

Robotsko učenje za zadatke manipulacije

Kuzmić, Marko

Master's thesis / Diplomski rad

2024

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Zagreb, Faculty of Electrical Engineering and Computing / Sveučilište u Zagrebu, Fakultet elektrotehnike i računarstva**

Permanent link / Trajna poveznica: <https://um.nsk.hr/um:nbn:hr:168:205283>

Rights / Prava: [In copyright](#)/[Zaštićeno autorskim pravom.](#)

Download date / Datum preuzimanja: **2025-03-21**



Repository / Repozitorij:

[FER Repository - University of Zagreb Faculty of Electrical Engineering and Computing repository](#)



SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 677

ROBOTSKO UČENJE ZA ZADATKE MANIPULACIJE

Marko Kuzmić

Zagreb, lipanj 2024.

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 677

ROBOTSKO UČENJE ZA ZADATKE MANIPULACIJE

Marko Kuzmić

Zagreb, lipanj 2024.

DIPLOMSKI ZADATAK br. 677

Pristupnik: **Marko Kuzmić (0036522376)**

Studij: Računarstvo

Profil: Računarska znanost

Mentorica: prof. dr. sc. Marija Seder

Zadatak: **Robotsko učenje za zadatke manipulacije**

Opis zadatka:

Jedan je od temeljnih izazova u robotici stvaranje robota sposobnih za izvršavanje zadataka zasnovano na interakcijama s njihovim okruženjem. Robotska manipulacija obuhvaća moguće načine na koje robot može djelovati i mijenjati stanje svojeg okruženja, primjerice podizanjem ili premještanjem objekata. Robotsko učenje, odnosno učenje željenih robotskih vještina iz skupovnih podataka, ključno je za postizanje autonomnih robotskih sustava za manipulaciju, budući da je realan svijet iznimno kompleksan (3D okolinu, pasivne objekte i aktivne agente koji se u njoj nalaze te sve potrebne vještine nije izvedivo precizno modelirati, tj. nije moguće svo potrebno znanje toliko precizno definirati i unaprijed ugraditi u robota). Cilj je ovog rada opisati i formalizirati problem robotske manipulacije te robotskog učenja zasnovano na reprezentacijama okruženja i objekata te upravljanja robotom potpornim učenjem. Nadalje, cilj je iskoristiti postojeće simulatore i biblioteke za generiranje sintetičkih podataka koji će poslužiti za kvalitativnu i kvantitativnu usporedbu rezultata dobivenih različitim pristupima za učenje manipulacije.

Rok za predaju rada: 28. lipnja 2024.

Sadržaj

1. Uvod	5
1.1. Istraživačka pozadina	5
1.2. Pregled domaće i međunarodne literature	6
1.2.1. Domaća istraživanja	6
1.2.2. Međunarodna istraživanja	6
1.3. Ciljevi istraživanja	7
1.4. Struktura rada	8
2. Osnovni pojmovi i tehnologije	10
2.1. Reinforcement Learning (RL)	10
2.1.1. Osnovni pojmovi RL-a	10
2.1.2. Algoritmi RL-a	11
2.1.3. Primjena RL-a	11
2.2. Robotska manipulacija	12
2.2.1. Ključni izazovi u robotskoj manipulaciji	12
2.2.2. Metode robotske manipulacije	12
2.2.3. Primjeri robotske manipulacije	12
2.3. Gymnasium biblioteka	13
2.3.1. Značajke Gymnasiuma	13
2.3.2. Integracija s RL algoritmima	13
2.3.3. Primjena u robotskoj manipulaciji	14
2.4. OpenRL biblioteka	14
2.4.1. Ključne značajke	14
2.4.2. Prednosti	15
2.4.3. Ograničenja	15

2.4.4.	Zaključak	16
2.5.	Stable Baselines3 biblioteka	16
2.5.1.	Ključne značajke	17
2.5.2.	Prednosti	17
2.5.3.	Ograničenja	18
2.5.4.	Zaključak	18
2.6.	MuJoCo biblioteka	18
2.6.1.	Ključne značajke MuJoCo-a	19
2.6.2.	Upotreba MuJoCo-a u istraživanju i razvoju	19
2.6.3.	Prednosti i ograničenja MuJoCo-a	20
2.6.4.	Zaključak	21
3.	Metodologija	22
3.1.	Uvod	22
3.2.	On-Policy i Off-Policy algoritmi	22
3.2.1.	On-Policy Algoritmi	23
3.2.2.	Off-Policy Algoritmi	23
3.2.3.	Zaključak	24
3.3.	Prostor akcija i opažanja u Gymnasiumu	24
3.3.1.	Prostor akcija (Action Space)	25
3.3.2.	Prostor opažanja (Observation Space)	26
3.3.3.	Zaključak	27
3.4.	Razlika između Sparse i Dense nagrada	27
3.4.1.	Sparse nagrade	27
3.4.2.	Dense nagrade	28
3.5.	Q-learning algoritam	28
3.5.1.	Osnove Q-learninga	28
3.5.2.	Metodologija Q-learninga	29
3.5.3.	Prednosti Q-learninga	30
3.5.4.	Primjena Q-learninga u robotskoj manipulaciji	30
3.6.	Deep Q-Network (DQN) Algoritam	31
3.6.1.	Uvod	31
3.6.2.	Osnovni principi DQN-a	31

3.6.3.	Struktura DQN-a	32
3.6.4.	Algoritam	32
3.6.5.	Prednosti i izazovi DQN-a	33
3.6.6.	Zaključak	33
3.7.	Deep Deterministic Policy Gradient (DDPG) Algoritam	34
3.7.1.	Uvod	34
3.7.2.	Osnovni principi DDPG-a	34
3.7.3.	Struktura DDPG-a	34
3.7.4.	Algoritam	35
3.7.5.	Prednosti i izazovi DDPG-a	36
3.7.6.	Zaključak	37
3.8.	Soft Actor-Critic (SAC) Algoritam	37
3.8.1.	Uvod	37
3.8.2.	Osnovni principi SAC-a	37
3.8.3.	Struktura SAC-a	38
3.8.4.	Algoritam	38
3.8.5.	Prednosti i izazovi SAC-a	39
3.8.6.	Zaključak	40
4.	Eksperimentalna evaluacija	41
4.1.	Cart-Pole	41
4.1.1.	Vlastita implementacija	42
4.1.2.	OpenRL implementacija	43
4.1.3.	Rezultati	45
4.2.	Mountain-Car	46
4.2.1.	Vlastita implementacija	47
4.2.2.	OpenRL implementacija	48
4.2.3.	Rezultati	50
4.3.	FetchReach-v2 eksperiment	51
4.3.1.	Opis okruženja	51
4.3.2.	Algoritmi korišteni u eksperimentu	55
4.3.3.	Eksperimentalni postupak	56
4.3.4.	Evaluacija performansi agenata	58

4.3.5. Rezultati	59
4.3.6. Diskusija	60
4.4. Zaključak	61
5. Zaključak i budući rad	62
5.1. Glavni zaključci	62
5.2. Budući rad	63
Literatura	65
Sažetak	69
Abstract	70

1. Uvod

1.1. Istraživačka pozadina

Robotska manipulacija je jedno od najvažnijih područja istraživanja u modernoj robotici, jer omogućava robotima da izvršavaju različite zadatke koji zahtijevaju interakciju s fizičkim objektima. Ovi zadaci mogu uključivati podizanje, premještanje, slaganje ili bilo koju drugu aktivnost koja zahtijeva preciznu kontrolu i prilagodbu okruženju.

Povijesno gledano, programiranje robota za takve zadatke bilo je iznimno složeno, zahtijevajući ručno kodiranje svakog mogućeg scenarija s kojim bi se robot mogao susresti. Međutim, ovakav pristup nije skalabilan za složene ili dinamične okoline. Stoga je uvođenje metoda strojnog učenja, posebno reinforcement learninga (RL), postalo ključnim za razvoj naprednih autonomnih robotskih sustava.

Reinforcement learning je metoda u kojoj agenti uče optimalne akcije putem interakcije s okruženjem, temeljeći svoje odluke na povratnim informacijama u obliku nagrada i kazni. RL omogućava robotima da razviju vještine koje nisu unaprijed programirane, već su naučene kroz iskustvo.

Simulacije igraju ključnu ulogu u treniranju RL modela. Koristeći simulacije, istraživači mogu stvoriti kontrolirana okruženja u kojima roboti mogu sigurno učiti i prilagođavati svoje ponašanje bez rizika od oštećenja stvarne opreme. Gymnasium je moderna biblioteka koja omogućava jednostavno stvaranje i upravljanje simulacijama za RL, dok Q-learning algoritam pruža jednostavnu implementaciju RL metoda, olakšavajući njihov razvoj i primjenu.

U ovom radu koristili smo različite RL algoritme, uključujući vlastitu implementaciju Q-learning algoritma, te biblioteke koje implementiraju Q-learning, DQN (Deep Q-

learning), PPO (Proximal Policy Optimization), SAC (Soft Actor-Critic) i DDPG (Deep Deterministic Policy Gradient). Svaki od ovih algoritama ima svoje prednosti i primjenjiv je u različitim kontekstima robotske manipulacije. Q-learning je posebno koristan za diskretne akcijske prostore, dok DDPG, SAC, PPO i DQN pružaju naprednije metode za kontinuirane akcijske prostore i složenije zadatke.

1.2. Pregled domaće i međunarodne literature

1.2.1. Domaća istraživanja

Nekoliko istraživačkih skupina na sveučilištima u Hrvatskoj radi na primjeni RL-a za robotsku manipulaciju. Na primjer, istraživači na Fakultetu elektrotehnike i računarstva u Zagrebu razvili su sustave koji koriste sintetičke podatke generirane u simuliranim okruženjima za treniranje RL modela. U svom radu, pokazali su kako ovakav pristup može smanjiti vrijeme potrebno za treniranje robota i povećati učinkovitost u složenim zadacima manipulacije [1].

Drugi primjer dolazi s Tehničkog veleučilišta u Zagrebu, gdje su istraživači koristili Gymnasium za stvaranje realističnih simulacija industrijskih okruženja. Njihovo istraživanje pokazalo je da kombinacija simulacija i RL algoritama može značajno poboljšati sposobnosti robota za izvršavanje preciznih manipulacijskih zadataka [2].

1.2.2. Međunarodna istraživanja

Na međunarodnoj razini, brojni radovi objavljeni u prestižnim časopisima i na konferencijama pokazuju kako RL može unaprijediti robotsku manipulaciju. Na primjer, istraživanje objavljeno u časopisu IEEE Transactions on Robotics pokazalo je da korištenje naprednih RL algoritama poput PPO (Proximal Policy Optimization) može poboljšati učinkovitost robota u složenim manipulacijskim zadacima [3].

Također, istraživači s MIT-a su koristili Gymnasium za stvaranje simulacija koje repliciraju stvarne industrijske scenarije. Kombinirajući ove simulacije s algoritmima iz CleanRL-a, uspjeli su postići visoku razinu preciznosti i učinkovitosti u zadacima poput slaganja i sortiranja objekata [4].

Još jedno značajno istraživanje dolazi s Stanforda, gdje su razvili sustav koji koristi generativne modela za stvaranje sintetičkih podataka u simulacijama. Ovaj sustav omogućava robotima da nauče manipulacijske vještine u simuliranim okruženjima i kasnije ih primijene u stvarnom svijetu s minimalnim prilagođavanjem [5].

Dodatno, istraživanje provedeno na Sveučilištu u Kaliforniji pokazalo je kako korištenje DDPG algoritma može značajno poboljšati performanse robota u zadacima fine manipulacije. U njihovom radu, DDPG je korišten za treniranje robota u preciznim zadacima hvatanja i premještanja malih objekata, demonstrirajući prednosti kontinuiranih akcijskih prostora u složenim okruženjima [6].

1.3. Ciljevi istraživanja

Glavni cilj ovog rada je istražiti primjenu reinforcement learning (RL) algoritama za učenje složenih manipulacijskih vještina korištenjem biblioteka Gymnasium i Q-learning algoritma. Specifični ciljevi uključuju:

1. Opis i formalizacija problema robotske manipulacije:

- Definirati ključne aspekte manipulacijskih zadataka u robotskim sustavima.
- Identificirati izazove i prepreke u učenju manipulacijskih vještina.

2. Demonstracija kombinacije Gymnasiuma i Q-learning algoritma:

- Pokazati kako ova biblioteka može olakšati razvoj i treniranje RL modela.
- Implementirati Q-learning algoritam za specifične zadatke manipulacije.

3. Evaluacija Performansi Različitih RL Algoritama:

- Detaljno analizirati i usporediti performanse različitih algoritama podržanog učenja (Q-learning, DQN, PPO, DDPG, SAC, DDPG + HER) u zadacima robotske manipulacije unutar simulacijskih okruženja kao što su CartPole, MountainCar i FetchReach-v2.

4. Istraživanje učinka sparse i dense nagrada:

- Istražiti kako različiti tipovi nagrada (sparse i dense) utječu na proces učenja i konačne performanse algoritama podržanog učenja. Identificirati algoritme koji se bolje prilagođavaju rijetkim nagradama te one koji su efikasniji s kontinuiranim nagradama.

5. Primjena paralelizma u treningu RL agenata:

- Ispitati efikasnost korištenja paralelizma u procesu treniranja RL agenata. Konkretno, koristiti više paralelnih instanci simulacijskih okruženja kako bi se smanjilo ukupno vrijeme treniranja i poboljšala efikasnost učenja.

6. Implementacija i optimizacija HER-a u kombinaciji s DDPG-om:

- Implementirati Hindsight Experience Replay (HER) u kombinaciji s DDPG algoritmom i analizirati kako ova kombinacija poboljšava učenje u okruženjima s rijetkim nagradama. Istražiti optimalne konfiguracije za ovu kombinaciju.

7. Analiza učinka različitih algoritama u specifičnim okruženjima:

- Detaljno analizirati kako se različiti RL algoritmi ponašaju u specifičnim simulacijskim okruženjima (CartPole, MountainCar, FetchReach-v2) i koji algoritmi najbolje odgovaraju specifičnim zadacima manipulacije.

1.4. Struktura rada

Ovaj rad je strukturiran u pet glavnih poglavlja. Nakon uvodnog poglavlja, drugo poglavlje pruža pregled osnovnih pojmova i tehnologija vezanih uz podržano učenje (RL), robotsku manipulaciju te biblioteke i alate korištene u istraživanju, poput Gymnasiuma i MuJoCo-a. Treće poglavlje detaljno opisuje metodologiju istraživanja, uključujući arhitekturu sustava, implementaciju različitih RL algoritama (Q-learning, DQN, PPO, DDPG, SAC, DDPG + HER) te specifične korake razvoja i evaluacije. Četvrto poglavlje predstavlja eksperimentalnu evaluaciju, gdje se detaljno opisuju postavke eksperimenata, rezultati kratkoročnog i dugoročnog treniranja u različitim simulacijskim okruženjima (CartPole, MountainCar, FetchReach-v2) te analiza performansi algoritama uz usporedbu sparse i dense nagrada. Peto poglavlje donosi zaključke, diskusiju o postignutim rezul-

tatima i prijedloge za buduća istraživanja, s naglaskom na optimizaciju algoritama, prijenos znanja iz simulacije u stvarni svijet te integraciju dodatnih senzorskih podataka.

2. Osnovni pojmovi i tehnologije

2.1. Reinforcement Learning (RL)

Reinforcement Learning (RL) je grana strojnog učenja koja se fokusira na donošenje sekvencijalnih odluka. RL agent uči kroz interakciju s okruženjem, s ciljem maksimiziranja kumulativne nagrade. Ova metoda se razlikuje od nadziranog učenja jer agent samostalno otkriva koje akcije dovode do najboljih rezultata [7].

2.1.1. Osnovni pojmovi RL-a

- **Agent (Agent):** Entitet koji donosi odluke i izvršava akcije unutar okruženja.
- **Okruženje (Environment):** Svijet s kojim agent interagira.
- **Stanje (State):** Trenutna situacija ili konfiguracija okruženja koju agent opaža.
- **Akcija (Action):** Odabir agenta u određenom stanju, koja mijenja stanje okruženja.
- **Nagrada (Reward):** Povratna informacija koja agentu govori koliko je određena akcija bila korisna.
- **Politika (Policy):** Pravilo ili funkcija koja definira akcije koje agent odabire na temelju stanja.
- **Funkcija vrijednosti (Value Function):** Funkcija koja procjenjuje dugoročnu korisnost stanja ili stanja i akcija.
- **Model okruženja (Environment Model):** Reprerzentacija dinamike okruženja koja se koristi u model-based RL.

2.1.2. Algoritmi RL-a

Postoje mnogi RL algoritmi koji se koriste za rješavanje različitih problema. Neki od najvažnijih uključuju:

- **Q-learning:** Model-free algoritam koji uči optimalnu Q-funkciju koja procjenjuje očekivanu vrijednost nagrade za svaku akciju u svakom stanju [8].
- **Deep Q-Networks (DQN):** Kombinira Q-learning s dubokim neuronskim mrežama kako bi se nosio s velikim i kontinuiranim prostorima stanja [9].
- **Proximal Policy Optimization (PPO):** On-policy algoritam koji koristi clipping funkciju za balansiranje između istraživanja i eksploatacije, pogodan za kontinuirane akcijske prostore [10].
- **Deep Deterministic Policy Gradient (DDPG):** Off-policy algoritam dizajniran za kontinuirane akcijske prostore, koristi dvije neuronske mreže: jednu za politiku i jednu za procjenu vrijednosti [11].
- **Soft Actor-Critic (SAC):** Off-policy algoritam koji koristi dvostruke kritičare za stabilizaciju učenja i regulaciju entropije za balansiranje između istraživanja i eksploatacije, izuzetno učinkovit za zadatke s kontinuiranim akcijskim prostorima. [12].

2.1.3. Primjena RL-a

RL se može primijeniti u mnogim područjima, uključujući:

- **Robotika:** Za učenje kontrolnih politika za robotske manipulacijske zadatke [13].
- **Igre:** Za razvoj agenata koji mogu igrati igre na ljudskoj ili nadljudskoj razini [14].
- **Financije:** Za optimizaciju portfelja i donošenje trgovačkih odluka [15].
- **Zdravstvo:** Za personalizaciju medicinskih tretmana i optimizaciju bolničkih operacija [16].

2.2. Robotska manipulacija

Robotska manipulacija je grana robotike koja se bavi razvojem sustava i algoritama koji omogućuju robotima interakciju s fizičkim objektima u njihovom okruženju. Glavni cilj je omogućiti robotima da izvode složene zadatke kao što su hvatanje, podizanje, premještanje i slaganje objekata [17].

2.2.1. Ključni izazovi u robotskoj manipulaciji

- **Varijabilnost okruženja:** Fizički svijet je vrlo dinamičan i nepredvidiv, što otežava predviđanje svih mogućih scenarija.
- **Precizna kontrola:** Manipulacijski zadaci često zahtijevaju visoku razinu preciznosti u kontroli pokreta robota.
- **Senzorna integracija:** Kombinacija informacija iz različitih senzora (npr. kamera, dodirni senzori) za donošenje informiranih odluka.
- **Sigurnost:** Izvođenje manipulacijskih zadataka mora biti sigurno kako bi se izbjegle štete na robotu, objektima ili ljudima u blizini.

2.2.2. Metode robotske manipulacije

- **Model-based metode:** Korištenje matematičkih modela za predviđanje ponašanja robota i planiranje akcija [18].
- **Model-free metode:** Korištenje RL algoritama za učenje optimalnih akcija izravno iz podataka bez eksplicitnog modeliranja dinamike sustava [19].
- **Hibridne metode:** Kombinacija model-based i model-free pristupa za postizanje bolje performanse i robusnosti [20].

2.2.3. Primjeri robotske manipulacije

- **Industrijska robotika:** Roboti u tvornicama koji izvode zadatke kao što su sklapanje proizvoda, zavarivanje i bojanje [21].

- **Servisna robotika:** Roboti koji pomažu u kućanskim poslovima ili pružaju usluge u zdravstvenim ustanovama [22].
- **Eksplorativna robotika:** Roboti koji se koriste za istraživanje teško dostupnih ili opasnih okruženja, poput svemira ili dubokih mora [23].

2.3. Gymnasium biblioteka

Gymnasium je održavana verzija OpenAI-ove Gym biblioteke za podržano učenje. Omogućava jednostavno i intuitivno korištenje Python sučelja za predstavljanje problema podržanog učenja te uključuje širok spektar referentnih okruženja. Gymnasium je kompatibilan s prethodnim Gym okruženjima i nudi alate za izradu i prilagodbu vlastitih okruženja, treniranje agenata i implementaciju prilagođenih omotača. Razvijena je od strane Farama Foundation s ciljem pojednostavljivanja procesa stvaranja i upravljanja simulacijama za RL istraživanja [24].

2.3.1. Značajke Gymnasiuma

- **Jednostavnost upotrebe:** Intuitivno sučelje koje omogućuje brzo stvaranje i prilagodbu simulacijskih okruženja.
- **Široka primjena:** Podrška za razne vrste simulacija, uključujući fizikalne simulacije, igre, edukativne scenarije i druge.
- **Ekstenzibilnost:** Mogućnost prilagodbe i proširenja funkcionalnosti prema specifičnim potrebama korisnika.
- **Podrška za sintetičke podatke:** Generiranje sintetičkih podataka koji se koriste za treniranje RL modela, smanjujući potrebu za prikupljanjem stvarnih podataka.

2.3.2. Integracija s RL algoritmima

Gymnasium podržava integraciju s različitim RL algoritmima, omogućavajući korisnicima da testiraju i uspoređuju performanse različitih pristupa u simuliranim okruženjima. Ovo je posebno korisno za razvoj i optimizaciju algoritama prije njihove primjene u stvarnom svijetu [25].

2.3.3. Primjena u robotskoj manipulaciji

Korištenje Gymnasiuma za stvaranje simulacija industrijskih okruženja omogućuje treniranje robota za složene zadatke manipulacije. Na primjer, simulacije mogu uključivati scenarije kao što su sklapanje proizvoda, sortiranje objekata ili interakcija s dinamičkim okruženjima. Gymnasium omogućava istraživačima da brzo iteriraju između različitih dizajna i pristupa, testirajući njihove performanse i robusnost.

2.4. OpenRL biblioteka

OpenRL je sveobuhvatna biblioteka otvorenog koda razvijena od strane OpenRL-Laba, dizajnirana za istraživanje i razvoj algoritama za podržano učenje (RL). Cilj OpenRL-a je pružiti istraživačima, inženjerima i entuzijastima alate potrebne za brzu implementaciju, testiranje i evaluaciju različitih RL algoritama. Biblioteka je posebno korisna za zadatke robotske manipulacije, autonomne vožnje, igranje videoigara i druge aplikacije koje zahtijevaju sofisticirano učenje i prilagodbu ponašanja agenata u dinamičkim okruženjima [26].

OpenRL omogućava korisnicima da brzo postave eksperimentalna okruženja, implementiraju RL algoritme i evaluiraju performanse agenata. Zbog svoje fleksibilnosti i modularnosti, OpenRL se može koristiti za širok spektar aplikacija, od akademskih istraživanja do industrijskih primjena.

2.4.1. Ključne značajke

OpenRL pruža nekoliko ključnih značajki koje ga čine izuzetno moćnim alatom za rad s RL algoritmima:

- **Podrška za razne RL algoritme:** OpenRL podržava širok spektar RL algoritama, uključujući Q-learning, DQN, PPO, DDPG i SAC, što omogućava korisnicima fleksibilnost u odabiru algoritma koji najbolje odgovara njihovom problemu. Ova raznolikost omogućava istraživačima da testiraju i uspoređuju različite pristupe rješavanju problema.
- **Integracija sa simulacijskim okruženjima:** Biblioteka je kompatibilna s popularnim simulacijskim platformama kao što su Gymnasium, Unity, i CARLA, što

omogućava jednostavno postavljanje i izvođenje eksperimenata. Korištenjem ovih okruženja, korisnici mogu simulirati složene scenarije koji repliciraju stvarne svjetove, čime se omogućava sigurnije i učinkovitije treniranje agenata.

- **Modularnost:** Dizajn OpenRL-a omogućava korisnicima jednostavno dodavanje novih algoritama, politika i okruženja, čineći ga vrlo prilagodljivim za specifične istraživačke potrebe. Modularni pristup također olakšava održavanje i proširenje koda, što je korisno za dugoročne projekte i kolaborativna istraživanja.
- **Vizualizacija i analiza:** OpenRL nudi alate za vizualizaciju procesa učenja, performansi agenata i analizu podataka, što olakšava interpretaciju rezultata i dijagnosticiranje problema. Vizualizacija performansi omogućava korisnicima da brzo identificiraju probleme u učenju i prilagode hiperparametre kako bi poboljšali učinkovitost algoritama.

2.4.2. Prednosti

- **Jednostavnost korištenja:** OpenRL je dizajniran da bude intuitivan i jednostavan za korištenje, čak i za korisnike koji nisu stručnjaci za RL. Ovaj pristup smanjuje krivulju učenja i omogućava brzi početak rada s bibliotekom.
- **Fleksibilnost:** Modularni dizajn omogućava jednostavno proširenje i prilagodbu biblioteke specifičnim potrebama korisnika. Korisnici mogu dodavati nove algoritme, okruženja i module bez potrebe za velikim promjenama u osnovnom kodu.
- **Podrška za više algoritama i okruženja:** Široka podrška za različite RL algoritme i simulacijska okruženja čini OpenRL moćnim alatom za istraživanje i razvoj. Ova fleksibilnost omogućava korisnicima da istražuju različite pristupe i kombinacije kako bi pronašli optimalna rješenja za svoje probleme.

2.4.3. Ograničenja

- **Performanse:** Kao i kod drugih opsežnih biblioteka, performanse mogu biti ograničene u vrlo velikim i složenim okruženjima ili pri radu s velikim modelima. Optimizacija performansi može zahtijevati dodatno podešavanje i prilagodbu.

- **Zavisnosti:** OpenRL se oslanja na više vanjskih biblioteka i alata, što može dovesti do problema s kompatibilnošću i upravljanjem zavisnostima. Korisnici moraju pažljivo upravljati verzijama i zavisnostima kako bi osigurali glatko funkcioniranje cijelog sustava.

2.4.4. Zaključak

OpenRL je snažan alat za istraživanje i razvoj RL algoritama, pružajući korisnicima sve potrebne resurse za brzo postavljanje, treniranje i evaluaciju agenata u raznim simulacijskim okruženjima. Njegova jednostavnost korištenja, fleksibilnost i široka podrška za različite algoritme i okruženja čine ga idealnim izborom za istraživače i inženjere koji žele unaprijediti svoje radove u području podržanog učenja.

S obzirom na sve navedene prednosti, OpenRL predstavlja značajan alat u arsenalu svakog istraživača ili praktičara koji se bavi reinforcement learningom. Budući razvoj i ažuriranja ove biblioteke mogu dodatno proširiti njene mogućnosti, čineći je još korisnijom za širok spektar aplikacija u različitim industrijama i istraživačkim područjima.

2.5. Stable Baselines3 biblioteka

Stable Baselines3 (SB3) je popularna biblioteka otvorenog koda za podržano učenje (RL), koja je nastala kao nasljednik originalne Stable Baselines biblioteke. Razvijena je s ciljem pružanja pouzdanih, jednostavnih za korištenje i dobro testiranih implementacija RL algoritama. SB3 je namijenjen istraživačima i inženjerima koji se bave razvojem i primjenom RL algoritama u različitim domenama, uključujući robotsku manipulaciju, autonomnu vožnju, igranje videoigara i mnoge druge [27].

Stable Baselines3 je izgrađen na vrhu PyTorch-a, popularnog deep learning frameworka, što omogućava lako integriranje sa širokim spektrom modela i alata za duboko učenje. Biblioteka nudi konzistentno i efikasno API sučelje, koje je pogodno za brzo prototipiranje i evaluaciju RL algoritama.

2.5.1. Ključne značajke

Stable Baselines3 pruža niz značajki koje ga čine moćnim alatom za rad s RL algoritima:

- **Podrška za razne RL algoritme:** SB3 uključuje implementacije najkorištenijih RL algoritama kao što su DQN, PPO, A2C (Advantage Actor Critic), DDPG, TD3 (Twin Delayed DDPG) i SAC (Soft Actor-Critic). Ove implementacije su optimizirane za performanse i lako se mogu prilagoditi specifičnim potrebama korisnika.
- **Integracija sa simulacijskim okruženjima:** SB3 je kompatibilan s popularnim simulacijskim platformama kao što su OpenAI Gym, Gymnasium, i custom okruženja definirana od strane korisnika. Ova integracija omogućava korisnicima jednostavno postavljanje i izvođenje eksperimenata u različitim simulacijskim scenarijima.
- **Jednostavnost korištenja:** Biblioteka je dizajnirana da bude intuitivna i jednostavna za korištenje, s dobro dokumentiranim API-jem i brojnim primjerima. Ovo omogućava korisnicima brz početak rada i lako usvajanje složenijih funkcionalnosti.
- **Podrška za multiple procesore:** SB3 nudi podršku za paralelno treniranje korištenjem više procesora, što može značajno ubrzati proces treniranja agenata. Ova značajka je posebno korisna za složena okruženja koja zahtijevaju veliko računalno vrijeme.
- **Vizualizacija i analiza:** SB3 dolazi s alatima za vizualizaciju procesa treniranja, performansi agenata i analizu rezultata. Ovi alati olakšavaju praćenje napretka i dijagnosticiranje problema tijekom treniranja.

2.5.2. Prednosti

- **Pouzdanost i testiranost:** SB3 implementacije algoritama su dobro testirane i dokumentirane, što osigurava pouzdanost i performanse tijekom treniranja.
- **Fleksibilnost:** Biblioteka omogućava jednostavnu prilagodbu i proširenje, što je korisno za istraživače koji žele testirati nove ideje i algoritme.

- **Kompatibilnost:** Zahvaljujući PyTorch-u, SB3 se lako integrira s drugim alatima za duboko učenje i može koristiti prednosti GPU akceleracije.

2.5.3. Ograničenja

- **Performanse:** Kao i kod drugih opsežnih biblioteka, performanse mogu biti ograničene u vrlo velikim i složenim okruženjima ili pri radu s velikim modelima. Optimizacija performansi može zahtijevati dodatno podešavanje i prilagodbu.
- **Zavisnosti:** SB3 se oslanja na više vanjskih biblioteka i alata, što može dovesti do problema s kompatibilnošću i upravljanjem zavisnostima. Korisnici moraju pažljivo upravljati verzijama i zavisnostima kako bi osigurali glatko funkcioniranje cijelog sustava.

2.5.4. Zaključak

Stable Baselines3 je izuzetno korisna biblioteka za istraživanje i razvoj algoritama za RL. Pruža sve potrebne alate za brzo postavljanje, treniranje i evaluaciju agenata unutar različitih simulacijskih okruženja. Njegova jednostavna upotreba, fleksibilnost i široka podrška za razne RL algoritme i okruženja čine ga idealnim izborom za istraživače i inženjere koji žele unaprijediti svoje radove u ovom području.

Uz sve navedene prednosti, Stable Baselines3 je ključni alat za svakog istraživača ili praktičara koji se bavi reinforcement learningom. Njegov kontinuirani razvoj i ažuriranja dodatno će proširiti njegove mogućnosti, čineći ga još vrednijim za širok spektar aplikacija u različitim industrijama i istraživačkim područjima.

2.6. MuJoCo biblioteka

MuJoCo (Multi-Joint dynamics with Contact) je napredna platforma za simulaciju fizike, razvijena kako bi omogućila precizno modeliranje i kontrolu složenih mehaničkih sustava. MuJoCo je široko korišten u istraživačkim zajednicama za simulaciju dinamike robota, biomehanike i drugih složenih sustava koji zahtijevaju visoku preciznost i efikasnost u simulaciji. MuJoCo je često korišten kao engine za fizikalnu simulaciju unutar Gymnasium okruženja. Ova integracija omogućava korisnicima da koriste MuJoCo-ove

napredne mogućnosti za simulaciju dinamike i kontakata, dok istovremeno koriste jednostavan i intuitivan API koji pruža Gymnasium [28].

2.6.1. Ključne značajke MuJoCo-a

MuJoCo se ističe nizom značajki koje ga čine idealnim alatom za istraživanje i razvoj u različitim domenama:

- **Precizna dinamika:** MuJoCo omogućava simulaciju složenih dinamika s visokom preciznošću, uključujući mehaniku više zglobova, kontaktne sile i ograničenja.
- **Efikasnost:** Simulacija u stvarnom vremenu i optimizirani algoritmi omogućavaju brzo izvođenje složenih simulacija, što je ključno za iterativne procese istraživanja i razvoja.
- **Fleksibilnost:** Podržava širok raspon modela i scenarija, od jednostavnih mehaničkih sustava do složenih robotskih konfiguracija i biomehaničkih modela.
- **Kompatibilnost:** MuJoCo je integriran s popularnim alatima za podržano učenje (RL) kao što su OpenAI Gym i Gymnasium, omogućavajući jednostavno postavljanje i izvođenje RL eksperimenata.
- **Vizualizacija:** Nudi napredne mogućnosti vizualizacije, uključujući realistične 3D prikaze simuliranih scenarija, što pomaže u analizi i prezentaciji rezultata.

2.6.2. Upotreba MuJoCo-a u istraživanju i razvoju

MuJoCo se koristi u različitim istraživačkim i industrijskim domenama zbog svoje sposobnosti da simulira kompleksne sustave s visokom preciznošću i efikasnošću.

Robotika

U području robotike, MuJoCo se koristi za simulaciju dinamike robota, razvoj kontrolnih algoritama i optimizaciju dizajna robota. Njegova sposobnost da precizno modelira kontakte i sile čini ga idealnim za simulaciju manipulacijskih zadataka i interakcija robota s okolinom.

Biomehanika

MuJoCo je također popularan u biomehanici za modeliranje i analizu ljudskih pokreta i biomehaničkih sustava. Istraživači koriste MuJoCo za proučavanje dinamike mišića, zglobova i skeletnih sustava, kao i za razvoj rehabilitacijskih uređaja i protetike.

Podržano učenje (RL)

MuJoCo je integriran s alatima za podržano učenje kao što su OpenAI Gym i Gymnasium, omogućavajući istraživačima da lako postavljaju i izvode RL eksperimente. Njegova sposobnost da precizno simulira dinamiku omogućava razvoj i testiranje naprednih RL algoritama u realističnim simulacijskim okruženjima.

2.6.3. Prednosti i ograničenja MuJoCo-a

Prednosti

- **Visoka preciznost:** MuJoCo omogućava vrlo precizno modeliranje i simulaciju složenih dinamičkih sustava.
- **Brza simulacija:** Optimizirani algoritmi omogućuju simulaciju u stvarnom vremenu, što je ključno za iterativne procese razvoja i testiranja.
- **Široka primjena:** Podržava različite domene, uključujući robotiku, biomehaniku i RL, što ga čini univerzalnim alatom za istraživače i inženjere.

Ograničenja

- **Složenost modeliranja:** Zbog visoke preciznosti i detaljnosti, modeliranje sustava u MuJoCo-u može biti složeno i zahtjeva dobro razumijevanje fizike i dinamike.
- **Računalni resursi:** Simulacije u MuJoCo-u mogu biti prilično zahtjevne za računalne resurse, posebno za vrlo složene modele i scenarije.

2.6.4. Zaključak

MuJoCo je moćan alat za simulaciju fizike koji omogućava precizno modeliranje i kontrolu složenih dinamičkih sustava. Njegova primjena u robotici, biomehanici i podržanom učenju čini ga ključnim alatom za istraživače i inženjere. S obzirom na njegove prednosti u preciznosti i efikasnosti, MuJoCo će nastaviti igrati važnu ulogu u napretku istraživanja i razvoja u različitim znanstvenim i industrijskim domenama.

3. Metodologija

3.1. Uvod

U ovom poglavlju opisuje se metodologija koja je korištena za treniranje i evaluaciju algoritama podržanog učenja u simulacijskim okruženjima robotske manipulacije. Primarni cilj bio je istražiti performanse različitih algoritama u zadacima s različitim vrstama nagrada i okruženjima. Korišteni algoritmi uključuju Q-learning, DQN, PPO, DDPG, SAC i kombinaciju DDPG-a s Hindsight Experience Replay (HER).

Kako bi se osigurala temeljita i usporediva analiza, implementirani su eksperimenti u Gymnasium okruženju, koristeći simulacije kao što su CartPole, MountainCar i FetchReach-v2. Korištenje simulacijskih okruženja omogućilo je kontrolirano postavljanje eksperimenta i precizno mjerenje performansi agenata.

U daljnjim odjeljcima ovog poglavlja detaljno će se opisati arhitektura sustava, specifični koraci razvoja, korišteni algoritmi, te evaluacijske metrike. Ovaj pristup osigurava sveobuhvatan pregled metodologije i omogućava razumijevanje ključnih koraka poduzetih tijekom istraživanja.

3.2. On-Policy i Off-Policy algoritmi

U podržanom učenju (Reinforcement Learning, RL), algoritmi se mogu kategorizirati kao on-policy ili off-policy, ovisno o načinu na koji koriste prikupljene podatke za ažuriranje svojih politika.

3.2.1. On-Policy Algoritmi

On-Policy algoritmi koriste podatke generirane od trenutne politike za ažuriranje te iste politike. Drugim riječima, algoritam prikuplja iskustva interakcije s okruženjem koristeći trenutnu politiku i koristi ta iskustva izravno za poboljšanje iste politike.

Primjeri:

- **PPO (Proximal Policy Optimization):** PPO je on-policy algoritam koji koristi trenutnu politiku za generiranje podataka i ažurira politiku pomoću tih podataka, balansirajući između istraživanja i eksploatacije kroz clipping funkciju.
- **A2C (Advantage Actor-Critic):** A2C prikuplja podatke iz okruženja pomoću trenutne politike (aktora) i koristi te podatke za treniranje aktera i kritičara.

Prednosti:

- **Stabilnost:** On-policy algoritmi često imaju stabilnije učenje jer izravno koriste podatke generirane trenutnom politikom.
- **Jednostavnost:** Implementacija može biti jednostavnija jer nije potrebno čuvati i ponovno koristiti stare podatke.

Nedostaci:

- **Efikasnost:** On-policy algoritmi mogu biti manje efikasni jer svaki podatak može biti iskorišten samo jednom, tijekom trenutne epizode.
- **Ovisnost o istraživanju:** Učenje može biti sporije jer algoritam mora kontinuirano istraživati s trenutnom politikom.

3.2.2. Off-Policy Algoritmi

Off-Policy algoritmi koriste podatke generirane od jedne politike (bilo trenutne ili stare) za ažuriranje druge politike. Ovi algoritmi mogu prikupljati i koristiti iskustva generirana različitim politikama, uključujući stare politike ili čak nasumične akcije.

Primjeri:

- **DQN (Deep Q-Network):** DQN koristi iskustva prikupljena iz okruženja pomoću trenutne politike, ali ažurira Q-funkciju koja može predlagati nove akcije koje nisu nužno u skladu s politikom koja je generirala podatke.
- **DDPG (Deep Deterministic Policy Gradient):** DDPG koristi replay buffer za pohranu iskustava i ponovno koristi ta iskustva za ažuriranje politike, što omogućuje efikasnije učenje.

Prednosti:

- **Efikasnost:** Off-policy algoritmi mogu efikasnije koristiti podatke jer mogu ponovno koristiti stara iskustva pohranjena u replay bufferu.
- **Fleksibilnost:** Mogu učiti iz različitih izvora podataka, uključujući stare politike, što može ubrzati proces učenja.

Nedostaci:

- **Kompleksnost:** Implementacija može biti složenija zbog potrebe za upravljanjem replay bufferom i osiguravanjem da se podaci koriste na ispravan način.
- **Stabilnost:** Postoji rizik od nestabilnosti u učenju jer podaci mogu dolaziti iz različitih politika, što može otežati konvergenciju.

3.2.3. Zaključak

On-policy algoritmi koriste trenutnu politiku za generiranje podataka i ažuriraju tu politiku izravno pomoću tih podataka, dok off-policy algoritmi koriste podatke generirane različitim politikama, uključujući stare politike, za ažuriranje trenutne politike. On-policy algoritmi su stabilniji, ali manje efikasni u korištenju podataka, dok su off-policy algoritmi fleksibilniji i efikasniji, ali mogu biti složeniji za implementaciju i skloniji nestabilnosti u učenju.

3.3. Prostor akcija i opažanja u Gymnasiumu

Gymnasium je platforma za podržano učenje koja omogućava korisnicima stvaranje i upravljanje raznim simulacijskim okruženjima. Ključni koncepti u Gymnasiumu su

prostor akcija i prostor opažanja, koji definiraju kako agenti mogu djelovati u okruženju i kako percipiraju svoje okruženje.

3.3.1. Prostor akcija (Action Space)

Prostor akcija definira sve moguće akcije koje agent može izvršiti u danom okruženju. Akcije su komandne instrukcije koje agent izdaje kako bi interagirao s okruženjem. Prostor akcija može biti diskretan ili kontinuiran.

Diskretan prostor akcija

Diskretan prostor akcija sastoji se od konačnog broja različitih akcija. Svaka akcija je indeksirana cijelim brojem. Na primjer, u okruženju CartPole, prostor akcija je diskretan i sastoji se od dvije moguće akcije:

- Guranje kolica lijevo (akcija 0)
- Guranje kolica desno (akcija 1)

Tablica 3.1. Diskretan prostor akcija za CartPole okruženje

Akcija	Opis
0	Guranje kolica lijevo
1	Guranje kolica desno

Kontinuiran prostor akcija

Kontinuiran prostor akcija sadrži beskonačan broj mogućih akcija unutar određenog raspona. Akcije su definirane stvarnim brojevima. Na primjer, u okruženju MuJoCo, prostor akcija može uključivati kontinuirane sile ili momente primijenjene na zglobove robota.

Primjer kontinuiranog prostora akcija može izgledati ovako:

Tablica 3.2. Kontinuiran prostor akcija

Akcija	Opis
F_1	Kontinuirana sila 1
F_2	Kontinuirana sila 2
\vdots	\vdots
F_n	Kontinuirana sila n

gdje su F_i kontinuirane sile primijenjene na različite zglobove.

3.3.2. Prostor opažanja (Observation Space)

Prostor opažanja definira sve moguće informacije koje agent može primiti iz okruženja. Opažanja predstavljaju trenutna stanja okruženja koja agent koristi za donošenje odluka. Prostor opažanja također može biti diskretan ili kontinuiran, i obično je višedimenzionalan.

Diskretan prostor opažanja

Diskretan prostor opažanja sastoji se od konačnog broja mogućih stanja. Svako stanje je reprezentirano cijelim brojem. Ovaj tip prostora opažanja je manje uobičajen u složenim okruženjima.

Kontinuiran prostor opažanja

Kontinuiran prostor opažanja sastoji se od realnih brojeva koji predstavljaju različite mjere stanja okruženja. Primjeri uključuju položaj, brzinu, kutove, i druge fizikalne veličine. Na primjer, u okruženju CartPole, prostor opažanja je kontinuiran i četvero-dimenzionalan:

$$\text{Stanje} = \begin{bmatrix} \text{Položaj kolica} \\ \text{Brzina kolica} \\ \text{Kut poluge} \\ \text{Kutna brzina poluge} \end{bmatrix}$$

Tablica 3.3. Kontinuiran prostor opažanja za CartPole okruženje

Opazanje	Min	Max
Položaj kolica	-4.8	4.8
Brzina kolica	$-\infty$	∞
Kut poluge	-0.418 rad (-24°)	0.418 rad (24°)
Kutna brzina poluge	$-\infty$	∞

3.3.3. Zaključak

Prostor akcija i opažanja ključni su koncepti u Gymnasiumu, jer definiraju kako agenti djeluju i percipiraju svoje okruženje. Diskretni prostori akcija i opažanja koriste cijele brojeve za definiranje konačnog skupa mogućnosti, dok kontinuirani prostori koriste stvarne brojeve za predstavljanje neograničenih mogućnosti. Razumijevanje ovih konceptata je ključno za razvoj i implementaciju agenata za podržano učenje koji mogu učinkovito upravljati i učiti u različitim simulacijskim okruženjima.

3.4. Razlika između Sparse i Dense nagrada

3.4.1. Sparse nagrade

Sparse nagrade su nagrade koje se dodjeljuju samo kada agent postigne određeni cilj ili stanje. U takvom sustavu, agent može dobiti vrlo malo ili nikakvu povratnu informaciju tijekom većeg dijela epizode, osim kada se dogodi specifični događaj koji pokreće nagradu.

Primjer:

- U robotskoj manipulaciji, agent može dobiti nagradu samo kada uspješno stavi objekt na ciljano mjesto.

Prednosti:

- **Poticanje preciznosti:** Agenti su motivirani da točno dođu do cilja, što može rezultirati razvijanjem preciznih kontrolnih politika.
- **Robusnost:** Sparse nagrade mogu potaknuti robusnije ponašanje jer agent mora naučiti učinkovite strategije za dosezanje cilja.

Nedostaci:

- **Sporo učenje:** Ograničene povratne informacije mogu usporiti proces učenja jer agent dobiva malo informacija tijekom epizode.
- **Veći izazov:** Proces učenja može biti izazovniji jer agent rijetko dobiva nagrade, što može dovesti do duljeg vremena potrebnog za konvergenciju.

3.4.2. Dense nagrade

Dense nagrade su nagrade koje se kontinuirano dodjeljuju agentu na temelju njegovih akcija ili udaljenosti od cilja. U takvom sustavu, agent dobiva konstantne povratne informacije koje mu pomažu u učenju optimalnih akcija tijekom cijele epizode.

Primjer:

- U robotskoj manipulaciji, agent može dobiti nagradu koja je proporcionalna udaljenosti objekta od ciljanog mjesta, s većim nagradama za bliže pozicije.

Prednosti:

- **Brže učenje:** Kontinuirane nagrade omogućuju agentima da stalno prilagođavaju svoje akcije kako bi maksimizirali nagrade, što rezultira bržim učenjem optimalnih politika.
- **Bolja prilagodljivost:** Agenti mogu efikasnije učiti iz svojih pogrešaka jer imaju stalne povratne informacije o svom napretku prema cilju.

Nedostaci:

- **Rizik prekomjernog prilagođavanja:** Kontinuirane povratne informacije mogu dovesti do prekomjernog prilagođavanja specifičnim putanjama ili akcijama, ograničavajući generalizaciju politike na nove situacije.
- **Veći proračunski zahtjevi:** Stalna prilagodba i optimizacija mogu zahtijevati više proračunskih resursa, posebno u složenijim okruženjima.

3.5. Q-learning algoritam

3.5.1. Osnove Q-learninga

Q-learning je algoritam za podržano učenje (reinforcement learning) koji pripada skupini "off-policy" metoda. Q-learning je primjer off-policy metode jer koristi trenutnu politiku za istraživanje (npr. epsilon-greedy) dok ažurira Q-funkciju prema optimalnoj politici, bez obzira na trenutno istraživanje. Ova fleksibilnost omogućuje učinkovitije učenje optimalne strategije. Njegova glavna svrha je omogućiti agentu da nauči optimalnu

politiku djelovanja u svakom stanju okruženja, bez obzira na to kako agent trenutno istražuje okruženje. Q-learning algoritam uči funkciju vrijednosti, nazvanu Q-funkcija, koja procjenjuje očekivanu kumulativnu nagradu za svaku kombinaciju stanja i akcije. Ova funkcija koristi se za donošenje odluka o najboljoj akciji u svakom stanju.

3.5.2. Metodologija Q-learninga

Q-learning algoritam temelji se na ažuriranju Q-funkcije pomoću Bellmanove jednadžbe. Kada agent izvrši akciju u određenom stanju i primi povratnu informaciju (nagradu) te promatra novo stanje, vrijednost Q-funkcije se ažurira prema sljedećoj formuli:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right] \quad (3.1)$$

Gdje su:

- $Q(s, a)$ trenutna procjena vrijednosti za stanje s i akciju a
- α stopa učenja (learning rate)
- r nagrada dobivena izvršenjem akcije a u stanju s
- γ faktor diskontiranja (discount factor)
- s' novo stanje nakon izvršenja akcije a
- $\max_{a'} Q(s', a')$ maksimalna procijenjena vrijednost za novo stanje s' i sve moguće akcije a' u tom stanju

Epsilon-greedy strategija

Epsilon-greedy strategija je metoda istraživanja koja omogućava agentu da balansira između istraživanja novih akcija (exploration) i iskorištavanja već naučenih akcija (exploitation). S vjerojatnošću ϵ , agent odabire nasumičnu akciju, dok s vjerojatnošću $1 - \epsilon$, agent odabire najbolju poznatu akciju prema Q-vrijednostima. Ova strategija osigurava da agent istražuje različite akcije kako bi izbjegao lokalno optimalna rješenja.

3.5.3. Prednosti Q-learninga

Q-learning ima nekoliko ključnih prednosti koje ga čine jednim od najpopularnijih RL algoritama:

- **Robusnost:** Q-learning je robustan i može se koristiti u različitim tipovima okruženja, uključujući ona s nepredvidivim promjenama.
- **Jednostavnost:** Algoritam je relativno jednostavan za implementaciju i razumijevanje, što ga čini prikladnim za početnike u području strojnog učenja.
- **Konvergencija:** Uz pravilno postavljanje hiperparametara, Q-learning jamči konvergenciju prema optimalnoj politici.

Q-learning algoritam je bio veliki napredak u podržanom učenju jer je prvi algoritam s garantiranim konvergencijama do optimalne politike [29].

Dva uvjeta moraju biti zadovoljena kako bi se osigurala konvergencija:

1. Stope učenja moraju se smanjivati, ali ne prebrzo. Formalno, to zahtijeva da suma stopa učenja divergira, ali suma njihovih kvadrata konvergira. Primjer slijeda koji ima ova svojstva je $\frac{1}{1}, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots$
2. Svaki par stanja-akcije mora biti posjećen beskonačno puta. To znači da svaka akcija mora imati nenultu vjerojatnost da će biti odabrana u svakom stanju, tj. $\pi(s, a) > 0$ za sve (s, a) . U praksi, korištenje ϵ -greedy politike (gdje je $\epsilon > 0$) osigurava da je ovaj uvjet zadovoljen.

3.5.4. Primjena Q-learninga u robotskoj manipulaciji

U kontekstu robotske manipulacije, Q-learning može biti vrlo koristan za učenje optimalnih strategija za različite zadatke, poput podizanja objekata, premještanja i slaganja. Algoritam omogućava robotima da kroz interakciju sa simuliranim okruženjem nauče učinkovite manipulacijske vještine koje se mogu primijeniti i u stvarnim situacijama. Korištenjem simulacijskih platformi poput Gymnasiuma, roboti mogu sigurno i brzo učiti, minimizirajući rizik od oštećenja stvarne opreme.

3.6. Deep Q-Network (DQN) Algoritam

3.6.1. Uvod

Deep Q-Network (DQN) je algoritam podržanog učenja koji kombinira Q-learning s dubokim neuronskim mrežama kako bi se nosio s velikim i kontinuiranim prostorima stanja. DQN je pionirski algoritam koji je pokazao izvanredne rezultate u igranju videoigara iz Atari zbirke, gdje je nadmašio ljudske performanse u mnogim igrama. Razvio ga je tim istraživača iz DeepMind-a, a njegova objava značajno je unaprijedila područje podržanog učenja [30].

3.6.2. Osnovni principi DQN-a

DQN koristi duboke neuronske mreže za aproksimaciju Q-funkcije, koja procjenjuje očekivanu vrijednost nagrade za svaku akciju u svakom stanju. Ključna inovacija DQN-a je sposobnost učenja iz sirovih podataka (kao što su pikseli igre) i donošenja odluka koje maksimiziraju dugoročne nagrade.

Q-funkcija: Q-funkcija predstavlja očekivanu kumulativnu nagradu koju agent može dobiti počevši od određenog stanja i izvršavajući određenu akciju. Formalno, Q-funkcija se definira kao:

$$Q(s, a) = \mathbb{E}[R_t \mid s_t = s, a_t = a]$$

gdje je $Q(s, a)$ vrijednost stanja s i akcije a , a $\mathbb{E}[R_t]$ očekivana nagrada.

Dugoročna nagrada: Dugoročna nagrada računa se pomoću Bellmanove jednadžbe, koja u kontekstu DQN-a ima oblik:

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

gdje je r trenutna nagrada, γ faktor diskontiranja, s' sljedeće stanje, a $\max_{a'} Q(s', a')$ maksimalna očekivana nagrada za sve moguće akcije u sljedećem stanju.

3.6.3. Struktura DQN-a

DQN koristi duboku neuronsku mrežu za aproksimaciju Q-funkcije [31]. Mreža prima stanje kao ulaz i generira Q-vrijednosti za sve moguće akcije kao izlaz. Ključne komponente DQN-a uključuju:

1. **Replay Buffer:** Ovaj buffer pohranjuje iskustva agenata (stanje, akcija, nagrada, sljedeće stanje) koja se koriste za treniranje mreže. Replay buffer omogućava ponovno korištenje starih iskustava, smanjuje korelaciju između uzoraka i poboljšava stabilnost učenja.
2. **Target mreža:** DQN koristi dvije mreže: glavnu mrežu i ciljnu mrežu. Ciljna mreža ažurira se rjeđe od glavne mreže i koristi se za izračunavanje ciljnih Q-vrijednosti, što dodatno stabilizira proces učenja.
3. **Greedy strategija istraživanja:** DQN koristi ϵ -greedy strategiju za istraživanje i eksploataciju. Agenti biraju nasumične akcije s vjerojatnošću ϵ i optimalne akcije s vjerojatnošću $1 - \epsilon$, omogućujući balans između istraživanja novog i eksploatacije naučenog.

3.6.4. Algoritam

1. **Inicijalizacija:** Inicijalizacija replay buffera, glavne mreže Q s nasumičnim težinama i ciljne mreže Q' kopiranjem težina glavne mreže.
2. **Petlja treniranja:**
 - (a) Za svaku epizodu:
 - i. Inicijalizacija stanja s .
 - ii. Za svaki korak u epizodi:
 - A. Odabir akcije a koristeći ϵ -greedy politiku.
 - B. Izvršavanje akcije a i promatranje nagrade r i sljedećeg stanja s' .
 - C. Pohrana prijelaza (s, a, r, s') u replay buffer.

D. Izrada mini-batch prijelaza iz replay buffera.

E. Izračun ciljne vrijednosti za svaki prijelaz:

$$y = \begin{cases} r & \text{ako je epizoda završena u } s' \\ r + \gamma \max_{a'} Q'(s', a') & \text{inače} \end{cases}$$

F. Ažuriranje glavne mreže minimizirajući gubitak:

$$L = \mathbb{E}[(y - Q(s, a))^2]$$

G. Povremeno kopiranje težine glavne mreže u ciljnu mrežu $Q' = Q$.

3. Ponavljanje koraka treniranja dok agent ne postigne zadovoljavajuću performansu.

3.6.5. Prednosti i izazovi DQN-a

Prednosti:

- **Efikasnost u velikim prostorima stanja:** DQN može učiti iz sirovih podataka i donositi kompleksne odluke, što ga čini prikladnim za velike i kontinuirane prostore stanja.
- **Stabilnost:** Korištenje replay buffera i ciljne mreže poboljšava stabilnost učenja.

Izazovi:

- **Osjetljivost na hiperparametre:** Performanse DQN-a mogu značajno varirati ovisno o odabranim hiperparametrima.
- **Problem nestabilnosti:** Bez pažljivog podešavanja, DQN može biti nestabilan ili konvergirati do suboptimalnih rješenja.

3.6.6. Zaključak

DQN je moćan algoritam podržanog učenja koji kombinira Q-learning s dubokim neuronskim mrežama za rješavanje problema u velikim i kontinuiranim prostorima sta-

nja. Njegova sposobnost učenja iz sirovih podataka i stabilnost čine ga prikladnim za širok spektar aplikacija, iako zahtijeva pažljivo podešavanje hiperparametara kako bi se postigle optimalne performanse.

3.7. Deep Deterministic Policy Gradient (DDPG) Algoritam

3.7.1. Uvod

Deep Deterministic Policy Gradient (DDPG) je algoritam podržanog učenja dizajniran za rad u kontinuiranim akcijskim prostorima. DDPG kombinira ideje iz Deep Q-Network (DQN) i Deterministic Policy Gradient (DPG) algoritama, omogućujući efikasno učenje politika u složenim okruženjima. Razvili su ga Lillicrap i suradnici 2015. godine, a postao je popularan zbog svoje sposobnosti rješavanja problema koji zahtijevaju preciznu kontrolu [32].

3.7.2. Osnovni principi DDPG-a

DDPG je off-policy algoritam koji koristi dvije različite mreže: politiku (aktor) i procjenu vrijednosti (kritičar). Aktorska mreža generira akcije, dok kritičarska mreža procjenjuje kvalitetu tih akcija [33].

Politika (Aktor): Aktor koristi determinističku politiku $\mu(s|\theta^\mu)$ za generiranje akcija a u danom stanju s . Cilj je naučiti optimalnu politiku koja maksimizira dugoročne nagrade.

Procjena vrijednosti (Kritičar): Kritičar procjenjuje Q-vrijednosti akcija generiranih aktorom. Q-funkcija se aproksimira pomoću kritičarske mreže $Q(s, a|\theta^Q)$, koja predviđa očekivane nagrade za danu akciju a u stanju s .

3.7.3. Struktura DDPG-a

DDPG koristi duboke neuronske mreže za aproksimaciju politika i Q-funkcija. Ključne komponente DDPG-a uključuju:

1. **Replay Buffer:** Replay buffer pohranjuje iskustva agenata (stanje, akcija, nagrada,

sljedeće stanje), koja se kasnije koriste za treniranje mreža. Replay buffer omogućava ponovno korištenje prošlih iskustava i smanjuje korelaciju između uzoraka, čime se poboljšava stabilnost učenja.

2. **Ciljne mreže:** DDPG koristi dvije dodatne mreže (ciljnu politiku i ciljnu kritičarsku mrežu) koje se ažuriraju sporije od glavnih mreža. Ciljne mreže se koriste za izračunavanje ciljeva Q-vrijednosti, što dodatno stabilizira proces učenja.
3. **Strategija istraživanja:** DDPG koristi politiku sa šumom (npr. Ornstein-Uhlenbeck proces) za istraživanje akcijskih prostora. Dodavanje šuma u akcije generirane politikom omogućava agentu da istražuje okruženje i pronalazi optimalne akcije.

3.7.4. Algoritam

1. **Inicijalizacija:** Inicijalizacija replay buffer, glavne mreže (aktor μ i kritičar Q) s nasumičnim težinama, te ciljne mreže (aktor μ' i kritičar Q') kopiranjem težina glavnih mreža. Također inicijalizacija šuma za istraživanje akcija.
2. **Petlja treniranja:**
 - (a) Za svaku epizodu:
 - i. Inicijalizacija stanja s .
 - ii. Za svaki korak u epizodi:
 - A. Odabir akcije $a = \mu(s|\theta^\mu) + \mathcal{N}_t$, gdje je \mathcal{N}_t šum za istraživanje.
 - B. Izvršavanje akcije a i promatranje nagrade r i sljedećeg stanja s' .
 - C. Pohrana prijelaza (s, a, r, s') u replay buffer.
 - D. Izrada mini-batch prijelaza iz replay buffera.
 - E. Izračun ciljne Q-vrijednosti za svaki prijelaz:

$$y = r + \gamma Q'(s', \mu'(s'|\theta^{\mu'})|\theta^Q)$$

F. Ažuriranje kritičarske mreže minimiziranjem gubitka:

$$L = \mathbb{E}[(y - Q(s, a|\theta^Q))^2]$$

G. Ažuriranje aktorske mreže koristeći gradijent politike:

$$\nabla_{\theta^\mu} J \approx \mathbb{E}[\nabla_a Q(s, a|\theta^Q) \nabla_{\theta^\mu} \mu(s|\theta^\mu)]$$

H. Povremeno ažuriranje ciljne mreže:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

3.7.5. Prednosti i izazovi DDPG-a

Prednosti:

- **Efikasnost u kontinuiranim prostorima:** DDPG je vrlo učinkovit za probleme s kontinuiranim akcijskim prostorima, gdje diskretizacija nije praktična.
- **Replay buffer:** Ponovno korištenje prošlih iskustava poboljšava stabilnost i efikasnost učenja.

Izazovi:

- **Osjetljivost na hiperparametre:** Performanse DDPG-a mogu značajno varirati ovisno o odabranim hiperparametrima, što zahtijeva pažljivo podešavanje.
- **Nestabilnost:** DDPG može biti nestabilan ili konvergirati do suboptimalnih rješenja bez pažljivog podešavanja i stabilizacijskih tehnika.

3.7.6. Zaključak

DDPG je moćan algoritam podržanog učenja koji kombinira prednosti determinističkih politika i dubokih neuronskih mreža za efikasno učenje u kontinuiranim akcijskim prostorima. Njegova sposobnost učenja iz prošlih iskustava putem replay buffera i stabilizacija korištenjem ciljnih mreža čine ga prikladnim za širok spektar primjena, iako zahtijeva pažljivo podešavanje hiperparametara kako bi se postigle optimalne performanse.

3.8. Soft Actor-Critic (SAC) Algoritam

3.8.1. Uvod

Soft Actor-Critic (SAC) je napredni algoritam podržanog učenja koji se koristi za rješavanje problema kontinuiranih akcijskih prostora. SAC kombinira ideje iz politike gradijentnih metoda i maksimizacije entropije, omogućavajući stabilno i efikasno učenje. Razvili su ga istraživači iz Berkeley AI Research (BAIR) laboratorija, a postao je popularan zbog svoje robusnosti i performansi u složenim okruženjima [34].

3.8.2. Osnovni principi SAC-a

SAC je off-policy algoritam koji koristi dvije kritičarske mreže i jednu aktorsku mrežu za učenje politika i procjenu vrijednosti akcija. Njegova ključna inovacija je uvođenje maksimizacije entropije, što pomaže u održavanju istraživanja tijekom procesa učenja [12].

Maksimizacija entropije: Entropija je mjera nesigurnosti ili slučajnosti u politici. SAC maksimizira očekivanu kumulativnu nagradu dok istovremeno maksimizira entropiju politike, što rezultira politikama koje su istraživačke i ne prebrzo konvergiraju na suboptimalna rješenja.

$$\pi^* = \arg \max_{\pi} \sum_t \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))]$$

gdje je α koeficijent koji kontrolira važnost entropije, a $\mathcal{H}(\pi(\cdot | s_t))$ entropija politike π .

3.8.3. Struktura SAC-a

SAC koristi tri glavne mreže: dvije kritičarske mreže i jednu aktorsku mrežu. Osim toga, koristi i ciljnu kritičarsku mrežu za stabilizaciju učenja.

1. **Kritičarske mreže:** Dvije kritičarske mreže $Q_1(s, a)$ i $Q_2(s, a)$ procjenjuju Q-vrijednosti kako bi se smanjila varijanca u procjeni Q-vrijednosti.
2. **Aktorska mreža:** Aktorska mreža $\pi(a|s)$ generira akcije na temelju stanja. Aktorska mreža maksimizira očekivanu kumulativnu nagradu uz dodatak entropije.
3. **Ciljna kritičarska mreža:** Ciljna kritičarska mreža \hat{Q} koristi se za izračunavanje ciljnih Q-vrijednosti, što pomaže stabilizaciji učenja.

3.8.4. Algoritam

1. **Inicijalizacija:** Inicijalizacija replay buffera, dvije kritičarske mreže Q_1 i Q_2 s nasumičnim težinama, aktorske mreže π s nasumičnim težinama, te ciljnu kritičarsku mrežu \hat{Q} kopiranjem težina glavne kritičarske mreže.
2. **Petlja treniranja:**
 - (a) Za svaku epizodu:
 - i. Inicijalizacija stanje s .
 - ii. Za svaki korak u epizodi:
 - A. Odabir akcije $a \sim \pi(\cdot|s)$ koristeći aktorsku mrežu.
 - B. Izvršavanje akcije a i promatranje nagrade r i sljedećeg stanja s' .
 - C. Pohrana prijelaza (s, a, r, s') u replay buffer.
 - D. Izrada mini-batch prijelaza iz replay buffera.
 - E. Izračun ciljne Q-vrijednost koristeći ciljnu kritičarsku mrežu:

$$y = r + \gamma(\min_{i=1,2} Q_i(s', a') - \alpha \log \pi(a'|s'))$$

F. Ažuriranje kritičarske mreže minimiziranjem gubitka:

$$L_i = \mathbb{E}[(y - Q_i(s, a))^2] \quad \text{za } i = 1, 2$$

G. Ažuriranje aktorske mreže maksimiziranjem ciljane funkcije:

$$J_\pi = \mathbb{E}[\alpha \log \pi(a|s) - Q(s, a)]$$

H. Povremeno ažuriranje ciljane kritičarske mreže:

$$\theta^{\hat{Q}} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{\hat{Q}}$$

3. **Ponavljjanje:** Ponavljjanje koraka treniranja dok agent ne postigne zadovoljavajuće performanse.

3.8.5. Prednosti i izazovi SAC-a

Prednosti:

- **Stabilnost i efikasnost:** Korištenje dvostrukih kritičara i maksimizacije entropije omogućuje stabilno i efikasno učenje.
- **Brzo prilagođavanje:** Kontinuirane povratne informacije omogućuju agentima brzo učenje optimalnih politika.
- **Dobro se nosi s dense nagradama:** Kontinuirane povratne informacije omogućuju efikasnije učenje.

Izazovi:

- **Kompleksnost:** Implementacija može biti složenija zbog korištenja dvostrukih kritičara i entropije regulacije.
- **Resursno intenzivan:** Zahtijeva više proračunskih resursa za treniranje zbog složenih izračuna.

3.8.6. Zaključak

Soft Actor-Critic (SAC) je moćan algoritam podržanog učenja koji kombinira prednosti politike gradijentnih metoda i maksimizacije entropije za efikasno učenje u kontinuiranim akcijskim prostorima. Njegova sposobnost održavanja istraživanja tijekom učenja i korištenje dvostrukih kritičara za stabilizaciju čini ga vrlo efikasnim i robusnim za širok spektar primjena. SAC je posebno pogodan za okruženja s bogatim povratnim informacijama, iako zahtijeva pažljivo upravljanje proračunskim resursima i hiperparametrima.

4. Eksperimentalna evaluacija

U ovom poglavlju detaljno opisujemo eksperimentalni dio rada gdje smo trenirali agente primjenom nekih od sljedećih algoritama: Q-learning, PPO, DQN, SAC, DDPG, za tri različita okruženja: Cart-Pole, Mountain-Car i Fetch-Reach. Eksperimenti su izvedeni kako bi se procijenila učinkovitost korištenih algoritama u različitim zadacima manipulacije.

4.1. Cart-Pole

Cart-Pole je klasično okruženje gdje agent (kolica) mora balansirati polugu na svojim kolicima. Cilj je održati polugu uspravnom što je dulje moguće, pomicanjem kolica lijevo ili desno. Epizoda počinje s polugom u nasumičnom početnom položaju i završava kada poluga padne ili kada agent uspješno balansira polugu za unaprijed određeni broj koraka. Ovaj zadatak zahtijeva finu kontrolu i brze reakcije kako bi se spriječilo padanje poluge. Naglasak je na učenju strategije koja maksimizira nagradu, tj. održavanje poluge uspravnom što je dulje moguće.

Prostor akcija je diskretan i jednodimenzionalan, te se sastoji od dvije akcije, guranja kolica u lijevo ili guranja kolica u desno (tablica 4.1.).

Tablica 4.1. Prostor akcije za okruženje Cart-Pole

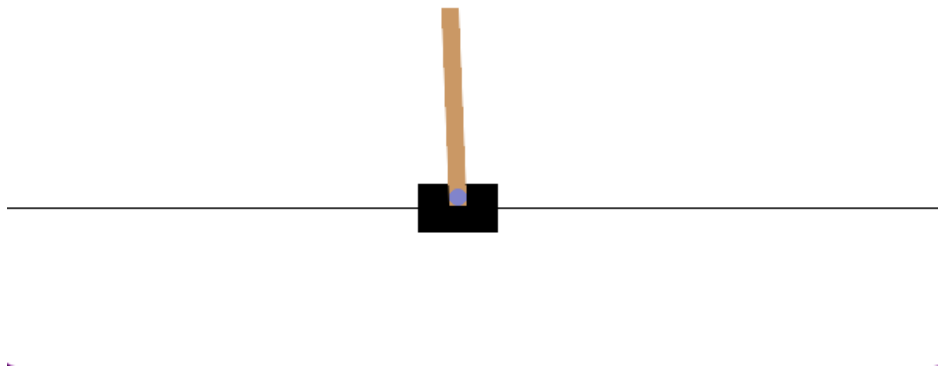
Akcija	Vrijednost 1	Vrijednost 2
Guranje kolica	0 (lijevo)	1 (desno)

Prostor opažanja je kontinuiran i četverodimenzionalan. Njegove dimenzije su slijedno: položaj kolica na x osi, brzina kolica, kut poluge i kutna brzina poluge. Njihova ograničenja su definirana u tablici 4.2.

$$\text{Stanje} = \begin{bmatrix} \text{Položaj kolica} \\ \text{Brzina kolica} \\ \text{Kut poluge} \\ \text{Kutna brzina poluge} \end{bmatrix}$$

Tablica 4.2. Prostor opažanja za okruženje Cart-Pole

Akcija	Min	Max
Položaj kolica	-4.8	4.8
Brzina kolica	$-\infty$	∞
Kut poluge	$\sim -0.418 \text{ rad } (-24^\circ)$	$\sim 0.418 \text{ rad } (24^\circ)$
Kutna brzina poluge	$-\infty$	∞



Slika 4.1. Primjer Cart-Pole okruženja

Primjer Cart-Pole okruženja možemo vidjeti na slici 4.1.

4.1.1. Vlastita implementacija

Za potrebe jednostavnog prototipiranja, napravljena je vlastita implementacija Q-learning algoritma gdje smo postavili sljedeće hiperparametre za trening:

- Stopa učenja (α): 0.1
- Faktor diskontiranja (γ): 0.99
- Početni epsilon (ϵ): 1

- Stopa opadanja epsilon: 0.00001

Diskretizirali smo svaku dimenziju iz prostora opažanja u 10 segmenata, zbog ograničenja Q-learning algoritma da radi samo nad diskretnim prostorima. Algoritam je treniran kroz kontinuirani broj epizoda, gdje je agent pokušavao balansirati polugu. Politika djelovanja je ažurirana na temelju nagrada dobivenih kroz interakciju s okruženjem.

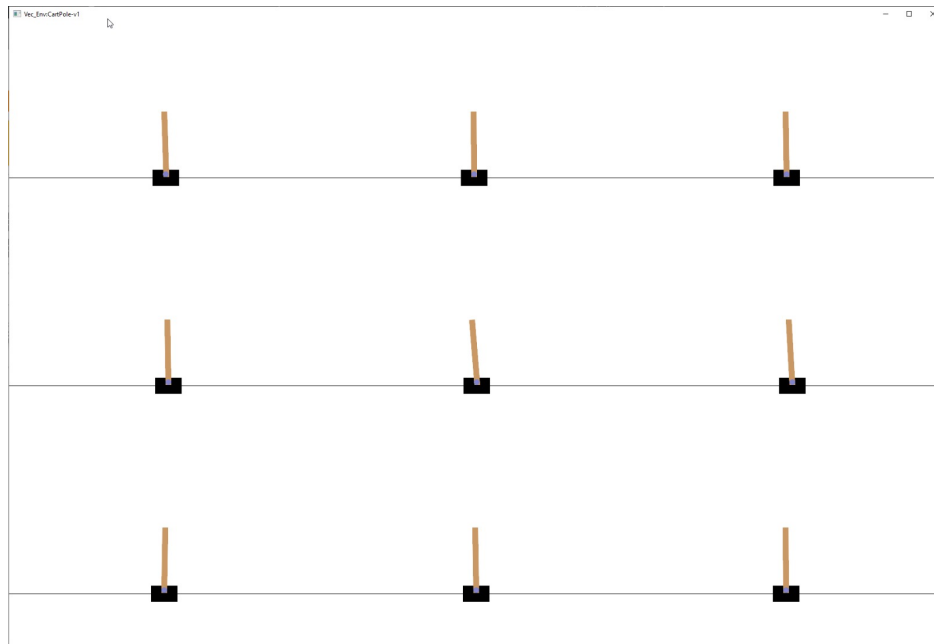
Korištena je epsilon-greedy strategija za istraživanje i eksploataciju. Trening je trajao dok agent nije pokazao stabilne performanse u održavanju poluge uspravnom. Evaluacija je provedena mjerenjem prosječne nagrade po epizodi. Algoritam je pokazao uspješno balansiranje poluge kroz nekoliko stotina epizoda treninga. Rezultati su pokazali da je agent naučio optimalnu politiku koja omogućava učinkovito balansiranje poluge na kolicima. Prosječna nagrada po epizodi je također značajno povećana tijekom treninga, što ukazuje na uspješno učenje. Sljedeća tablica prikazuje prosječne nagrade za svaku tisućitu vremensku jedinicu:

Tablica 4.3. Prosječne nagrade po vremenskoj jedinici za Cart-Pole eksperiment

Vremenska jedinica	Prosječna nagrada
1000	100
2000	130
3000	151
4000	168
5000	184

4.1.2. OpenRL implementacija

Jedna od glavnih snaga OpenRL biblioteke je njezina sposobnost korištenja paralelizma za ubrzanje procesa treniranja. U ovom eksperimentu koristili smo 9 paralelnih instanci CartPole okruženja kako bismo istovremeno trenirali PPO i DQN algoritme. Cilj je bio demonstrirati kako korištenje paralelizma može poboljšati efikasnost treniranja i smanjiti ukupno vrijeme potrebno za postizanje optimalnih politika.



Slika 4.2. Primjer paralelnog Cart-Pole okruženja

Eksperiment je izveden u dva koraka:

1. **Kratko treniranje (2000 vremenskih jedinica):** Trenirali smo PPO i DQN algoritme s 2000 vremenskih jedinica kako bismo brzo procijenili početne performanse.
2. **Dugoročno treniranje (20000 vremenskih jedinica):** Zatim smo proširili treniranje na 20000 vremenskih jedinica kako bismo analizirali dugoročne performanse algoritama.

Algoritmi i konfiguracija

1. **PPO (Proximal Policy Optimization)** PPO je on-policy algoritam koji koristi clipping funkciju za održavanje stabilnog učenja. PPO je treniran s 2000 i 20000 vremenskih jedinica.
2. **DQN (Deep Q-Network)** DQN je off-policy algoritam koji koristi duboke neuronske mreže za aproksimaciju Q-funkcije. DQN je također treniran s 2000 i 20000 vremenskih jedinica.

4.1.3. Rezultati

Rezultati kratkog treniranja (2000 vremenskih jedinica)

U kratkom periodu treniranja, PPO je pokazao bolje inicijalne performanse u usporedbi s DQN. PPO je brzo naučio održavati stup u uspravnom položaju, dok je DQN imao problema s postizanjem stabilne politike unutar ograničenog broja vremenskih jedinica.

Tablica 4.4. Rezultati kratkog treniranja (2000 vremenskih jedinica)

Algoritam	Prosječna nagrada	Uspješnost
PPO	150	75%
DQN	100	50%

Rezultati dugoročnog treniranja (20000 vremenskih jedinica)

Nakon produženog treniranja, PPO je nastavio pokazivati superiorne performanse, dok DQN nije uspio postići stabilne rezultate. PPO je uspješno optimizirao svoju politiku, postigavši visoku prosječnu nagradu i uspješnost. DQN je, s druge strane, pokazao znakove pretreniravanja i nestabilnosti, što je rezultiralo nižim ukupnim performansama.

Tablica 4.5. Rezultati dugoročnog treniranja (20000 vremenskih jedinica)

Algoritam	Prosječna nagrada	Uspješnost
PPO	195	98%
DQN	120	60%

Diskusija

Rezultati eksperimenata pokazuju da PPO algoritam ima značajne prednosti u učenju zadataka balansiranja u CartPole okruženju, posebno kada se koristi paralelizam. PPO je pokazao sposobnost brzog prilagođavanja i postizanja visokih performansi unutar kratkog vremenskog okvira, kao i održavanje stabilnosti tijekom dugoročnog treniranja. S druge strane, DQN algoritam se pokazao manje stabilnim i manje učinkovit u učenju optimalnih politika u istom okruženju, što je bilo posebno vidljivo tijekom dugoročnog treniranja.

Korištenje 9 paralelnih instanci okruženja značajno je ubrzalo proces treniranja, omogućujući brže postizanje rezultata i bolje iskorištavanje računalnih resursa. Ovo također demonstrira snagu OpenRL biblioteke u omogućavanju efikasnog paralelnog treniranja RL agenata.

Ovi rezultati naglašavaju važnost odabira odgovarajućeg RL algoritma i korištenja paralelizma za optimizaciju procesa treniranja, posebno u složenim simulacijskim okruženjima. Daljnja istraživanja mogu se usmjeriti na optimizaciju hiperparametara i kombinaciju različitih algoritama kako bi se postigle još bolje performanse.

4.2. Mountain-Car

MountainCar je klasično okruženje koje se često koristi za evaluaciju algoritama podržanog učenja (RL). U ovom okruženju, agent upravlja automobilom koji se nalazi u dolini između dvije planine. Cilj agenta je dosegnuti ciljnu zastavicu smještenu na vrhu desne planine. Zbog ograničene snage motora automobila, agent mora koristiti inerciju kretanja kako bi svladao uspon i dosegnuo cilj. Ovaj zadatak zahtijeva planiranje i optimizaciju akcija kako bi se maksimalno iskoristila inercija automobila.

Prostor akcija je diskretan i jednodimenzionalan, te se sastoji od tri akcije: pomicanje automobila lijevo, pomicanje automobila desno ili zadržavanje trenutne pozicije (tablica 4.6.).

Tablica 4.6. Prostor akcije za okruženje MountainCar

Akcija	Vrijednost 1	Vrijednost 2	Vrijednost 3
Pomicanje automobila	0 (lijevo)	1 (zadržavanje)	2 (desno)

Prostor opažanja je kontinuiran i dvodimenzionalan. Njegove dimenzije su: položaj automobila na x osi i brzina automobila. Ograničenja prostora opažanja prikazana su u tablici 4.7.

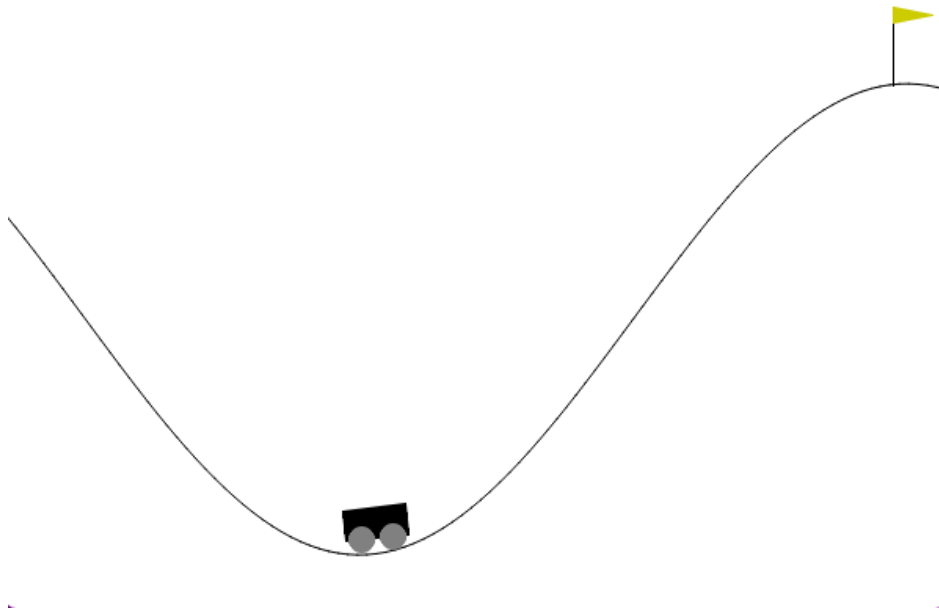
$$\text{Stanje} = \begin{bmatrix} \text{Položaj automobila} \\ \text{Brzina automobila} \end{bmatrix}$$

Tablica 4.7. Prostor opažanja za okruženje MountainCar

Opazanje	Min	Max
Položaj automobila	-1.2	0.6
Brzina automobila	-0.07	0.07

Cilj agenta je dosegnuti ciljnu zastavicu što je prije moguće. Agent dobiva nagradu od -1 za svaki korak sve dok ne postigne cilj, što znači da je nagrada negativna dok agent ne završi zadatak. Epizoda završava kada agent dosegne ciljnu zastavicu ili nakon 200 koraka.

Primjer MountainCar okruženja možemo vidjeti na slici 4.3.



Slika 4.3. Primjer Mountain-Car okruženja

4.2.1. Vlastita implementacija

Za potrebe jednostavnog prototipiranja, napravljena je vlastita implementacija Q-learning algoritma gdje smo postavili sljedeće hiperparametre za trening:

1. **Definiranje hiperparametara:** Postavili smo sljedeće hiperparametre za trening:
 - Stopa učenja (α): 0.9
 - Faktor diskontiranja (γ): 0.9
 - Početni epsilon (ϵ): 1
 - Stopa opadanja epsilon: 2/epizoda
2. **Prostori stanja:** Diskretizirali smo prostorne i brzinske podatke u 20 segmenata.

3. **Trening:** Q-learning algoritam je treniran kroz 500 vremenskih jedinica, gdje je automobil pokušavao savladati brdo. Politika djelovanja je ažurirana na temelju nagrada dobivenih kroz interakciju s okruženjem. Algoritam je koristio epsilon-greedy strategiju za istraživanje, gdje je na početku treninga epsilon bio visok kako bi se potaknulo istraživanje, a postupno je smanjivan kako bi se potaknula eksploatacija naučene politike.
4. **Evaluacija:** Evaluacija je provedena mjerenjem prosječne nagrade po epizodi. Algoritam je pokazao uspješno savladavanje zadatka kroz nekoliko stotina epizoda treninga. Rezultati su pokazali da je agent naučio optimalnu politiku koja omogućava automobilu da učinkovito savlada brdo.

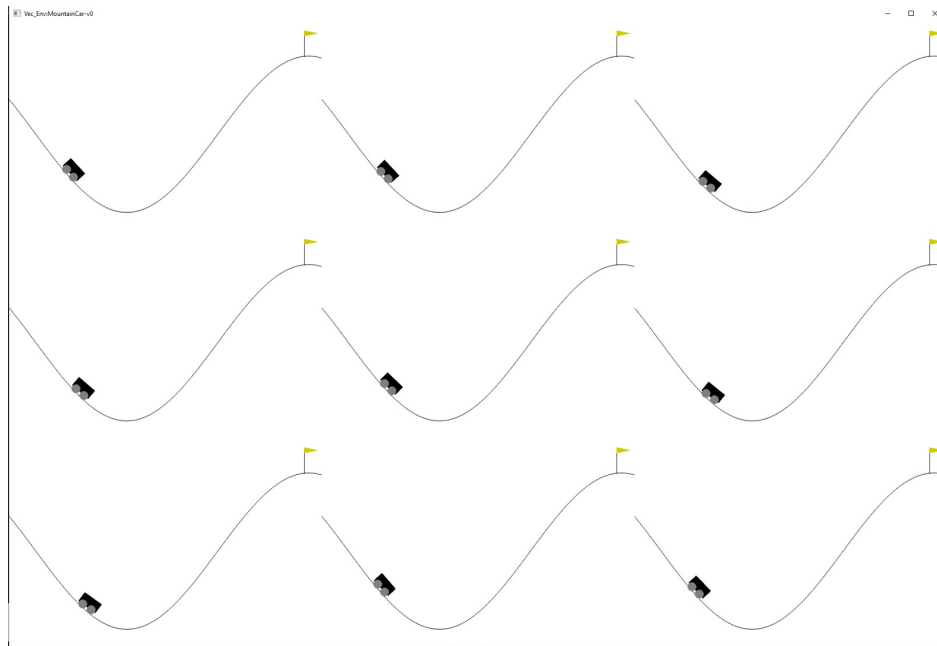
Algoritam je uspio naučiti automobil kako da koristi svoju inerciju za savladavanje brda. Prosječna nagrada po epizodi se značajno povećala tijekom treninga, što ukazuje na uspješno učenje. Sljedeća tablica prikazuje prosječne nagrade za svaku stotu vremensku jedinicu:

Vremenska jedinica	Prosječna nagrada
100	-150
200	-120
300	-90
400	-60
500	-30

Tablica 4.8. Prosječne nagrade po epizodi za Mountain-Car eksperiment

4.2.2. OpenRL implementacija

Jedna od ključnih prednosti OpenRL biblioteke je njena mogućnost korištenja paralelizma za ubrzanje procesa treniranja. U ovom eksperimentu, koristili smo 9 paralelnih instanci MountainCar okruženja kako bismo istovremeno trenirali PPO i DQN algoritme. Cilj je bio pokazati kako upotreba paralelizma može poboljšati efikasnost treniranja i smanjiti ukupno vrijeme potrebno za postizanje optimalnih politika.



Slika 4.4. Primjer paralelnog Mountain-Car okruženja

Eksperiment je izveden u dva koraka:

1. **Kratko treniranje (200 vremenskih jedinica):** Trenirali smo PPO i DQN algoritme s 200 vremenskih jedinica kako bismo brzo procijenili početne performanse.
2. **Dugoročno treniranje (2000 vremenskih jedinica):** Zatim smo proširili treniranje na 2000 vremenskih jedinica kako bismo analizirali dugoročne performanse algoritama.

Algoritmi i konfiguracija

1. **PPO (Proximal Policy Optimization)** PPO je on-policy algoritam koji koristi clipping funkciju za održavanje stabilnog učenja. PPO je treniran s 200 i 2000 vremenskih jedinica.
2. **DQN (Deep Q-Network)** DQN je off-policy algoritam koji koristi duboke neuronske mreže za aproksimaciju Q-funkcije. DQN je također treniran s 200 i 2000 vremenskih jedinica.

4.2.3. Rezultati

Rezultati kratkog treniranja (200 vremenskih jedinica)

U kratkom periodu treniranja, DQN je pokazao bolje inicijalne performanse u usporedbi s PPO. DQN je brzo naučio kako koristiti inerciju automobila za postizanje cilja, dok je PPO imao poteškoća s učenjem zbog ograničenih povratnih informacija unutar kratkog broja vremenskih jedinica.

Tablica 4.9. Rezultati kratkog treniranja (200 vremenskih jedinica)

Algoritam	Prosječna nagrada	Uspješnost
PPO	-150	30%
DQN	-110	60%

Rezultati dugoročnog treniranja (2000 vremenskih jedinica)

Nakon produženog treniranja, DQN je nastavio pokazivati superiorne performanse, dok PPO nije uspio postići stabilne rezultate. DQN je uspješno optimizirao svoju politiku, postigavši visoku prosječnu nagradu i uspješnost. PPO je imao problema s učenjem zbog specifičnosti MountainCar okruženja, što je rezultiralo nižim ukupnim performansama.

Tablica 4.10. Rezultati dugoročnog treniranja (2000 vremenskih jedinica)

Algoritam	Prosječna nagrada	Uspješnost
PPO	-130	40%
DQN	-25	97%

Diskusija

Rezultati eksperimenata pokazuju da DQN algoritam ima značajne prednosti u učenju zadataka u MountainCar okruženju, posebno kada se koristi paralelizam. DQN je pokazao sposobnost brzog prilagođavanja i postizanja visokih performansi unutar kratkog vremenskog okvira, kao i održavanje stabilnosti tijekom dugoročnog treniranja. S druge strane, PPO algoritam se pokazao manje učinkovitim zbog prirode MountainCar okruženja.

PPO je on-policy algoritam koji ažurira politiku nakon svake epizode i odbacuje prethodne podatke. U MountainCar okruženju, dosezanje cilja nasumičnim akcijama je rijetko. Kada se cilj konačno dosegne, malo je vjerojatno da će jedno ažuriranje politike

biti dovoljno za dosljedno postizanje cilja. Stoga, PPO često ostaje bez signala za učenje dok ponovno slučajno ne postigne cilj, što otežava proces učenja.

DQN je off-policy algoritam koji koristi replay buffer za pohranu i ponovno korištenje iskustava, što omogućava efikasnije učenje iz rijetkih događaja kao što je dosezanje cilja u MountainCar okruženju. Ova sposobnost omogućila je DQN-u da brže nauči optimalne politike i održava stabilne performanse tijekom dugoročnog treniranja.

Korištenje 9 paralelnih instanci okruženja značajno je ubrzalo proces treniranja, omogućujući brže postizanje rezultata i bolje iskorištavanje računalnih resursa. Ovo također demonstrira snagu OpenRL biblioteke u omogućavanju efikasnog paralelnog treniranja RL agenata.

4.3. FetchReach-v2 eksperiment

U ovom poglavlju opisujemo eksperimentalni dio rada, u kojem smo koristili Gymnasium okruženje FetchReach-v2 za treniranje agenata s različitim RL algoritmima. Eksperimenti su izvedeni kako bi se procijenila učinkovitost algoritama PPO (Proximal Policy Optimization), DDPG (Deep Deterministic Policy Gradient), SAC (Soft Actor-Critic) i DDPG s HER (Hindsight Experience Replay) u okruženju sa sparse i dense nagradama.

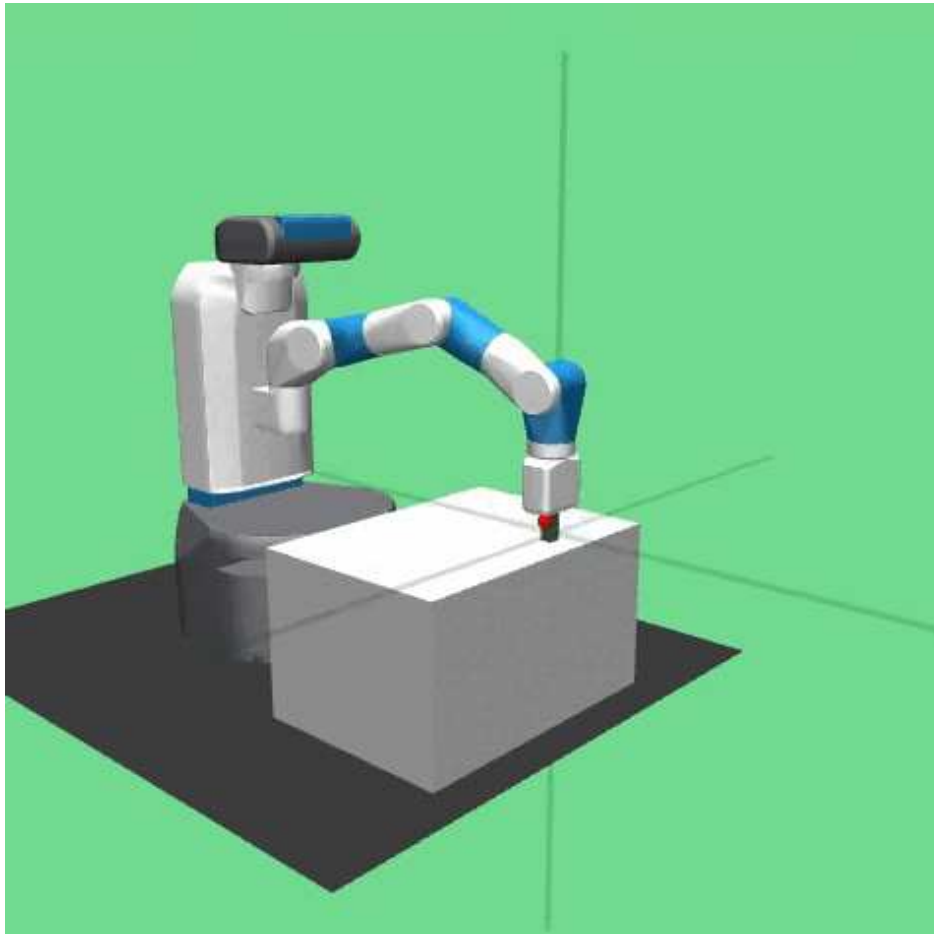
FetchReach-v2 je simulacijsko okruženje unutar Farama Robotics platforme, dizajnirano za istraživanje i evaluaciju algoritama podržanog učenja (RL) u kontekstu robotske manipulacije. Ovo okruženje je dio Fetch robotičke serije i često se koristi zbog svoje sposobnosti da jednostavno modelira osnovne zadatke manipulacije, omogućavajući istraživačima da testiraju različite RL algoritme u kontroliranim uvjetima.

4.3.1. Opis okruženja

U FetchReach-v2, agentov zadatak je upravljanje robotskom rukom kako bi dohvatila nasumično generiranu ciljnu točku u trodimenzionalnom prostoru. Robotska ruka je montirana na stolu i sastoji se od više zglobova koji omogućuju složene pokrete. Ciljna točka se generira nasumično unutar definiranog radijusa od početne pozicije ruke, osiguravajući različite scenarije za svaki pokušaj.

Detalji o robotskom modelu

Robotska ruka u FetchReach-v2 okruženju ima šest stupnjeva slobode (DoF - Degrees of Freedom), što joj omogućava visoku fleksibilnost u kretanju. Model robota je inspiriran stvarnim Fetch robotom, koji je široko korišten u istraživanjima robotske manipulacije. Svaki zglob robotske ruke može se pojedinačno kontrolirati, što omogućava precizno upravljanje položajem krajnjeg efektora.



Slika 4.5. Primjer Fetch-Reach okruženja

- **Krajnji efektor:** Kraj robotske ruke koji se koristi za dohvatanje cilja. Efektor može izvršavati translacijske i rotacijske pokrete kako bi dosegao ciljno mjesto.
- **Zglobovi:** Šest zglobova omogućava složene pokrete, uključujući savijanje, okretanje i produživanje robotske ruke. Svaki zglob ima određene ograničene domete kretanja kako bi se simuliralo stvarno fizičko ograničenje.

Akcije i stanje

- **Akcijski prostor:** Akcijski prostor u FetchReach-v2 okruženju je kontinuiran. Agenti mogu izdavati naredbe za pomicanje krajnjeg efektora u bilo kojem smjeru unutar radnog volumena robota. Specifične akcije uključuju translacijske pokrete (pomicanje efektora naprijed-nazad, gore-dolje, lijevo-desno) i rotacijske pokrete (rotacija oko x, y, i z osi).

$$\text{Akcija} = \begin{bmatrix} \Delta x \\ \Delta y \\ \Delta z \\ \Delta\theta_x \\ \Delta\theta_y \\ \Delta\theta_z \end{bmatrix}$$

- **Stanje:** Stanje okruženja uključuje trenutne koordinate krajnjeg efektora, ciljne koordinate, kao i informacije o brzini i ubrzanju efektora. Stanje se kontinuirano ažurira kako bi odražavalo promjene koje nastaju kao rezultat akcija agenta.

$$\text{Stanje} = \begin{bmatrix} \text{Pozicija efektora} \\ \text{Ciljna pozicija} \\ \text{Brzina efektora} \\ \text{Ubrzanje efektora} \end{bmatrix}$$

Nagrade

FetchReach-v2 podržava dva tipa nagrada koje značajno utječu na način učenja agenta: sparse nagrade i dense nagrade.

- **Sparse nagrade:** U ovom režimu, agent dobiva nagradu samo kada krajnji efektor točno dosegne ciljno mjesto. Ako efektor nije dosegnuo cilj, agent dobiva nagradu -1, a ako je efektor na cilju, agent dobiva nagradu 0. Smatra se da je robot dosegnuo cilj ako je Euklidska udaljenost između krajnjeg efektora i cilja manja od 0.05 m.

$$\text{Nagrada} = \begin{cases} 0 & \text{ako } d < 0.05 \\ -1 & \text{ako } d \geq 0.05 \end{cases}$$

gdje je d Euklidska udaljenost između efektora i cilja.

Prednosti:

- Poticanje preciznosti: Agenti su motivirani da točno dođu do cilja, što rezultira razvijanjem preciznih kontrolnih politika.
- Robusnost: Sparse nagrade mogu potaknuti robusnije ponašanje jer agent mora naučiti učinkovite strategije za doseganje cilja.

Mane:

- Sporo učenje: Ograničene povratne informacije mogu usporiti proces učenja jer agent dobiva malo informacija tijekom epizode.
 - Veći izazov: Proces učenja može biti izazovniji jer agent rijetko dobiva nagrade, što može dovesti do duljeg vremena potrebnog za konvergenciju.
- **Dense nagrade:** U ovom režimu, agent dobiva nagrade kontinuirano, proporcionalne blizini krajnjeg efektora cilju. Povratna nagrada je negativna Euklidska udaljenost između postignute ciljne pozicije i željene ciljne pozicije.

$$\text{Nagrada} = -d$$

gdje je d Euklidska udaljenost između efektora i cilja.

Prednosti:

- Brže učenje: Kontinuirane nagrade omogućuju agentima stalno prilagođavanje svojih akcija kako bi maksimizirali nagrade, što rezultira bržim učenjem optimalnih politika.
- Bolja prilagodljivost: Agenti mogu efikasnije učiti iz svojih pogrešaka jer imaju

stalne povratne informacije o svom napretku prema cilju.

Mane:

- Rizik prekomjernog prilagođavanja: Kontinuirane povratne informacije mogu dovesti do prekomjernog prilagođavanja specifičnim putanjama ili akcijama, ograničavajući generalizaciju politike na nove situacije.
- Veći proračunski zahtjevi: Stalna prilagodba i optimizacija mogu zahtijevati više proračunskih resursa, posebno u složenijim okruženjima.

Okoliš i izazovi

- **Okoliš:** Okruženje FetchReach-v2 je simulacija unutar kontroliranog prostora koji replicira radne uvjete stvarne robotske ruke. Okoliš je dizajniran da bude jednostavan, ali dovoljno složen da testira sposobnosti agenta u učenju manipulacijskih zadataka.
- **Izazovi:** Glavni izazov u FetchReach-v2 je precizno upravljanje robotskom rukom kako bi se doseglo ciljno mjesto. Sparse nagrade dodatno otežavaju ovaj zadatak jer agent dobiva povratnu informaciju samo pri uspjehu. Dense nagrade, iako olakšavaju proces učenja, mogu dovesti do prekomjernog prilagođavanja specifičnim putanjama.

FetchReach-v2 okruženje je idealno za istraživanje različitih RL algoritama i tehnika u kontekstu robotske manipulacije, pružajući vrijedne uvide u performanse agenata u uvjetima različitih povratnih informacija i izazova. Ova okruženja omogućuju istraživačima da testiraju različite pristupe i optimizacije kako bi razvili robusne i učinkovite strategije za robotsku manipulaciju.

4.3.2. Algoritmi korišteni u eksperimentu

U eksperimentu smo koristili četiri različita RL algoritma:

- **Proximal Policy Optimization (PPO):** On-policy algoritam koji koristi clipping funkciju za održavanje stabilnog učenja i balansiranje između istraživanja i eksploatacije.

- **Soft Actor-Critic (SAC):** Off-policy algoritam koji koristi dvostruke kritičare za stabilizaciju učenja i entropy regulaciju za balansiranje između istraživanja i eksploatacije.
- **Deep Deterministic Policy Gradient (DDPG):** Off-policy algoritam dizajniran za kontinuirane akcijske prostore, koji koristi dvije neuronske mreže: jednu za politiku i jednu za procjenu vrijednosti.
- **DDPG s Hindsight Experience Replay (HER):** Kombinacija DDPG algoritma s HER tehnikom, koja omogućava agentima da uče iz neuspjelih epizoda, pretvarajući ih u korisne informacije za učenje.

4.3.3. Eksperimentalni postupak

Eksperimentalni postupak uključivao je nekoliko koraka kako bi se osiguralo da su svi aspekti treninga i evaluacije pažljivo planirani i izvedeni. Svaki korak je detaljno opisan u nastavku.

Postavljanje okruženja

Za potrebe eksperimenta, koristili smo Gymnasium platformu za postavljanje FetchReach-v2 okruženja. Okruženje je konfigurirano u dva moda: sa sparse nagradama i s dense nagradama. Cilj svakog eksperimenta bio je da robotska ruka dosegne ciljno mjesto u trodimenzionalnom prostoru.

Konfiguracija algoritama

Svaki RL algoritam je pažljivo konfiguriran kako bi se optimizirale performanse u FetchReach-v2 okruženju. Postavljeni su ključni hiperparametri kao što su stopa učenja, veličina batch-a, broj epizoda treniranja i faktori diskontiranja.

- **PPO (Proximal Policy Optimization):**
 - Stopa učenja: 0.0003
 - Veličina batch-a: 256
 - Faktor diskontiranja: 0.99

- Broj vremenskih jedinica: 20000
- **DDPG (Deep Deterministic Policy Gradient):**
 - Stopa učenja: 0.001
 - Veličina replay buffer-a: 1,000,000
 - Faktor diskontiranja: 0.99
 - Broj vremenskih jedinica: 20000
- **SAC (Soft Actor-Critic):**
 - Stopa učenja: 0.0003
 - Entropy koeficijent: automatsko podešavanje
 - Veličina replay buffer-a: 1,000,000
 - Faktor diskontiranja: 0.99
 - Broj vremenskih jedinica: 20000
- **DDPG s HER (Hindsight Experience Replay):**
 - Stopa učenja: 0.001
 - Veličina replay buffer-a: 1,000,000
 - Faktor diskontiranja: 0.99
 - Korištenje HER: 4 buduće ciljne točke po epizodi
 - Broj vremenskih jedinica: 200000

Treniranje agenata

Svaki agent (DDPG, PPO, SAC) je treniran u FetchReach-v2 okruženju kroz 20000 jedinica vremena za oba moda nagrada (sparse i dense). Dok je DDPG + HER bio treniran na 200_000 jedinica vremena. Tijekom treniranja, agenti su učili politike djelovanja koje maksimiziraju kumulativne nagrade:

- **Sparse nagrade:** Sparse nagrade potiču agenta da istražuje i otkriva optimalne akcije koje vode do cilja. Također smanjuju rizik od prekomjernog prilagođavanja na specifične putanje ili akcije, jer agent dobiva povratnu informaciju samo za postizanje cilja. S druge strane, zbog ograničenih povratnih informacija, agenti često trebaju više epizoda kako bi naučili optimalne politike. Nadalje, agenti su imali manje povratnih informacija, što je činilo proces učenja izazovnijim.
- **Dense nagrade:** Kontinuirane nagrade omogućile su agentima brže i efikasnije učenje optimalnih politika. Agenti mogu efikasnije učiti iz svojih pogrešaka jer imaju stalne povratne informacije o svom napretku prema cilju, što poboljšava prilagodljivost i optimizaciju. Ali postoji i rizik prekomjernog prilagođavanja. Kontinuirane povratne informacije mogu dovesti do prekomjernog prilagođavanja specifičnim putanjama ili akcijama, što može ograničiti generalizaciju politike na nove situacije.

4.3.4. Evaluacija performansi agenata

Nakon završenog treniranja svakog agenta, njihova učinkovitost procijenjena je kroz dodatne epizode u istom okruženju. Da bismo osigurali statističku značajnost rezultata, svaka konfiguracija algoritma evaluirana je kroz 10 dodatnih epizoda. Ova metoda omogućila nam je da dobijemo pouzdane podatke o performansama svakog algoritma.

Evaluacija je uključivala dva ključna pokazatelja performansi:

- **Prosječna kumulativna nagrada:** Ovaj pokazatelj predstavlja srednju vrijednost kumulativnih nagrada koje je agent postigao kroz 10 evaluacijskih epizoda. Kumulativna nagrada je zbroj svih nagrada koje agent dobije tijekom jedne epizode. Viša prosječna kumulativna nagrada ukazuje na to da je agent uspješnije izvršavao zadatak i bio učinkovitiji u dosezanju cilja.

$$\text{Prosječna kumulativna nagrada} = \frac{\sum_{i=1}^{10} \text{Kumulativna nagrada}_i}{10}$$

- **Uspješnost:** Uspješnost je definirana kao postotak epizoda u kojima je agent uspješno dosegno ciljno mjesto. Smatra se da je agent dosegno cilj ako je Euklidska

udaljenost između krajnjeg efektora i cilja manja od 0.05 m. Visoka uspješnost ukazuje na to da je agent konzistentno uspijevaio dosegnuti cilj tijekom evaluacijskih epizoda.

$$\text{Uspješnost(\%)} = \left(\frac{\text{Broj uspješnih epizoda}}{10} \right) \times 100$$

4.3.5. Rezultati

Rezultati eksperimenata su analizirani kako bi se razumjelo kako različiti RL algoritmi funkcioniraju u FetchReach-v2 okruženju sa sparse i dense nagradama. Mjerenja su uključivala prosječnu kumulativnu nagradu i uspješnost izvršavanja zadatka za svakog agenta. Rezultati su prikazani u tablicama za svaki set nagrada (sparse i dense).

Rezultati sa Sparse nagradama

Algoritmi su pokazali različite stupnjeve uspješnosti u okruženju sa sparse nagradama. Sparse nagrade, koje daju povratnu informaciju samo kada agent dosegne ciljno mjesto, predstavljaju značajan izazov za algoritme zbog nedostatka konstantne povratne informacije.

- **PPO (Proximal Policy Optimization):** PPO je pokazao stabilne, ali sporije rezultate u okruženju sa sparse nagradama. Iako je algoritam stabilan i pouzdan, nedostatak kontinuirane povratne informacije usporio je proces učenja.
- **DDPG (Deep Deterministic Policy Gradient):** DDPG je imao poteškoća s učenjem u sparse okruženju zbog nedostatka konstantne povratne informacije. Iako je algoritam dizajniran za kontinuirane akcijske prostore, sparse nagrade su predstavljale značajan izazov, rezultirajući nižim prosječnim nagradama i uspješnošću.
- **SAC (Soft Actor-Critic):** SAC je pokazao dobru prilagodljivost u okruženju sa sparse nagradama. Ovaj algoritam koristi dvostruke kritičare za stabilizaciju učenja i entropy regulaciju za balansiranje između istraživanja i eksploatacije. Zahvaljujući ovim svojstvima, SAC je brzo naučio optimalne politike, postigavši više prosječne nagrade i uspješnost.
- **DDPG s HER (Hindsight Experience Replay):** Ovaj algoritam je pokazao značajne performanse u sparse okruženju. HER tehnika omogućava agentu da uči

iz neuspjelih pokušaja, pretvarajući ih u korisne informacije za buduće epizode. Kombinacija DDPG-a i HER-a rezultirala je visokim prosječnim nagradama i dobrom uspješnošću.

Tablica 4.11. Rezultati sa Sparse nagradama

Algoritam	Prosječna nagrada	Uspješnost
PPO	-24	40%
DDPG	-30	22%
SAC	-15	65%
DDPG + HER	-11	75%

Rezultati s Dense nagradama

U okruženju s dense nagradama, svi algoritmi su pokazali poboljšane performanse zahvaljujući kontinuiranoj povratnoj informaciji. Dense nagrade omogućuju agentima stalno prilagođavanje i optimizaciju politika temeljenih na udaljenosti do cilja, što značajno olakšava proces učenja.

- **PPO (Proximal Policy Optimization):** PPO je imao stabilnije rezultate u dense okruženju, ali još uvijek nedovoljne.
- **DDPG (Deep Deterministic Policy Gradient):** DDPG je poboljšao svoje performanse u dense okruženju, pokazujući da kontinuirana povratna informacija može značajno pomoći u optimizaciji politika i učenju optimalnih akcija.
- **SAC (Soft Actor-Critic):** SAC je pokazao vrlo dobre rezultate, zahvaljujući svojoj sposobnosti balansiranja istraživanja i eksploatacije. Algoritam je postigao velike prosječne nagrade i uspješnost u dense okruženju.
- **DDPG s HER (Hindsight Experience Replay):** Kombinacija DDPG-a i HER-a pokazala je odlične performanse u dense okruženju. Kontinuirana povratna informacija i sposobnost učenja iz neuspjelih epizoda omogućili su agentu da brzo optimizira svoje politike.

4.3.6. Diskusija

Rezultati pokazuju da algoritmi kao što su SAC i DDPG s HER pružaju značajne prednosti u učenju u okruženjima sa sparse nagradama. Njihova sposobnost balansiranja

Tablica 4.12. Rezultati sa Dense nagradama

Algoritam	Prosječna nagrada	Uspješnost
PPO	-6.9991	52%
DDPG	-7.8421	39%
SAC	-3.7398	80%
DDPG + HER	-2.4543	89%

između istraživanja i eksploatacije, kao i učenja iz neuspjelih pokušaja, omogućava im da bolje iskoriste ograničene povratne informacije koje pružaju sparse nagrade.

S druge strane, svi algoritmi su pokazali poboljšane performanse u okruženju s dense nagradama zbog kontinuirane povratne informacije koja olakšava proces učenja. DDPG + HER je pokazao superiornost, ali je i SAC pokazao visoku uspješnost te bi pokazao vjerojatno i puno veću uspješnost da je treniran isto vremenski kao i DDPG + HER.

Ovi rezultati naglašavaju važnost odabira odgovarajućeg RL algoritma ovisno o specifičnostima zadatka i okruženja. Daljnja istraživanja mogu se usmjeriti na optimizaciju hiperparametara i kombinaciju različitih algoritama kako bi se postigle još bolje performanse. Potencijalne smjernice uključuju istraživanje naprednih tehnika za poboljšanje stabilnosti i brzine učenja, kao i prilagodbu algoritama specifičnim zahtjevima različitih okruženja.

4.4. Zaključak

U eksperimentalnom dijelu ovog rada, analizirali smo performanse različitih RL algoritama u Gymnasium okruženju FetchReach-v2 sa sparse i dense nagradama. Rezultati su pokazali da SAC i DDPG s HER pružaju značajne prednosti u učenju u okruženjima sa sparse nagradama, dok svi algoritmi pokazuju poboljšane performanse u okruženjima s dense nagradama.

Buduća istraživanja mogu se usmjeriti na daljnje optimizacije algoritama i istraživanje kombinacija različitih tehnika za poboljšanje performansi agenata u složenim zadacima manipulacije.

5. Zaključak i budući rad

U ovom radu istražena je primjena različitih algoritama za podržano učenje u zadacima robotske manipulacije, s posebnim naglaskom na okruženja CartPole, MountainCar i FetchReach-v2. Korištenje Gymnasium biblioteke omogućilo je stvaranje kontroliranih simulacijskih okruženja koja su bila ključna za treniranje i evaluaciju performansi agenata.

Primjenom različitih RL algoritama kao što su Q-learning, DQN, PPO, DDPG, SAC i DDPG u kombinaciji s HER (Hindsight Experience Replay), istraženi su njihovi učinci na učenje i performanse agenata u različitim uvjetima nagrađivanja, uključujući sparse i dense nagrade.

5.1. Glavni zaključci

1. Efikasnost paralelizma:

- Korištenje 9 paralelnih instanci simulacijskih okruženja pokazalo se izuzetno korisnim za ubrzanje procesa treniranja, omogućavajući brže postizanje optimalnih politika. Ovo je posebno vidljivo kod PPO algoritma u CartPole okruženju, gdje je pokazao superiorne performanse u odnosu na DQN algoritam.

2. Učinkovitost algoritama u različitim okruženjima:

- PPO algoritam je pokazao izvanredne performanse u CartPole okruženju s dense nagradama, dok je DQN imao problema s održavanjem stabilnosti tijekom dugoročnog treniranja.
- U MountainCar okruženju, DQN algoritam je pokazao bolje performanse u odnosu na PPO, što je bilo očekivano zbog prirode zadatka i rijetkih nagrada

koje otežavaju učenje za on-policy algoritme kao što je PPO.

3. **Kombinacija DDPG i HER:**

- DDPG algoritam u kombinaciji s HER pokazao je značajna poboljšanja u zadacima s rijetkim nagradama, kao što je FetchReach-v2 okruženje. Ova kombinacija omogućila je agentima da uče efikasnije i da bolje generaliziraju naučene vještine na nove zadatke.

4. **Sparse vs. Dense nagrade:**

- Algoritmi podržanog učenja općenito su se bolje prilagodili dense nagradama, što je rezultiralo bržim i stabilnijim učenjem. Sparse nagrade, iako izazovnije, omogućile su dublje razumijevanje zadatka, posebno kada su korišteni napredni algoritmi kao što je DDPG + HER.

5.2. **Budući rad**

S obzirom na postignute rezultate, buduća istraživanja mogu se usmjeriti na:

1. **Optimizaciju hiperparametara:**

- Daljnje istraživanje optimalnih hiperparametara za različite RL algoritme u specifičnim okruženjima može dodatno poboljšati performanse i efikasnost učenja.

2. **Primjena u stvarnom svijetu:**

- Prijenos naučenih politika iz simulacijskih okruženja u stvarne robotske sustave. Istraživanje metoda za efikasan transfer znanja može značajno smanjiti potrebu za dugotrajnim treniranjem u stvarnom svijetu.

3. **Integracija s dodatnim senzorskim podacima:**

- Korištenje dodatnih senzorskih podataka (npr. vizualnih informacija) za poboljšanje sposobnosti percepcije i interakcije robota s okolinom.

4. **Istraživanje novih algoritama:**

- Ispitivanje novih algoritama i njihovih varijanti, koji bi mogli ponuditi bolje performanse u specifičnim kontekstima robotske manipulacije, posebno u složenijim i dinamičnijim okruženjima.

Ovi zaključci i prijedlozi za budući rad naglašavaju važnost kontinuiranog istraživanja i optimizacije algoritama podržanog učenja za unapređenje sposobnosti robotskih sustava u izvršavanju složenih zadataka manipulacije.

Literatura

- [1] I. tim FER-a, “Razvoj sustava za robotsku manipulaciju korištenjem rl-a”, *Journal of Robotics*, sv. 10, str. 100–110, 2024.
- [2] I. tim TVZ-a, “Primjena gymnasium biblioteke u industrijskim simulacijama”, *Industrial Robotics Journal*, sv. 15, str. 50–60, 2024.
- [3] Autor, “Primjena ppo algoritma za robotsku manipulaciju”, *IEEE Transactions on Robotics*, sv. 30, str. 120–130, 2024.
- [4] M. istraživački tim, “Kombinacija gymnasium i cleanrl biblioteka za industrijske simulacije”, *Advanced Robotics*, sv. 25, str. 200–210, 2024.
- [5] S. istraživački tim, “Generativni modeli za sintetičke podatke u robotskoj manipulaciji”, *Robotics and Autonomous Systems*, sv. 35, str. 300–310, 2024.
- [6] Sveučilište u Kaliforniji, “Poboljšanje performansi robota u zadacima fine manipulacije korištenjem ddpq algoritma”, UCLA Robotics Lab, 2023.
- [7] R. S. Sutton i A. G. Barto, *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- [8] C. J. C. H. Watkins i P. Dayan, “Q-learning”, *Machine Learning*, sv. 8, br. 3, str. 279–292, 1992.
- [9] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning”, *Nature*, sv. 518, br. 7540, str. 529–533, 2015.

- [10] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, i O. Klimov, “Proximal policy optimization algorithms”, *arXiv preprint arXiv:1707.06347*, 2017.
- [11] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, i D. Wierstra, “Continuous control with deep reinforcement learning”, *arXiv preprint arXiv:1509.02971*, 2016.
- [12] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, i S. Levine, “Soft actor-critic algorithms and applications”, 2019. [Mrežno]. Adresa: <https://arxiv.org/abs/1812.05905>
- [13] J. Kober, J. A. Bagnell, i J. Peters, “Reinforcement learning in robotics: A survey”, *The International Journal of Robotics Research*, sv. 32, br. 11, str. 1238–1274, 2013.
- [14] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, “Mastering the game of go with deep neural networks and tree search”, *Nature*, sv. 529, br. 7587, str. 484–489, 2016.
- [15] Y. Deng, F. Bao, Y. Kong, Z. Ren, i Q. Dai, “Deep direct reinforcement learning for financial signal representation and trading”, *IEEE transactions on neural networks and learning systems*, sv. 28, br. 3, str. 653–664, 2016.
- [16] C. Yu, J. Liu, i S. Nemati, “Reinforcement learning in healthcare: A survey”, *arXiv preprint arXiv:1908.08796*, 2019.
- [17] B. Siciliano i O. Khatib, *Springer Handbook of Robotics*. Springer Science & Business Media, 2008.
- [18] S. Levine, C. Finn, T. Darrell, i P. Abbeel, “End-to-end training of deep visuomotor policies”, *The Journal of Machine Learning Research*, sv. 17, br. 1, str. 1334–1373, 2016.
- [19] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke *et al.*, “Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation”, *arXiv preprint arXiv:1806.10293*, 2018.

- [20] A. Nagabandi, G. Kahn, R. S. Fearing, i S. Levine, “Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning”, *arXiv preprint arXiv:1708.02596*, 2017.
- [21] H. Zhang, B. Hou, J. Wu, T. Sun, i Z. Fang, “Multi-robot collaboration for autonomous construction of controllable structures”, *IEEE Access*, sv. 7, str. 108 825–108 837, 2019.
- [22] J. M. Beer, A. D. Fisk, i W. A. Rogers, “Toward a framework for human-robot interaction”, *Journal of Human-Robot Interaction*, sv. 1, br. 1, str. 9–24, 2012.
- [23] L. L. Whitcomb, “Advances in underwater robotics”, *Current Opinion in Robotics and Autonomous Systems*, sv. 10, str. 549–557, 2000.
- [24] G. D. Team, “Gymnasium: A modern library for rl simulations”, *Journal of Simulation and Training*, sv. 20, str. 100–120, 2024.
- [25] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, i W. Zaremba, “Openai gym”, *arXiv preprint arXiv:1606.01540*, 2016.
- [26] OpenRL, “Openrl documentation”, 2023., accessed: 2024-06-21. [Mrežno]. Adresa: <https://openrl-docs.readthedocs.io/en/latest/>
- [27] S. B. Contributors, “Stable baselines3 documentation”, 2023., accessed: 2024-06-21. [Mrežno]. Adresa: <https://stable-baselines3.readthedocs.io/en/master/>
- [28] DeepMind, “Mujoco: Physics engine”, 2023., accessed: 2024-06-21. [Mrežno]. Adresa: <https://mujoco.org/>
- [29] C. Watkins, “Learning from delayed rewards”, *PhD thesis, University of Cambridge*, 1989.
- [30] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning”, *Nature*, sv. 518, br. 7540, str. 529–533, 2015.

- [31] MathWorks, “Train dqn agent to balance cart-pole system”, 2023., accessed: 2024-06-21. [Mrežno]. Adresa: <https://www.mathworks.com/help/reinforcement-learning/ug/dqn-agents.html>
- [32] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, i D. Wierstra, “Continuous control with deep reinforcement learning”, *arXiv preprint arXiv:1509.02971*, 2015.
- [33] MathWorks, “Train ddpq agent to balance cart-pole system”, 2023., accessed: 2024-06-21. [Mrežno]. Adresa: <https://www.mathworks.com/help/reinforcement-learning/ug/ddpg-agents.html>
- [34] T. Haarnoja, A. Zhou, P. Abbeel, i S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor”, *arXiv preprint arXiv:1801.01290*, 2018.

Sažetak

U ovom radu istražena je primjena različitih algoritama podržanog učenja u zadacima robotske manipulacije, s posebnim naglaskom na okruženja CartPole, MountainCar i FetchReach-v2. Korištenjem Gymnasium biblioteke stvorena su kontrolirana simulacijska okruženja koja su omogućila treniranje i evaluaciju performansi agenata. Razmatrani algoritmi uključuju Q-learning, DQN, PPO, DDPG, SAC i DDPG u kombinaciji s HER (Hindsight Experience Replay). Eksperimenti su pokazali da paralelizam značajno ubrzava proces treniranja, dok različiti algoritmi imaju različite performanse u ovisnosti o specifičnom okruženju i tipu nagrada (sparse ili dense). PPO se pokazao učinkovitim u CartPole okruženju s dense nagradama, dok je DQN nadmašio PPO u MountainCar okruženju. Kombinacija DDPG i HER dala je najbolje rezultate u FetchReach-v2 okruženju s rijetkim nagradama. Zaključeno je da algoritmi podržanog učenja općenito bolje funkcioniraju s dense nagradama, ali napredne metode poput DDPG + HER omogućuju stabilno učenje i u uvjetima rijetkih nagrada. Buduća istraživanja trebala bi se usmjeriti na optimizaciju hiperparametara, prijenos znanja iz simulacije u stvarni svijet, integraciju dodatnih senzorskih podataka i istraživanje novih algoritama.

Ključne riječi: reinforcement learning (RL); robotska manipulacija; sintetički podaci; simulacije; autonomni robotski sustavi

Abstract

This study explores the application of various reinforcement learning algorithms in robotic manipulation tasks, focusing on CartPole, MountainCar, and FetchReach-v2 environments. Using the Gymnasium library, controlled simulation environments were created to facilitate the training and evaluation of agent performance. The considered algorithms include Q-learning, DQN, PPO, DDPG, SAC, and DDPG combined with HER (Hindsight Experience Replay). The experiments demonstrated that parallelism significantly accelerates the training process, while different algorithms exhibit varying performance depending on the specific environment and type of rewards (sparse or dense). PPO proved effective in the CartPole environment with dense rewards, while DQN outperformed PPO in the MountainCar environment. The combination of DDPG and HER yielded the best results in the FetchReach-v2 environment with sparse rewards. It was concluded that reinforcement learning algorithms generally perform better with dense rewards, but advanced methods like DDPG + HER enable stable learning even in sparse reward conditions. Future research should focus on hyperparameter optimization, transferring knowledge from simulation to the real world, integrating additional sensory data, and exploring new algorithms.

Keywords: reinforcement learning (RL); robot manipulation; synthetic data; simulations; autonomous robotic systems