



PRIMENA UČENJA USLOVLJAVANJEM NA OBUČAVANJE AGENTA ZA AUTONOMNO KRETANJE U OKRUŽENJU *CarRacing-v0*

USING REINFORCEMENT LEARNING TO TRAIN AN AGENT FOR THE *CarRacing-v0* ENVIRONMENT

Novica Šarenac, *Fakultet tehničkih nauka, Novi Sad*

Oblast – RAČUNARSTVO I AUTOMATIKA

Kratak sadržaj – *U radu je opisan trening i evaluacija agenta za autonomno kretanje u OpenAI Gym okruženju CarRacing-v0. Okruženje predstavlja stazu za vožnju automobila iz ptičije perspektive. Za treniranje agenta korištene su tehnike učenja uslovljavanjem. Poređene su osobine algoritama u pogledu rezultata postignutih u okruženju, dužine treniranja i implementacionih detalja. Implementirani i poređeni algoritmi su: Deep Q-Network (DQN), Advantage Actor Critic (A2C) i Asynchronous Advantage Actor Critic (A3C).*

Ključne reči: učenje uslovljavanjem, *CarRacing-v0*, Deep Q-Network, Advantage Actor Critic, Asynchronous Advantage Actor Critic

Abstract – *This paper presents training and evaluation of the agent for autonomous driving in OpenAI Gym environment CarRacing-v0. Environment is a top-down view of racing track. Agent is trained using reinforcement learning techniques. Algorithms are compared in terms of results achieved in the environment, training time and implementation details. Algorithms that are implemented and evaluated are: Deep Q-Network (DQN), Advantage Actor Critic (A2C) i Asynchronous Advantage Actor Critic (A3C).*

Keywords: reinforcement learning, *CarRacing-v0*, Deep Q-Network, Advantage Actor Critic, Asynchronous Advantage Actor Critic

1. UVOD

U današnje vreme prisutna je velika ekspanzija veštačke inteligencije u svim sferama računarskih nauka. Sve više se teži razvoju softvera koji se na neki način može okarakterisati kao intelligentan. Veliku ulogu u toj ekspanziji ima razvoj hardvera, koji omogućava treniranje složenih modela.

Najčešća upotreba veštačke inteligencije je u detektovanju i prepoznavanju objekata na slikama, razumevanju i klasifikaciji teksta, sistemima za preporučivanje, odgovaranju na pitanja, i slično.

Sve češće veštačka inteligencija se upotrebljava za obučavanje intelligentnih agenata koji fleksibilno i samostalno izvšavaju neke akcije u okruženju.

NAPOMENA:

Ovaj rad proistekao je iz master rada čiji mentor je bio dr Aleksandar Kovačević, vanr. prof.

Tako u auto-industriji mnoge kompanije rade na razvoju autonomnih vozila, u privredi mnogi procesi teže ka automatizaciji, u softveru se razvijaju agensi čiji je cilj ostvarenje što boljih rezultata u simuliranim okruženjima (poput video igara).

U ovom radu opisan je trening i evaluacija agenta za autonomno kretanje u simuliranom okruženju. Okruženje predstavlja vožnju automobila iz ptičije perspektive. Za treniranje agenta korišćena su i upoređena tri algoritma učenja uslovljavanjem: Deep Q-Network (DQN) [1], Advantage Actor Critic (A2C) i Asynchronous Advantage Actor Critic (A3C) [2].

2. METODOLOGIJA

Svaki od pristupa za treniranje agenta koristi različite korake preprocesiranja stanja (observacija) koja daje okruženje, kao i različite metode dubokog učenja uslovljavanjem.

2.1. Okruženje *CarRacing-v0*

CarRacing-v0 [3] je okruženje sa OpenAI Gym platforme. Okruženje podrazumeva 2D simulaciju upravljanja automobilom iz ptičije perspective.

Stanje koje okruženje vraća prilikom interakcije je slika dimenzija 96x96 piksela, sa tri kanala (RGB).

Nagrade se računaju za svaki frejm (korak u okruženju). Nagrada je -0.1 za svaki frejm i 1000/N za svaku posećenu pločicu na stazi, gde je N ukupan broj pločica na stazi. Kako se staza generiše na slučajan način za svaku igru, ukupan broj pločica na stazi varira između 230 i 320. Epizoda se završava kada se posete sve pločice na stazi ili kada prođe 1000 frejmova.

Jedna od najvažnijih karakteristika ovog okruženja je što ima kontinualne akcije, za razliku od većine okruženja, kao što su *Atari* igre koje imaju diskretne akcije. Svaka akcija sastoji se od tri kontinualne vrednosti, koje odgovaraju: smeru (levo ili desno), ubrzanju i kočenju. Upravljanje smerom je u rangu između -1.0 i 1.0, gde je -1.0 potpuno skretanje levo i 1.0 potpuno skretanje desno. Ubrzavanje i kočenje je u rangu od 0 do 1.

U donjem delu prozora nalaze se grafički predstavljeni indikatori: brzina, ABS senzori, pozicija volana i žiroskop. Ove informacije nisu dostupne u brojevnom obliku. Smatra se da je okruženje rešeno ako je srednja vrednost nagrade u 100 uzastopnih pokušaja 900.

2.2. Primena DQN metode

Jedna od tri primenjene metode je DQN metoda. Da bi se agent mogao uspešno obučavati potrebno je stanje koje se dobija od okruženja prilagoditi DQN metodi. Takođe je potrebno odabrat optimalnu arhitekturu neuronske mreže, koja je pogodna za rad sa takvim stanjem.

2.2.1 Stanje u DQN metodi

Okruženje *CarRacing-v0* reprezentuje svako stanje sa 96x96x3 RGB matricom. Iz razloga kao što su brža konvergencija i trening vrši se niz koraka preprocesiranja originalnog stanja.

Stanje se prvo prebacuje u *grayscale*, čime se dobija matrica dimenzija 96x96. S obzirom na to da se indikatori sa donjeg dela prozora ne mogu iskoristiti u brojevnom obliku, oni se uklanjuju iz stanja. Iseca se prostor dimenzija 12x96 sa donjeg dela svake slike. Takođe se iseca prostor dimenzija 96x6 piksela sa leve i desne strane svake slike. Nakon izvršenih koraka preprocesiranja dobija se slika dimenzija 84x84 piksela.

Takođe se primenjuje tehnika *frame-skipping*. Agent vidi i odabira akciju za svako k -to stanje, umesto za svako stanje. Poslednja odabrana akcija se ponavlja za preskočena stanja. Pravljenje koraka u okruženju je računski veoma jeftino i brzo, za razliku od odabiranja akcija, pa ova tehnika omogućava agentu da odigra k puta više igara bez velikog povećanja vremena treninga. Vrednost korištena za parametar k je 4.

2.2.2 Arhitektura DQN mreže

Stanje koje se koristi za obučavanje agenta je slika ekrana okruženja, pa je odabrana arhitektura koja predstavlja konvolutivnu mrežu. Arhitektura se sastoji od konvolutivnih, *batch normalization* i potpuno povezanih slojeva. Ulaz u neuronsku mrežu je preprocesirana slika dimenzija 84x84.

Prvi sloj je 2D konvolutivni sloj. Broj kanala na izlazu tog sloja je 16. Koristi se filter dimenzija 5x5 i korak (*stride*) pomeranja filtera 2. Sledeći sloj je *batch normalization* sloj.

Drugi konvolutivni sloj na ulazu očekuje 16 kanala, a broj kanala na izlazu je 32. Dimenzije filtera u ovom sloju su 5x5, a korak pomeranja 2. I nakon ovog sloja vrši se *batch normalization*.

Broj kanala na ulazu i izlazu trećeg konvolutivnog sloja je 32. Koristi se filter dimenzija 5x5 sa korakom pomeranja 2. Sledeći sloj je takođe *batch normalization*. Izlaz se transformiše u jednodimenzioni vektor, koji je pogodan za potpuno povezane slojeve.

Poslednja dva sloja su potpuno povezani slojevi. Prvi potpuno povezan sloj na ulazu ima dimenziju 1568, a na izlazu 256. Poslednji sloj na ulazu ima dimenziju 256, a na izlazu 4, što se poklapa sa brojem akcija. Izlaz neuronske mreže predstavlja Q -vrednosti za svaku akciju.

2.2.3 Implementacija DQN metode

Obučavanje agenta *DQN* metodom koristi okruženje i korake preprocesiranja opisane u sekciji 2.2.1. U

implementaciji se koriste dve mreže: *policy* mreža i *target* mreža. *Policy* mreža koristi se za odabiranje akcija, a *target* mreža služi za računanje očekivanih Q vrednosti.

U *replay* memoriji se čuvaju parovi: stanje, akcija, nagrada, sledeće stanje i da li je epizoda završena. Kapacitet *replay* memorije koji je korišten je 10 000 primera. Tokom treniranja izvršava se glavna trening petlja. Broj izvršenih koraka u trening petljama je 150 000.

Prvi korak u trening petlji je računanje Q -vrednosti za akcije u trenutnom stanju. Za ovo računanje koristi se *policy* mreža. Akcije su diskretizovane, tako da postoje četiri moguće akcije: skretanje levo, skretanje desno, gas i kočnica. Tokom treninga, za odabir akcije se koristi ϵ -greedy, da bi agent mogao istraživati okruženje. Ukoliko je slučajno odabran broj manji od ϵ agent na slučajan način odabira akciju. U suprotnom uzima se akcija za koju je Q -vrednost maksimalna. Vrednost ϵ se menja tokom treninga (1).

$$\epsilon = \epsilon_{final} + (\epsilon_{start} - \epsilon_{final}) * e^{(-1 * step / \epsilon_{step})} \quad (1)$$

Vrednost ϵ_{start} je 0.9, ϵ_{final} je 0.05 i ϵ_{step} je 30 000. Parametar *step* je korak do koga je trening petlja stigla.

Sledeći korak je izvršavanje akcije u okruženju. Nakon toga se par koji se sastoji od stanja, akcije, novog stanja (dobjegenog iz okruženja) i statusa epizode čuva u *replay* memoriji.

Update policy mreže je sledeći korak obučavanja. Prvo se na slučajan način uzima *batch* sakupljenih primera iz *replay* memorije. Veličina *batch*-a je 32. Računaju se Q -vrednosti za urađene akcije za sva stanja u *batch*-u, koristeći *policy* mrežu. Zatim se računaju očekivane Q -vrednosti za sledeća stanja iz *batch*-a koristeći *target* mrežu. Nakon toga se računa loss i vrši *update policy* mreže.

Update target mreže vrši se povremeno, na svakih 10 000 koraka. *Update* se vrši tako što *target* mreža preuzme parametre *policy* mreže. Algoritam za optimizaciju i *update* neuronske mreže korišten u implementaciji je *RMSprop* [4].

2.3. Primena actor-critic metoda

Actor-critic metode koje su implementirane su A3C i A2C. Kao i kod DQN metode, i za *actor-critic* metode je potrebno prilagoditi stanje koje se dobija od okruženja i pronaći optimalnu arhitekturu neuronske mreže. Za obe implementirane metode se koriste isti koraci preprocesiranja stanja, kao i ista arhitektura neuronske mreže.

2.3.1. Stanje u actor-critic metodama

Preprocesiranje stanja (96x96x3 RGB matrice) koje se dobija od *CarRacing-v0* okruženja sastoji se od niza koraka preprocesiranja.

Prvi korak preprocesiranja je prebacivanje slike u *grayscale*, čime se dobija stanje dimenzija 96x96 piksela.

Sledeći korak preprocesiranja je centriranje piksela oko 0. To se vrši oduzimanjem srednje vrednosti od svakog piksela na slici. Kao i kod DQN metode, indikatori sa

donjeg dela prozora se uklanjaju iz stanja i iseca se prostor dimenzija 12×96 piksela sa donjeg dela svakog stanja. Takođe se iseca prostor dimenzija 96×6 piksela sa leve i desne strane stanja. Nakon preprocesiranja, dimenzija stanja je 84×84 piksela.

Krajnje stanje, koje se koristi u *actor-critic* metodama formira se od nekoliko poslednjih observacija (preprocesiranih 84×84 slika) koje daje okruženje. Observacije se čuvaju i stanje koje je ulaz u neuronsku mrežu za A3C i A2C algoritme je dimenzija $84 \times 84 \times N$. N je broj observacija koji se čuva. Vrednost N koja je korištena je 5.

2.3.2 Arhitektura mreže kod A2C i A3C metoda

Arhitektura neuronske mreže korištene kod A2C i A3C metoda je konvolutivna neuronska mreža. Arhitektura se sastoji od konvolutivnih i potpuno povezanih slojeva. Ulaz u neuronsku mrežu opisan je u sekciji 2.3.1 i ima dimenziju $84 \times 84 \times 5$.

Prvi sloj je 2D konvolutivni sloj. Broj kanala na izlazu tog sloja je 16. Koristi se filter dimenzija 8×8 i korak pomeranja filtera 4.

Dруги конвolutивни слој на улазу има 16 канала, а број излазних канала је 32. Димензије филтера у овом слоју су 3×3 , а корак померања филтера 2.

Izlaz другог конвolutивног слоја трансформише се у једнодимензиона вектор, који је погодан за потпuno пoveзane slojeve. Следећи слој је потпuno пoveзан и на улазу има димензију 2592, а на излазу 256.

Neuronska mreža има два излаза. Један представља функцију вредности, а други политику. Функција вредности и политика се разликују само у последњем слоју.

Улазна димензија за последњи слој функције вредности је 256, а излазна је 1. Излаз из tog слоја је вредност trenutnog stanja. Исто као код функције вредности, и код последnjeg слоја политике улазна димензија је 256. Излазна димензија политике је 4 и поклапа се са бројем акција.

2.3.4 Implementacija A3C metode

A3C метода користи окruženje и кораке preprocesiranja opisane u sekciji 2.3.1. Како је A3C метода заснована на multiprocesingu prvi korak u implementaciji је да се inicijalizuje više процеса, од којих сваки има своју trening petlju. Broj процеса који је korišten je 8. Postoji jedan globalni model, чији update vrše svi процеси.

Svaki процес има локалнуinstancu modela, instancu okruženja, као и referencu ka globalnom modelu. Inicijalizuje се i memorija која чува парове: korak, vrednost stanja, nagrada, logaritam verovatnoće akcije, entropija i da li je epizoda završena. Подаци сачувани у тој memoriji користе се за један update мреже, па је kapacitet memorije jednak броју корака који су потребни за један update. Broj koraka korišten u implementaciji je 40. Takođe se izvršava главна trening petlja. Broj izvršenih koraka u trening petljи u jedном процесу је 150 000.

Prvi korak u trening petljи je update локалне мреже i preuzimanje parametara глобалне мреже. Следећи корак је petlja која izvršava broj akcija potreban за један update

globalne mreže. Prvo se računaju verovatnoće akcija, logaritmi tih verovatnoća i vrednost trenutnog stanja. Затим се, користећи verovatnoće, бира akcija која ће се izvršiti u okruženju. Akcije су kontinualне и добијају се тако што се прво odabere akcija која има максималну verovatnoću. Moguće akcije су skretanje levo, skretanje desno, gas i kočnica. Затим се тај вектор помножи са verovatnoćom i та akcija се искористи као akcija u okruženju. Nakon odabira akcije računa се entropija, која подстиче istraživanje agenta.

Следећи корак у trening petlji је izvršavanje akcije u okruženju, чиме се добија ново stanje i nagrada. Корак, vrednost stanja, nagrada, logaritam verovatnoće akcije, entropija i da ли је epizoda završena чувају се у memoriji.

Када се izvrši potreban број корака, računa се loss i vrši update globalne mreže. Prvo se računaju kumulativне nagrade. Затим, користећи kumulativne nagrade и vrednosti stanja računa се advantage i loss за funkciju vrednosti.

Следећи корак је računanje loss-a за политику, који користи izračunati advantage i logaritme verovatnoće akcija iz memorije. Користећи loss функцију vrednosti, loss функцију политике и entropiju računa се krajnja vrednost loss-a. Vrednost koeficijenta entropije korištena u implementaciji je 0,01.

Када се izračuna loss vrši se update globalnog modela. Алгоритам за optimizaciju i update neuronske mreže korišтен у implementaciji је Adam [5].

2.3.5 Implementacija A2C metode

A2C метода користи окruženje и кораке preprocesiranja opisane u sekciji 2.3.1. I A2C метода заснована је на паралелним agentima, који istražuju više окruženja. Agenti су sinhroni i koriste isti globalni model, чији update se vrši тек када сви agenti izvrše potreban број корака. Agenti imaju zajedničку trening petlju.

Inicijalizују се процеси од којих сваки има по једну instance окruženja. Broj korištenih процеса је 8. Када се uradi reset или izvrši korak u okruženjima, чека се да сви процеси urade reset ili korak. Тakođe, као и код A3C методе, inicijalizује се memorija, у којој се чувају парови: korak, vrednosti stanja, nagrade, logaritmi verovatnoće akcije, entropije i da ли су epizode zavrшene. Подаци сачувани у memoriji služe за update neuronske mreže, па је kapacitet jednak броју корака који се izvrše за један update po процесу. Broj koraka korišten u implementaciji је 5. Broj koraka који се izvršava u glavnoј trening petljи је 125 000.

Prvi korak u trening petljи је izvršavanje akcija potrebnih за један update мреже. Izvršава се 5 akcija po процесу. Prvo se računaju verovatnoće akcija, logaritmi tih verovatnoća i vrednost trenutnog stanja. Od dobijenih verovatnoće бира се akcija која ће се izvršiti u okruženju. Akcije су, као и код A3C методе kontinualне и biraju се на исти начин.

Nakon odabira akcija računa се entropija. Затим се одabrane akcije izvršavaju у свим окruženjima у исто време. Nakon izvršavanja, од сваког процеса се добија ново stanje i nagrada.

Koraci, vrednosti stanja, nagrade, logaritmi verovatnoća akcija, entropije i da li su epizode završene čuvaju se u memoriji. Kada se izvrši potreban broj koraka, vrši se *update* mreže. Prvo se računaju kumulativne nagrade. Pomoću njih i vrednosti stanja se računa *advantage*, a zatim i *loss* za funkciju vrednosti. Za računanje *loss-a* za politiku koriste se izračunati *advantage* i logaritmi verovatnoća akcija, koji su sačuvani u memoriji.

Krajnji *loss* se računa pomoću *loss-a* funkcije vrednosti, *loss-a* politike i entropije. Koeficijent entropije koji je korišten je 0,1. Nakon toga vrši se *update* modela. Algoritam za optimizaciju i *update* mreže korišten u implementaciji je Adam.

3. REZULTATI

Okruženje *CarRacing-v0* smatra se rešenim ukoliko agent ostvari prosečnu vrednost nagrade veću od 900 u 100 uzastopnih pokušaja. Ti rezultati su daleko iznad rezultata koje može ostvariti čovek. Mera tačnosti korištena za evaluaciju sva tri modela je prosečna ostvarena nagrada u 100 uzastopnih epizoda.

Kod DQN metode broj koraka prilikom treninga koji se pokazao kao optimalan je 150 000. Agent i sa manjim brojem koraka (oko 100 000) obučavanja nauči da se kreće unapred, ali ne može da savlada krivine. Tek sa većim brojem koraka agent uči da skreće u krivinama. DQN agent može savladati čak i velike krivine, sa kojima ostali modeli imaju poteškoća. Prosečna osvojena nagrada u 100 uzastopnih epizoda ovog modela je 178.

Optimalan broj koraka prilikom treninga A3C agenta je 150 000. Agent jako brzo (za manje od 40 000 koraka) nauči da prati stazu i kreće se, međutim, kao i DQN agentu, potreban mu je veći broj koraka obučavanja da bi naučio da skreće u krivinama. U 100 uzastopnih epizoda A3C agent ostvaruje prosečnu nagradu 241.

Kod A2C metode optimalan broj koraka prilikom treninga je 125 000. Slično A3C agentu, vrlo brzo nauči da se kreće unapred i prati stazu, ali tek sa velikim brojem koraka obučavanja počinje pravilno skretati u krivinama. Prosečna osvojena nagrada u 100 uzastopnih epizoda ovog modela je 263.

Iz tabele 1. vidi se da A2C metoda ostvaruje najveću prosečnu nagradu u 100 uzastopnih epizoda.

Broj koraka treniranja	Metoda	Prosečna osvojena nagrada u 100 epizoda
Slučajne akcije	-	-54
150 000	DQN	178
150 000	A3C	241
125 000	A2C	263

Tabela 1: Poređenje rezultata svih implementiranih metoda

4. ZAKLJUČAK

U ovom radu opisano je obučavanje agenta za autonomno kretanje u simuliranom okruženju *CarRacing-v0*. Agent je obučavan koristeći metode učenja uslovljavanjem. Implementirane metode su: *Deep Q-Network* (DQN), *Advantage Actor Critic* (A2C) i *Asynchronous Advantage Actor Critic* (A3C).

Na osnovu rezultata može se zaključiti da agenti trenirani metodama učenja uslovljavanjem mogu naučiti da se snalaze u okruženjima poput *CarRacing-v0*. Metode korištene u radu su se pokazale kao veoma uspešne u učenju direknto iz vizuelnih podataka. Takođe, sve metode su veoma osetljive na promene hiperparametara. Performanse svih modela u velikoj meri zavise od načina preprocesiranja slike i formiranja stanja.

Actor-critic metodama potrebno je mnogo manje vremena da počnu da konvergiraju i uče. Kod njih se učenje može primetiti već nakon 40 000 koraka, gde agent počinje da prati stazu, ali ima probleme prilikom skretanja. Kod DQN metode učenje se može primetiti nešto kasnije, nakon 100 000 koraka obučavanja. Takođe vreme potrebno za trening *actor-critic* metoda je značajno manje nego kod DQN metode, što je rezultat korištenja multiprocesinga.

Metode A2C i A3C imaju slične karakteristike upravljanja automobilom. Agenti obučavani korištenjem A2C i A3C metoda imaju karakteristike brze vožnje i na taj način pokušavaju osvojiti veće nagrade. Ti agenti imaju pravilno skretanje u većini krivina. Agent obučavan korištenjem DQN metode ima karakteristike sporije vožnje, ali pravilnog skretanja u krivinama i na taj način pokušava osvojiti što veću nagradu. Oba načina vožnje imaju prednosti i mane. Problem koji se pojavljuje kod agenata obučavanih A2C i A3C metodama je skretanje u velikim krivinama, dok problem agenta obučenog DQN metodom je povratak na stazu nakon silaska sa staze.

5. LITERATURA

- [1] V. Mnih, K. Kavukcouoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, D. Hassabis „Human-level control through deep reinforcement learning“
- [2] V. Mnih, A.P. Badia, M. Mirza, A. Graves, T. Harley, T.P. Lillicrap, D. Silver, K. Kavukcouoglu „Asynchronous Methods for Deep Reinforcement Learning“
- [3] <https://gym.openai.com/envs/CarRacing-v0/> [pristupljeno 7.9.2019.]
- [4] S. Ruder „An overview of gradient descent optimization algorithms“
- [5] D. P. Kingma, J. L. Ba „Adam: A Method for Stochastic Optimization“

Kratka biografija:



Novica Šarenac rođen je u Nevesinju 1994. godine. Osnovne akademske studije iz oblasti Računarstvo i automatika, završio je na Fakultetu tehničkih nauka, 2017. godine., kada upisuje i master akademske studije iz iste oblasti. Položio je sve ispite predložene planom i programom.