

# Osnove podržanog učenja

---

Ozretić, Petar

**Master's thesis / Diplomski rad**

**2022**

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Split, University of Split, Faculty of science / Sveučilište u Splitu, Prirodoslovno-matematički fakultet**

Permanent link / Trajna poveznica: <https://urn.nsk.hr/urn:nbn:hr:166:303398>

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

Download date / Datum preuzimanja: **2024-05-06**

Repository / Repozitorij:

[Repository of Faculty of Science](#)



PRIRODOSLOVNO-MATEMATIČKI FAKULTET  
SVEUČILIŠTA U SPLITU

PETAR OZRETIĆ

## **OSNOVE PODRŽANOG UČENJA**

DIPLOMSKI RAD

Split, rujan 2022.

PRIRODOSLOVNO-MATEMATIČKI FAKULTET  
SVEUČILIŠTA U SPLITU

ODJEL ZA MATEMATIKU

## OSNOVE PODRŽANOG UČENJA

DIPLOMSKI RAD

Student:  
Petar Ozretić

Mentor:  
izv. prof. dr. sc.  
Saša Mladenović

Split, rujan 2022.

# Uvod

Učenje uz pomoć okoline prati ljudska bića od prvih dana, kad se beba igra, maše rukama, promatra svijet itd., nema izričitog učitelja, ali ima direktnu senzomotornu vezu sa svojom okolinom. Ova veza je bogat izvor informacija o uzrocima i posljedicama, akcijama i reakcijama, i o tome što treba napraviti da bi došli do svog cilja, e.g. *ako zaplačem, promijenit će mi pelenu.*

Interakcija s okolinom je forma učenja koja prati čovjeka tijekom cijelog života - kada npr. čovjek uči voziti automobil, mora biti izrazito pozoran na svoju okolinu i kako ona reagira na ono što ona čini - *ako vozači na cesti blicaju i trube, vjerojatno sam ja taj koji vozi krivom stranom autocestu.*

**Podržano učenje** (PU) je učenje postupaka koji će maksimizirati nekakvu numeričku nagradu. Učeniku nitko ne govori što bi trebalo učiniti - sam mora isprobati radnje i otkriti koja od njih mu donosi najveću nagradu. Najzanimljiviji i najteži su slučajevi kada neka radnja ne utječe samo na trenutnu nagradu, nego i na iduću situaciju, a time i na sve buduće nagrade. Ova dva svojstva - *učenje iz pokušaja i pogrešaka te odgođena nagrada*, odnosno dalekosežne posljedice radnji - su razlikovne karakteristike podržanog učenja u odnosu na druge paradigmе korištene u umjetnoj inteligenciji.

Dok je razlika između podržanog i nenadziranog učenja - koje pokušava pronaći skrivenu strukturu u skupu podataka - očita, granicu između podržanog i nadziranog učenja je puno teže razlučiti. Nadzirano učenje je učenje iz

skupa za treniranje koji sadrži označene primjere koje dobijamo od svez-najućeg "nadziratelja". Svaki primjer je dan kao opis situacije zajedno s ispravnom radnjom koju bi trebalo poduzeti u toj situaciji, najčešće je to određivanje kojoj kategoriji ta situacija pripada. Cilj ove vrste učenja je ge-neralizacija danih odgovora kako bi to kasnije mogli primijeniti na situacije koje se ne nalaze u skupu za treniranje. Iako je nadzirano učenje jako važno, nije dostatno da obuhvati učenje iz interakcija s okolinom jer je jako često nepraktično izraditi toliko primjera koji bi pokrili apsolutno sve situacije. Na ovom neistraženom teritoriju treba biti sposoban učiti iz vlastitog iskustva.

Još jedna razlikovna karakteristika koju se ne susreće kod drugih vrsta učenja je kompromis između istraživanja i iskorištavanja trenutnog znanja. Da bi dobio što veću nagradu učenik mora koristiti one radnje koje je već nekad u prošlosti isprobao i za njih dobio dobru nagradu. Ali da bi uopće otkrio koje su to radnje mora isprobati radnje koje do sada nije. Učenik mora iskorištavati svoje trenutno znanje da bi dobio nagradu, ali mora i nastaviti istraživati da bi u budućnosti mogao još bolje birati. Vrijedi uočiti da iako je dilema između istraživanja i iskorištavanja desetljećima bila tema brojnih matematičkih istraživanja, ipak i dalje nije razriješena.

Još jedno bitno svojstvo podržanog učenja koje je za rad važno spomenuti - za razliku od drugih vrsta učenja, i općenito prakse u programiranju, gdje se veći problem razbija na više manjih - podržano učenje promatra problem isključivo u cjelini. *Nebitno je ako na vozačkom ispit u cijelo vrijeme vozi savršeno, ako se na kraju vozač zabije u zid.*

# Sadržaj

<b>Uvod</b>	<b>iii</b>
<b>Sadržaj</b>	<b>v</b>
<b>1 Osnovni pojmovi podržanog učenja</b>	<b>1</b>
1.1 Agent i okolina . . . . .	1
1.2 Markovljev proces odlučivanja . . . . .	3
1.3 Funkcija vrijednosti . . . . .	6
1.4 Primjer: Mrežni svijet . . . . .	8
1.4.1 Optimalna strategija . . . . .	10
<b>2 Dinamičko programiranje</b>	<b>14</b>
2.1 Vrednovanje strategije (predviđanje) . . . . .	15
2.2 Poboljšavanje strategije (kontrola) . . . . .	19
2.3 Iteracija strategija . . . . .	21
2.4 Iteracija vrijednosti . . . . .	23
2.5 Napomene . . . . .	25
2.5.1 Asinkrono DP . . . . .	25
2.5.2 Generalizirana iteracija strategije . . . . .	25
<b>3 Predviđanje bez modela</b>	<b>27</b>

3.1	Monte-Carlo Učenje . . . . .	28
3.1.1	Primjer: Blackjack . . . . .	29
3.2	Učenje s vremenskom razlikom . . . . .	32
3.2.1	Primjer: vožnja kućí . . . . .	34
3.3	MC vs. TD . . . . .	36
3.3.1	Primjer: Petar predviđa . . . . .	38
<b>4</b>	<b>Kontrola bez modela</b>	<b>40</b>
4.1	Monte-Carlo kontrola . . . . .	41
4.2	Kontrola s vremenskom razlikom . . . . .	43
4.2.1	Primjer: vjetroviti Mrežni svijet [1] . . . . .	46
<b>5</b>	<b>Praktični primjer: Blackjack</b>	<b>49</b>
<b>6</b>	<b>Zaključak</b>	<b>61</b>
	<b>Literatura</b>	<b>63</b>

# Poglavlje 1

## Osnovni pojmovi podržanog učenja

### 1.1 Agent i okolina

Da bi se prilagodili standardnoj nomenklaturi PU-a, učenika, tj. program koji donosi odluke nazivamo *agent*, a *okolina* je sve ono izvan našeg agenta s čime agent vrši interakciju. Radnju ili odluku koju agent vrši ili donosi nazivamo *akcija*<sup>1</sup>. Osim ovih, još su četiri podeljena PU-a: strategija, nagrada/nagradni signal, vrijednosna funkcija (en. *value function*), te model okoline.

*Strategija* definira kako se agent ponaša u određenoj situaciji. Ugrubo govoreći može se reći da je *strategija* preslikavanje iz stanja okoliša u skup svih akcija koje se mogu poduzeti u tim stanjima. Strategija može biti i stohastička - u smislu da poručuje samo kolika je vjerojatnost da će agent poduzeti neku akciju u određenoj situaciji.

---

<sup>1</sup>U literaturi se za pojmove *agent*, *okolina* i *akcija* znaju koristiti i više inženjerski pojmovi *kontroler*, *kontrolirani sustav* i *kontrolni signal* redom.

### 1.1. Agent i okolina

Pomoću *nagradnog signala* definira se cilj problema PU-a. U svakom koraku okolina šalje agentu broj koji se naziva *nagrada*. Agentov je jedini cilj maksimizirati ukupnu nagradu - nagradnim signalom se onda može definirati koja su ponašanja poželjna, a koja ne. Nagradni signal je temeljni element prilagođavanja strategije agenta - ako neka akcija daje malenu nagradu može se prilagoditi strategiju da drugi put u takvoj situaciji agent isproba neku drugu akciju.

Nagradni signal govori što je dobro u datom trenutku, dok *vrijednosna funkcija* kaže što je dobro dugoročno gledano. Drugim riječima, *vrijednost* nekog stanja u kojem se agent nalazi je ukupna količina nagrade koju agent može prikupiti ako se nalazi u tom stanju - znači uzima se u obzir trenutna nagrada koja se dobije u tom stanju i stanja do kojih se može doći iz njega, zajedno sa njihovim nagradama. Upravo se **vrijednost** najviše uzima u obzir kada se vrši procjena odluke. Poželjne su one akcije koje će rezultirati stanjem s najvećom vrijednosti, a ne samo najvećom nagradom jer vrijednost stanja može dugoročno donijeti veću nagradu. Nažalost, vrijednosti je puno teže procijeniti. Dok nagrada proizlazi neposredno iz okoliša, vrijednosti je potrebno procjenjivati iz niza opservacija koje agent učini tokom svog rada. Efikasno procjenjivanje vrijednosne funkcije je najvažniji dio gotovo svih algoritama PU-a.

*Model okoline* nije nužan element PU-a. Model govori kako će okolina reagirati na akciju agenta - u kojem će se stanju naći i koliku će nagradu dobiti. Modeli zapravo služe za planiranje budućih akcija.

## 1.2. Markovljev proces odlučivanja

### 1.2 Markovljev proces odlučivanja

Markovljev proces odlučivanja (MDP) je klasična formalizacija slijednog donošenja odluka kada akcija ne utječe samo na trenutnu nagradu, već i sve sljedeće situacije ili stanja i s njima povezane nagrade.

**Definicija 1.1 (MDP)** *MDP je uređena trojka  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}_0)$  gdje je*

*$\mathcal{S}$  (prebrojiv) neprazan skup čije elemente nazivamo stanja;*

*$\mathcal{A}$  (prebrojiv) neprazan skup čije elemente nazivamo akcije;*

*$\mathcal{P}_0$  matrica prijelaza, koja svakom paru  $(s, a) \in \mathcal{S} \times \mathcal{A}$  pridružuje vjerojatnosnu mjeru povrh  $\mathcal{S} \times \mathbb{R}$  koju označavamo sa  $\mathcal{P}_0(\cdot|s, a)$ .*<sup>2</sup>

Kao što je već spomenuto MDP je alat za modeliranje problema slijednog donošenja odluka. Ako je dan MDP  $\mathcal{M}$ , interakcija izgleda ovako: Neka je  $t \in \mathbb{N}$  trenutno vrijeme (razina, trenutak), neka su  $S_t \in \mathcal{S}$  i  $A_t \in \mathcal{A}$  nasumično odabранo stanje sustava i akcija koju će agent poduzeti u vremenu  $t$  redom. U idućem koraku, kao posljedicu njegove akcije, agent dobija numeričku nagradu  $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$ , i nalazi se u novom stanju  $S_{t+1}$ . Ovim postupkom MDP i agent nam daju tzv. putanju:  $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$

U Markovljevom procesu odlučivanja vjerojatnosti koje daje matrica prijelaza  $\mathcal{P}_0$  u potpunosti opisuju dinamiku okoline - tj. vjerojatnosti za sve vrijednosti  $S_t$  i  $R_t$  ovise o neposredno prethodnom stanju  $S_{t-1}$  i akciji  $A_{t-1}$ <sup>3</sup>, i samo o njima, uopće se ne obazire na stanja i akcije koje im prethode. Ovo je najbolje promatrati kao svojstvo stanja, a ne cjelokupnog procesa.

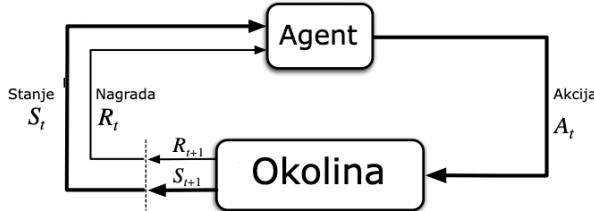
---

<sup>2</sup>što znači da za skup  $\mathcal{U} \subset \mathcal{S} \times \mathbb{R}$ ,  $\mathcal{P}_0(\mathcal{U}|s, a)$  govori kolika je vjerojatnost da će iduće stanje i njemu pripadajuća nagrada biti u skupu  $\mathcal{U}$ , ako se agent trenutno nalazi u stanju  $s$  i poduzme akciju  $a$ .

<sup>3</sup>Skup  $\mathcal{U}$  iz prethodne napomene je u ovom slučaju jednočlan, tj.  $\mathcal{U} = \{S_t\} \times \{R_t\}$ , pa pišemo  $\mathcal{P}_0(S_t, R_t | S_{t-1}, A_{t-1})$ .

## 1.2. Markovljev proces odlučivanja

Stanje u sebi mora sadržavati sve potrebne informacije o prošlim interakcijama agenta i sustava koje imaju utjecaj na budućnost - u tom slučaju kaže se da stanje ima *Markovljevo svojstvo*.



Slika 1.1: Interakcija agenta i okoline u MDP-u [2]

Iz funkcije  $\mathcal{P}_0$  može se izračunati sve o okolini što bi moglo biti korisno ili zanimljivo, kao na primjer *matricu prijelaza stanja*  $\mathcal{P} : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$

$$\mathcal{P}(s'|s, a) = \mathbb{P}[S_t = s' | S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} \mathcal{P}_0(s', r | s, a)$$

Također je moguće i *očekivanu nagradu* izraziti kao funkciju stanja i akcije  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ :

$$r(s, a) = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} \mathcal{P}_0(s', r | s, a)$$

Kako je već spomenuto, agentov je cilj maksimizirati sveukupnu nagradu - znači ne trenutnu, već kumulativnu, dugoročnu nagradu. Iako se ovo možda čini limitirajuće, u praksi se pokazalo kao fleksibilno i široko primjenjivo. Na primjer ako agent treba naučiti kako pobjeći iz labirinta, a "nagrada" je -1 za svaki korak koji agent napravi prije nego li uspije izaći - na ovaj način agent će nastojati pobjeći iz labirinta što prije.

Budući da agent uvijek nalazi način kako maksimizirati svoju nagradu, a cilj je da agent izvrši određeni zadatak, onda mu za taj zadatak mora biti ponuđena nagrada na način da maksimizirajući nagrade agent usput postigne

## 1.2. Markovljev proces odlučivanja

i željeni cilj. Od iznimne je važnosti da se nagrade uistinu podudaraju s onim što se želi postići - npr. agent koji uči igrati šah bi trebao biti nagrađen samo za pobjeđivanje, ne i za postizanje određenih podciljeva kao što su uzimanje protivničkih figura ili zauzimanje sredine ploče. Nagrađivanjem podciljeva može doći do situacije da agent nađe način da maksimizira svoju nagradu bez postizanja želenog cilja. Npr. može uzeti neprijateljsku figuru pod cijenu gubljenja partije.

Ako je niz nagrada koje se dobiju nakon koraka  $t$  označen s  $R_{t+1}, R_{t+2}, R_{t+3}\dots$ , koji se točno dio ovoga niza želi maksimizirati? Općenito agent želi maksimizirati *očekivanu dobit* - eng. *expected return*, gdje se dobit, u oznaci  $G_t$ , definira kao funkciju niza nagrada. Najjednostavnije je uzeti sumu nagrada  $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$ , gdje je  $T$  zadnji korak. Ovo ima smisla u primjerima koji u sebi po prirodi stvari podrazumijevaju neki zadnji korak, tj. kada se interakciju agenta i okoline može na prirodan način podijeliti u tzv. *epizode*: partija igre, jedan prolaz kroz labirint, bilo koja vrsta ponavljane interakcije. Svaka epizoda završava u posebnom stanju koje nazivamo *terminalno stanje*, nakon čega se vraća na početno stanje (ili na neko od početnih stanja ako ih ima više). Svaka epizoda je neovisna od prethodne, bez obzira kako je ona završila - pobjedom ili porazom.

Iako ovakvi epizodni scenariji pokrivaju puno slučajeva, od posebnog su interesa i kontinuirane interakcije koje se nastavljaju bez jasnog kraja ili zadnjeg koraka. U tom slučaju je  $T = \infty$  pa se uvodi koncept *korekcije*, eng. *discounting*. Agent bira akcije na način da suma nagrada s korekcijom bude maksimalna, tj. bira  $A_t$  da bi maksimizirao korigiranu dobit

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1.1)$$

parametar  $\gamma$ ,  $0 \leq \gamma \leq 1$ , nazivamo *korekcijski faktor*.

### 1.3. Funkcija vrijednosti

Korekcijski faktor određuje trenutnu vrijednost budućih nagrada: nagrada koju može dobiti za  $k$  koraka je vrijedna za faktor  $\gamma^{k-1}$  manje nego da je primi sada. Ako je  $\gamma < 1$  i skup  $\mathcal{R}$  omeđen, onda je beskonačna suma u (1.1) isto omeđena. Ako je  $\gamma = 0$ , onda je taj agent "kratkovidan" i zanima ga samo maksimiziranje trenutne nagrade: zanima ga samo kako odabratи  $A_t$  kako bi maksimizirao samo  $R_{t+1}$ . Općenito, ovo će spriječiti agenta od dobijanja nekih budućih odgođenih nagrada i smanjiti njegovu ukupnu dobit. Koliko se  $\gamma$  približava 1, to taj agent sve više uzima u obzir buduće nagrade - postaje "dalekovidniji".

## 1.3 Funkcija vrijednosti

Gotovo svi algoritmi PU-a se temelje na određivanju funkcije vrijednosti - funkcija koja nam govori "koliko je dobro" za agenta da bude u određenom stanju. Vrijednosne funkcije definiramo s obzirom na ponašanje agenta, odnosno *strategije* koju agent slijedi. Formalno, *strategija* je mapiranje koje svakom stanju pridruži distribuciju na skupu akcija<sup>4</sup>. Ako agent slijedi strategiju  $\pi$  u trenutku  $t$ , onda nam  $\pi(a|s)$  kaže kolika je vjerojatnost da je  $A_t = a$  ako je  $S_t = s$ , tj. ako se u trenutku  $t$  agent nalazi u stanju  $s$ , kolika je vjerojatnost da će poduzeti akciju  $a$ .

Vrijednosna funkcija stanja  $s$  u strategiji  $\pi$ , u oznaci  $v_\pi(s)$ , je očekivana dobit agenta koji počevši od stanja  $s$  slijedi strategiju  $\pi$ . Formalno se za MDP-jeve  $v_\pi$  može definirati sa

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right],$$

za sve  $s \in \mathcal{S}$ . Funkciju  $v_\pi$  se naziva *vrijednosna funkcija stanja* za *strategiju*  $\pi$ .

---

<sup>4</sup>tj. za svaku se akciju dobije kolika je vjerojatnost da će biti odabrana.

### 1.3. Funkcija vrijednosti

Slično, ako se u stanju  $s$  odabere akcija  $a$  slijedeći strategiju  $\pi$  definira se vrijednost

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

Funkcija  $q_\pi$  zove se *vrijednosna funkcija akcije za strategiju  $\pi$* .

Vrijednosne funkcije  $v_\pi$  i  $q_\pi$  mogu se procijeniti iz iskustva - ako agent slijedi strategiju  $\pi$  i za svako stanje prati ukupnu dobit koja slijedi nakon tog stanja, onda ti prosjeci konvergiraju k vrijednosti toga stanja  $v_\pi(s)$  kada broj posjeta tom stanju teži prema beskonačno. Ako se uz to odvojeno prati prosjeke dobiti za svaku akciju poduzetu iz svakog stanja na analogan način procjenjuje se  $q_\pi(s, a)$ . Ovakve tehnike procjene nazivaju se *Monte Carlo metode* jer se temelje na prosjeku mnogih slučajnih uzoraka stvarnih dobiti.

Iz otprije navedene definicije *očekivane dobiti* može se primijetiti da ju je moguće izraziti preko rekurzivne relacije

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma[R_{t+2} + \gamma^1 R_{t+3} + \gamma^2 R_{t+4} + \dots] \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned} \tag{1.2}$$

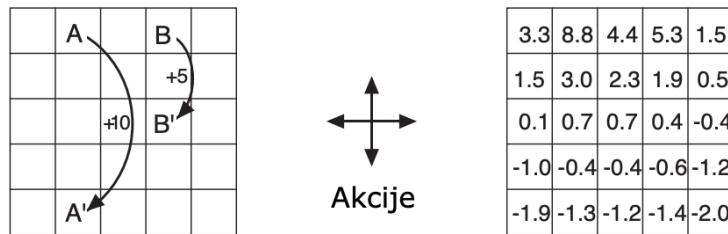
Slično vrijedi i za vrijednosne funkcije. Za proizvoljnu strategiju  $\pi$  i stanje  $s$  vrijedi

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r \mathcal{P}_0(s', r | s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']] \\ &= \sum_a \pi(a|s) \sum_{s', r} \mathcal{P}_0(s', r | s, a) [r + \gamma v_\pi(s')], \forall s \in \mathcal{S} \end{aligned} \tag{1.3}$$

#### 1.4. Primjer: Mrežni svijet

Jednadžbu (1.3) nazivamo *Bellmanova jednadžba* za  $v_\pi$  - ona izražava vezu između vrijednosti trenutnog stanja i vrijednosti stanja koje slijedi nakon njega. Bellmanova jednadžba uzima u obzir sva stanja koja okolina može proizvesti, pridajući svakom vjerojatnost hoće li se dogoditi. Lako se pokaže da je vrijednosna funkcija  $v_\pi$  jedinstveno rješenje pripadne Bellmanove jednadžbe. Kroz tri sljedeća poglavlja koristit će se Bellmanove jednadžbe kao temelj za izračun, odnosno procjenu,  $v_\pi$ .

## 1.4 Primjer: Mrežni svijet



Slika 1.2: *Mrežni svijet: dinamika nagrada (lijevo) i vrijednosna funkcija stanja (desno) pri korištenju strategije s jednakom vjerojatnošću kretanja u sva četiri smijera [1]*

Slika (1.2) predstavlja model jednostavnog konačnog MDP-a. Svako polje na ploči predstavlja jedno od stanja okoline. Na svakom polju agent može birati između četiri akcije: gore, dolje, lijevo i desno. Svaka akcija deterministički pomiče agenta u tom smjeru na ploči. Akcija koja bi pomakla agenta sa ploče, prelazak ruba ploče, ne mijenja položaj agenta, ali nosi sa sobom nagradu -1. Sve druge akcije donose nagradu 0, osim akcija koje se poduzmu u posebnim stanjima  $A$  i  $B$ . U stanju  $A$  sve četiri akcije donose

#### 1.4. Primjer: Mrežni svijet

agentu nagradu +10 i prebacuju ga u stanje  $A'$ . U stanju  $B$  sve četiri akcije donose agentu nagradu +5 i prebacuju ga u stanje  $B'$ .

Uz pretpostavku da agent svaku od akcija bira s jednakom vjerojatnošću na slici (1.2) s desne strane mogu se vidjeti vrijednosti stanja za danu strategiju s korekcijom  $\gamma = 0.9$ . Vrijednosna funkcija je izračunata rješavanjem sustava jednadžbi (1.3) i lako se provjeri vrijede li, npr. za vrijednost 0.7 u samom središtu tablice vrijedi

$$\begin{aligned} & 1/4 * (0 + 0.9 * 2.3) + 1/4 * (0 + 0.9 * 0.4) + \\ & 1/4 * (0 - 0.9 * 0.4) + 1/4 * (0 + 0.9 * 0.7) = \\ & 1/4 * 0.9[2.3 + 0.4 - 0.4 + 0.7] = \\ & 0.67499 \approx 0.7 \end{aligned}$$

1/4 je vjerojatnost odabira svake akcije (pomoću nje se izračunava očekivanje), 0 u svakoj od zagrada je izravna nagrada za odabir te akcije, 0.9 korekcija svake vrijednosti idućeg stanja i na kraju se vrijednost zaokružuje na jednu decimalu.

Negativne vrijednosti stanja na donjem rubu su posljedica velike vjerojatnosti da će se prijeći rub ako se biraju podjednako sve akcije. Stanje  $A$  je najbolje stanje koju god strategiju odabrali - ali je njegova vrijednost manja od nagrade koja se sigurno dobije u njemu jer se iz stanja  $A$  prelazi direktno u stanje  $A'$  u kojem je visoka vjerojatnost da će agent prijeći rub ploče. S druge strane, vrijednost stanja  $B$  je veća nego nagrada koja se u njemu dobije jer ono vodi u stanje  $B'$  čija je vrijednost također pozitivna - vjerojatnost kazne (negativna nagrada) za prelazak granice iz  $B'$  niža je od vjerojatnosti da će agent opet nabasati na stanja  $A$  ili  $B$ .

U ovom su primjeru dani pozitivni skalari kao nagrade za postizanje cilja, odnosno negativni skalari za prelazak ruba kao svojevrsna kazna i nula u svim ostalim slučajevima. Međutim, je li bitan predznak nagrada koje dajemo ili

#### 1.4. Primjer: Mrežni svijet

samo razlika između njih? Ako se na svaku nagradu u ovom primjeru doda neka konstanta  $c$  uvrštavanjem u (1.1) dobije se:

$$G_t = \sum_{k=0}^{\infty} \gamma^k [R_{t+k+1} + c] = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + c \sum_{k=0}^{\infty} \gamma^k \quad (1.4)$$

Kako se u pravilu uzima  $\gamma$  manja od 1, to je desni pribrojnik konvergentni geometrijski red. Dakle dodavanje konstante  $c$  svakoj nagradi pomiče vrijednost svakog stanja za  $v_c = c \sum_{k=0}^{\infty} \gamma^k$  što ne utječe na odnos vrijednosti stanja - ovo vrijedi za sve strategije.

##### 1.4.1 Optimalna strategija

Riješiti problem PU, ugrubo govoreći, znači pronaći strategiju koja dugoročno prikupi puno nagrada - što nasumično pogađanje i pretpostavljanje koje se do sada koristilo definitivno nije postiglo. U slučaju konačnih MDP-ova optimalnu strategiju može se i strogo definirati. Prvo na skupu svih strategija definira se uređaj uz pomoć njihovih vrijednosnih funkcija.

**Definicija 1.2** Neka su za MDP  $\mathcal{M}$  sa skupom stanja  $\mathcal{S}$  dane strategije  $\pi$  i  $\pi'$  sa vrijednosnim funkcijama  $v_\pi$  i  $v_{\pi'}$  redom. Za strategiju  $\pi$  kažemo da je bolja ili jednaka strategiji  $\pi'$ , i pišemo  $\pi \geq \pi'$ , ako za svaki  $s \in \mathcal{S}$  vrijedi  $v_\pi(s) \geq v_{\pi'}(s)$ .

Važno je uočiti da uvijek postoji barem jedna strategija koja je bolja ili jednaka od svih ostalih<sup>5</sup>. Iako možda postoji više takvih strategija sve ih se označava sa  $\pi_*$  i nazivaju se *optimalna strategija*.

---

<sup>5</sup>Za svako stanje  $s \in \mathcal{S}$  postoji neka strategija  $\pi_s$  za koju je vrijednost stanja  $s$  najveća - tj. veća ili jednaka od vrijednosti tog stanja za bilo koju drugu strategiju. Sada se jednostavno konstruira nova strategija  $\pi_*$  koja će se u svakom stanju  $s \in \mathcal{S}$  ponašati kao strategija  $\pi_s$ . Lako se vidi da je strategija  $\pi_*$  bolja ili jednaka od svih drugih.

#### 1.4. Primjer: Mrežni svijet

Sve optimalne strategije imaju istu vrijednosnu funkciju čiji je naziv *optimalna funkcija stanja* i definira se sa

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

za sve  $s \in \mathcal{S}$ .

Uz to sve optimalne strategije imaju i zajedničku *optimalnu vrijednosnu funkciju akcije* definiranu sa

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

za sve  $s \in \mathcal{S}$  i sve  $a \in \mathcal{A}$ .

U kontekstu funkcije  $v_*$  lako se uočiti da vrijedi

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]$$

Jer je  $v_*$  vrijednosna funkcija za neku strategiju, ona mora zadovoljavati uvjete dane Bellmanovim jednadžbama (1.3). Ali budući da je  $v_*$  optimalna vrijednosna funkcija ti se uvjeti mogu zapisati na način da se ne referira na iti jednu specifičnu strategiju - ovaj se zapis zove *Bellmanova jednadžba optimalnosti*:

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi*}(s, a) \quad (1.5)$$

$$\begin{aligned} &= \max_a \mathbb{E}_{\pi*}[G_t | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi*}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \end{aligned} \quad (1.6)$$

$$= \max_a \sum_{s', r} \mathcal{P}_0(s', r | s, a) [r + \gamma v_*(s')], \quad (1.7)$$

#### 1.4. Primjer: Mrežni svijet

za sve  $s \in \mathcal{S}$

Bellmanova jednadžba optimalnosti na intuitivnoj razini govori da vrijednost stanja u optimalnoj strategiji mora biti jednaka očekivanoj dobiti za odabir najbolje akcije u tom stanju

(1.6) i (1.7) su dva oblika Bellmanove jednadžbe optimalnosti za  $v_*$ . Za  $q_*$  Bellmanova jednadžba optimalnosti je

$$\begin{aligned} q_*(s, a) &= \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a] \\ &= \sum_{s', r} \mathcal{P}_0(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')] \end{aligned} \quad (1.8)$$

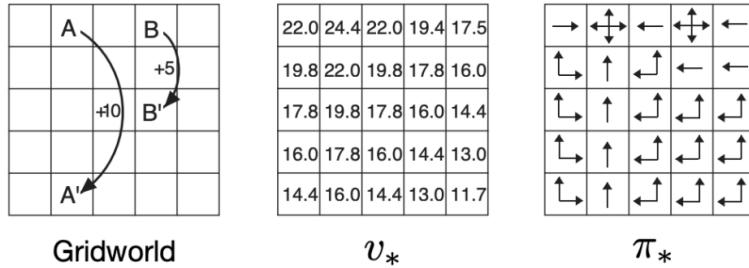
Za konačne MDP-ove Bellmanova jednadžba optimalnosti (1.7) za  $v_*$  ima jedinstveno rješenje - radi se zapravo o sustavu jednadžbi, po jedna za svako stanje. Ako je poznata dinamika okoliša, tj. matrica prijelaza  $\mathcal{P}_0$ , sustav je općenito rješiv bilo kojom metodom za rješavanje sustava nelinearnih jednadžbi. Kada se radi izračun za  $v_*$ , lako je odrediti optimalnu strategiju - za svako stanje  $s$  postojat će barem jedna akcija koja daje maksimum postignut u Bellmanovoj jednadžbi optimalnosti. Bilo koja strategija koja sa vjerojatnošću većom od nula bira te i samo te akcije bit će optimalna. Analogno zaključivanje bi bilo da se radi o *pohlepnom*<sup>6</sup> ponašanju s obzirom na optimalnu funkciju vrijednosti  $v_*$ .

Kada se ovo primjeni na spomenuti primjer Mrežnog svijeta rješavanjem Bellmanove jednadžbe za  $v_*$ : može se vidjeti na slici (1.3), u sredini su prikazane optimalne vrijednosti za svako stanje te desno optimalna strategija za ovaj primjer.

---

<sup>6</sup>Pohlepno donošenje odluka bi značilo biranje opcija uzimajući u obzir samo lokalna tj. trenutna saznanja bez obzira na dalekosežne posljedice - općenito znači odabir **trenutno** najboljih opcija. Ljepota  $v_*$  leži u tome da kada se koristi za određivanje trenutnih posljedica akcija, pohlepna je strategija zapravo optimalna jer  $v_*$  sadrži u sebi i sve moguće buduće nagrade.

#### 1.4. Primjer: Mrežni svijet



Slika 1.3: Mrežni svijet: optimalna rješenja [1]

Iako rješavanje Bellmanove optimalne jednadžbe daje optimalnu strategiju za problem PU-a, ova metoda je sama po sebi rijetko korisna. Da bi ju se uopće razmotrilo trebaju biti ispunjena tri uvjeta : (1) potpuno poznavanje okoline, (2) dovoljan broj računalnih resursa za izračun; (3) sva stanja imaju Markovljevo svojstvo;

U stvarnom svijetu rijetko je slučaj da su sva tri uvjeta ispunjena, a najčešće nije nijedan. Npr. za partiju igre Backgammon, isto analogno vrijedi i za Go ili Šah, iako vrijede uvjeti (1) i (3), uvjet (2) je nepremostiva prepreka: igra ima  $10^{20}$  različitih stanja i najmoćnijim današnjim računalima bi trebale tisuće godina za rješavanje pripadnih Bellmanovih jednadžbi.

U podržanom učenju je stoga potrebno zadovoljiti se približnim rješenjima čime će se pozabaviti poglavljia koja slijede.

## Poglavlje 2

# Dinamičko programiranje

Dinamičko programiranje (DP) se odnosi na niz algoritama koji se koriste za izračun optimalne strategije u slučaju da postoji savršen model okoliša kao MDP-a. Iako od klasičnih algoritama DP-a nema velike koristi u praksi - zbog zahtjeva za savršenim modelom okoline ili prevelikih računalnih zahtjeva - ipak su važni iz teoretske perspektive jer sve ostale metode, koje će biti naknadno spomenute, u nekoj se mjeri temelje na DP-u: pokušavaju doći do istih rezultata kao DP samo s manje računanja ili bez savršenog modela okoline. U ovom poglavlju, a i nadalje, pretpostavka je da se radi s konačnim MDP-om, tj. da su svi skupovi MDP-a konačni.

Općenito govoreći, osnovna je ideja DP-a, i PU-a korištenje vrijednosnih funkcija kako bi se olakšala potraga za optimalnom strategijom. Fokus ovog poglavlja je prikaz načina korištenja DP-a za izračun vrijednosnih funkcija definiranih u prethodnom poglavlju. Već je otprije poznato, kada su jednom utvrđene optimalne vrijednosne funkcije  $v_*$  ili  $q_*$  lako se dolazi do optimalne strategije. Da bi se došlo do optimalnih vrijednosnih funkcija DP se koristi Bellmanovim optimalnim jednadžbama (1.5) i (1.8) kao pravilom ažuriranja kojim poboljšava aproksimacije funkcije koju računa.

## 2.1. Vrednovanje strategije (predviđanje)

### 2.1 Vrednovanje strategije (predviđanje)

Prvo je vrijedno razmotriti kako izračunati vrijednosnu funkciju stanja  $v_\pi$  za proizvoljnu strategiju  $\pi$ . U literaturi se ovo zove *vrednovanje strategije* ili *problem predviđanja*. U prošlom je poglavlju pokazano za sve  $s \in \mathcal{S}$

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s',r} \mathcal{P}_0(s',r|s,a) [r + \gamma v_\pi(s')], \end{aligned} \quad (2.1)$$

gdje je  $\pi(a|s)$  vjerojatnost da će agent poduzeti akciju  $a$  ako se nalazi u stanju  $s$  i slijedi strategiju  $\pi$ . Očekivanja imaju index  $\pi$  da bi naznačili uvjetovanost o  $\pi$  - tj. da vrijedi ako se slijedi strategiju  $\pi$ .

Funkcija  $v_\pi$  postoji i jedinstvena je ako je  $\gamma < 1$  ili agent sigurno dolazi do terminalnog stanja - do eventualnog cilja ili kraja istraživanja - odakle god počeli, slijedeći strategiju  $\pi$ .

Ako je dinamika okoliša u potpunosti poznata, onda vrijedi (2.2) sustav od  $|\mathcal{S}|$  jednadžbi s  $|\mathcal{S}|$  nepoznanica. Iako rješavanje sustava može biti računalno zahtjevno, ono je općenito dosta jednostavno s teoretske strane. U ovom su slučaju najprikladnije iterativne metode: Promatra se niz aproksimacija vrijednosnih funkcija  $v_0, v_1, v_2, \dots$ , gdje svaka ide iz  $\mathcal{S}$  u  $\mathbf{R}$ . Prva se aproksimacija uzima proizvoljno (uz pretpostavku da je vrijednost terminalnih stanja, ako ih ima, 0), a sve sljedeće se dobiju koristeći Bellmanove jednadžbe za  $v_\pi$  (2.2) kao pravilo ažuriranja:

## 2.1. Vrednovanje strategije (predviđanje)

$$\begin{aligned} v_{k+1}(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_k(S_{t+1})|S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma v_k(s')], \end{aligned} \quad (2.2)$$

, za sve  $s \in \mathcal{S}$ .

Očito je  $v_k = v_\pi$  fiksna točka ovog pravila ažuriranja jer Bellmanove jednadžbe za  $v_\pi$  osiguravaju jednakost u tom slučaju. Općenito, može se pokazati da niz  $v_k$  konvergira ka  $v_\pi$  kada  $k \rightarrow \infty$  pod istim uvjetima koji garantiraju egzistenciju  $v_\pi$ . Ovaj se algoritam naziva *iterativno vrednovanje strategije*.

---

### Algoritam 1: Iterativno vrednovanje strategije

---

**Podatci:** strategija  $\pi$  koju treba vrednovati

**Odredi:** mali  $\epsilon > 0$  kojim određujemo željena preciznost

Za sve  $s \in \mathcal{S}$  proizvoljno postavi  $V(s)$  uz uvjet da je

$V(\text{terminal}) = 0$ ;

$\Delta \leftarrow \epsilon + 1$ ;

**dok**  $\Delta > \epsilon$  **čini**

$\Delta \leftarrow 0$ ;

**za**  $s \in \mathcal{S}$  **čini**

$v \leftarrow V(s)$ ;

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma V(s')]$ ;

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$ ;

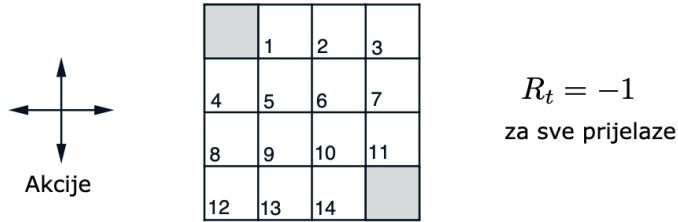
**kraj**

**kraj**

---

**Primjer 2.1 (Mrežni svijet: 4x4)** U mrežnom svijetu sa slike ([1]) vrijedi: Neterminalna stanja su  $\mathcal{S} = 1, 2, \dots, 14$ , za svako stanje moguće su 4

## 2.1. Vrednovanje strategije (predviđanje)



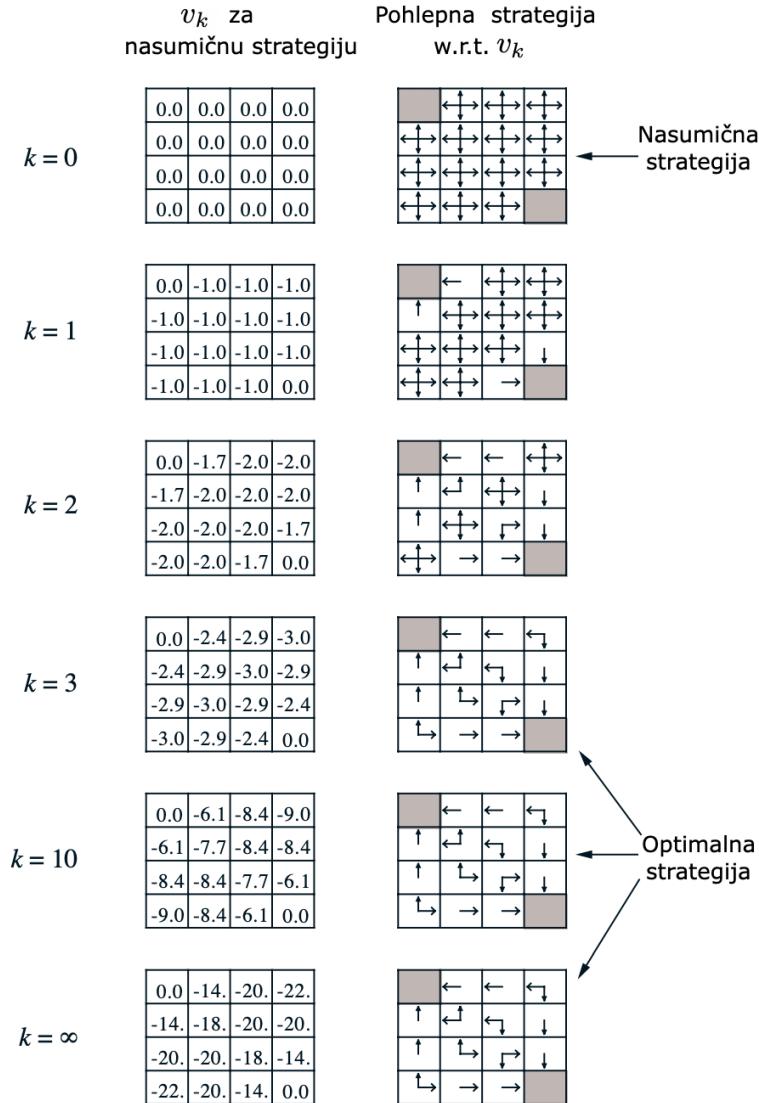
akcije  $\mathcal{A} = \{gore, dolje, lijevo, desno\}$  koje deterministički prebacuju agenta u pripadno stanje, osim ako bi prešao rub ploče, u tom slučaju stanje ostaje nepromijenjeno.

Nagrada iznosi  $-1$  za sve prijelaze dok agent ne dođe do terminalnog stanja (siva polja na ploči).

Dakle, vrijedi npr.  $\mathcal{P}_0(6, -1|5, desno) = 1$ ,  $\mathcal{P}_0(7, -1|7, desno) = 1$  i  $\mathcal{P}_0(10, r|5, desno) = 0$  za sve  $r \in \mathcal{R} \subset \mathbf{R}$ .

Na lijevoj strani slike (2.1), na stranici 18, vidi se niz vrijednosnih funkcija  $\{v_k\}$  dobivenih iterativnim vrednovanjem nasumične strategije - agent s jednakom vjerojatnosti, 0.25, bira svaku od akcija.

## 2.1. Vrednovanje strategije (predviđanje)



Slika 2.1: Konvergencija iterativnog vrednovanja strategije. Lijevi stupac prikazuje niz aproksimacija vrijednosne funkcije stanja za jednakovjerojatnu nasumičnu strategiju. U desnom stupcu vidi se niz pohlepnih strategija temeljenih na izračunatim procjenama vrijednosnih funkcija. Strategija dobivena na kraju je garantirano bolja od nasumične, odnosno one na kojoj su temeljene procjene funkcija, ne nužno i općenito optimalna - ali u ovom slučaju je. [1]

## 2.2. Poboljšavanje strategije (kontrola)

### 2.2 Poboljšavanje strategije (kontrola)

Razlog zašto uopće računati vrijednosnu funkciju neke strategije je da bi se pomoću nje moglo pronaći bolje strategije, tzv. *kontrola*. Uz pretpostavku da je izračunata vrijednosna funkcija  $v_\pi$  neke determinističke strategije  $\pi$ , za neko stanje  $s$  vrijedi znati isplati li se promijeniti strategiju da se u tom stanju izabere akciju  $a \neq \pi(s)$ . Ono što je moguće učiniti da bi se odgovorilo na ovo pitanje je izabirati akciju  $a$  u stanju  $s$ , a u svim ostalim slučajevima slijediti strategiju  $\pi$ . Vrijednost ovakvog načina ponašanja je

$$q_\pi(s, a) = \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] = \sum_{s', r} \mathcal{P}_0(s', r | s, a)[r + \gamma v_\pi(s')],$$

Pitanje je je li ova vrijednost veća ili manja od  $v_\pi(s)$ . U slučaju da je veća, može se zaključiti da je novonastala strategija bolja od prijašnje. Točnost takvog zaključivanja je posljedica općenitijeg rezultata poznatog kao *Teorem poboljšanja strategije*.

**Teorem 2.2 (poboljšanja strategije)** *Neka su  $\pi$  i  $\pi'$  strategije takve da za njih vrijedi*

$$\forall s \in \mathcal{S}, q_\pi(s, \pi'(s)) \geq v_\pi(s).$$

*Tada je strategija  $\pi'$  barem jednako dobra kao strategija  $\pi$ , tj. za svako stanje  $s \in \mathcal{S}$  vrijedi*

$$v_{\pi'}(s) \geq v_\pi(s).$$

**Dokaz.** Neka vrijedi (2.2), tj  $q_\pi(s, \pi'(s)) \geq v_\pi(s)$ , onda vrijedi i sljedeće

## 2.2. Poboljšanje strategije (kontrola)

$$\begin{aligned}
v_\pi(s) &\leq q_\pi(s, \pi'(s)) && \text{(primjeni (2.2))} \\
&= \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = \pi'(s)] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \\
&\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) | S_t = s] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}[R_{t+2} + \gamma v_\pi(S_{t+2}) | S_{t+1}, A_{t+1} = \pi'(S_{t+1})] | S_t = s] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_\pi(S_{t+2}) | S_t = s] \\
&\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_\pi(S_{t+3}) | S_t = s] \\
&\vdots \\
&\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots | S_t = s] \\
&= v_{\pi'}(s).
\end{aligned}$$

Sada se primjenjuje ista logika na *sva* stanja: u svakom stanju agent izabire akciju koja je najbolja po  $q_\pi(s, a)$ . Tj. formira se novu *pohlepnu* strategiju  $\pi'$  definiranu s

$$\begin{aligned}
\pi'(s) &= \operatorname{argmax}_a q_\pi(s, a) && \text{(2.3)} \\
&= \operatorname{argmax}_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\
&= \operatorname{argmax}_a \sum_{s', r} \mathcal{P}_0(s', r | s, a) [r + \gamma v_\pi(s')]
\end{aligned}$$

gdje  $\operatorname{argmax}_a$  označava vrijednost za koju je izraz maksimalan. Konstrukcija ovakve pohlepne strategije zadovoljava uvjete Teorema (2.2) tako da je jasno da je ova strategija barem jednako dobra kao početna strategija - možda čak i bolja.

### 2.3. Iteracija strategija

Očigledno, ako vrijedi pretpostavka da je nova pohlepna strategija  $\pi'$  jednako dobra kao i izvorna  $\pi$  i da **nije** bolja od nje, onda vrijedi  $v_\pi = v_{\pi'}$  te iz (2.3) slijedi:

$$\begin{aligned} v_{\pi'}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma v_\pi(s')], \forall s \in \mathcal{S}. \end{aligned} \quad (2.4)$$

Budući da je ovo isto što i Bellmanova jednadžba optimalnosti (1.5), onda  $v_{\pi'}$  mora biti i  $v_*$ . Pa su i  $\pi$  i  $\pi'$  optimalne strategije. Proces poboljšanja strategije će uvijek dati bolju strategiju od polazne, osim ako je polazna strategija već bila optimalna<sup>1</sup>.

## 2.3 Iteracija strategija

U ovoj se fazi iskoristio  $v_\pi$  da bi se poboljšalo strategiju  $\pi$  i došlo do strategije  $\pi'$ , pa se može nastaviti analogno, tj. izračunati vrijednosnu funkciju  $v_{\pi'}$  i iskoristiti je da bi došli do još bolje strategije  $\pi''$ . Na ovaj način dobije se niz monotono rastućih strategija (s obzirom na uređaj koji je definiran) i njima pripadajućih vrijednosnih funkcija

$$\pi_0 \xrightarrow{\mathbb{E}} v_{\pi_0} \xrightarrow{\mathbb{I}} \pi_1 \xrightarrow{\mathbb{E}} v_{\pi_1} \xrightarrow{\mathbb{I}} \pi_2 \xrightarrow{\mathbb{E}} \dots \xrightarrow{\mathbb{I}} \pi_* \xrightarrow{\mathbb{E}} v_{\pi_*}$$

gdje  $\xrightarrow{\mathbb{E}}$  predstavlja *vrednovanje* strategije, a  $\xrightarrow{\mathbb{I}}$  *poboljšanje* strategije - eng. *improvement*. Svaka iduća strategija je izričito bolja od prethodne (osim ako je ova već bila optimalna). Jer konačan MDP ima konačno mnogo determinističkih strategija, to ovaj proces sigurno konvergira ka optimalnoj strategiji i optimalnoj funkciji vrijednosti u konačno mnogo iteracija. Ovakav način pronalaska optimalne strategije naziva se *iteracija strategija*.

---

<sup>1</sup>Slično zaključivanje i dokazi se mogu primjeniti i za stohastičke strategije.

### 2.3. Iteracija strategija

---

**Algoritam 2:** Iteracija strategija za procjenu  $\pi \approx \pi_*$

---

**Inicijaliziraj:** za sve  $s \in \mathcal{S}$  (nasumično) postavi  $V(s) \in \mathbf{R}$  i

$$\pi(s) \in \mathcal{A}(s), V(\text{terminal}) = 0$$

*Vrednovanje strategije:*

$$\Delta \leftarrow \epsilon + 1;$$

**dok**  $\Delta > \epsilon$  **čini**

$$\Delta \leftarrow 0;$$

**za**  $s \in \mathcal{S}$  **čini**

$$v \leftarrow V(s);$$

$$V(s) \leftarrow \sum_{s',r} \mathcal{P}_0(s',r|s,\pi(s))[r + \gamma V(s')];$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|);$$

**kraj**

**kraj**

*Poboljšanje strategije:*

$$\text{strategija-stabilna} \leftarrow \text{true};$$

**za**  $s \in \mathcal{S}$  **čini**

$$\text{stara-akcija} \leftarrow \pi(s);$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma V(s')];$$

**ako**  $\text{stara-akcija} \neq \pi(s)$  **onda**

$$\text{strategija-stabilna} \leftarrow \text{false};$$

**kraj**

**kraj**

**ako**  $\text{strategija-stabilna}$  **onda**

$$\quad \quad \quad \text{zaustavi algoritam i vrati } V \approx v_* \text{ i } \pi \approx \pi_*;$$

**kraj**

**inače**

$$\quad \quad \quad \text{idi na } \text{Vrednovanje strategije};$$

**kraj**

## 2.4. Iteracija vrijednosti

Negativna strana iteracije strategije je to što svaka iteracija uključuje vrednovanje strategije - što može biti zahtjevan račun. Ako se vrednovanje strategije radi iterativno, onda se konvergencija ka optimalnoj funkciji vrijednosti  $v_*$  postiže tek na kraju, tj. u graničnoj vrijednosti. U primjeru na slici (2.1) vidljivo je da svako vrednovanje strategije nakon treće iteracije nema nikakav utjecaj na dobivenu pohlepnu strategiju - ovo sugerira da bi se neka vrednovanja strategija moglo preskočiti.

Štoviše, pokaže se da postoji nekoliko kriterija po kojima je moguće preskočiti vrednovanje strategije uz zadržavanje garantirane konvergencije ka optimalnoj strategiji.

Posebno je zanimljiv slučaj kada vrednovanje strategije staje nakon samo jednog prolaza po svim stanjima - ovaj se algoritam zove *iteracija vrijednosti*. Opisan je jednostavnim pravilom ažuriranja koje kombinira poboljšanje strategije i preskakanje koraka vrednovanja strategije:

$$\begin{aligned} v_{k+1}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma v_k(s')], \forall s \in \mathcal{S}. \end{aligned} \quad (2.5)$$

Za proizvoljno odabrani  $v_0$  niz  $\{v_k\}$  konvergira ka  $v_*$  po istim uvjetima koji garantiraju postojanje  $v_*$ . Algoritam se izvodi iz Bellmanovih jednadžbi (1.5) - iteracija se dobije tako što se Bellanova jednadžba pretvorí u pravilo ažuriranja.

Vrijedi spomenuti i način zaustavljanja algoritma. Kao i kod vrednovanja strategije, iteracije vrijednosti formalno trebaju beskonačno mnogo iteracija da bi konvergirale ka  $v_*$ . U praksi se zaustavljaju kada se vrijednosna funkcija u jednom prijelazu promijeni za dovoljno malu vrijednost. Precizniji zapisi su vidljivi u Algoritmu (3).

## 2.4. Iteracija vrijednosti

---

**Algoritam 3:** Iteracija vrijednosti za procjenu  $\pi \approx \pi_*$

---

**Inicijaliziraj:** za sve  $s \in \mathcal{S}$  (nasumično) postavi  $V(s) \in \mathbf{R}$ ,

$V(\text{terminal}) = 0$  te odaberi mali  $\epsilon > 0$  kojim se

određuje željenu preciznost

$\Delta \leftarrow \epsilon + 1;$

**dok**  $\Delta > \epsilon$  **čini**

$\Delta \leftarrow 0;$

**za**  $s \in \mathcal{S}$  **čini**

$v \leftarrow V(s);$

$V(s) \leftarrow \max_a \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma V(s')];$

$\Delta \leftarrow \max(\Delta, |v - V(s)|);$

**kraj**

**kraj**

Vrati determinističku strategiju  $\pi \approx \pi_*$  za koju je

$\pi(s) = \operatorname{argmax}_a \sum_{s',r} \mathcal{P}_0(s',r|s,a)[r + \gamma V(s')];$

---

## 2.5. Napomene

Iteracija vrijednosti zapravo kombinira, u svakom prolazu kroz stanja, prolaz vrednovanja strategije i prolaz poboljšanja strategije.

## 2.5 Napomene

### 2.5.1 Asinkrono DP

U slučaju da je skup stanja jako velik čak i jedan prolaz skupom može bit računalno jako "skup." Asinkroni DP algoritmi ne ovise o sistematskom ažuriranju **svih** stanja - ažuriraju se stanja proizvoljnim redom, koristeći pri tome dostupne vrijednosti stanja. Može se dogoditi da je vrijednost nekog stanja ažurirana više puta dok nekog drugog nije niti jednom. Međutim, izbjegavanje punog prolaza po stanjima ne znači da će se nužno proći jeftinije jer sve vrijednosti stanja treba prije ili kasnije ažurirati - ne mogu se neke jednostavno ignorirati - ali nije potrebno čekati da se izvrte sve vrijednosti da bi mogli unaprijediti tu strategiju.

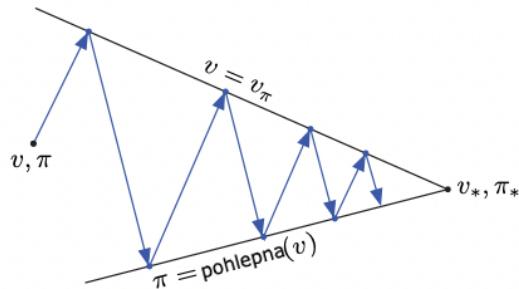
### 2.5.2 Generalizirana iteracija strategije

Iteracija strategije se sastoji od dva isprepletena procesa: izračun vrijednosne funkcije da odgovara trenutnoj strategiji (vrednovanje strategije) i formiranje pohlepne strategije s obzirom na trenutnu vrijednosnu funkciju (poboljšanje strategije). Ovi procesi alterniraju, da bi idući mogao započeti prethodni mora završiti - što zapravo i nije nužno. Kod asinkronih DP metoda ovi su procesi isprepleteni na finiji način i dokle god proces ažurira sva stanja konačni rezultati bi trebali biti isti - konvergencija ka optimalnoj vrijednosnoj funkciji i optimalnoj strategiji.

Pod nazivom *generalizirana iteracija strategija* (GPI) misli se na općenitu

## 2.5. Napomene

ideju međudjelovanja *vrednovanja strategije* i *poboljšanja strategije* neovisno o granularnosti i detaljima tih procesa. Skoro sve metode PU-a može se opisati kao GPI - sve imaju strategije i vrijednosne funkcije, s tim da se strategije poboljšavaju s obzirom na vrijednosne funkcije, a funkcije se računaju po strategiji.



Slika 2.2: GPI [2]

Jednom kad se oba procesa stabiliziraju, tj. nema više promjena, znači da je postignuta optimalna vrijednosna funkcija i optimalna strategija - a to se događa samo ako je vrijednosna funkcija konzistentna s trenutnom strategijom i trenutna je strategija pohlepna s obzirom na vrijednosnu funkciju.

Odnos vrednovanja i poboljšanje u GPI-u može se promatrati kao suparništvo i suradnja. Suparništvo jer svaki proces na neki način povlači na svoju stranu - formiranje pohlepne strategije s obzirom na trenutnu funkciju vrijednosti čini funkciju netočnom za novu strategiju, a izračun nove funkcije čini da strategija nije više pohlepna u odnosu na nju. Međutim, dugoročno ova dva procesa dolaze do jednog zajedničkog rješenja: optimalne vrijednosne funkcije i optimalne strategije.

# Poglavlje 3

## Predviđanje bez modela

Metode koje su dosad spomenute su jako korisne i vrijedne, no imaju jako veliki nedostatak koji je bio prisutan iako ne i istaknut - ove se metode mogu koristiti samo ako je dinamika okoliša u potpunosti poznata. Ako je smjer iz kojeg će vjetar zapuhati nepoznat, ili nisu poznati svi parametri koji utječu na produktivnost tvornice, onda se ne mogu ni koristiti. Upravo su zato korisne *metode predviđanja bez modela*, eng. model-free prediction, za procjenu vrijednosne funkcije i otkrivanje optimalne strategije. Postoje dvije glavne klase metoda predviđanja bez modela:

*Monte-Carlo učenje* - metode koji prijeđu cijelu putanju agenta i procjenjuju vrijednost iz dobiti uzoraka;

*Učenje s vremenskom razlikom*, eng. *temporal-difference learning* - metode koje gledaju jedan korak unaprijed i procjenjuju dobit nakon tog jednog koraka.

Dakle, odbacuje se ogromnu pretpostavku ponašanja okoliša - koja je ne-realna za većinu zanimljivih problema - i radi se s metodama koje koriste *iskustvo* agenta, tj. njegove interakcije s okolišem. Kao i do sada, metoda traženja optimalne strategije podijeljena je na dva dijela: vrednovanje strate-

### 3.1. Monte-Carlo Učenje

gije, tj. nalaženje vrijednosne funkcije za danu strategiju, i poboljšanje dane strategije. Predviđanje se bavi prvim dijelom - ako je zadana strategija, vrši se procjena koliko je dobra, tj. traži se njenu vrijednosnu funkciju.

## 3.1 Monte-Carlo Učenje

Pojam "Monte-Carlo" se općenito odnosi na metode procjene koje u sebi sadrže neki nasumični dio koji ih u značajnoj mjeri određuje. Ovdje se taj pojam koristi za metode koje računaju prosjek cijelokupne dobiti - za razliku od metoda koje promatraju djelomičnu dobit kojima se bavi iduće poglavlje.

MC koristi najjednostavniju ideju za slučaj da modela nema: da bi odredio vrijednost nekog stanja računa prosjek dobiti postignutih iz tog stanja. Ako su dana primjerice tri uzorka koja počinju iz istog stanja sa dobitima 7,3 i 8, prosjek tih vrijednosti je 6 pa je pretpostavka da je vrijednost toga stanja 6. Važno je uočiti da epizoda treba završiti kako bi se mogla izračunati dobit, tj. agent treba doći do terminalnog stanja. Kako broj uzoraka raste, aproksimacija teži stvarnoj vrijednosti tog stanja.

Cilj je procjeniti  $v_\pi(s)$ , vrijednost stanja  $s$  ako se slijedi strategiju  $\pi$ , a poznat je skup epizoda u kojima je agent slijedio strategiju  $\pi$  i prošao kroz stanje  $s$ . Svako pojavljivanje stanja  $s$  u epizodi naziva se *posjet* stanju  $s$ . Kako agent u jednoj epizodi može posjetiti isto stanje više puta:

*MC metoda prvog posjeta* procjenjuje vrijednost  $v_\pi(s)$  kao prosjek<sup>1</sup> dobiti od prvog posjeta stanju  $s$ , dok *MC metoda svakog posjeta* računa prosjek

---

<sup>1</sup> Inkrementalni izračun aritmetičke sredine (vidi algoritam (4)):

Arit. sredine  $\mu_1, \mu_2, \dots$  niza  $x_1, x_2, \dots$  može se računati inkrementalno:

$$\mu_k = \frac{1}{k} \sum_{j=1}^k x_j = \frac{1}{k} (x_k + \sum_{j=1}^{k-1} x_j) = \frac{1}{k} (x_k + (k-1)\mu_{k-1}) = \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

### 3.1. Monte-Carlo Učenje

nakon svakog posjeta stanju  $s$ . Iako su ove dvije MC metode jako slične među njima ipak postoje male teoretske razlike.

---

**Algoritam 4:** *MC metoda prvog posjeta za procjenu vrijednosti  $V \approx v_\pi$*

---

Zadana je strategija  $\pi$  koju treba procjeniti;

**Inicijaliziraj:** za sve  $s \in \mathcal{S}$  (nasumično) postavi  $V(s) \in \mathbf{R}$  ;

*Dobiti( $s$ )*  $\leftarrow$  prazna lista, za sve  $s \in \mathcal{S}$ ;

**dok** *true* **čini**

Generiraj epizodu sljedeći strategiju  $\pi$ :

$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ ;

$G \leftarrow 0$ ;

**za** *svaki korak*  $t \in [T - 1, T - 2, \dots, 0]$  **čini**

$G \leftarrow \gamma G + R_{t+1}$ ;

**ako** Stanje  $S_t$  se **ne** nalazi u nizu  $S_0, S_1, \dots, S_{t-1}$  **onda**

Dodaj  $G$  na kraj liste *Dobiti( $S_t$ )*;

$V(S_t) \leftarrow \text{prosjek}(\text{Dobiti}(S_t))$  ;

**kraj**

**kraj**

**kraj**

---

#### 3.1.1 Primjer: Blackjack

Cilj popularne casino igre je skupiti karte čija je suma što veća a da ne premašuje sumu od 21. Sve karte sa slikom vrijede 10, karte s brojem vrijede koliko pokazuju, osim asa koji vrijedi 1 ili 11 ovisno o tome što više odgovara igraču. Igrač igra protiv djelitelja (kuće, casina). Igra započinje tako što oboje dobiju po dvije karte s tim da je vidljiva samo jedna djeliteljeva karta,

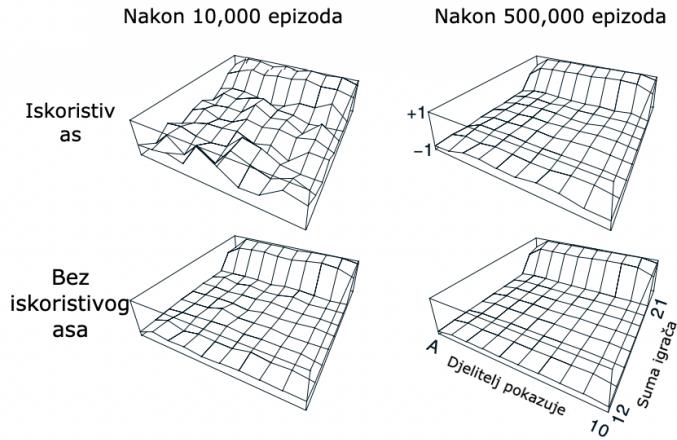
### 3.1. Monte-Carlo Učenje

dok drugu vidi samo djelitelj. Ako igrač skupi 21 automatski pobjeđuje, osim ako i djelitelj ima 21, u tom slučaju je izjednačeno. Ako je zbroj manji od 21, može se zatražiti još karata jednu po jednu dok ili igrač ne odluči stati ili zbroj na kartama premaši 21. Ako je 21 premašen, igrač gubi tu partiju igre. Ako igrač odluči stati, red je na djelitelju. Djelitelj uvijek igra istu strategiju: vuče karte dok ne postigne sumu 17 ili više. Ako djelitelj prebací 21, igrač pobjeđuje, inače pobjeđuje onaj koji je bliže 21.

Nagrade su sljedeće: +1, -1 i 0 za pobjedu u igri, poraz i izjednačeno redom. Sve akcije unutar jedne partije daju nagradu 0 i ne radi se korekcija ( $\gamma = 0$ ) tako da nagrade u terminalnom stanju zapravo daju dobit. U svakom potezu, tj. stanju, igrač može birati uzima li još jednu kartu ili staje - stanja ovise o kartama koje ima igrač i karti djelitelja koju svi vidimo. Pretpostavljamo da djelitelj izvlači karte iz beskonačnog špila, tj. da nema smisla brojati karte. Kada igrač ima asa kojeg može brojati kao 11 bez da prijeđe 21, onda se kaže da je taj as *iskoristiv* - u tom slučaju se broji kao 11, kad bi ga brojali kao 1 igrač bi imao manje od 11 pa je poznato koja je njegova iduća akcija. Sigurno će uzeti još jednu kartu jer nema što izgubiti. Znači, igrač donosi odluku temeljenu na tri činjenice: njegova trenutna suma (12-21), karta koju pokazuje djelitelj (as-10) i ima li iskoristivog asa. Sveukupno 200 stanja.

Neka je strategija agenta da stane ako ima sumu 20 ili 21, inače uvijek uzima još jednu kartu. Da bi otkrili vrijednosnu funkciju stanja ove strategije MC metodom simulira se velik broj partija igre koristeći ovu strategiju i računa se prosjek dobiti za svako stanje. Na ovaj način dobije se procjenu vrijednosne funkcije stanja prikazanog na slici (3.1.1) - procjena vrijednosti stanja kada igrač ima iskoristivog asa je manje sigurna jer je to stanje puno rjeđe, ali nakon 500000 partija vrijednosna funkcija je dosta precizno proci-

### 3.1. Monte-Carlo Učenje



Slika 3.1: Procjena vrijednosne funkcije stanja MC metodom za strategiju koja staje samo na 20 i 21 [2]

jenjena.

Iako je u ovom slučaju znanje o okolini potpuno, ipak nije jednostavno primijeniti algoritme DP-a na njega. DP metode pretpostavljaju poznatu distribuciju vjerojatnosti za idući događaj, tj. da je poznata matrica prijelaza  $\mathcal{P}_0$  - a nju nije lako odrediti za Blackjack. Sve moguće vjerojatnosti događaja bi trebalo izračunati prije nego li se primijene algoritmi DP-a, a ti izračuni su često komplikirani i podložni greškama. S druge strane, generiranje partija potrebnih za MC metode je iznimno lako i baš je ova činjenica velika prednost MC metode čak i ako je dinamika okoline u potpunosti poznata.

Dodatna prednost MC metoda je njena procjena vrijednosti za svako stanje koja je neovisna od ostalih - kao što je slučaj kod DP-a. Dalje, računalna cijena procjene jednog stanja **ne** ovisi o ukupnom broju stanja MDP-a - zbog ovoga su MC metode poželjne kada je bitna samo vrijednost jednog ili manjeg podskupa stanja.

### 3.2. Učenje s vremenskom razlikom

## 3.2 Učenje s vremenskom razlikom

Ideja koja se pokazala kao centralna za PU je učenje s vremenskom razlikom, eng. temporal-difference (TD) learning. TD je kombinacija ideja MC učenja i DP-a. Kao i MC, TD metode mogu učiti direktno iz iskustva bez modela koji opisuje dinamiku okoline. Poput DP-a, TD metode ažuriraju procjene dijelom na temelju unaprijed određenih procjena, ne čekajući na konačan rezultat. Odnos između TD-a, DP-a i MC-s je stalna tema diskusije u teoriji PU-a.

I TD i MC metode koriste iskustvo da bi riješili problem predviđanja. Ako je dostupno iskustvo dobiveno praćenjem strategije  $\pi$ , obje metode ažuriraju svoje pretpostavke o procjeni  $V$  vrijednosti  $v_\pi$  za neterminalno stanje  $S_t$  koje se pojavilo u tom iskustvu. MC metode čekaju do kraja kako bi postigle *dobit* koju zatim upotrijebe za poboljšanje  $V(S_t)$ . MC metodu za nestacionarne probleme opisuje se pravilom ažuriranja <sup>2</sup>

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

gdje je konstanta  $\alpha$  tzv. *parametar veličine koraka*. Dok MC metode čekaju do kraja epizode da bi odredile inkrement za  $V(S_t)$  (znači tek kada je poznat  $G_t$ ), TD metode čekaju samo do idućeg koraka. U trenutku  $t + 1$  odmah uz pomoć dobivene nagrade  $R_{t+1}$  i procjene za  $V(S_{t+1})$  rade pretpostavku. Najjednostavniji oblik TD-a koristi pravilo ažuriranja

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

čim prijeđe u stanje  $S_{t+1}$  i primi nagradu  $R_{t+1}$ . Znači, dok MC čeka na ukupnu dobit  $G_t$ , TD koristi procjenu  $R_{t+1} + \gamma V(S_{t+1})$ . Ova TD metoda poznata je kao TD *jednog koraka* ili TD(0).

---

<sup>2</sup>Primijenimo inkrementalni izračun aritmetičke sredine iz fusnote (1)

### 3.2. Učenje s vremenskom razlikom

---

**Algoritam 5:**  $TD(0)$  za procjenu vrijednosti  $V \approx v_\pi$

---

Dana je strategija  $\pi$  koju treba procijeniti

**Inicijaliziraj:** za sve  $s \in \mathcal{S}$  (nasumično) postavi  $V(s) \in \mathbf{R}$ ;

Odaberi parametar algoritma  $\alpha \in \langle 0, 1 \rangle$ ;

**dok** *true* (*za svaku epizodu*) **čini**

    Inicijaliziraj  $S$ :

**za** *svaki korak epizode*: **čini**

$A \leftarrow$  akcija dana od  $\pi$  za stanje  $S$ ;  
        poduzmi akciju  $A$  i promatraj  $R, S'$ ;  
         $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ ;  
         $S \leftarrow S'$ ;  
        dok ne dođe do terminalnog  $S$ ;

**kraj**

**kraj**

---

### 3.2. Učenje s vremenskom razlikom

Korisno je uočiti izraz unutar zagrada u pravilu ažuriranja TD(0) koji je moguće protumačiti kao neku vrstu pogreške: mjeri se razliku između procijenjene vrijednosti za  $S_t$  i bolje procjene  $R_{t+1} + V(S_{t+1})$ . Taj izraz, tzv. *pogreška TD-a*

$$\delta_t = R_{t+1} + V(S_{t+1}) - V(S_t)$$

se u raznim oblicima pojavljuje u svim područjima PU-a.

#### 3.2.1 Primjer: vožnja kući

Svatko tko svakodnevno vozi kući s posla pokušava predvidjeti koliko je vremena potrebno za prelaženje tog puta. U trenutku odlaska s radnog mjesta, zabilježe se vremenske prilike, dan u tjednu, prognoza vremena i sve što bi moglo biti važno. Na primjer, Petar je prošli petak krenuo točno u 6 popodne i pretpostavlja da će mu trebati 30 minuta do kuće. Sjeda u auto u 6:05 i primjećuje da počinje padati kiša. Budući da je vožnja u kišnim uvjetima otežana, Petar predviđa da će mu od tog trenutka do kuće trebati 35 minuta, tj. da će cijeli put trajati ukupno 40 minuta. Međutim, 15 minuta kasnije dok je izlazio s brze ceste pokaže se da je taj dio puta prošao prilično glatko pa smanjuje predviđeno ukupno trajanje puta na 35 minuta. Nažalost, ispred njega vozi spori kamion na jednotračnoj cesti sa zabranom pretjecanja i Petar mora voziti iza njega sve dok ne skrene u svoju ulicu u 6:40. Za tri minute stiže kući. Niz stanja, vremena i pretpostavki vide se u tablici [2].

### 3.2. Učenje s vremenskom razlikom

<i>Stanje</i>	<i>Proteklo vrijeme</i> <i>(u minutama)</i>	<i>Predviđeno vr.</i> <i>preostalo</i>	<i>Predviđeno vr.</i> <i>ukupno</i>
Petar kreće iz ureda, petak u 6	0	30	30
Petar dolazi do auta, pada kiša	5	35	40
Petar silazi s brze ceste	20	15	35
jedna traka, iza kamiona	30	10	40
Petar ulazi u svoju ulicu	40	3	43
Petar dolazi kući	43	0	43

Nagrade u ovome primjeru su proteklo vrijeme za svaki dio puta<sup>3</sup> Budući da je izostavljena korekcija ( $\gamma = 1$ ) dobit svakog stanja je stvarno vrijeme koje je preostalo iz njega. Vrijednost svakog stanja je predviđeno preostalo vrijeme. U trećem stupcu su trenutne procjene za svako od stanja. Najjednostavniji način prikaza MC metode je graf predviđenog ukupnog vremena (zadnji stupac) na nizu stanja, kao što je prikazano na Slici 3.2.1 (lijevo) - crvene strelice pokazuju promjenu u predviđanju MC metode, one naine točno odgovaraju pogrešci između procijenjene vrijednosti (procijenjenog preostalog vremena) u svakom stanju i stvarne dobiti (stvarnog preostalog vremena). Na primjer, kod izlaska s brze ceste Petar je mislio da će mu trebati još 15 minuta, a zapravo mu je trebalo 23. Ako se primjeni pravilo  $V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$  na ovaj trenutak povećat će se vrijednost ovoga stanja, tj predviđenog preostalog vremena. Stvarna greška  $G_t - V(S_t)$  u ovome stanju iznosi 8 minuta. Pod pretpostavkom da se koristi parametar veličine koraka  $\alpha = 0.5$ , onda se predviđeno vrijeme u tom stanju ažurira povećanjem za 4 minute kao posljedica ovoga iskustva. Ovo je možda pre-

---

<sup>3</sup>Da se radi o optimizaciji cilj bi bio minimizirati vrijeme putovanja i u tom slučaju bi nagrade bile negativne za količinu proteklog vremena. S obzirom na to da u ovom primjeru fokus nije na optimizaciji, radi jednostavnosti prikaza koriste se pozitivni brojevi.

### 3.3. MC vs. TD

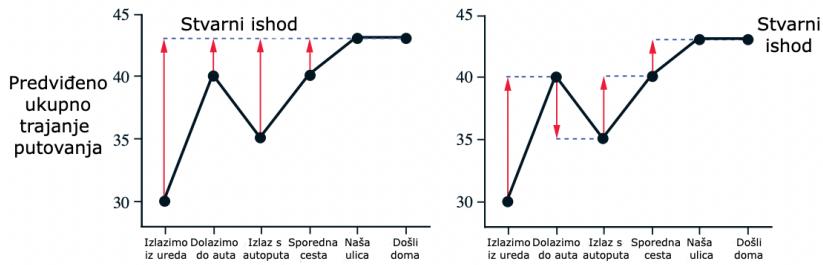
velika promjena jer se najvjerojatnije neće često ponavljati stanje sa sporim kamionom. U svakom se slučaju ažuriranje vrši tek na kraju epizode, tj. pri dolasku kući - tek je tada poznata stvarna dobit od svakog stanja. Je li nužno čekati kraj epizode da bi se postigao trenutak učenja iz novog stanja? Recimo da neki drugi dan Petar napušta ured i predviđa da će mu trebati 30 minuta, međutim taj dan zapne u ogromnoj gužvi i 25 minuta nakon napuštanja uredu još uvijek se sporo kreće u koloni. U tom trenutku Petar predviđa da mu preostaje još bar 25 minuta putovanja, znači ukupno 50. Dok čeka u koloni već zna da je njegovo početno predviđanje od 30 minuta bilo preoptimistično. Po MC treba čekati do kraja jer stvarna dobit još uvijek nije poznata.

Po TD-u se odmah uči i mijenja početno predviđanje sa 30 na 50 minuta. Zapravo se svako predviđanje pomiče. U prethodnom primjeru vožnje od petka, na slici 3.2.1 (desno) se vide promjene predviđanja koje daje TD pravilo  $V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$  (za  $\alpha = 1$ ). Osim što se na ovaj način Petar ima čime zabavljati dok čeka u koloni, postoji još nekoliko računalnih prednosti učenja baziranog na trenutnim pretpostavkama umjesto čekanja do terminalnog stanja da bi se saznala stvarna dobit - neke od tih prednosti su navedene u dalnjem tekstu.

## 3.3 MC vs. TD

Očita prednost TD metoda nad MC metodama je da su implementirane na inkrementalan način. Kod MC metoda potrebno je dočekati kraj epizode, jer se tek onda sazna je kolika je dobit, dok se kod TD-a čeka samo jedan korak. U mnogo slučajeva se ovo pokazalo kao odlučujući faktor jer u brojnim primjenama epizode traju jako dugo, a kod kontinuiranih slučaj uopće nema

### 3.3. MC vs. TD



Slika 3.2: Preporučene promjene predviđanja u primjeru "vožnja kući" od MC (lijevo) i TD (desno) metoda [1]

epizoda. Uz to neke MC metode moraju ignorirati ili značajno korigirati epizode u kojima su korištene eksperimentalne akcije što značajno usporava učenje - TD metode nisu podložne ovim problemima jer one uče prilikom svakog pomaka bez obzira koje su akcije poduzete kasnije.

Jako je zgodno učiti o jednoj pretpostavci temeljem druge pretpostavke, bez čekanja na stvarni rezultat, postavlja se pitanje može li se u tom slučaju garantirati konvergeniciju ka točnom odgovoru. Iako neće biti ovdje dokazano, na svu sreću odgovor je da! Za bilo koju fiksnu strategiju  $\pi$  TD(0) će konvergirati ka  $v_\pi$  ako je parametar veličine koraka dovoljno malen.

Ako i TD i MC metode asimptotski konvergiraju ka ispravnim predviđanjima, sljedeće logično pitanje je koja od njih dođe prije do rezultata, tj. koja metoda brže uči.

Ovo je još uvijek otvoreno pitanje, u smislu da ne postoji formalni dokaz da jedna metoda konvergira brže od druge. Ali u praksi se pokazalo da TD metode u pravilu konvergiraju brže od MC metoda na stohastičkim zadatcima.

U slučaju rada s konačnom količinom iskustva, recimo 10 epizoda ili 100 vremenskih koraka, ustaljena je praksa iterirati po tom iskustvu dok se ne ko-

### 3.3. MC vs. TD

nvergira do odgovora - pokazalo se da metode TD(0) i MC obe konvergiraju, ali prema različitim odgovorima!

#### 3.3.1 Primjer: Petar predviđa

Petar koji je došao kući u ovom je primjeru prediktor koji promatra rezultate za dane akcije:

A, 0, B, 0	B, 1
B, 1	B, 1
B, 1	B, 1
B, 1	B, 0

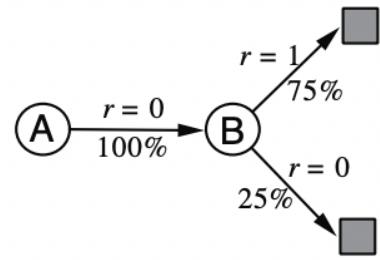
Prva epizoda počinje u stanju A, uz nagradu 0 prelazi u stanje B i onda uz nagradu 0 prelazi u terminalno stanje. Idućih šest epizoda direktno iz stanja B prelaze u terminalno stanje uz nagradu 1 i zadnja epizoda terminira uz nagradu 0. Iz danog skupa podataka koje su pretpostavljene vrijednosti stanja tj. koliki su  $V(A)$  i  $V(B)$ ?

Nema nikakvih značajnih dilema oko pretpostavke  $V(B) = \frac{6}{8} = \frac{3}{4}$ . Međutim, što reći o  $V(A)$  za dani skup podataka?

Postoje 2 razumna odgovora. Jedan je kada Petar primijeti da u 100% slučajeva iz stanja A prelazi u stanje B (uz nagradu 0), a budući da je  $V(B) = \frac{3}{4}$ , to je onda i  $V(A) = \frac{3}{4}$ . Ovaj se odgovor može promatrati kao da se prvo modelira MDP, grafika (3.3.1), i onda se računaju procjene za taj model. Ovo je i odgovor koji bi dao TD(0).

Drugi razuman odgovor je da svako promatrano pojavljivanje A završava s dobiti 0, dakle  $V(A) = 0$ . Ovo je odgovor koji bi dala MC metoda - i odgovor koji daje minimum srednje kvadratne greške na skupu za treniranje. Dapače, greška u ovom slučaju iznosi 0 ali ipak se prvi odgovor čini smisleniji. Ako je proces koji Petar promatra Markovljev, onda će prvi odgovor dati

### 3.3. MC vs. TD



Slika 3.3: Jednostavan model MDP-a [1]

manju grešku na budućim podatcima - iako je MC odgovor bolji za postojeće podatke.

# Poglavlje 4

## Kontrola bez modela

Sve prethodno napisano bilo je svojevrstan uvod u ovo poglavlje koje je zapravo centralni dio ovog rada. Konačno će se moći izvući zaključci o načinu na koji neki robot, bačen u nepoznato okruženje, otkriva što treba napraviti - kako maksimizirati svoju nagradu, tj. birati akcije koje će dugoročno donijeti najveću dobit bez da unaprijed zna išta o okolišu u kojem se nalazi. Sve druge tehnike PU, koje zbog sadržajnog ograničenja nisu spomenute u ovom radu, temelje se na ovome.

U prethodnom je poglavlju predstavljen set potrebnih alata- predviđanje bez modela, tj. kako vrednovati danu strategiju. Pomoću toga je moguće procijeniti vrijednosnu funkciju strategije u nepoznatom MDP-u. Taj osnovni alat će koristiti za metode kontrole, tj. unapređenju kako bi došli do optimalne vrijednosne funkcije i optimalne strategije.

Većina zanimljivih problema PU-a kao npr. botovi koji igraju računalne igrice, savijanje proteina, upravljanje fuzijskom reakcijom itd. imaju MDP, tj. mogu se opisati kao MDP-ovi - imaju stanja okoline i neka pravila/logiku po kojoj se ponašaju. Međutim, ili je taj MDP nepoznat, nepoznata su pravila ponašanja okoline pa se oslanja na metode učenja bez modela, ili je

#### 4.1. Monte-Carlo kontrola

MDP poznat ali je toliko složen da ga je nepraktično koristiti za bilo što osim uzorkovanja - u tom slučaju se opet radi o području učenja bez modela kod kojeg se uči iz iskustva pokušavajući razne stvari. Dobar primjer potonjeg je stvarni svijet - i kad bi znali točna pravila stvarnog svijeta (vjerojatno neka teorija kaosa), neizmjerno je jednostavnije uzeti uzorak - napraviti eksperiment - i vidjeti što se dogodilo nego računati ponašanje svakog atoma.

## 4.1 Monte-Carlo kontrola

Ako se na tren opet obrati pažnja na generaliziranu iteraciju strategija (GPI), njezina je osnovna ideja alterniranje dva procesa: *vrednovanje strategije* i poboljšanje strategije. Ovaj proces konvergira ka optimalnoj vrijednosnoj funkciji i optimalnoj strategiji. Postavlja se pitanje što je potrebno poduzeti kako bi se upotreboom ta dva procesa dobila kontrola bez modela.

Najjednostavnije bi bilo iskoristiti MC predviđanje za vrednovanje strategije i pohlepno poboljšanje strategije - međutim pojavljuje se nekoliko poteškoća. MC predviđanje uzme određen broj uzoraka i računa aritmetičke sredine dobiti kako bi aproksimirali vrijednosnu funkciju strategije, zatim se generira nova strategija koja se ponaša pohlepno u odnosu na procijenjenu vrijednosnu funkciju. U ovom su se slučaju pojavila dva problema:

*Bez modela:* kako je moguće biti bez modela ako se koristi funkciju  $v_\pi$ . Ako je poznata samo vrijednosna funkcija stanja i koristi se pohlepni algoritam s obzirom na nju, iz (2.3) tj.

$$\pi'(s) = \operatorname{argmax}_a \sum_{s',r} \mathcal{P}_0(s',r|s,a) [r + \gamma v_\pi(s')]$$

vidi se da je i dalje potreban model okoline - tj. matrica prijelaza  $\mathcal{P}_0$ , koje nema! Alternativa je korištenje vrijednosne funkcije akcije  $q_\pi$  jer u tom

#### 4.1. Monte-Carlo kontrola

slučaju je dovoljno tražiti maksimum po  $Q$  vrijednostima

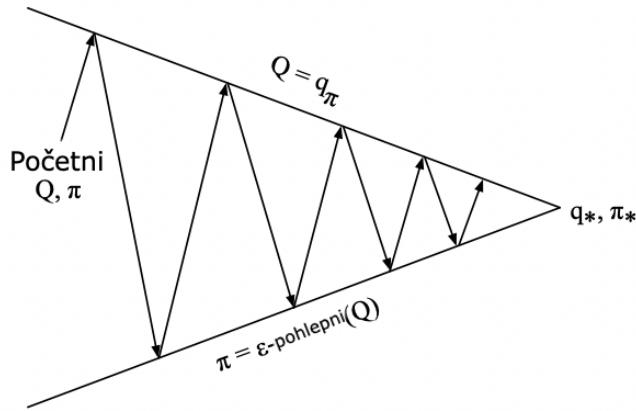
$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$$

Ovo omogućuje kontrolu u slučaju da nema informacije o modelu - ako se zna informacija  $q$  i želi se raditi pohlepno poboljšanje strategije u svakom stanju bira se akciju koja maksimizira  $q$  vrijednosti.

*Problem istraživanja:* ako je dana strategija koja se slijedi pohlepna, onda nema nikakve garancije da putanje kojima se prolazi pokrivaju cijeli prostor stanja - lako se može dogoditi da postoje dijelovi prostora stanja koji imaju veći potencijal, daju bolje nagrade, no do kojih agent nikad neće ni doći. Kratak primjer će poslužiti kao ilustracija: Petar se nalazi ispred dvaju vrata iza kojih se nalaze različite **stohastičke** nagrade. On otvorit će lijeva vrata i dobije nagradu 0, tj.  $V(ljevo) = 0$ . Nulom nije baš oduševljen pa isproba drugi put desna vrata i dobije nagradu +1, tj.  $V(desno) = +1$ . Dalje se ponaša *pohlepno* pa opet otvara desna vrata i dobije nagradu +3 pa je  $V(desno) = +2$  itd. On otvara desna vrata i dobija neku pozitivnu nagradu koja varira između 1 i 3. Međutim, je li stvarno odabrao najbolja vrata? - problem je što mu zapravo nije poznata vrijednost lijevih vrata jer ih je probao otvoriti samo jedan put. Dakle, treba nastaviti istraživati sve kako bi bio siguran da zna prave vrijednosti svih opcija koje mu se nude. Ako prestane istraživati određene akcije može kasnije donijeti krive odluke i zapeti u nekom lokalnom optimumu. Kako bi se ovaj problem rješio, uz garanciju da se uvijek istražuje sva stanja i sve akcije, koristi se dosta jednostavna ideja - iako ima mnogo sofisticiranih metoda tako je teško postići bolje rezultate od nje. Koristi se tzv.  $\epsilon$  – *pohlepno istraživanje*. Svaki put kada agent bira akciju sa vjerojatnošću  $\epsilon$  izabrat će neku nasumičnu akciju, a sa vjerojatnošću  $1 - \epsilon$  bira pohlepnu akciju. Dakle, uglavnom se prati pohlepnu strategiju ali s malom vjerojatnošću se vrše i druge akcije. Možda se čini

#### 4.2. Kontrola s vremenskom razlikom

naivno ali ima neka dobra svojstva. Garantira se istraživanje svih akcija i poboljšavanje strategije. Lako se pokaže (vidi [1]) da  $\epsilon$ -*pohlepno istraživanje* daje poboljšanje strategije i može se koristiti u spomenutoj verziji GPI-a<sup>1</sup>.



Slika 4.1: Monte-Carlo iteracija strategija [1]:

*Vrednovanje strategije:* MC vrednovanje,  $Q = q_\pi$

*Poboljšanje strategije:*  $\epsilon$  – *pohlepno*

## 4.2 Kontrola s vremenskom razlikom

Kod učenja MC metode nadovezala se na nju metoda vremenske razlike, TD za razliku od MC-a može se koristiti bez čekanja na kraj epizode ili za kontinuirane slučajeve koji nikad ne dodu do terminalnog stanja. Isprobavanje

---

<sup>1</sup>Kao i kod DP-a nije nužno do kraja vrednovati strategiju - nekad već nakon nekoliko koraka ima dovoljno informacija za konstruiranje puno bolje strategije, bez potrebe za dodatnim iteracijama kako bi se vrijednosnu funkciju izračunalo do kraja. U ekstremnom slučaju nakon svake epizode (završene putanje) može se raditi iteraciju poboljšanja strategije.

#### 4.2. Kontrola s vremenskom razlikom

iste ideje kao kod MC-a se zapravo samo po sebi nameće, radi se o generaliziranoj iteraciji strategije. Iz već prethodno navedenih razloga, potrebno je koristiti vrijednosnu funkciju akcije  $q_\pi$  i  $\epsilon - pohlepno$  poboljšanje strategije. Na ovaj se način automatski dobije jedan od najboljih algoritama PU-a - jedino poboljšanje koje će se primijeniti je dodatno povećanje frekvencije iteracija. Kako TD ne treba čekati niti do kraja epizode da bi ažurirali vrijednosnu funkciju, dapače može se ažurirati i nakon samo jednog koraka, a uvijek je poželjno koristiti najnovije dostupne vrijednosti kod biranja akcija, onda se može povećati broj iteracija primjenjujući ih na svakom vremenskom koraku.

Generalno pravilo TD učenja je tzv. **SARSA**

$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A))$$

Dakle ako je tražena vrijednost za uređeni par stanja i akcije  $(S, A)$ . Iz stanja  $S$  agent vrši akciju  $A$  i okolina dodjeljuje nagradu  $R$  te agent završava u stanju  $S'$ . Zatim slijedeći odabranu strategiju agent vrši akciju  $A'$  tako da se u dva koraka dobije  $(S, A) \xrightarrow{R} (S', A')$  odakle naziv SARSA. Počinje se od odabrane  $Q$  vrijednosti i pomiče se u smjeru koji daje nagrada uvećana za razliku korigirane vrijednosti idućeg stanja i vrijednosti trenutnog stanja.

## 4.2. Kontrola s vremenskom razlikom

---

### Algoritam 6: SARSA za procjenu $Q \approx q_*$

---

**Podatci:** parametri algoritma: veličina koraka  $\alpha \in \langle 0, 1 \rangle$ , maleni

$$\epsilon > 0$$

**Inicijaliziraj:**  $Q(s, a)$  za sve  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$  po volji, osim

$$Q(\text{terminal}, \cdot) = 0$$

**za svaku epizodu čini**

Inicijaliziraj  $\mathcal{S}$ ;

Odaberi  $A$  u stanju  $S$  koristeći strategiju izvedenu iz  $Q$  (npr.

$$\epsilon - \text{pohlepno});$$

**dok**  $S$  nije terminalan, za svaki korak epizode čini

Poduzmi akciju  $A$  i promotri  $R$  i  $S'$ ;

Odaberi  $A'$  u stanju  $S'$  koristeći strategiju izvedenu iz  $Q$  (npr.

$$\epsilon - \text{pohlepno});$$

$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A));$$

$S \leftarrow S'$ ;

$A \leftarrow A'$ ;

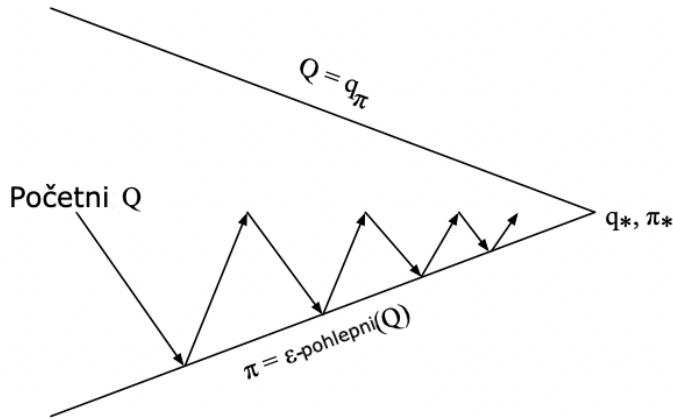
**kraj**

**kraj**

---

Drugim riječima: U nekom stanju se razmatra o poduzimanju određene akcije. Ako se ta akcija poduzme i promotri se dobivena nagrada i vrijednost iduće akcije koju će agent poduzeti, taj postupak daje procjenu vrijednosti strategije koja se slijedi. Ta se procjena koristi za ažuriranje vrijednosti para stanja i akcije iz kojih se započelo - ova generalna ideja je zapravo već viđena nekoliko puta i trebala bi odmah podsjetiti na Bellmanove jednadžbe.

## 4.2. Kontrola s vremenskom razlikom



Slika 4.2: TD iteracija strategija (nakon **svakog** vremenskog koraka) [1]:

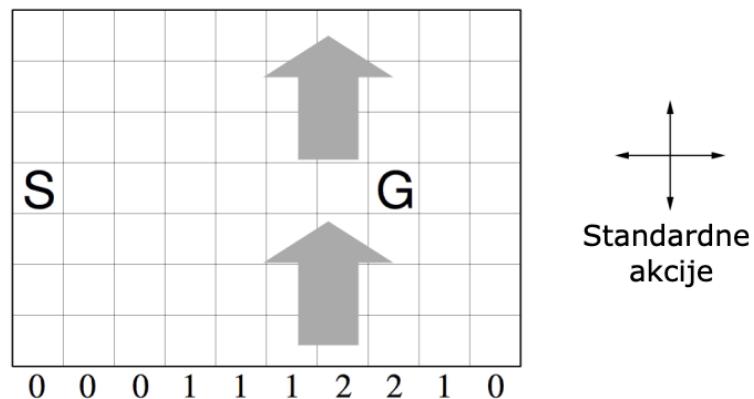
*Vrednovanje strategije:* SARSA  $Q_\pi$

*Poboljšanje strategije:*  $\epsilon$  – pohlepno

### 4.2.1 Primjer: vjetroviti Mrežni svijet [1]

Na slici (4.2.1) je primjer Mrežnog svijeta s naznačenim početnim (S) i konačnim (G) stanjem te jednim dodatkom: postoji ”vjetar” koji puše od dolje prema gore sredinom tablice. Akcije agenta su standardne (gore, dole, lijevo i desno), ali ako prolazi sredinom tablice vjetar, tj. okolina, ga pomiče prema gore što utječe na stanje u kojem će završiti nakon akcije – jačina vjetra ovisi o stupcu a može se vidjeti ispod tablice. Na primjer, ako se agent nalazi u polju desno od cilja i kreće prema lijevo, završit će u polju neposredno iznad cilja. Za svaki korak nagrada je -1 sve dok agent ne dođe do cilja. Na grafu (4.2.1) se vidi rezultat primjene  $\epsilon$  – pohlepne SARSA-e na ovaj primjer sa  $\epsilon = 0.1$  i  $\alpha = 0.5$  te početnim vrijednostima  $Q(s, a) = 0$  za sve  $s$  i  $a$ . Iz grafa se iščitava da je za završetak prve epizode bilo potrebno više od 2000 koraka. Međutim, nakon prve epizode druga je puno brža – dok se dođe do 8000 učinjenih vremenskih koraka već je odavno otkrivena

#### 4.2. Kontrola s vremenskom razlikom

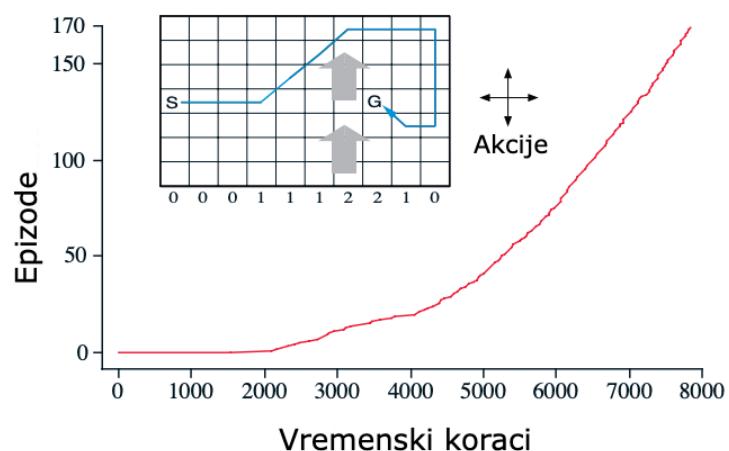


Slika 4.3: Vjetroviti Mrežni svijet

optimalna strategija (vidi se na slici) ali zbog  $\epsilon - pohlepnog$  istraživanja je prosječna duljina epizode 17 koraka - za dva više od optimalnih 15.

MC metode ne koriste u ovome slučaju jer ne postoji garancija za strategiju da će doći do kraja epizode, npr. kad bi strategija uputila agenta da stalno stoji na istom mjestu. SARSA nema ovih problema jer uči tijekom same epizode pa mijenja danu strategiju s drugačijom svaki put kad nauči nešto novo.

#### 4.2. Kontrola s vremenskom razlikom



Slika 4.4: Rješenje vjetrovitog Mrežnog svijeta & Graf potrošenih vremenskih koraka i prijeđenih epizoda [1]

# Poglavlje 5

## Praktični primjer: Blackjack

Kako bi se nadišlo puko teoretiziranje i navođenje primjera iz knjige, primjenjivost metode je iskušana u ”stvarnom svijetu” primjenom par linija koda. Kod pažljivijeg iščitavanja rada, oštro oko je sigurno primijetilo nedostatak pratećeg primjera u poglavlju 4.1 **Monte-Carlo kontrola** - ovo će poglavlje to ispraviti. Kao što i sam naslov poglavlja sugerira, traži se optimalna strategija za igranje Blackjacka. Budući da su pravila i algoritam MC kontrole već poznati, odmah ide modeliranje problema.

Program je u svrhu ovog rada napisan u *vanilla* JavaScript-u i pokretan je uz pomoć Node.js-a v16.13.1 , no očigledno je toliko jednostavan da se lako da implementirati u bilo kojem programskom jeziku i pokrenuti u većini okolina. Program je radi preglednosti podijeljen u dvije datoteke: *pomocni.js* - modul koji sadrži definicije svih funkcija koje se koriste i *program.js* - koji sadrži lako čitljivu implementaciju algoritma. Napokon evo i koda:

```
1 \\                                         10 J Q K
2 const KARTE = "A,2,3,4,5,6,7,8,9,10,10,10,10".split(',');
3 const izvuciKartu = () => KARTE[Math.floor(Math.random() *
    KARTE.length)];
```

```

4
5 const inkrementalnoIzracunajProsjek = (stariProsjek,
    novaVrijednost, n) => stariProsjek + (novaVrijednost -
    stariProsjek) / (n + 1);
6
7 const resetirajNovuTablicuVrijednosti = () => ({
8     sAsom: [...Array(10)].map((el, i) =>
9         [...Array(10)].map((elInner) => ({
10            h: { val: 0, brojPonavljanja: 0 },
11            s: { val: 0, brojPonavljanja: 0 }
12        }))
13    ),
14    bezAsa: [...Array(10)].map((el, i) =>
15        [...Array(10)].map((elInner) => ({
16            h: { val: 0, brojPonavljanja: 0 },
17            s: { val: 0, brojPonavljanja: 0 }
18        }))
19    )
20 });

```

Listing 5.1: pomocni.js 1/3

Većina koda je samoobjašnjiva, *resetirajNovuTablicuVrijednosti* vraća objekt sa dvije 10x10 tablice u koje će se spremati prosjeci ostvarenih dobiti iz danog stanja za trenutnu strategiju, tj. vrijednosti tih stanja.

```

1 const zapocniPartiju = () => {
2     let iskoristivAS = false;
3     let suma = 0;
4
5     while (suma < 12) {
6         const izvucenaKarta = izvuciKartu();
7         if (izvucenaKarta === "A" && !iskoristivAS) {
8             suma += 10;

```

```

9     iskoristivAS = true;
10    } else if (izvucenaKarta === "A" && iskoristivAS) {
11        suma += 1;
12    } else {
13        suma += +izvucenaKarta;
14    }
15}
16
17 return {
18     suma,
19     iskoristivAS,
20     putanja: [],
21     dobit: undefined,
22     protivnikPokazuje: izvuciKartu()
23 };
24};
25
26 const poduzimamoAkciju = (partija, izvucenaKarta) => {
27     partija.putanja.push({
28         suma: partija.suma,
29         odabranaAkcija: "h",
30         iskoristivAS: partija.iskoristivAS
31     });
32     if (izvucenaKarta === "A") {
33         partija.suma += 1;
34     } else {
35         partija.suma += +izvucenaKarta;
36     }
37
38     if (partija.suma > 21 && partija.iskoristivAS) {
39         partija.suma -= 9;
40         partija.iskoristivAS = false;
41     }

```

```

42 };
43
44 const odigrajEpsilonPohlepono = (partija, epsilon,
45   staraTablicaVrijednosti) => {
46   let odabranaAkcija = undefined;
47
48   while (odabranaAkcija !== "s" && partija.suma < 22) {
49     const trenutnoStanje = (partija.iskoristivAS
50       ? staraTablicaVrijednosti.sAsom
51       : staraTablicaVrijednosti.bezAsa)[partija.suma - 12][
52         partija.protivnikPokazuje === "A" ? 0 : +partija.
53           protivnikPokazuje - 1
54     ];
55
56     const [preferiranaAkcija, nePreferiranaAkcija] =
57       trenutnoStanje === 'h' ? ["h", "s"] : ["s", "h"];
58
59     odabranaAkcija =
60       (Math.random() > epsilon && preferiranaAkcija) ||
61       nePreferiranaAkcija;
62
63     if (odabranaAkcija === "h") {
64       poduzimamoAkciju(partija, izvuciKartu());
65     } else {
66       partija.putanja.push({
67         suma: partija.suma,
68         odabranaAkcija,
69         iskoristivAS: partija.iskoristivAS
70       });
71     }
72   }
73 }
74
75 }
76
77 }
78
79 }
80
81 }
82
83 }
84
85 }
86
87 }
88
89 }
90
91 }
92
93 }
94
95 }
96
97 }
98
99 }
100
101 }
102
103 }
104
105 }
106
107 }
108
109 }
110
111 }
112
113 }
114
115 }
116
117 }
118
119 }
120
121 }
122
123 }
124
125 }
126
127 }
128
129 }
130
131 }
132
133 }
134
135 }
136
137 }
138
139 }
140
141 }
142
143 }
144
145 }
146
147 }
148
149 }
150
151 }
152
153 }
154
155 }
156
157 }
158
159 }
160
161 }
162
163 }
164
165 }
166
167 }
168
169 }
170
171 }
172
173 }
174
175 }
176
177 }
178
179 }
180
181 }
182
183 }
184
185 }
186
187 }
188
189 }
190
191 }
192
193 }
194
195 }
196
197 }
198
199 }
200
201 }
202
203 }
204
205 }
206
207 }
208
209 }
210
211 }
212
213 }
214
215 }
216
217 }
218
219 }
220
221 }
222
223 }
224
225 }
226
227 }
228
229 }
230
231 }
232
233 }
234
235 }
236
237 }
238
239 }
240
241 }
242
243 }
244
245 }
246
247 }
248
249 }
250
251 }
252
253 }
254
255 }
256
257 }
258
259 }
259
260 }
261
262 }
263
264 }
265
266 }
267
268 }
269
270 }
271
272 }
273
274 }
275
276 }
277
278 }
279
280 }
281
282 }
283
284 }
285
286 }
287
288 }
289
290 }
291
292 }
293
294 }
295
296 }
297
298 }
299
300 }
301
302 }
303
304 }
305
306 }
307
308 }
309
310 }
311
312 }
313
314 }
315
316 }
317
318 }
319
320 }
321
322 }
323
324 }
325
326 }
327
328 }
329
330 }
331
332 }
333
334 }
335
336 }
337
338 }
339
340 }
341
342 }
343
344 }
345
346 }
347
348 }
349
350 }
351
352 }
353
354 }
355
356 }
357
358 }
359
360 }
361
362 }
363
364 }
365
366 }
367
368 }
369
370 }
371
372 }
373
374 }
375
376 }
377
378 }
379
380 }
381
382 }
383
384 }
385
386 }
387
388 }
389
390 }
391
392 }
393
394 }
395
396 }
397
398 }
399
400 }
401
402 }
403
404 }
405
406 }
407
408 }
409
410 }
411
412 }
413
414 }
415
416 }
417
418 }
419
420 }
421
422 }
423
424 }
425
426 }
427
428 }
429
430 }
431
432 }
433
434 }
435
436 }
437
438 }
439
440 }
441
442 }
443
444 }
445
446 }
447
448 }
449
450 }
451
452 }
453
454 }
455
456 }
457
458 }
459
460 }
461
462 }
463
464 }
465
466 }
467
468 }
469
470 }
471
472 }
473
474 }
475
476 }
477
478 }
479
480 }
481
482 }
483
484 }
485
486 }
487
488 }
489
490 }
491
492 }
493
494 }
495
496 }
497
498 }
499
500 }
501
502 }
503
504 }
505
506 }
507
508 }
509
510 }
511
512 }
513
514 }
515
516 }
517
518 }
519
520 }
521
522 }
523
524 }
525
526 }
527
528 }
529
530 }
531
532 }
533
534 }
535
536 }
537
538 }
539
540 }
541
542 }
543
544 }
545
546 }
547
548 }
549
550 }
551
552 }
553
554 }
555
556 }
557
558 }
559
560 }
561
562 }
563
564 }
565
566 }
567
568 }
569
570 }
571
572 }
573
574 }
575
576 }
577
578 }
579
580 }
581
582 }
583
584 }
585
586 }
587
588 }
589
590 }
591
592 }
593
594 }
595
596 }
597
598 }
599
600 }
601
602 }
603
604 }
605
606 }
607
608 }
609
610 }
611
612 }
613
614 }
615
616 }
617
618 }
619
620 }
621
622 }
623
624 }
625
626 }
627
628 }
629
630 }
631
632 }
633
634 }
635
636 }
637
638 }
639
640 }
641
642 }
643
644 }
645
646 }
647
648 }
649
650 }
651
652 }
653
654 }
655
656 }
657
658 }
659
660 }
661
662 }
663
664 }
665
666 }
667
668 }
669
670 }
671
672 }
673
674 }
675
676 }
677
678 }
679
680 }
681
682 }
683
684 }
685
686 }
687
688 }
689
690 }
691
692 }
693
694 }
695
696 }
697
698 }
699
700 }
701
702 }
703
704 }
705
706 }
707
708 }
709
710 }
711
712 }
713
714 }
715
716 }
717
718 }
719
720 }
721
722 }
723
724 }
725
726 }
727
728 }
729
730 }
731
732 }
733
734 }
735
736 }
737
738 }
739
740 }
741
742 }
743
744 }
745
746 }
747
748 }
749
750 }
751
752 }
753
754 }
755
756 }
757
758 }
759
759
760 }
761
762 }
763
764 }
765
766 }
767
768 }
769
769
770 }
771
772 }
773
774 }
775
776 }
777
778 }
779
779
780 }
781
782 }
783
784 }
785
786 }
787
788 }
789
789
790 }
791
792 }
793
794 }
795
796 }
797
798 }
799
799
800 }
801
802 }
803
804 }
805
806 }
807
808 }
809
809
810 }
811
812 }
813
814 }
815
816 }
817
818 }
819
819
820 }
821
822 }
823
824 }
825
826 }
827
828 }
829
829
830 }
831
832 }
833
834 }
835
836 }
837
838 }
839
839
840 }
841
842 }
843
844 }
845
846 }
847
848 }
849
849
850 }
851
852 }
853
854 }
855
856 }
857
858 }
859
859
860 }
861
862 }
863
864 }
865
866 }
867
868 }
869
869
870 }
871
872 }
873
874 }
875
876 }
877
878 }
879
879
880 }
881
882 }
883
884 }
885
886 }
887
888 }
889
889
890 }
891
892 }
893
894 }
895
896 }
897
898 }
899
899
900 }
901
902 }
903
904 }
905
906 }
907
908 }
909
909
910 }
911
912 }
913
914 }
915
916 }
917
918 }
919
919
920 }
921
922 }
923
924 }
925
926 }
927
928 }
929
929
930 }
931
932 }
933
934 }
935
936 }
937
938 }
939
939
940 }
941
942 }
943
944 }
945
946 }
947
948 }
949
949
950 }
951
952 }
953
954 }
955
956 }
957
958 }
959
959
960 }
961
962 }
963
964 }
965
966 }
967
968 }
969
969
970 }
971
972 }
973
974 }
975
976 }
977
978 }
979
979
980 }
981
982 }
983
984 }
985
986 }
987
988 }
989
989
990 }
991
992 }
993
994 }
995
996 }
997
998 }
999
999
1000 }
1001
1002 }
1003
1004 }
1005
1006 }
1007
1008 }
1009
1009
1010 }
1011
1012 }
1013
1014 }
1015
1016 }
1017
1018 }
1019
1019
1020 }
1021
1022 }
1023
1024 }
1025
1026 }
1027
1028 }
1029
1029
1030 }
1031
1032 }
1033
1034 }
1035
1036 }
1037
1038 }
1039
1039
1040 }
1041
1042 }
1043
1044 }
1045
1046 }
1047
1048 }
1049
1049
1050 }
1051
1052 }
1053
1054 }
1055
1056 }
1057
1058 }
1059
1059
1060 }
1061
1062 }
1063
1064 }
1065
1066 }
1067
1068 }
1069
1069
1070 }
1071
1072 }
1073
1074 }
1075
1076 }
1077
1078 }
1079
1079
1080 }
1081
1082 }
1083
1084 }
1085
1086 }
1087
1088 }
1089
1089
1090 }
1091
1092 }
1093
1094 }
1095
1096 }
1097
1098 }
1099
1099
1100 }
1101
1102 }
1103
1104 }
1105
1106 }
1107
1108 }
1109
1109
1110 }
1111
1112 }
1113
1114 }
1115
1116 }
1117
1118 }
1119
1119
1120 }
1121
1122 }
1123
1124 }
1125
1126 }
1127
1128 }
1129
1129
1130 }
1131
1132 }
1133
1134 }
1135
1136 }
1137
1138 }
1139
1139
1140 }
1141
1142 }
1143
1144 }
1145
1146 }
1147
1148 }
1149
1149
1150 }
1151
1152 }
1153
1154 }
1155
1156 }
1157
1158 }
1159
1159
1160 }
1161
1162 }
1163
1164 }
1165
1166 }
1167
1168 }
1169
1169
1170 }
1171
1172 }
1173
1174 }
1175
1176 }
1177
1178 }
1179
1179
1180 }
1181
1182 }
1183
1184 }
1185
1186 }
1187
1188 }
1189
1189
1190 }
1191
1192 }
1193
1194 }
1195
1196 }
1197
1198 }
1199
1199
1200 }
1201
1202 }
1203
1204 }
1205
1206 }
1207
1208 }
1209
1209
1210 }
1211
1212 }
1213
1214 }
1215
1216 }
1217
1218 }
1219
1219
1220 }
1221
1222 }
1223
1224 }
1225
1226 }
1227
1228 }
1229
1229
1230 }
1231
1232 }
1233
1234 }
1235
1236 }
1237
1238 }
1239
1239
1240 }
1241
1242 }
1243
1244 }
1245
1246 }
1247
1248 }
1249
1249
1250 }
1251
1252 }
1253
1254 }
1255
1256 }
1257
1258 }
1259
1259
1260 }
1261
1262 }
1263
1264 }
1265
1266 }
1267
1268 }
1269
1269
1270 }
1271
1272 }
1273
1274 }
1275
1276 }
1277
1278 }
1279
1279
1280 }
1281
1282 }
1283
1284 }
1285
1286 }
1287
1288 }
1289
1289
1290 }
1291
1292 }
1293
1294 }
1295
1296 }
1297
1298 }
1299
1299
1300 }
1301
1302 }
1303
1304 }
1305
1306 }
1307
1308 }
1309
1309
1310 }
1311
1312 }
1313
1314 }
1315
1316 }
1317
1318 }
1319
1319
1320 }
1321
1322 }
1323
1324 }
1325
1326 }
1327
1328 }
1329
1329
1330 }
1331
1332 }
1333
1334 }
1335
1336 }
1337
1338 }
1339
1339
1340 }
1341
1342 }
1343
1344 }
1345
1346 }
1347
1348 }
1349
1349
1350 }
1351
1352 }
1353
1354 }
1355
1356 }
1357
1358 }
1359
1359
1360 }
1361
1362 }
1363
1364 }
1365
1366 }
1367
1368 }
1369
1369
1370 }
1371
1372 }
1373
1374 }
1375
1376 }
1377
1378 }
1379
1379
1380 }
1381
1382 }
1383
1384 }
1385
1386 }
1387
1388 }
1389
1389
1390 }
1391
1392 }
1393
1394 }
1395
1396 }
1397
1398 }
1399
1399
1400 }
1401
1402 }
1403
1404 }
1405
1406 }
1407
1408 }
1409
1409
1410 }
1411
1412 }
1413
1414 }
1415
1416 }
1417
1418 }
1419
1419
1420 }
1421
1422 }
1423
1424 }
1425
1426 }
1427
1428 }
1429
1429
1430 }
1431
1432 }
1433
1434 }
1435
1436 }
1437
1438 }
1439
1439
1440 }
1441
1442 }
1443
1444 }
1445
1446 }
1447
1448 }
1449
1449
1450 }
1451
1452 }
1453
1454 }
1455
1456 }
1457
1458 }
1459
1459
1460 }
1461
1462 }
1463
1464 }
1465
1466 }
1467
1468 }
1469
1469
1470 }
1471
1472 }
1473
1474 }
1475
1476 }
1477
1478 }
1479
1479
1480 }
1481
1482 }
1483
1484 }
1485
1486 }
1487
1488 }
1489
1489
1490 }
1491
1492 }
1493
1494 }
1495
1496 }
1497
1498 }
1499
1499
1500 }
1501
1502 }
1503
1504 }
1505
1506 }
1507
1508 }
1509
1509
1510 }
1511
1512 }
1513
1514 }
1515
1516 }
1517
1518 }
1519
1519
1520 }
1521
1522 }
1523
1524 }
1525
1526 }
1527
1528 }
1529
1529
1530 }
1531
1532 }
1533
1534 }
1535
1536 }
1537
1538 }
1539
1539
1540 }
1541
1542 }
1543
1544 }
1545
1546 }
1547
1548 }
1549
1549
1550 }
1551
1552 }
1553
1554 }
1555
1556 }
1557
1558 }
1559
1559
1560 }
1561
1562 }
1563
1564 }
1565
1566 }
1567
1568 }
1569
1569
1570 }
1571
1572 }
1573
1574 }
1575
1576 }
1577
1578 }
1579
1579
1580 }
1581
1582 }
1583
1584 }
1585
1586 }
1587
1588 }
1589
1589
1590 }
1591
1592 }
1593
1594 }
1595
1596 }
1597
1598 }
1599
1599
1600 }
1601
1602 }
1603
1604 }
1605
1606 }
1607
1608 }
1609
1609
1610 }
1611
1612 }
1613
1614 }
1615
1616 }
1617
1618 }
1619
1619
1620 }
1621
1622 }
1623
1624 }
1625
1626 }
1627
1628 }
1629
1629
1630 }
1631
1632 }
1633
1634 }
1635
1636 }
1637
1638 }
1639
1639
1640 }
1641
1642 }
1643
1644 }
1645
1646 }
1647
1648 }
1649
1649
1650 }
1651
1652 }
1653
1654 }
1655
1656 }
1657
1658 }
1659
1659
1660 }
1661
1662 }
1663
1664 }
1665
1666 }
1667
1668 }
1669
1669
1670 }
1671
1672 }
1673
1674 }
1675
1676 }
1677
1678 }
1679
1679
1680 }
1681
1682 }
1683
1684 }
1685
1686 }
1687
1688 }
1689
1689
1690 }
1691
1692 }
1693
1694 }
1695
1696 }
1697
1698 }
1699
1699
1700 }
1701
1702 }
1703
1704 }
1705
1706 }
1707
1708 }
1709
1709
1710 }
1711
1712 }
1713
1714 }
1715
1716 }
1717
1718 }
1719
1719
1720 }
1721
1722 }
1723
1724 }
1725
1726 }
1727
1728 }
1729
1729
1730 }
1731
1732 }
1733
1734 }
1735
1736 }
1737
1738 }
1739
1739
1740 }
1741
1742 }
1743
1744 }
1745
1746 }
1747
1748 }
1749
1749
1750 }
1751
1752 }
1753
1754 }
1755
1756 }
1757
1758 }
1759
1759
1760 }
1761
1762 }
1763
1764 }
1765
1766 }
1767
1768 }
1769
1769
1770 }
1771
1772 }
1773
1774 }
1775
1776 }
1777
1778 }
1779
1779
1780 }
1781
1782 }
1783
1784 }
1785
1786 }
1787
1788 }
1789
1789
1790 }
1791
1792 }
1793
1794 }
1795
1796 }
1797
1798 }
1799
1799
1800 }
1801
1802 }
1803
1804 }
1805
1806 }
1807
1808 }
1809
1809
1810 }
1811
1812 }
1813
1814 }
1815
1816 }
1817
1818 }
1819
1819
1820 }
1821
1822 }
1823
1824 }
1825
1826 }
1827
1828 }
1829
1829
1830 }
1831
1832 }
1833
1834 }
1835
1836 }
1837
1838 }
1839
1839
1840 }
1841
1842 }
1843
1844 }
1845
1846 }
1847
1848 }
1849
1849
1850 }
1851
1852 }
1853
1854 }
1855
1856 }
1857
1858 }
1859
1859
1860 }
1861
1862 }
1863
1864 }
1865
1866 }
1867
1868 }
1869
1869
1870 }
1871
1872 }
1873
1874 }
1875
1876 }
1877
1878 }
1879
1879
1880 }
1881
1882 }
1883
1884 }
1885
1886 }
1887
1888 }
1889
1889
1890 }
1891
1892 }
1893
1894 }
1895
1896 }
1897
1898 }
1899
1899
1900 }
1901
1902 }
1903
1904 }
1905
1906 }
1907
1908 }
1909
1909
1910 }
1911
1912 }
1913
1914 }
1915
1916 }
1917
1918 }
1919
1919
1920 }
1921
1922 }
1923
1924 }
1925
1926 }
1927
1928 }
1929
1929
1930 }
1931
1932 }
1933
1934 }
1935
1936 }
1937
1938 }
1939
1939
1940 }
1941
1942 }
1943
1944 }
1945
1946 }
1947
1948 }
1949
1949
1950 }
1951
1952 }
1953
1954 }
1955
1956 }
1957
1958 }
1959
1959
1960 }
1961
1962 }
1963
1964 }
1965
1966 }
1967
1968 }
1969
1969
1970 }
1971
1972 }
1973
1974 }
1975
1976 }
1977
1978 }
1979
1979
1980 }
1981
1982 }
1983
1984 }
1985
1986 }
1987
1988 }
1989
1989
1990 }
1991
1992 }
1993
1994 }
1995
1996 }
1997
1998 }
1999
1999
2000 }
2001
2002 }
2003
2004 }
2005
2006 }
2007
2008 }
2009
2009
2010 }
2011
2012 }
2013
2014 }
2015
2016 }
2017
2018 }
2019
2019
2020 }
2021
2022 }
2023
2024 }
2025
2026 }
2027
2028 }
2029
2029
2030 }
2031
2032 }
2033
2034 }
2035
2036 }
2037
2038 }
2039
2039
2040 }
2041
2042 }
2043
2044 }
2045
2046 }
2047
2048 }
2049
2049
2050 }
2051
2052 }
2053
2054 }
2055
2056 }
2057
2058 }
2059
2059
2060 }
2061
2062 }
2063
2064 }
2065
2066 }
2067
2068 }
2069
2069
2070 }
2071
2072 }
2073
2074 }
2075
2076 }
2077
2078 }
2079
2079
2080 }
2081
2082 }
2083
2084 }
2085
2086 }
2087
2088 }
2089
2089
2090 }
2091
2092 }
2093
2094 }
2095
2096 }
2097
2098 }
2099
2099
2100 }
2101
2102 }
2103
2104 }
2105
2106 }
2107
2108 }
2109
2109
2110 }
2111
2112 }
2113
2114 }
2115
2116 }
2117
2118 }
2119
2119
2120 }
2121
2122 }
2123
2124 }
2125
2126 }
2127
2128 }
2129
2129
2130 }
2131
2132 }
2133
2134 }
2135
2136 }
2137
2138 }
2139
2139
2140 }
2141
2142 }
2143
2144 }
2145
2146 }
2147
2148 }
2149
2149
2150 }
2151
2152 }
2153
2154 }
2155
2156 }
2157
2158 }
2159
2159
2160 }
2161
2162 }
2163
2164 }
2165
2166 }
2167
2168 }
2169
2169
2170 }
2171
2172 }
2173
2174 }
2175
2176 }
2177
2178 }
2179
2179
2180 }
2181
2182 }
2183
2184 }
2185
2186 }
2187
2188 }
2189
2189
2190 }
2191
2192 }
2193
2194 }
2195
2196 }
2197
2198 }
2199
2199
2200 }
2201
2202 }
2203
2204 }
2205
2206 }
2207
2208 }
2209
2209
2210 }
2211
2212 }
2213
2214 }
2215
2216 }
2217
2218 }
2219
2219
2220 }
2221
2222 }
2223
2224 }
2225
2226 }
2227
2228 }
2229
2229
2230 }
2231
2232 }
2233
2234 }
2235
2236 }
2237
2238 }
2239
2239
2240 }
2241
2242 }
2243
2244 }
2245
2246 }
2247
2248 }
2249
2249
2250 }
2251
2252 }
2253
2254 }
2255
2256 }
2257
2258 }
2259
2259
2260 }
2261
2262 }
2263
2264 }
2265
2266 }
2267
2268
```

Izvlače se karate sve dok nije postignuta suma od barem 12 jer ako je suma 11 ili manje nema se što izgubiti, tj. nema nikakve dileme da je optimalna strategija uzeti još jednu kartu. Kuća ima samo jednu kartu koju igrač vidi, jer karta koju ne vidi nema nikakav utjecaj na njegove odluke pa je po potrebi izvlači tek nakon što igrač završi.

```
1 const kucaIgra = (partija) => {
2   if (partija.suma > 21) {
3     partija.dobit = -1;
4     return;
5   }
6   let kucaImaIskoristivogAsa = partija.protivnikPokazuje ===
7     "A";
8   let sumaKuce =
9     partija.protivnikPokazuje === "A" ? 10 : +partija.
10    protivnikPokazuje;
11
12  while (sumaKuce < 17) {
13    const izvucenaKarta = izvuciKartu();
14    if (izvucenaKarta === "A") {
15      sumaKuce += kucaImaIskoristivogAsa ? 1 : 10;
16      kucaImaIskoristivogAsa = kucaImaIskoristivogAsa || true
17      ;
18    } else {
19      sumaKuce += +izvucenaKarta;
20    }
21
22    if (sumaKuce > 21 && kucaImaIskoristivogAsa) {
23      sumaKuce -= 9;
24      kucaImaIskoristivogAsa=false;
25    }
26  }
```

```

24 if (sumaKuce > 21) {
25     partija.dobit = 1;
26     return;
27 }
28
29 if (sumaKuce === partija.suma) {
30     partija.dobit = 0;
31     return;
32 }
33 if (partija.suma > sumaKuce) {
34     partija.dobit = 1;
35     return;
36 } else {
37     partija.dobit = -1;
38     return;
39 }
40 };
41
42 const azurirajStanjaNoveTablice = (novaTablicaVrijednosti,
43                                     partija) => {
44     partija.putanja.forEach((stanje) => {
45         const stanjeUTablici = (stanje.iskoristivAS
46             ? novaTablicaVrijednosti.sAsom
47             : novaTablicaVrijednosti.bezAsa)[stanje.suma - 12][
48                 partija.protivnikPokazuje === "A" ? 0 : +partija.
49                     protivnikPokazuje - 1
50             ][stanje.odabranaAkcija];
51         stanjeUTablici.val = inkrementalnoIzracunajProsjek(
52             stanjeUTablici.val,
53             partija.dobit,
54             stanjeUTablici.brojPonavljanja
55         );
56         stanjeUTablici.brojPonavljanja += 1;

```

```

55    });
56 };
57
58 const formirajNovuTablicu = (novaTablica) => {
59     const formirana = {};
60     formirana.sAsom = novaTablica.sAsom.map((redak) =>
61         redak.map(({ h, s }) => h.val >= s.val ? 'h': 's' )
62     );
63
64     formirana.bezAsa = novaTablica.bezAsa.map((redak) =>
65         redak.map(({ h, s }) => h.val >= s.val ? 'h': 's' )
66     );
67     return formirana;
68 };
69
70 const ispirintajTablicu = (tablica) => {
71     console.log('s Asom:');
72     console.log('      A   2   3   4   5   6   7   8   9   10');
73     tablica.sAsom.forEach((redak, i) => console.log(i+12, JSON.
74         stringify(redak)))
75     console.log('=====');
76     console.log('bez Asa:');
77     console.log('      A   2   3   4   5   6   7   8   9   10');
78     tablica.bezAsa.forEach((redak, i) => console.log(i+12, JSON
79         .stringify(redak)))
80 }
81
80 module.exports = {
81     resetirajNovuTablicuVrijednosti ,
82     ispirintajTablicu ,
83     azurirajStanjaNoveTablice ,
84     zapocniPartiju ,
85     poduzimamoAkciju ,

```

```

86     odigrajEpsilonPohlepno ,
87     kucaIgra ,
88     formirajNovuTablicu
89 }
```

Listing 5.3: pomocni.js 3/3

- Funkcija *kucaIgra* odradi izvlačenje karata za kuću, nastavlja izvlačiti dok ne postigne barem 17 i određuje dobiti koje su poznate čim kuća završi.
- *azuriranjeStanjaNoveTablice* računa i postavlja nove procjene vrijednosti posjećenih stanja u danoj partiji (kako broj posjeta tih stanja raste, procjena će težiti stvarnoj vrijednosti tih stanja).
- *formirajNovuTablicu* iz tablice vrijednosti formira tablicu koja govori kako se ponaša pohlepna strategija s obzirom na te vrijednosti.

```

1 const {
2   resetirajNovuTablicuVrijednosti ,
3   ispirintajTablicu ,
4   formirajNovuTablicu ,
5   zapocniPartiju ,
6   odigrajEpsilonPohlepno ,
7   kucaIgra ,
8   azurirajStanjaNoveTablice
9 } = require('./pomocni')
10
11 const EPSILON = 0.1;
12
13 let staraTablicaVrijednosti = {
14   sAsom: [...Array(10)].map((_, i) => [...Array(10)].map(() =>
```

```

15 bezAsa: [...Array(10)].map(_, i) => [...Array(10)].map(() =
    > i + 12 < 17 ? 'h': 's'))
16 };
17 let novaTablicaVrijednosti = resetirajNovuTablicuVrijednosti(
    );
18 let partija;
19 let brojac = 0;
20
21 do {
22
23 if(brojac !== 0) staraTablicaVrijednosti =
    formirajNovuTablicu(novaTablicaVrijednosti);
24 novaTablicaVrijednosti = resetirajNovuTablicuVrijednosti();
25
26 for (let i = 0; i < 10000000; i++) {
27     partija = zapocniPartiju();
28     odigrajEpsilonPohlepono(partija, EPSILON,
        staraTablicaVrijednosti);
29     kucaIgra(partija);
30     azurirajStanjaNoveTablice(novaTablicaVrijednosti, partija
        );
31 }
32
33 console.log('brojac:', brojac)
34 brojac++
35 } while(JSON.stringify(staraTablicaVrijednosti) !== JSON.
    stringify(formirajNovuTablicu(novaTablicaVrijednosti)))
36
37 ispirintajTablicu(formirajNovuTablicu(novaTablicaVrijednosti)
    )

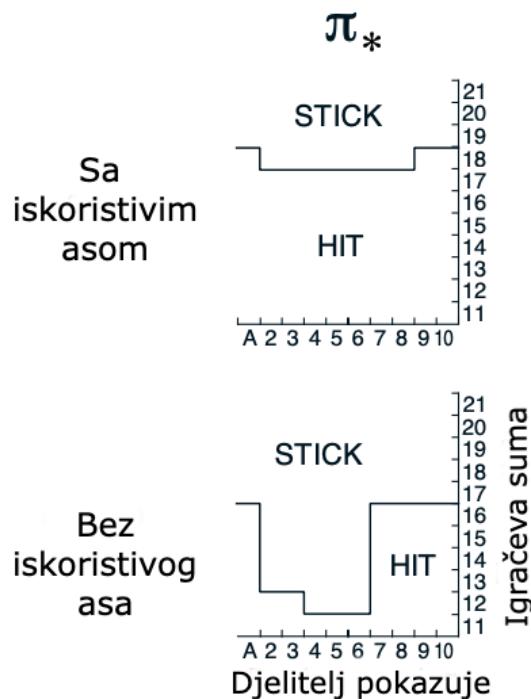
```

Listing 5.4: program.js

Dakle, import-aju se sve funkcije koje će se koristiti iz *pomocni.js* i postav-

lja se vrijednost epsilona (u ovom konkrentno slučaju 0.1, tj. u prosjeku svaku desetu akciju se iskoristi za istraživanje novih stanja). Zatim se definira *staru Tablicu Vrijednosti* kojom se utvrđuje početnu strategiju (u ovom je primjeru za početnu strategiju odabrana ista strategija koju koristi kuća, a iz teorije je poznato da od koje god strategije počeli uvijek će se na kraju doći do optimalne). U *do-while* petlji može se vidjeti implementaciju GPI-a gdje *for* petlja služi za vrednovanje strategije.

Rezultati strategije su sljedeći:



Slika 5.1: Rezultat kojeg su dobili Sutton i Barto [2]

Kod usporedbe primjera iz knjige i rezultata ovog jednostavnog programa predstavljenog u ovom radu, rezultati su prilično zadovoljavajući budući da su ostvarili prvotno zamišljenu svrhu. Kod stanja s iskoristivim asom dobije

se identičnu strategiju, a kod stanja bez iskoristivog asa se razlikuju u 3 (od 100) stanja - ovo je posljedica činjenice da se partije stvarno odvijaju nasumično. 10 milijuna odigranih partija je velik broj, ali i dalje je značajno manji od beskonačno.

Program je kod testiranja pokrenut oko 100 puta i u pravilu se rezultat razlikovao u jednom ili dva stanja od ovoga prikazanog na slici, a nerijetko je rezultat bila i optimalna strategija koju su postigli Sutton i Barto. Ono što se značajno razlikuje od izvođenja do izvođenja je broj potrebnih iteracija: od nekoliko iteracija do 100+. Ali kad se oslabi uvjet na *do-while* petlji tako da stane kada se strategije razlikuju npr. manje ili jednako 2 stanja, program u pravilu stane u manje od 5 iteracija - pomoćne funkcije za implementaciju ovoga pristupa mogu se vidjeti na sljedećem GitHub profilu <https://github.com/0zra/Diplomski>. Ovim je programom riješeno pitanje nalaženja optimalne strategije u nepoznatoj okolini pomoći e-pohlepnog istraživanja, poboljšanja strategije i metodama vremenske razlike.

```

[petarozretic@MacBook-Pro-2 Diplomski % node program.js
brojac: 0
brojac: 1
brojac: 2
brojac: 3
brojac: 4
brojac: 5
brojac: 6
brojac: 7
brojac: 8
brojac: 9
brojac: 10
brojac: 11
brojac: 12
brojac: 13
brojac: 14
s Asom:
      A   2   3   4   5   6   7   8   9   10
21  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
20  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
19  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
18  [ "h", "s", "s", "s", "s", "s", "s", "h", "h" ]
17  [ "h", "h", "h", "h", "h", "h", "h", "h", "h" ]
16  [ "h", "h", "h", "h", "h", "h", "h", "h", "h" ]
15  [ "h", "h", "h", "h", "h", "h", "h", "h", "h" ]
14  [ "h", "h", "h", "h", "h", "h", "h", "h", "h" ]
13  [ "h", "h", "h", "h", "h", "h", "h", "h", "h" ]
12  [ "h", "h", "h", "h", "h", "h", "h", "h", "h" ]
=====
bez Asa:
      A   2   3   4   5   6   7   8   9   10
21  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
20  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
19  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
18  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
17  [ "s", "s", "s", "s", "s", "s", "s", "s", "s" ]
16  [ "h", "s", "s", "s", "s", "h", "h", "h", "s" ]
15  [ "h", "s", "s", "s", "s", "h", "h", "h", "h" ]
14  [ "h", "s", "s", "s", "s", "h", "h", "h", "h" ]
13  [ "h", "s", "s", "s", "s", "h", "h", "h", "h" ]
12  [ "h", "s", "s", "s", "s", "h", "h", "h", "h" ]
petarozretic@MacBook-Pro-2 Diplomski %

```

Slika 5.2: Rezultat izvođenja programa: nakon 15 iteracija rezultat je optimalna strategija. U redu iznad matrica se nalazi karta koju pokazuje kuća, a s lijeve strane suma koju igrač ima.

# Poglavlje 6

## Zaključak

Podržano učenje zadnjih godina puni naslovnice medija: od pobjeđivanja pravaka u Go ili profesionalaca u Starcraftu II do slaganja proteina i kontroliranja plazme u fuzijskom rekatoru. U ovome radu, zbog sadržajnih ograničenja, obrađene su samo osnove PU-a, ali u dovoljnoj mjeri da bi se mogao naslutiti daljnji razvoj, iako nije izrijekom spomenut. Naime PU u kombinaciji s dubokim neuronskim mrežama, tzv. *Duboko podržano učenje*, tek otključava svoj pravi potencijal. Ovaj je rad zbog svojih unaprijed određenih okvira, postavio skroman temelj koji otvara mogućnosti dalnjem razvoju i nalaženju novih načina primjene podržanog učenja. Samo stotinjak linija koda postiglo je isti rezultat kao E. Thorp koji je svojom knjigom izmijenio kockarsku industriju. Dakle, prostora za rast i napredovanje svakako ima. No, ovaj je rad postigao prvotnu namjeru, doskočiti klasičnim poteškoćama i preprekama koje se pojavljuju kod korištenja različitih metoda PU. Njihovom kombinacijom napisani je program prilično uspješno došao do optimalne vrijednosne funkcije i optimalne strategije za igru Blackjacka. Kroz taj praktični primjer, osim navedenih rješivih prepreka, moglo se naslutiti i otvoreno pitanje odnosa iskorištavanja znanja i istraživanja - možda bi povećanjem epsilona sa 0.1 na

0.2 smanjili prosječan broj iteracija: je li 10 000 000 partija bilo pretjerivanje za vrednovanje strategije; Sutton i Barto su u svom primjeru bili zadovoljni sa 500 000. I u konačnici, ako se želi poboljšati rezultate koji bi se dobili isključivo podržanim učenjem, je li odgovor u kombinaciji podržanog učenja i dubokih neuronskih mreža, odnosno u dubokom PU-u.

# Literatura

- [1] Silver, D., *Introduction to reinforcement learning*  
<https://www.youtube.com/playlist?list=PLqYmG7hTraZDM-OYHWgPebj2MfCFzF0bQ>  
<https://www.davidsilver.uk/teaching/>
- [2] Sutton, R.S., Barto, A.G., *Reinforcement learning*  
<http://incompleteideas.net/book/RLbook2020.pdf>
- [3] Szepesvári, C., *Algorithms for Reinforcement Learning*  
<https://sites.ualberta.ca/~szepesva/papers/RLAlgsInMDPs.pdf>
- [4] Blažević, L., *Podržano učenje i Q-učenje*  
<https://www.mathos.unios.hr/~mdjumic/uploads/diplomski/BLA42.pdf>
- [5] Čupić, M., *Umjetna inteligencija: Podržano učenje*  
<http://java.zemris.fer.hr/nastava/ui/rl/rl-20200401.pdf>
- [6] Kaparthy, A., *Deep Reinforcement Learning: Pong from Pixels*  
<https://karpathy.github.io/2016/05/31/rl/>
- [7] *AlphaGo - The Movie — Full award-winning documentary*  
<https://www.youtube.com/watch?v=WXuK6gekU1Y>

## Literatura

- [8] *DeepMind StarCraft II Demonstration*  
<https://www.youtube.com/watch?v=cUTMhmVh1qs>
- [9] *Accelerating fusion science through learned plasma control*  
<https://deepmind.com/blog/article/Accelerating-fusion-science-through-learned-plasma-control>
- [10] *AlphaFold: The making of a scientific breakthrough*  
<https://www.youtube.com/watch?v=gg7WjuFs8F4>
- [11] *DeepMind x UCL — Deep Learning Lecture Series 2021*  
<https://www.youtube.com/playlist?list=PLqYmG7hTraZDVH599EIt1EWsU0sJbAodm>

## TEMELJNA DOKUMENTACIJSKA KARTICA

### PRIRODOSLOVNO-MATEMATIČKI FAKULTET SVEUČILIŠTA U SPLITU ODJEL ZA MATEMATIKU

### DIPLOMSKI RAD **OSNOVE PODRŽANOG UČENJA**

Petar Ozretić

#### **Sažetak:**

Ovaj se rad bavi osnovama podržanog učenja iznoseći sve poznate metode korištene u PU poput Bellmanovih jednadžbi, Dinamičkog programiranja, teorema poboljšanja strategije, Monte-Carlo učenja, učenja s vremenskom razlikom itd. Osnovni problem je bio nalaženje optimalne strategije uz pomoć optimalne vrijednosne funkcije, i to u okolini bez modela. Kombinacijom navedenih metoda, posebno MC iteracije strategija, napisan je kod koji je uspješno nalazio optimalne strategije na primjeru igre Blackjack.

#### **Ključne riječi:**

agent, akcija, strategija, vrijednosna funkcija, model okoline, Markovljevo svojstvo, Bellmanove jednadžbe, Monte Carlo kontrola, učenje s vremenskom razlikom, dinamičko programiranje, pohlepno istraživanje, SARSA vrednovanje strategije

#### **Podatci o radu:**

62 stranice, 15 slika i 1 tablica, 11 literaturnih navoda, pisano na hrvatskom jeziku)

**Mentor:** izv.prof. dr. sc. Saša Mladenović

#### **Članovi povjerenstva:**

izv.prof. dr. sc. Saša Mladenović

TEMELJNA DOKUMENTACIJSKA KARTICA

*doc.dr.sc. Vesna Gotovac Đogaš*

*dr.sc. Ana Perišić*

Povjerenstvo za diplomski rad je prihvatio ovaj rad *UPISATI DATUM  
ODOBRENJA OBRANE*

TEMELJNA DOKUMENTACIJSKA KARTICA

FACULTY OF SCIENCE, UNIVERSITY OF SPLIT  
DEPARTMENT OF MATHEMATICS

MASTER'S THESIS  
**BASICS OF REINFORCEMENT  
LEARNING**

Petar Ozretic

**Abstract:**

*This paper deals with the basics of reinforcement learning, presenting all known methods used in PU such as Bellman equations, Dynamic programming, strategy improvement theorem, MC learning, temporal difference learning, etc. The basic problem was finding an optimal strategy that converges to an optimal value function, and in a model-free environment. By combining the above methods, especially MC strategy iteration, a code was written that successfully found optimal strategies on the example of the Blackjack game.*

**Key words:**

*agent, action, strategy, value function, environment model, Markov property, Bellman equations, Monte Carlo control, temporal difference learning, dynamic programming, greedy strategy, SARSA strategy*

**Specifications:**

*62 pages, 15 figures and 1 table, 11 references, written in Croatian)*

**Mentor:** Saša Mladenović, PhD, Associate Professor

**Committee:**

*Saša Mladenović, PhD, Associate Professor*

*Vesna Gotovac Dogaš, PhD, Assistant professor*

*Ana Perišić, PhD, Postdoctoral Researcher*

TEMELJNA DOKUMENTACIJSKA KARTICA

This thesis was approved by a Thesis commettee on *UPISATI DATUM  
ODOBRENJA OBRANE*