



MATEMATICKO-FYZIKÁLNÍ
FAKULTA
Univerzita Karlova

BAKALÁŘSKÁ PRÁCE

Marek Sádovský

Použitie strojového učenia pre deskovú hru Azul

Katedra softwaru a výuky informatiky (201. • 32-KSVI)

Vedoucí bakalářské práce: RNDr. Tomáš Holan, Ph.D.

Studijní program: Programování a vývoj software

Praha 2025

Prohlašuji, že jsem tuto bakalářskou práci vypracoval(a) samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů. Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V dne

Podpis autora

Poděkování.

Název práce: Použitie strojového učenia pre deskovú hru Azul

Autor: Marek Sádovský

Katedra: Katedra softwaru a výuky informatiky (201. • 32-KSVI)

Vedoucí bakalářské práce: RNDr. Tomáš Holan, Ph.D., katedra

Abstrakt: Abstrakt.

Klíčová slova: klíčové slovo, složitější fráze

Title: Usage of the machine learning for the board game Azul

Author: Marek Sádovský

Department: Name of the department

Supervisor: RNDr. Tomáš Holan, Ph.D., department

Abstract: Abstract.

Keywords: key, words

Obsah

Úvod	6
1 Predstavenie hry	7
1.1 Príprava hry	7
1.2 Priebeh hry	7
1.2.1 Ponuka výloh	7
1.2.2 Obkladanie stien	7
1.3 Parametre hry	8
1.4 Ako hrať	8
2 Teoretické pozadie	9
2.1 Štruktúra knižnice	9
2.2 Reinforcement learning	9
2.3 Umelá inteligencia bez strojového učenia	9
3 Implemetácia proximal policy optimalizácie PPO	10
3.1 Prostredie a agent	10
3.2 Architektúra modelu	10
3.3 Výhodnosť a optimalizačný krok PPO	11
3.4 Technické výzvy	12
4 Konfigurácie PPO	13
4.1 Hodnotiace funkcie	13
4.1.1 Maximalizácia okamžitého zisku a ignorovanie steny . . .	13
4.1.2 Maximalizácia okamžitého zisku s rešpektom ku stene . .	13
Závěr	15
Seznam obrázků	16
Seznam tabulek	17
Seznam použitých zkratk	18
A Přílohy	19
A.1 První příloha	19

Úvod

Následuje několik ukázkových kapitol, které doporučují, jak by se měla závěrečná práce sázet. Primárně popisují použití T_EXové šablony, ale obecné rady poslouží dobře i uživatelům jiných systémů.

1 Predstavenie hry

V hre Azul je cieľom čo najlepšie a čo najefektívnejšie uložiť dlaždice na stenu. Pri čom musia dlaždice odpovedať nejakému predom určenému vzoru.

1.1 Príprava hry

Každý hráč má na začiatku hernú dosku (ďalej len doska). Na pravej strane dosky môže každý hráč vidieť stenu pozostávajúcu z 5 krát 5 políček kam sa budú časom dostávať dlaždice (**odkaz na kapitolu**). v ľavej časti sú zas zásobníky, do ktorých jednotliví hráči vkladajú dlaždice počas hry. Vľavo hore na doske sa nachádza skóre daného hráča. Na záver na spodku dosky sa nachádza podlaha, miesto kam padajú dlaždice čo sa nám nezmestili do zásobníku (**odkaz**).

Okrem dosiek jednotlivých hráčov sú na hracej ploche aj výlohy s ponukami dlaždíc (**odkaz na postup pri hre**). Počet výloh je závislý na počte hráčov (**odkaz**). Okrem základných výloh existuje aj špeciálna centrálna výloha (ďalej len stred).

1.2 Priebeh hry

Samotná hra sa delí na dve fázy, ktoré sa periodicky striedajú kým nenastane koniec hry.

1.2.1 Ponuka výloh

Na začiatku tejto fázy sa naplnia všetky výlohy okrem stredu z banku (**odkaz**). Do stredu sa presunie dlaždica prvého hráča (**odkaz**). Ďalej sa hráči striedajú vo výbere výlohy a druhu dlaždice v nej. Keď si vyberú výlohu a v nej dlaždicu, zoberú z výlohy všetky dlaždice daného druhu (kliknutím na jednu z dlaždíc vo výlohe) a uložia ju do jedného z možných zásobníkov (kliknutím na zásobník). Zásobník musí vždy obsahovať maximálne jeden druh dlaždíc, ak pridá viac dlaždíc, ako je kapacita, tak nadbytočné dlaždice skončia na podlahe (**odkaz**).

Po úspešnom presunutí dlaždíc do zásobníku sa ostatné dlaždice z danej výlohy presunú do stredu. Nasleduje ťah ďalšieho hráča.

V momente keď sa prvý hráč rozhodne brať zo stredu tak si okrem vybraného druhu dlaždíc berie aj dlaždicu prvého hráča, ktorú si uloží na podlahu (**odkaz**).

V momente keď príde hráč na ťah a nemá si odiaľ brať, všetky výlohy vrátane stredu sú prázdne, hra prechádza do druhej fázy.

1.2.2 Obkladanie stien

V tejto časti hráči presúvajú dlaždicu z plných zásobníkov do odpovedajúcich riadkov, pripisujú si okamžité body, následne vyprázdňujú podlahu a odčítavajú si za ňu bodu (**odkaz**).

Po tom, ako všetci hráči urobili vyššie popísaný krok, skontroluje sa či nenastal koniec hry, ten nastane ak niektorý z hráčov vyplnil celý riadok dlaždicami.

V prípade konca hry si hráči spočítajú bonusové body na základe nasledujúceho diagramu. - za každý celý riadok +2 body - za každý celý stĺpec +7 bodov - za 5 dlaždíc rovnakej farby na ploche +10 bodov

V prípade, že nenastal koniec hry opakuje sa prvá fáza hry

1.3 Parametre hry

Počet výloh je závislý od počtu hráčov: - pre 2 hráčov 5 výloh - pre 3 hráčov 7 výloh - pre 4 hráčov 9 výloh

V hre sa nachádza 100 unikátnych dlaždíc rovnomerne rozdelených na 5 druhov takže 20 dlaždíc každého druhu.

Doska obsahuje stenu o veľkosti 5 krát 5, 5 zásobníkov vo veľkostiach 1 až 5. Podlaha na spodku dosky má celkovú veľkosť 7. Ak by mal počet dlaždíc na podlahe prekročiť daný počet 7, je každá nadbytočná dlaždica odstránená.

Hodnotenie pri dopĺňaní na stenu je nasledovné. Ak doplnená dlaždica nesusedí s žiadnou inou dlaždicou tak dostaneme 1 bod, ak susedí aspoň s jednou ďalšou dlaždicou v rámci riadku alebo stĺpcu tak je počet bodov vypočítaný, ako súčet susediacich v riadku plus počet susediacich v stĺpci. ([odkaz](#)).

1.4 Ako hrať

Po spustení hry sa hráčovi/hráčom načíta menu kde uvidia tlačítka "Play", ktorým spustia hru, pod ním sa nachádza nastavenie na počet hráčov a ich druh, počet sa upravuje kliknutím na plus či mínus, a druh jednotlivých hráčov kliknutím na ikonu konkrétneho hráča, kde sú k dispozícii dve možnosti, AI (bot) alebo H (človek). Okrem toho kliknutím na meno hráča je možné ho zmeniť.

Po nastavení hráčov sa stlačením *Play* presunie hra do herného okna. Tu sa striedajú hráči v hraní, ako je popísané v priebehu hry, a to tak, že hráč na ťahu v rámci fázy jedna, Ponuka výloh, kliknutím na dlaždicu v rámci výloh, vyberie dlaždice z danej výlohy, daného druhu. Hráč ich môže ďalej položiť kliknutím pravého tlačidla myši, alebo ich môže uložiť do zásobníku kliknutím na zásobník, v prípade, že nemá žiadny zásobník môže ich uložiť na podlahu taktiež kliknutím na ňu.

V rámci druhej fázy je vidno preddefinovanú stenu hráča, takže hráč len klikne a automaticky sa presunie dlaždica z prvého plného zásobníku na stenu na správne miesto, toto iteruje kým má nejaký plný zásobník. V prípade, že nastane koniec hry po doplnení všetkých hráčov sa spočítajú bonusy a ukáže sa tabuľka s výsledkami.

2 Teoretické pozadie

V predchádzajúcej kapitole sme popísali pravidlá hry, skor ako sa pustíme priamo k práci treba si ešte predstaviť ako vyzerá daná knižnica aby bolo jenznačné ako s ňou operovať.

- Podporované platformy: Linux, Windows - Jazyk: C s .net 8 - Knižnice: Unity knižnice (Unity, UnityEngine atď.)

2.1 Štruktúra knižnice

Knižnica má čisto backend-ové rozhranie, čo znamená, že neobsahuje žiadnu možnosť ovládania aplikácie, na to však máme vedľajšie riešenie. V princípe na hranie nie je potreba žiadna iná trieda takže konkrétne informácie o jednotlivých triedach sú dohľadateľné v programátorskej dokumntácií. **TODO odkaz**

Board nám dáva možnosť zachytávať čo sa deje v hre pomocou eventov ktoré sa volajú keď hra očakáva nejakú interakciu z vonku. Z vonku sa volajú len dve metódy a to Move(), pre špecifikovanie výlohy dlaždice a zásobníka, a Calculate() na uloženie na stenu. Pre trénovanie umelej inteligencie môže byť ešte zaujímavá metóda EncodeBoardState() ktorá vráti stav ako číselný vektor.

2.2 Reinforcement learning

Reinforcement learning (učenie formou odmeňovania). V rámci tejto metódy strojového učenia je agent (**odkaz**), nasadený do prostredia a učí sa na základe interakcií s prostredím.

Agentovi sú vopred určené pravidlá ako sa môže správať a taktiež odmeňovacia funkcia, ktorá agentovi určuje či je jeho rozhodnutie pre neho prospešné alebo nie. Pomocou iteratívneho postupu sa agent metódou pokus omyl naučí ako sa správať v prostredí.

Prostredie nad ktorým budeme pracovať je hra Azul, spomínaná vyššie. Knižnica nám ponúka metódu EncodeBoardState, ktorá nám vracia vektor kde sú postupne zakódované stavy výloh, doska hráča čo je na ťahu a dosky ostatných hráčov. Pre zjednodušenie momentálny agent si vyberá len v rámci prvej fázy hry, kde existuje potenciálne 300 roznych akcií. Druhá fáza sa vie líšiť od štýlu hry, v rámci základnej hry, je deterministicky dané správanie.

Predpokladajme, že hrajú 4 hráči, to znamená, že máme 9 výloh plus stred, nech sa ťah skladá z troch čísel nejaké id výlohy, id dlaždice a id zásobníka, v tom prípade máme 10 výloh, 5 druhov dlaždíc, 5 zásobníkov a podlaha. $10 \times 5 \times (5 + 1) = 300$.

2.3 Umelá inteligencia bez strojového učenia

Okrem strojového učenia existujú aj iné možnosti, umelej inteligencie, vo forme heruistik alebo pomocou stavových stromov (minmax). Tieto algoritmi **TODO bla bla**

3 Implementácia proximal policy optimalizácie PPO

Táto kapitola detailne popisuje implementáciu algoritmu Proximal Policy Optimization (PPO) v jazyku C# pre potreby trénovania umelej inteligencie v hre Azul. Algoritmus bol implementovaný od základov bez použitia externých ML knižníc, čo umožnilo presnú kontrolu nad všetkými komponentmi systému.

3.1 Prostredie a agent

Prostredie predstavuje samotnú hru Azul, pričom každý stav hry je reprezentovaný, ako množina čísel opisujúcich aktuálne rozloženie dlaždíc **TODO ref** v rámci výloh, stav dosky hráča a stavy ostatných dosiek hráčov. Pre dosiahnutie menšieho vektora reprezentujúceho konkrétny stav, ho pred predaním kompresujeme na menšiu veľkosť. Agent získava tieto informácie vo forme vstupného vektora, ktorý je ďalej spracovaný jeho modelom.

Priestor akcií majme definovaný, ako množinu všetkých možných pohybov existujúcich v hre. To preto, že v každom stave existuje inak veľká množina legálnych ťahov čo by viedlo k zbytočným komplikáciám. Z tohoto dôvodu používa agent v rámci ťahu maskovanie ilegálnych možností aby nebral do úvahy ťahy ktoré sú neplatné.

Agent, ako už spomínané v predchádzajúcej kapitole, využíva rozhranie IBot na komunikáciu s herným prostredím. Pri každom ťahu agent vykoná nasledovné kroky:

1. Zakóduje aktuálny stav hry do numerického vektora.
2. Pomocou policy siete získa pravdepodobnosti možných akcií.
3. Vynuluje pravdepodobnosti neplatných akcií a opätovne normalizuje zvyšné.
4. Na základe pravdepodobností náhodne vzorkuje akciu.
5. Zaznamená stav, akciu, pravdepodobnosti a priebežnú odmenu pre neskorší tréning.

3.2 Architektúra modelu

Politika (policy) a aj hodnotová funkcia (value function), používajú implementáciu jednoduchej doprednej neurónovej siete (feedforward network) so zdieľanými vrstvami.

Architektúra našej neurónovej siete pozostáva z týchto vrstiev:

- Vstupná vrstva: počet neurónov zodpovedá veľkosti vektoru stavu.
- Skrytá vrstva 1: 256 neurónov, aktivačná funkcia ReLU

- Skrytá vrstva 2: 256 neurónov, aktivačná funkcia ReLU **TODO: vysvetliť prečo tato veľkosť**
- Výstupná vetva politiky: softmax nad počtom možných akcií
- Výstupná vetva hodnoty: jeden neurón (skalárna hodnota stavu)

Implementácia siete prebieha v C# prostredníctvom vlastnej triedy NeuralNetwork, ktorá podporuje základné operácie ako dopredný prechod, spätnú propagáciu a optimalizáciu parametrov.

3.3 Výhodnosť a optimalizačný krok PPO

Agent počas hrania zbiera trajektóriu stavov, akcií a odmien v každej epizóde. Odmeny r_t sú nasledovne diskontované pomocou diskontného faktoru $\gamma \in [0,1)$ následovne:

$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

Kde R_t je tzv. discounted return, diskontovaná kumulatívna odmena od času t .

Na rozdiel od klasického policy gradientu, PPO zavádza tzv. clipped surrogate objective, ktorý bráni príliš veľkým zmenám v politike medzi iteráciami. Základný tvar straty (loss) pre policy sieť je:

$$L^{CLIP}(\theta) = \mathbb{E}[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)]$$

Kde:

- θ sú parametre aktuálnej politiky
- $r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{OLD}}(a_t|s_t)}$ je pomer pravdepodobností akcie pred a po aktualizácii politiky.
- ϵ je hiperparameter, ktorý určuje o koľko sa môže politika zmeniť
- funkcia clip() zaručuje že aktualizácia politiky nepresiahne povolený rozsah

Táto funkcionálna zabezpečuje, že aktualizácia sa bude riadiť výhodnosťou, ale zároveň nedovolí politike „skočiť“ príliš ďaleko, čím sa zvyšuje stabilita tréningu.

Na podporu exploračie počas tréningu sa do stratovej funkcie pridáva aj entropia výstupu politiky:

$$L^{total} = L^{CLIP}(\theta) + c \cdot H[\pi_{\theta}(\cdot|s_t)]$$

Kde:

-

$$H[\pi_\theta] = - \sum_i p_i \cdot \log(p_i)$$

- je entropia politiky

- c je váhový koeficient určujúci vplyv entropie na celkovú stratu

Po výpočte straty parametre politiky aktualizujeme pomocou gradientného zostupu (gradient descent) v rámci spätnej propagácie (backpropagation). Ide o štandardný prístup pri optimalizácii parametrov siete, ktorý spočíva v iteratívnej aktualizácii váh a biasov tak, aby sa minimalizovala hodnota straty (loss funkcie).

3.4 Technické výzvy

Implementácia PPO priniesla viacero technických výziev. Medzi najväčšie patrila absencia pokročilých ML knižníc, keďže pracujeme v jazyku C#, čo si vyžadovalo ručné implementovanie všetkých komponentov: neurónových sietí, stratových funkcií, optimalizácie a numerických operácií. Z počiatku sa ukazoval problém so stabilitou, hodnoty v rámci výpočtov v neurónovej sieti často viedli k veľkým číslam umožňujúcim pretečenie alebo naopak v log funkciach k nule kde hrozil nedefinovaný stav. Ako riešenie tohoto problému sme použili clip hodnoty na obmedzenie hodnôt a epsilon na zabránenie $\log(0)$.

4 Konfigurácie PPO

Vzhľadom na komplexitu hry Azul nie je jednoduché určiť aká konfigurácia parametrov je najideálnejšia, či sa už jedná o veľkosť neurónovej siete alebo o rýchlosť učenia.

4.1 Hodnotiace funkcie

Hra Azul je náhodná a má veľké spektrum stratégií, preto si vytvoríme viaceré možnosti na porovnanie.

4.1.1 Maximalizácia okamžitého zisku a ignorovanie steny

Algoritmus 1 V hodnotení podľa maximalizácie okamžitého zisku sa ignorujú odložené hodnotenia a tým sa zanedbáva dlhodobější stratégia.

```
1: function VYHODNOTSTAV(tah, stav)
2:   reward = 0  $\leftarrow$  Nastavíme základnú hodnotu
3:   if tah == TahNaPodlahu then
4:     Vrátime -10 ▷ ukladať na podlahu je strata
5:   end if
6:   novyPocetNaDoske = tah.Pocet + stav.nasaDoska.pocetVRiadku
7:   if novyPocetNaDoske >= stav.nasaDoska.kapacitaRiadku then ▷
   Znamená, že dotaneme body v tomto kole reward+ = 5 ▷ 5 je maximálna veľkosť riadku
   ▷ Nechceme pridávať na podlahu
   reward- = novyPocetNaDoske - stav.nasaDoska.kapacitaRiadku
8:   else reward+ = tah.Pocet
9:   end if
10: end function
```

Toto riešenie je jednoduché a javý sa účine ale len voči náhodnému oponentovi. Pokiaľ oponent používa nejakú stratégiu nemá problém poraziť agenta s touto hodnotiacou funkciou. **TODO viac graf**

4.1.2 Maximalizácia okamžitého zisku s rešpektom ku stene

Na rozdiel od predchádzajúceho riešenia teraz, sa bude brať do úvahy aj to ako to čo berieme interaguje už uloženými vecami na stene.

Algoritmus 2 V hodnotení podľa maximalizácie okamžitého zisku sa ignorujú odložené hodnotenia a tým sa zanedbáva dlhodobejšia stratégia.

```

1: function VYHODNOTSTAV(tah, stav)
2:   reward = 0  $\leftarrow$  Nastavíme základnú hodnotu
3:   if tah == TahNaPodlahu then
4:     Vrátime -10 ▷ ukladať na podlahu je strata
5:   end if
6:   novyPocetNaDoske = tah.Pocet + stav.nasaDoska.pocetVRiadku
7:   if novyPocetNaDoske  $\geq$  stav.nasaDoska.kapacitaRiadku then ▷
     Znamená, že dotaneme body v tomto kole reward + = 5 ▷ 5 je maximálna
     veľkosť riadku
     ▷ Nechceme pridávať na podlahu
     reward - = novyPocetNaDoske - stav.nasaDoska.kapacitaRiadku
8:   else reward + = tah.Pocet
9:   end if
10: end function

```

Závěr

Seznam obrázků

Seznam tabulek

Seznam použitých zkratek

A Přílohy

A.1 První příloha