

UNIVERZA V LJUBLJANI
FAKULTETA ZA MATEMATIKO IN FIZIKO

Finančna matematika – 1. stopnja

Tim Kalan

Spodbujevalno učenje pri igranju namiznih iger

Delo diplomskega seminarja

Mentorica: izr. prof. dr. Marjetka Knez

Ljubljana, 2021

KAZALO

| | |
|--|---|
| 1. Uvod | 4 |
| 1.1. Motivacija | 4 |
| 1.2. Strojno učenje | 4 |
| 1.3. Struktura naloge | 4 |
| 2. Spodbujevalno učenje | 5 |
| 2.1. Osnovni koncepti | 5 |
| 2.2. Korak spodbujevalnega učenja | 7 |
| 2.3. Markovski proces odločanja | 8 |
| 2.4. Algoritmi | 8 |
| 2.5. Izboljšave | 8 |
| 3. Namizne igre | 8 |
| 3.1. Pregled konceptov teorije iger | 8 |
| 3.2. Kompleksnost iger | 9 |
| 3.3. Morda kaj o optimal board representationu? | 9 |
| 3.4. Pride še kaj v poštev tu? | 9 |
| 4. Spodbujevalno učenje pri namiznih igrah | 9 |
| 4.1. Parcialni model - »po-stanja« | 9 |
| 4.2. Učenje | 9 |
| 4.3. Kombinacija z iskanjem - generalised policy eval? | 9 |
| 4.4. Algoritem - zaključena celota | 9 |
| 5. Empirični rezultati | 9 |
| 5.1. m,n,k -igra | 9 |
| 6. Uspehi glede na velikost | 9 |
| 7. Primerjava, grafi | 9 |
| Literatura | 9 |

Spodbujevalno učenje pri igranju namiznih iger

POVZETEK

V povzetku na kratko opišite vsebinske rezultate dela. Sem ne sodi razlaga organizacije dela – v katerem poglavju/razdelku je kaj, pač pa le opis vsebine.

Reinforcement learning in board games

ABSTRACT

Prevod zgornjega povzetka v angleščino.

Math. Subj. Class. (2010): navedite vsaj eno klasifikacijsko oznako – dostopne so na www.ams.org/mathscinet/msc/msc2010.html

Ključne besede: Spodbujevalno učenje

Keywords: Reinforcement learning

1. UVOD

Namizne igre ljudje igramo že od prazgodovine. Na Kitajskem je bila igra Go znana kot ena izmed štirih umetnosti Kitajskega učenjaka poleg igranja inštrumenta s strunami, kaligrafije in slikanja. Spremljajo nas že zelo dolgo časa, zato je naravno, da jih želimo ljudje čim bolje igrati.

Z adventom računalnika in računalništva je bil ta problem postavljen v novi luči. Vprašanje ni bilo več samo, kako dobro lahko človek igra igro sam, temveč tudi do kakšnega nivoja lahko spravi računalnik. Izkazalo se je, da nam pri tem problemu (in mnogih drugih) zelo dobro koristi »umetna inteligenca« oz. metode strojnega učenja (SU). Eno izmed vej SU bomo predstavili v tem delu in pogledali, kako nam lahko pomaga pri igranju namiznih iger.

Ideja, da bi nek stroj igral igre ni nova, in kompleksnosti takega stroja so se zavedali ljudje že pred obstojem računalnika. Za konec uvodnega dela morda zabeležimo še citat iz eseja ameriškega pisatelja in pesnika Edgarja Allana Poea, ki govori o mehaničnem igralcu šaha:

»Če prej omenjenemu [igralcu šaha] rečemo čisti stroj, moramo biti pripravljeni priznati, da je zunaj vseh primerjav, najbolj čudovit izum človeštva.

1.1. Motivacija. Spodbujevalno učenje ima zelo lepo motivacijo, in sicer izhaja iz psihologije. Znana psihologa Thorndike in Skinner, sta na živalih izvajala eksperimente; postavila sta jih v neko novo situacijo, kjer je lahko žival naredila akcijo, ki je rezultirala v neki nagradi. Ko je bila žival ponovno postavljena v to situacijo, je hitreje ugotovila, katero akcijo mora storiti, da pride do nagrade.

Koncept, ki je opisan v zgornjem odstavku, se imenuje instrumentalno pogojevanje. Z njim se srečamo tudi ljudje; tako se namreč učijo otroci, odrasli ljudje pa se bolj zanesejo na logično razmišljanje. Vseeno pa je to motiviralo utemeljitelje spodbujevanega učenja

1.2. Strojno učenje. To relativno novo raziskovalno področje se deli na tri glavne veje:

- **Nadzorovano učenje** se ukvarja s tem, kako iz nekih označenih podatkov naučimo računalnik, da prepozna razne signale (slike, govor, tekst, ...) in to znanje uporabi za razpoznavo novih, neoznačenih podatkov.
- **Nenadzorovano učenje** odstrani označevanje iz podatkov in v njih probava odkriti skrite vzorce.
- **Spodbujevalno učenje** se ukvarja z »učenjem iz izkušenj.«

1.3. Struktura naloge. Naloga je razdeljena na štiri glavne dele. Na začetku so predstavljeni osnovni koncepti spodbujevalnega učenja in nekateri glavni algoritmi s tega področja. Potem se fokus obrne na namizne igre in ob nekaj malega teorije iger povzame osnovne koncepte, na katere naletimo tam. V naslednjem odseku potem združimo znanje iz prejšnjih dveh in predstavimo, kako nam teorija iger pripomore pri spodbujevalnem učenju v tem kontekstu. Na koncu pa so predstavljeni nekateri empirični rezultati, ki so posledica zgoraj navedene teorije.

2. SPODBUJEVALNO UČENJE

Spodbujevalno učenje se ukvarja s ti. učenjem iz interakcije oz. izkušenj. Čeprav se to na prvi pogled ne zdi kot računska metoda, pač pa stvar psihologije, bomo kmalu dognali, kako prevesti to idejo v računalniku razumljiv jezik.

2.1. Osnovni koncepti. V osnovi nas zanima precej preprosta stvar: kako preslikati neko opazovano situacijo v akcijo na tak način, da maksimiriziramo neko numerično nagrado. Pri tem ne obstaja opazovalec, ki bi nam povedal ali pa namignil, katere akcije so dobre, to moramo ugotoviti sami, s poskušanjem in napakami. V tem dejstvu se skriva bistvena razlika med spodbujevalnim učenjem in ostalimi vejami strojnega učenja.

Pomembna razlika tiči tudi v pomembnosti časa pri spodbujevalnem učenju. Pri drugih oblikah strojnega učenja se ponavadi ukvarjamo s tabelaričnimi podatki, tu pa ponavadi modeliramo nek dinamičen proces, zato je naravno, da je pomemben čas. Čeprav se ga da v tem kontekstu modelirati zvezno, je za naše namene dovolj, da ga jemljemo kot diskretne točke $t \in 1, \dots, T$, kjer T označuje nek končni čas (v splošnem je lahko seveda $T = \infty$).

2.1.1. Nagrada. Prvi pomemben koncept pri spodbujevanem učenju je nagrada. Kot smo že zgoraj omenili, to za nas pomeni neko numerično vrednost, pozitivno število indicira pozitivno nagrado, negativno pa »kazen«. S pomočjo tega koncepta formaliziramo *cilj* učenja. Edini cilj učenca je maksimizacija te nagrade, pri čemer je vredno omeniti, da na nagrado učenec lahko vpliva samo s svojimi akcijami (ne more recimo spremeniti načina, na katerega dobi nagrado).

Posebaj pomembno je na tem mestu poudariti, da akcije nimajo nujno neposredne nagrade. Le-te lahko pridejo v poljubnem kasnejšem časovnem obdobju. To je smiselno, če pomislimo z vidika namiznih iger: pri šahu ne razmišljamo samo o neposrednih akcijah, temveč razvijamo neko dolgoročno strategijo ki nas na koncu nagradi z zmago.

Zgled 2.1 (Križci in krožci). *Pri tej znani otroški igri (in pri mnogo drugih namiznih igrah) modeliramo nagrado na preprost način: če zmagamo, prejmemo nagrado 1, če izgubimo pa -1 . V vseh ostalih situacijah, torej za izenačenje in po vsaki potezi, prejmemo nagrado 0.*

Zavedati se moramo tudi potencialnih omejitev oz. pomanjkljivosti takega modela. Razmislimo malo o:

Definicija 2.2 (Hipoteza o nagradi). Vse cilje je mogoče opisati kot maksimizacijo neke kumulativne numerične nagrade.

Zgled 2.3 (Protiprimera hipotezi o nagradi). *Problem je, da hipoteza dovoljuje samo enodimezionalnost:*

- Ko kupujemo hamburger, nam je pomemben okus in cena; kaj nam več pomeni?
- Država želi med epidemijo ohraniti življenja in gospodarstvo; v kakšni meri naj prioritizira ti dve kategoriji?

Na tem mestu poudarimo še, da se da tudi v takih situacijah modelirati nagrado na zgoraj opisani način in da je ta koncept vseeno dovolj splošen, da zajame zelo velik razred problemov.

2.1.2. *Okolje*. Okolje predstavlja del našega sistema, na katerega učenec nima nobenega vpliva. Funkcija okolja je, da učencu pokaže **stanje** (angl. *state*) in mu da nagrado glede na **akcijo**, ki jo prejme od njega. Če se ponovno osredotočimo na namizne igre, bi lahko rekli, da je okolje igralna plošča pri šahu in nasprotnik - tudi nanj namreč nimamo vpliva. Okolje nam služi tudi kot sodnik akcij oz. stanj. V kontekstu programa za igranje iger torej okolje izbira nasprotnikove akcije, odloča katero stanje pomeni zmago in dodeljuje nagrade.

Zgled 2.4 (Križci in krožci). *Okolje za nas pomeni 3×3 igralno polje in našega nasprotnika.*

2.1.3. *Agent*. Zgoraj omenjenemu »učencu« v spodbujevanem učenju formalno rečemo *agent*. Njegov cilj je torej maksimizacija numerične nagrade, to težnjo pa dosega s pomočjo **strategije** (angl. *policy*), ki mu pove, katero akcijo naj izbere v določenem stanju. Za ocenjevanje stanja, si pomaga z **vrednostno funkcijo** (angl. *value function*). Kot ime implicira, je to funkcija, ki določa vrednosti stanjem (in akcijam).

Nagrada nam pove takojšnjo vrednost stanja, vrednostna funkcija pa to vrednost gleda na dolgi rok. Je izpeljanka nagrade, a veliko bolj primerna za maksimizacijo kot nagrada, saj upošteva, da so tudi stanja, ki ne prinesejo takojšnje nagrade, lahko veliko vredna.

Poleg tega je v splošnem lažje učenje prek vrednostnih funkcij kot prek strategij neposredno, saj je ponavlja stanj mnogokrat manj kot možnih strategij agenta.

Formalno gledano:

Definicija 2.5. Naj R_t, S_t, A_t zaporedoma označujejo nagrado, stanje in akcijo ob času t . Definiramo naslednje pojme:

- Agentova **strategija** je takšna preslikava $\pi : S \rightarrow A$ da velja

$$a = \pi(s) \text{ oz.}$$

$$\pi(a|s) = P(A_t = a | S_t = s).$$

Pri čemer prva formula definira *deterministično* strategijo, druga pa *stohastično*. a in s sta realizaciji akcije in stanja v času t (to sta načeloma slučajni spremenljivki.)

- Naj bodo R_{t+1}, \dots, R_T nagrade, ji jih bomo prejeli od trenutka t do končnega časa. **Povračilo** (angl. *return*) G_t definiramo kot

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T,$$

kjer je $\gamma \in [0, 1]$ *diskontni faktor*. Predstavlja dejstvo, da imamo raje nagrade, ki bodo prišle prej. Formalno gledano, je cilj učenja maksimizacija pričakovanega povračila

- Naj bo π dana strategija agenta. **Vrednostna funkcija stanja** (angl. *state value function*) glede na strategijo $v_\pi(s)$ je

$$v_\pi(s) = \mathbb{E}[G_t | S_t = s].$$

Predstavlja torej pričakovani izplen, če se vedemo skladno s strategijo π .

- Naj bo π še vedno dana strategija agenta. **Vrednostna funkcija akcije** (tudi stanja-akcije) (angl. *action value function*) glede na strategijo $q_\pi(s, a)$ definiramo

$$q_\pi(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a].$$

Pove nam pričakovani izplen, če ob času t naredimo akcijo a , nato pa se vedemo skladno s strategijo π .

Zgled 2.6 (Križci in krožci). *Agent je v tem primeru računalniški program, ki prejme igralno ploščo, nasprotnikove poteze in nagrade, vrne pa optimalno strategijo (to si želimo).*

2.1.4. *Model.* Model je nenujen del našega sistema. Predstavlja znanje, ki ga ima agent o svojem okolju. Če imamo model, da lahko uporabimo, da napovemo, kako se bo vedlo okolje in s tem premaknemo agentovo učenje iz čistih poskusov in napak na *načrtovanje* (angl. *planning*). Model je torej poleg strategije in vrednostne funkcije še tretja komponenta agenta. Na podlagi modela lahko agent »preračuna« smiselnost svojih akcij, brez da bi dejansko karkoli storil.

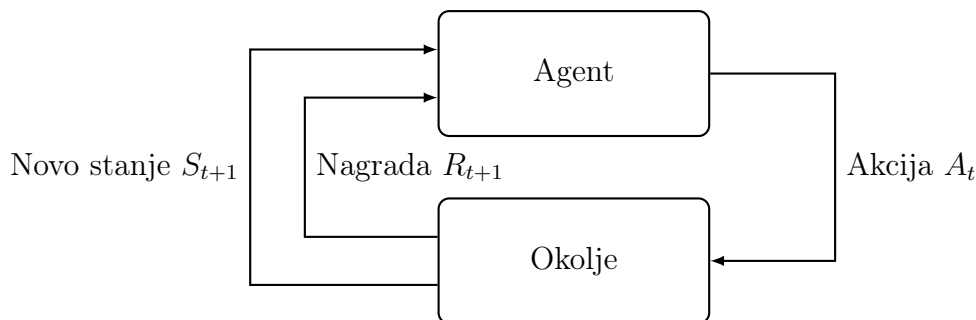
Prisotnost modela je glavna ločnica med dvema velikima, a zelo različnima vejama spodbujevalnega učenja.

2.2. **Korak spodbujevalnega učenja.** Spodbujevalno učenje se pogosto ukvarja s procesi, ki naravno razpadejo v ti **epizode**. Tak proces so recimo namizne igre, kjer so epizode precej naravno posamezne igre. Ni pa nujno, da je delitev tako naravna (ali pa sploh možna oz. smiselna). Za namene te diplomske naloge lahko privzamemo, da taka delitev obstaja.

Ideja učenja je, da agenta spustimo v okolje in mu dovolimo, da doživi (igra) mnogo epizod. Nato na nek način (bo razjasnjeno kasneje) ob nekih določenih časih (npr. po koncu epizode) posodobi svojo strategijo (in/ali vrednostno funkcijo).

Dejanski korak (npr. poteza v namizni igri) v epizodi pa formalno gledano opredelimo:

- Agent naredi akcijo A_t ob prejetem stanju S_t in prejme nagrado R_t .
- Okolje prejme akcijo A_t , posreduje agentu stanje S_{t+1} in nagrado R_{t+1}



2.2.1. *Raziskovanje in izkoriščanje.* Eden izmed glavnih problemov, s katerim se srečamo pri spodbujevalnem učenju je problem raziskovanja in izkoriščanja. Ko se agent uči, začne dojemati katere akcije ali pa kombinacije akcij mu pripeljejo nagrado. Ko to ugotovi, seveda lahko začne te akcije *izkoriščati* in prejemativso nagrado, jim pripada. Pri tem pa naletimo na problem. Agent lahko izkorišča te akcije in nikoli ne ugotovi, da neka druga akcija prinese še višjo nagrado; tega ne izve, ker ne *raziskuje*. Če pa samo raziskuje pa nikoli ne izkoristi potencialnih nagrad, ki jih sreča to je, ničesar se ne nauči.

Uravnoteženje raziskovanja in izkoriščanja je pomemben problem, a se izkaže, da ima dokaj enostavno rešitev (ki deluje dovolj dobro). Spoznali jo bomo v kratkem.

Morda se nekaterim bralcem zdi, da smo zaenkrat preceč »mahali z rokami«, to je zato, ker želimo, da se do te točke razvije intuicija o predstavljenih pojmi. V nadaljevanju bomo do sedaj opisane stvari bolj formalizirali.

2.3. Markovski proces odločanja. Spomnimo se najprej procesa spodbujevalnega učenja in ga poskusimo opisati bolj formalno: imamo zaporedje časovnih korakov $t = 0, 1, 2, \dots$, ob katerih med sabo interaktirata agent in okolje. Ob koraku t agent prejme od okolja stanje (oz. reprezentacijo stanja) $S_t \in \mathcal{S}$, kjer \mathcal{S} označuje množico vseh stanj. Na podlagi stanja in strategije, ki jo ima, izbere akcijo $A_t \in \mathcal{A}(S_t)$, kjer $\mathcal{A}(S_t)$ predstavlja množico akcij, ki jih ima agent na voljo v stanju S_t . Rezultat te akcije je nagrada $R_{t+1} \in \mathcal{R}$ in novo stanje S_{t+1} .

Čeprav se da vse opisane koncepte posplošiti na števne in celo neštevne množice stanj in akcij, se bomo mi omejili na končne množice. To je glede na problem, s katerim se ukvarajmo, dovolj.

2.3.1. Markovska veriga. Dogajanje pri spodbujevalnem učenju lahko v grobem opišemo s slučajnim procesom stanj $(S_t)_{t=0}^T$. Zato je pomembno, da si natančno pogledamo nekaj lastnosti, ki jih lahko pričakujemo.

Definicija 2.7 (Markovska veriga). Slučajni proces $(S_t)_{t=0}^T$ na končnem verjetnostnem prostoru (Ω, P) je **Markovska veriga**, če zanj velja Markovska lastnost

$$P(S_{t+1} = s_{t+1} | S_t = s_t, \dots, S_0 = s_0) = P(S_{t+1} = s_{t+1} | S_t = s_t).$$

Na kratko to *prehodno verjetnost* označimo $p_{ss'} := P(S_{t+1} = s' | S_t = s)$. Opazimo, da lahko te verjetnosti zložimo v matriko $\mathcal{P} := [p_{ss'}]_{s, s' \in \mathcal{S}}$.

Zdaj Markovsko verigo predstavimo še na alternativni način: kot dvojico $(\mathcal{S}, \mathcal{P})$, kjer je \mathcal{P} zgoraj definirana matrika.

Markovska lastnost pomeni, da je prihodnost neodvisna od preteklosti, če poznamo sedanost. Spodbujevalno učenje se ukvarja predvsem s problemi, kjer to dejstvo drži. Tudi pri našem ciljnim problemu to načeloma velja: če pogledamo igralno ploščo na katerikoli točki pogosto izvemo enako o trenutnem stanju, kot če bi opazovali igro od začetka.

2.4. Algoritmi.

2.4.1. *Dinamično programiranje.*

2.4.2. *Monte Carlo.*

2.4.3. *TD(0).*

2.4.4. *TD(λ).*

2.5. Izboljšave.

2.5.1. *Nevronske mreže.*

3. NAMIZNE IGRE

3.1. Pregled konceptov teorije iger.

3.1.1. *Nashevo ravnotežje.*

3.1.2. *Igre z vsoto nič.*

3.1.3. *Ekstenzivne igre.*

3.2. Kompleksnost iger.

3.2.1. *Game tree* ...

3.3. Morda kaj o optimal board representationu?

3.4. Pride še kaj v poštev tu?

4. SPODBUJEVALNO UČENJE PRI NAMIZNIH IGRAH

4.1. Parcialni model - »po-stanja«.

4.2. Učenje.

4.2.1. *Samoigra*.

4.2.2. *Igre iz podatkovnih baz*.

4.2.3. *Naključni nasprotnik*.

4.3. Kombinacija z iskanjem - generalised policy eval?

4.4. Algoritem - zaključena celota.

4.4.1. *Opomba: deluje, tudi ko ni vsota 0*.

5. EMPIRIČNI REZULTATI

5.1. **m,n,k-igra**.

5.1.1. *Kompleksnost m,n,k-igre*.

6. USPEHI GLEDE NA VELIKOST

6.0.1. *Morda tudi kaj z bolj modificiranimi mnk igrami*.

7. PRIMERJAVA, GRAFI

LITERATURA

- [1] Imran Ghory. Reinforcement learning in board games. 2004.
- [2] David Silver. Introduction to reinforcement learning. <https://deepmind.com/learning-resources/-introduction-reinforcement-learning-david-silver>, 2015.
- [3] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An introduction*. The MIT Press, Cambridge, Massachusetts, 2 edition, 2015.
- [4] Csaba Szepesvari. *Algorithms for Reinforcement Learning*. Morgan & Claypool Publishers, Alberta, Canada, 2009.