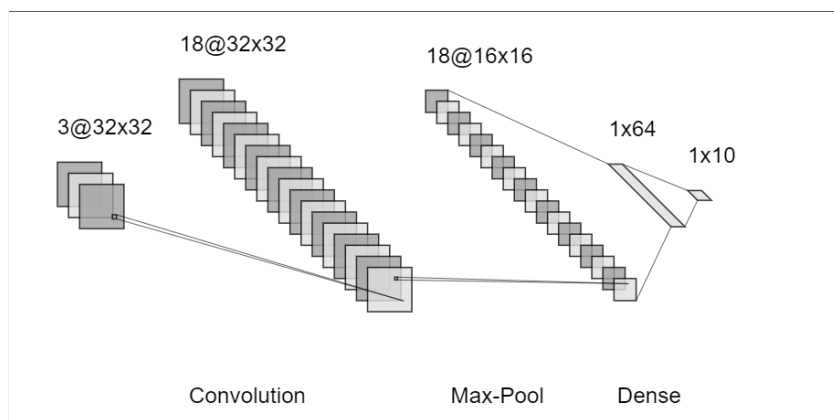
Konvolucijski modeli (engl. *Convolutional neural networks*) su modeli koji uz potpuno povezane slojeve imaju najmanje jedan konvolucijski sloj (engl. *convolution layer*). Osim spomenutih slojeva, konvolucijski modeli sadrže i slojeve sažimanja (engl. *pooling layers*), slojeve u kojima provodimo nelinearnost podataka te sloj koji višedimenzionalni ulaz pretvara u jednodimenzionalni i pritom priprema podatke za obradu u potpuno povezanim slojevima (engl. *flatten layer*).

Operacija konvolucije provodi se nad ulazom i jezgrom (engl. *kernel*) (slobodnim parametrima koje učimo) gdje kao rezultat dobivamo mapu značajki koja pokazuje gdje se koja značajka nalazi u ulaznim podacima (npr. slici). Dimenzije mapa značajki i njihovu robusnost korigiramo korištenjem atributa koraka konvolucije (engl. *stride*) i nadopunjavanja ulaznih podataka (engl. *padding*). Slično konvolucijskom sloju, sloj sažimanja odgovoran je za smanjenje prostora značajki (smanjenje dimenzionalnosti podataka) i dodatno za izdvajanje dominantnih značajki. Razlikujemo dvije vrste sažimanja: sažimanje maksimalnom vrijednosti (engl. *max pooling*) i sažimanje srednjom vrijednosti (engl. *average pooling*). Prilikom sažimanja značajki maksimalnom vrijednošću u obzir uzimamo samo značajku najveće vrijednosti te na taj način uklanjamo šum ulaza i potencijalno biramo značajku najveće važnosti.

Korištenje konvolucijskih modela biti će nam iznimno potrebno u situacijama kada su ulazni podatci u formi slike, odnosno kada su nam važne lokalne interakcije između ulaznih podataka (piksela) te njihova vremenska i prostorna ovisnost.

Slika 2.3 [6] prikazuje jednostavni konvolucijski model koji se sastoji od konvolucijskog sloja, sloja sažimanja maksimalnom vrijednosti, sloja koji 3-dimenzionalne

podatke pretvara u 1-dimenzionalne, te dva potpuno povezana sloja. Ulaz u konvolucijski sloj predstavlja RGB slika (3 kanala) dimenzije 32×32 . Primjenom konvolucije (veličina jezgre 3, korak 3, nadopuna 1) izvlačimo 18 kanala značajki dimenzije 32×32 (dimenzija se nije promijenila iz razloga što nadopunjavamo ulaz). Primjenom sažimanja maksimumom (veličina jezgre 2, korak 2) smanjujemo broj značajki na dimenziju 16×16 . Prvi potpuno povezani sloj na svoj ulaz dobije vektor dimenzije 4608 kojeg pretvara u vektor izlaza dimenzije 64. Posljednji potpuno povezani sloj koji je ujedno i posljednji sloj u ovom konvolucijskom modelu za izlaz predaje vektor dimenzije 10. Isječak koda koji prikazuje implementaciju jednostavnog konvolucijskog modela koristeći biblioteku *PyTorch* prikazan je kodom 2.



Slika 2.3: Arhitektura konvolucijskog modela

```

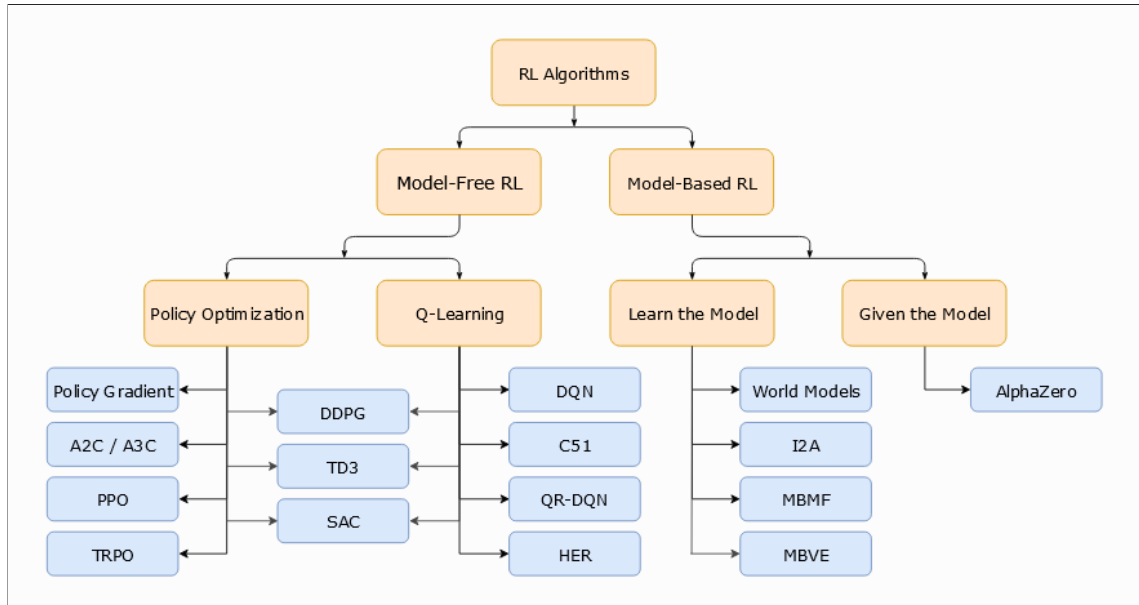
1 import torch.nn as nn
2
3 model = nn.Sequential(
4     nn.Conv2d(in_channels=3, out_channels=18, kernel_size=(3, 3),
5               ↪ stride=(1, 1), padding=1),
6     nn.MaxPool2d(kernel_size=2, stride=2),
7     nn.Flatten(),
8     nn.Linear(in_features=18 * 16 * 16, out_features=64),
9     nn.Linear(in_features=64, out_features=10)
10 )

```

Kôd 2: Implementacija konvolucijskog modela na slici 2.3 koristeći biblioteku *PyTorch*

2.3. Algoritmi podržanog učenja

TODO dodati dio oko value functions, advantage functions... kako bi bilo kasnije lakše objasniti algoritme, funkcije gubitka, što koja mreža predstavlja...



Slika 2.4: Podjela algoritama podržanog učenja [2]

Svaki algoritam strojnog učenja definiran modelom, gubitkom i metodom optimizacije. Model je postupak obrade (odnosno skup funkcija) sa slobodnim parametrima koji za određen ulaz daje pripadajući izlaz. Gubitak je mjera koja na formaliziran način vrednuje slobodne parametre modela, odnosno pokazuje u kojoj mjeri se mi ne slažemo s onim što je model predstavio kao izlaz. Metoda optimizacije (optimizacijski postupak) jest način na koji pronalazimo optimalne parametre koji su važni kako bi minimizirali prethodno navedenu komponentu - gubitak. Navedene tri glavne komponente biti će važno napomenuti pri svakom predstavljanju algoritma jer su to glavne odrednice pri analizi algoritama strojnog učenja.

2.3.1. Deep Q Learning

TODO <https://github.com/udacity/deep-learning/blob/master/reinforcement/nature14236.pdf>

2.3.2. Double Deep Q Learning

TODO <https://arxiv.org/pdf/1509.06461.pdf>

2.3.3. Actor Critic

TODO https://www.tensorflow.org/tutorials/reinforcement_learning/actor_critic

3. OpenAI Gym

OpenAI Gym je Python biblioteka (engl. *library*) otvorenog koda (engl. *open source*) koja služi za razvijanje i usporedbu agenata u odabranim okolinama. Iznimno je popularna u sferi simpatizera i programera koji se bave razvijanjem modela podržanog učenja zbog jednostavnosti korištenja, velikog broja dostupnih okolina i jednostavnog stvaranja novih okolina, te jednostavne interakcije agenta i okoline. OpenAI Gym biblioteka se redovito održava i trenutno je na verziji 0.24.1.

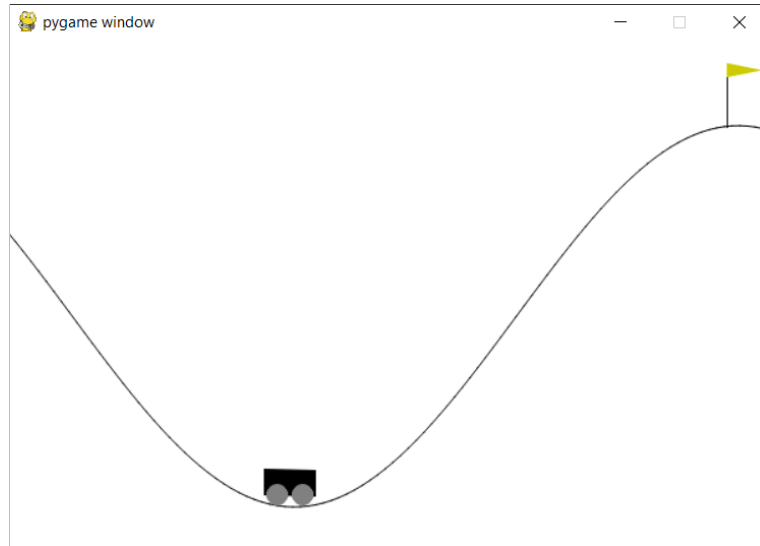
3.1. Struktura

Interakcija agenta i okoline podijeljena je na epizode. Na početku svake epizode, početno stanje se nasumično uzorkuje iz distribucije, i interakcija se nastavlja sve dok se okolina ne nađe u terminalnom stanju [4].

```
1 import gym
2
3 env = gym.make("MountainCar-v0")
4
5 observation = env.reset()
6
7 done = False
8 while not done:
9     action = env.action_space.sample()
10    observation, reward, done, info = env.step(action)
11
12    env.render()
13
14 env.close()
```

Kôd 3: Jednostavan primjer integracije agenta i Gym okoline (1 epizoda)

Kod 3 prikazuje potpunu implementaciju jednostavne interakcije agenta i okoline. Agent u ovom jednostavnom slučaju nasumično odabere akciju iz skupa svih dostupnih akcija za tu okolinu (linija 9). Osnovni kostur sastoji se od koraka specifikacije okoline (linija 3), inicijalizacije okoline (linija 5) te interakcija okoline i agenta - agent predaje okolini odabranu akciju, okolina vraća povratnu informaciju (linije 7 - 14).



Slika 3.1: Rezultat pokretanja koda 3

3.1.1. Okolina

Temelj oko kojeg se zasniva OpenAI Gym biblioteka jest razred (engl. *class*) `Env` koji u suštini implementira simulator koji pokreće okruženje u kojem naš agent može interagirati s okolinom. Točnije rečeno, enkapsulira sva potrebna ponašanja i metode koje su potrebne za jednostavnu interakciju. Objekt tipa `Env` stvara se pozivanjem funkcije `gym.make()` (kod 3 linija 3) kojoj se predaje identifikator okoline (`id`) zajedno s opcionalnim argumentima (metapodacima).

3.1.2. Interakcija s okolinom

Kao što je vidljivo iz koda 3 osnovne metode koje se pozivaju nad instancom razreda `Env` su `reset` i `step`. Funkcija `reset` postavlja okruženje u početno stanje i vraća njegovu vrijednost (kod 3. linija 5). S druge strane, funkciji `step` (kod 3. linija 10) predaje se jedna od ispravnih akcija koja inicira prijelaz okoline iz jednog stanja u drugo. Funkcija vraća 4 vrijednosti: vrijednost prostora stanja (engl. *observation*), iznos nagrade kao rezultat poduzimanja određene akcije (engl. *reward*), zastavicu koja

signalizira jesmo li došli u završno stanje okoline (engl. *done*), te neke dodatne informacije.

Još jedna vrlo često korištena funkcija jest `render` (kod 3. linija 12) koja služi kako bi se u određenom formatu prikazala okolina. Dostupni formati su: `human` (otvara se skočni prozor sa slikom stanja okoline - slika 3.1), `rgb_array` (numpy array RGB vrijednosti) i `ansi` (string reprezentacija okoline).

3.1.3. Prostor akcija i prostor stanja

Osnovna struktura okruženja opisana je atributima `observation_space` i `action_space` koji su dio razreda `Env` i čija se vrijednost može razlikovati zavisno o okolini. Atribut `action_space` opisuje numeričku strukturu svih legitimnih akcija koje se mogu izvesti nad određenom okolinom. S druge strane, atribut `observation_space` definira strukturu objekta koje predstavlja opis stanja u kojem se okolina nalazi.

Format validnih akcija i stanja okoline, odnosno struktura tih podataka, definirana je razredima `Box`, `Discrete`, `MultiBinary` i `MultiDiscrete`. Svi navedeni razredi nasljeđuju i implementiraju glavne metode nadrazreda `Space`.

Razred `Box` predstavlja strukturu podataka u kontinuiranom n -dimenzionalnom prostoru. Prostor i njegove validne vrijednosti omeđene su gornjim i donjim granicama koje se jednostavno postavljaju pri inicijalizaciji strukture pridruživanjem željenih vrijednosti atributima `high` i `low`. Kod 4 prikazuje inicijalizaciju `Box` strukture podataka koja je sastavljena od 3-dimenzionalnog vektora čije su vrijednosti omeđene odozdo i odozgo vrijednostima -1 i -2 . Metoda `sample(self)` nasumično uzorkuje element iz prostora koristeći različite distribucije ovisno o ograničenjima prostora.

```
1 >>> import numpy as np
2 >>> from gym.spaces import Box
3 >>> b = Box(low=-1.0, high=2.0, shape=(3,), dtype=np.float32)
4 >>> b.sample()
5 array([-0.3791686 , -0.35007873,  0.8138365 ], dtype=float32)
```

Kôd 4: Primjer korištenja strukture kontinuiranog prostora `Box`

Razred `Discrete` s druge strane, predstavlja strukturu podataka u diskretnom n -dimenzionalnom prostoru gdje su validne vrijednosti sve cjelobrojne vrijednosti unutar intervala $[0, n - 1]$ (početna vrijednost se može specificirati). Kod 5 prikazuje inicijalizaciju `Discrete` strukture podataka ovisno o specificiranoj početnoj vrijednosti.

```

1 >>> from gym.spaces import Discrete
2 >>> d = Discrete(3)                # {0, 1, 2}
3 >>> d = Discrete(3, start=-1)     # {-1, 0, 1}
4 >>> d.sample()
5 1

```

Kôd 5: Primjer korištenja strukture diskretnog prostora `Discrete`

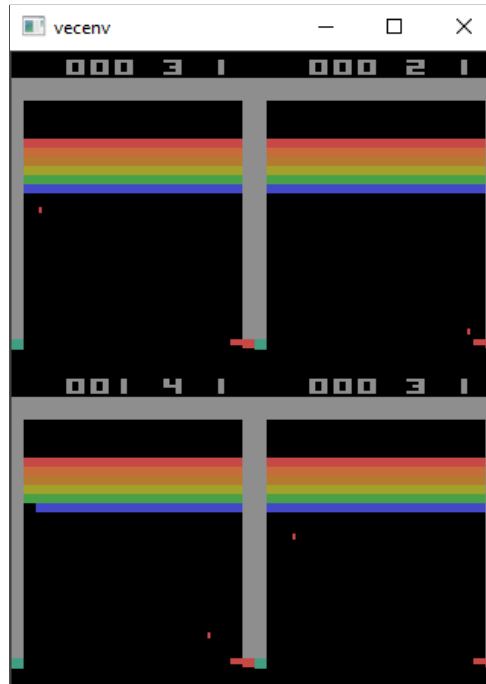
3.1.4. Omotači

Omotači (engl. *Wrappers*) su prikladne strukture koje omogućavaju izmjenu elemenata postojećeg okruženja bez potrebe za mijenjanjem originalnog koda. Omotači omogućavaju modularnost, mogu se implementirati prema vlastitim potrebama i ulančavati. Ova funkcionalnost vrlo se često koristi u situacijama kada pri treniranju modela želimo normalizirati ulaze (skalirati vrijednosti slikovnih jedinica), provesti regularizaciju (podrezivanje vrijednosti nagrade), transformirati ulaze u *PyTorch* dimenzije, implementirati metodu preskakanja slikovnih okvira... Navedene funkcionalnosti moguće je postići tako da definiramo vlastiti omotač koji će nasljeđivati ili obični `Wrapper` nadrazred ili specifičnije razrede poput `ObservationWrapper`, `RewardWrapper`, `ActionWrapper`...

3.1.5. Vektorizirana okruženja

Koristeći standardne metode stvaranja i interakcije s Gym okruženjem, pokrećemo samo jednu instancu okruženja i na taj način ne iskorištavamo računalnu snagu koja nam je dostupna u potpunosti. Vektorizirana okruženja (engl. *Vectorized environments*) su okruženja koja paralelno pokreću više kopija istog okruženja u svrhu poboljšanja učinkovitosti i ubrzanja procesa učenja agenta i njegove interakcije s kopijama okolina.

Python biblioteka *stable-baselines* [1] pruža gotove omotače koji omogućavaju pokretanje n paralelnih neovisnih instanci okolina. Paralelna obrada interno je implementirana korištenjem *multiprocessing* paketa. Zbog paralelne interakcije agenta s više okolina potrebno je prilagoditi strukturu objekata. Agent predaje omotanoj okolini niz od n akcija i dobiva listu povratnih informacija (stanje pojedinačne okoline, nagrada i zastavica terminalnog stanja) u formi liste n vrijednosti.



Slika 3.2: Vektorizirano okruženje sa 4 paralelne asinkrone instance

3.2. Determinizam atari okruženja

Upravljanje agentima i njihovo izvođenje akcija u određenoj atari okolini interno je implementirano koristeći objektno orijentiranu razvojnu cjelinu (engl. *framework*) *The Arcade Learning Environment* (skraćeno *ALE*). Na taj način razdvajaju se slojevi emulacije okruženja (koristeći Atari 2600 emulator *Stella*), sloj koji kontrolira samog agenta (*ALE*) i sloj koji pruža jednostavnu programsku implementaciju i interakciju (*OpenAI Gym*) [3].

Originalna Atari 2600 konzola nije imala izvor entropije za generiranje pseudoslučajnih brojeva i iz tog razloga okruženje je bilo u potpunosti determinističko - svaka igra počinje u istom stanju, a ishodi su u potpunosti određeni stanjem i akcijom [7]. Iz tog razloga vrlo je jednostavno postići visoku uspješnost agenta u okolini pamćenjem dobrog niza akcija. Nama to ne odgovara jer želimo postići da agent nauči donositi dobre odluke. pristupi kojima se nastoji uvesti određen stupanj stohastike su *sticky actions*, *frame skipping*, *initial no-ops*, te *random action noise*.

Frame skipping pristup pri svakom koraku okoline uzima zadnju akciju koju je agent odabrao i ponavlja ju n broj puta (kroz n slikovnih okvira). Ovisno o verzijama okoline, broj ponavljanja može se specificirati tipom `int` - cijelobrojnom fiksnom vrijednošću (determinizam), ili tipom `tuple(2)` - bira se slučajno odabrana vrijednost

unutar specificiranog intervala (stohastika). Osim potencijalnog uvođenja stohastike, ovaj pristup pojednostavljuje problem podržanog učenja i ubrzava izvođenje [7]. Tehnika *initial no-ops* označava da na početku epizode okolina ignorira akcije agenta slučajan broj puta iz intervala $[0, n]$ - izvršava akciju NOOP. Osim što na početku agent ne zna kada će okolina početi koristiti njegove akcije, ostatak interakcije s okolinom je i dalje deterministički [7]. Sličan pristup bio bi i da potičemo raznolikost početnih stanja objekta s kojim upravljamo. To bi postigli nasumičnim poduzimanjem niza akcija kojima pomičemo objekt prije samog aktiviranja posebne akcije FIRE [8].

Tehnika *random action noise* zamjenjuje odabranu akciju agenta nasumičnom akcijom samo ako je vjerojatnost zamjene manja od specificirane vrijednosti. Ovo uspješno uvođenje stohastike dolazi s negativnim posljedicama poduzimanja slučajne akcije koja može značajno poremetiti politiku agenta i smanjiti učinak [7]. Odabir slučajne akcije bolje je zamijeniti pristupom ponavljanja prethodno odabrane akcije uzimajući u obzir parametar vjerojatnosti zamjene akcije - *stickiness*. Riječ je o *sticky actions* pristupu koji uvodi stohastiku, te ne ometa odabir akcije agenta koji može biti siguran da njegove akcije neće biti pogubne za ostatak iteracije s okolinom.

Pristupi *sticky actions* i *frame skipping* dio su *OpenAI Gym* biblioteke i mogu se specificirati prilikom instanciranja okoline, dok je *random action noise* i *initial no-ops* tehnike potrebno manualno implementirati ili se osloniti na izvedbu u poznatim *Python* bibliotekama *stable-baselines* i *stable-baselines3*.

3.3. Okruženja

U *OpenAI Gym* ekosustavu dostupno je puno okruženja koja omogućuju interakciju s agentom. Neki od njih su: *Atari* - skup od Atari 2600 okolina, *MuJoCo* (punim nazivom *Multi-Joint dynamics with Contact*) - skup okolina za provođenje istraživanja i razvoja u robotici, biomehanici, grafici i drugim područjima gdje je potrebna brzina i točna simulacija, *Classic Control* - skup okolina koje opisuju poznate fizikalne eksperimente. Opisat ćemo neke od najkorištenijih okolina.

3.3.1. Okruženje CartPole

Ovo okruženje modelira fizikalni problem održavanja ravnoteže. Inačica je sličnog fizikalnog problema pod nazivom *obrnuto njihalo* (engl. *inverted pendulum*). Za pomična kolica zakvačen je stupić. Njegovo težište nalazi se iznad središta mase i na taj način osigurava da sustav nije stabilan. Zglob, odnosno dodirna točka između stupića

i kolica nema trenja niti drugih gubitaka. Također, kolica koja se kreću vodoravno po putanji u 2 smjera nemaju trenja niti drugih gubitaka. Cilj ovog fizikalnog problema jest uravnotežiti stup primjenom sila i pomicanjem kolica u lijevom ili desnom smjeru.

Za svaki poduzeti korak okolina dodjeljuje nagradu u vrijednosti $+1$. Struktura valjanih akcija koje agent može poduzeti (`action_space`) instanca je razreda `Discrete(2)` - skup akcija je diskretan i u svakom koraku je moguće odabrati 1 od maksimalno 2 dostupne akcije. Opis značenja svake akcije prikazan je u tablici 3.1. S druge strane, objekt koji predstavlja strukturu stanja okoline u određenom vremenskom trenutku (`observation_space`) instanca je razreda `Box(4)` - stanje se sastoji od 4 kontinuirane vrijednosti od kojih su neke ograničene i odozdo i odozgo. Točan opis i granice vrijednosti predloženi su u tablici 3.2.

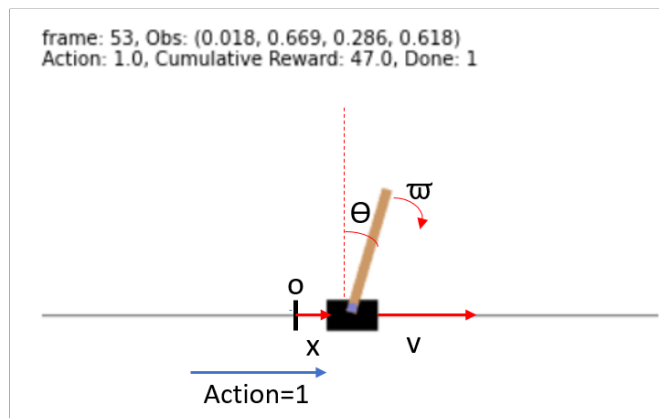
Tablica 3.1: Opis valjanih akcija okoline CartPole - atribut `action_space`

Akcija	Opis akcije
0	Pomak kolica ulijevo
1	Pomak kolica udesno

Tablica 3.2: Opis strukture okoline okoline CartPole - atribut `observation_space`

Indeks	Opis	Donja granica	Gornja granica
0	Pozicija kolica	-4.8	4.8
1	Brzina kolica	$-\infty$	∞
2	Nagib štapića i kolica	-0.418rad	0.418rad
3	Brzina štapića na vrhu	$-\infty$	∞

Početno stanje okoline inicijalizira se pozivom metode `reset()` slučajnim vrijednostima iz uniformne razdiobe na intervalu $[-0.05, 0.05]$. Okolina podržava 3 uvjeta zaustavljanja (uvjeti koji označuju da je riječ o terminalnom stanju): nagib štapića i kolica je izvan intervala $[-0.2095, 0.2095]\text{rad}$, pozicija sredine kolica je izvan intervala $[-2.4, 2.4]$ (sredina kolica dotiče rub vidljivog prostora) i duljina epizode veća je od 500 koraka. Na slici 3.3 prikazana je okolina CartPole zajedno s vrijednostima okoline.



Slika 3.3: Izgled i vrijednosti CartPole okoline [10]

3.3.2. Okruženje Breakout

Ova okolina simulira poznatu Atari 2600 igru u kojoj je cilj sakupiti što više bodova pomičući platformu i održavajući lopticu na ekranu. Platforma je postavljena na dnu ekrana, na fiksnoj visini i moguće ju je pomicati u dva smjera. Loptica se odbija između zidova, platforme i 6 razina *ciglenih* blokova čijim razbijanjem se sakupljaju bodovi. Ako loptica padne ispod platforme koju igrač kontrolira, gubi se život. Igra završava kada igrač potroši 5 života, odnosno kada 5 puta loptica padne ispod platforme.

Skup validnih akcija je instanca razreda `Discrete(4)` - skup akcija je diskretan i u svakom koraku je moguće odabrati 1 od maksimalno 4 dostupne akcije. Opis značenja svake akcije prikazan je u tablici 3.3. Kao opis trenutnog stanja okoline moguće je dobiti RGB vrijednosti svakog piksela slike (slike dimenzije 210×160) ili vrijednosti radne memorije *ALE* okoline (128 bajta) - što je korisno jer možemo preskočiti korak učenja reprezentacije okoline (preskačemo dio gdje algoritmi učenja moraju iz piksela slike naučiti reprezentaciju). Razlike u strukturi objekta okoline prikazane su u tablici 3.4.

Tablica 3.3: Opis valjanih akcija okoline Breakout - atribut `action_space`

Akcija	Opis akcije	Detaljniji opis akcije
0	NOOP	Ne poduzima se nikakva akcija
1	FIRE	Akcija koja pokreće igru
2	RIGHT	Platforma se pomiče udesno
3	LEFT	Platforma se pomiče ulijevo

Tablica 3.4: Opis strukture objekta okoline Breakout - atribut `observation_space`

Indeks	Struktura
RAM vrijednosti	<code>Box(0, 255, (128,), uint8)</code>
Vrijednosti RGB slike	<code>Box(0, 255, (210, 160, 3), uint8)</code>

Nagrada dolazi u obliku bodova koji se dobivaju uništavajući *ciglene* blokove. Vrijednost nagrade ovisi o boji cigle. Izgled same okoline prikazan je na slici 3.4.



Slika 3.4: Okolina Breakout

4. Implementacija

4.1. Okolina CartPole

4.1.1. Implementacija Deep Q Learning algoritma

https://github.com/PauloSankovic/retro-ai/blob/master/models/DeepQNetworkAgent-state_env%3Dcartpole_gamma%3D0.99_lr%3D5e-4_bs%3D32_es%3D1_ee%3D0.02.pickle

4.1.2. Implementacija Double Deep Q Learning algoritma

4.1.3. Implementacija Actor Critic algoritma

https://github.com/PauloSankovic/retro-ai/blob/master/models/ActorCriticAgent_env%3Dcarpole_gamma%3D0.99_lr%3D3e-2.pickle
video igranja: <https://raw.githubusercontent.com/PauloSankovic/retro-ai/master/playground/video/cart-pole/openaigym.video.0.9768.video000000.mp4>

4.1.4. Usporedba

4.2. Okolina Breakout

4.2.1. Implementacija Deep Q Learning algoritma

https://github.com/PauloSankovic/retro-ai/blob/master/models/DeepQNetworkAgent-state_env%3Dbreakout_step%3D1080000.pickle

4.2.2. Implementacija Double Deep Q Learning algoritma

<https://github.com/PauloSankovic/retro-ai/blob/master/models/>

DoubleDeepQNetworkAgent-state_env%3Dbreakout_step%3D1400000_
v%3D2.pickle

video igranja: <https://github.com/PauloSankovic/retro-ai/tree/master/playground/video/breakout-ddqn>

4.2.3. Implementacija Actor Critic algoritma

4.2.4. Usporedba

4.3. Usporedba algoritama

5. Moguća poboljšanja

Korištenje drugih optimizacijskih postupaka osim Adam-a kao npr. SGD, SGD s momentom, SGD s restartom, RMSProp, RMSProp s momentom Korištenje restarta kod promjene stope učenja... igranje s hiperparametrima

- Implementacija ostalih algoritama podržanog učenja

- Isprobati više verzija arhitekture dubokog modela i procijeniti benefite svakog.

- Naučiti agente i provjeriti kako se ponašaju na novim vrstama okolina i kako se suočavaju s novim problemima.

6. Zaključak

LITERATURA

- [1] Vectorized environments. URL https://stable-baselines.readthedocs.io/en/master/guide/vec_envs.html.
- [2] Part 2: Kinds of rl algorithms. URL https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html.
- [3] M. G. Bellemare, Y. Naddaf, J. Veness, i M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, jun 2013.
- [4] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, i Wojciech Zaremba. Openai gym, 2016.
- [5] Ian Goodfellow, Yoshua Bengio, i Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [6] Alexander LeNail. Nn-svg: Publication-ready neural network architecture schematics. *Journal of Open Source Software*, 4(33):747, 2019. doi: 10.21105/joss.00747. URL <https://doi.org/10.21105/joss.00747>.
- [7] Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew J. Hausknecht, i Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *CoRR*, abs/1709.06009, 2017. URL <http://arxiv.org/abs/1709.06009>.
- [8] Nicolas Maquaire. Are the space invaders deterministic or stochastic?, Sep 2020. URL <https://towardsdatascience.com/are-the-space-invaders-deterministic-or-stochastic-595a30becae2>.
- [9] Wikipedia. Reinforcement learning — Wikipedia, the free encyclopedia. <http://en.wikipedia.org/w/index.php?title=>

Reinforcement%20learning&oldid=1090221087, 2022. [Online; accessed 25-June-2022].

- [10] Jian Xu. How to beat the cartpole game in 5 lines, Feb 2021. URL <https://towardsdatascience.com/how-to-beat-the-cartpole-game-in-5-lines-5ab4e738c93f>.
- [11] Marko Čupić. *Umjetna Inteligencija - Uvod u strojno učenje*. 2022. <http://java.zemris.fer.hr/nastava/ui/ml/ml-20220430.pdf>.

Automatizirani razvoj agenata za različita okruženja

Sažetak

Ovaj rad proučava problematiku podržanog učenja gdje se bez znanja o pravilima i funkcioniranju specifične okoline, želi konstruirati agent kojemu je cilj pronaći optimalnu strategiju koja maksimizira očekivanu dobit u određenom vremenskom okviru. Osim opisa podržanog učenja, dubokih modela i algoritama koji se baziraju na dubokim modelima, detaljno je prezentirana struktura *OpenAI Gym* biblioteke i njenih okolina, te napravljena usporedba u kojoj se analizira uspješnost naučenih agenata.

Programsko rješenje implementirano je u programskom jeziku *Python*, primarno koristeći *PyTorch* radni okvir za dizajn i učenje dubokih neuronskih mreža, te *OpenAI Gym* biblioteka za simulaciju i testiranje ponašanja agenata.

Ključne riječi: Podržano učenje, duboko učenje, neuronske mreže, OpenAI Gym, algoritmi neovisni o modelu

Automated design of agents for various environments

Abstract

This paper studies the issue of reinforcement learning where without knowing the rules and internal dynamics of environment, the goal is to build an agent whose aim is to find the optimal strategy that maximizes the expected reward in a given time frame. Aside from describing the basics of reinforcement learning, deep learning models and model free algorithms, the key concepts of *OpenAI Gym* library and its environments are broken down in detail, and a comparison is made between agents.

Software implementation is written in *Python*, primarily using the *PyTorch* framework for designing and learning deep neural networks, and the *OpenAI Gym* library for simulating and testing agent behavior.

Keywords: Reinforcement learning, deep learning, neural networks, OpenAI Gym, Model-Free RL algorithms

POPIS SLIKA

2.1. Ciklus interakcije agenta s okolinom [9]	4
2.2. Arhitektura potpuno povezanog modela	5
2.3. Arhitektura konvolucijskog modela	7
2.4. Podjela algoritama podržanog učenja [2]	8
3.1. Rezultat pokretanja koda 3	11
3.2. Vektorizirano okruženje sa 4 paralelne asinkrone instance	14
3.3. Izgled i vrijednosti CartPole okoline [10]	17
3.4. Okolina Breakout	18

POPIS TABLICA

3.1. Opis valjanih akcija okoline CartPole - atribut <code>action_space</code> . . .	16
3.2. Opis strukture okoline okoline CartPole - atribut <code>observation_space</code>	16
3.3. Opis valjanih akcija okoline Breakout - atribut <code>action_space</code> . . .	17
3.4. Opis strukture objekta okoline Breakout - atribut <code>observation_space</code>	18

POPIS PROGRAMSKIH ISJEČAKA

1.	Implementacija potpuno povezanog modela na slici 2.4 koristeći biblioteku <i>PyTorch</i>	6
2.	Implementacija konvolucijskog modela na slici 2.3 koristeći biblioteku <i>PyTorch</i>	7
3.	Jednostavan primjer integracije agenta i Gym okoline (1 epizoda) . . .	10
4.	Primjer korištenja strukture kontinuiranog prostora <code>Box</code>	12
5.	Primjer korištenja strukture diskretnog prostora <code>Discrete</code>	13