

Reinforcement learning pentru jocul de sah

Cristian Ivan

1 Introducere

Aceasta lucrare de licență are ca subiect aplicarea algoritmilor de reinforcement learning pentru diferite variante ale jocului de sah. Obiectivul general al lucrării este dezvoltarea unui algoritm de reinforcement learning și apoi testarea acestuia cu ajutorul unui engine de sah - Fairy Stockfish, un engine derivat din Stockfish, cel mai cunoscut și performant engine de sah actual. În scopul acestei lucrări de licență am studiat diferite capitole de rețele neuronale, reinforcement learning, algoritmi de căutare minimax cu alpha-beta pruning, cât și alte metode care sunt folosite în mod uzual în inteligența artificială pentru sah. Motivația alegerii acestui subiect, cât și alte abordări existente vor fi acoperite în subcapitole viitoare.

În afara regulilor standard de sah, vom folosi încă trei variante ale jocului de sah:

1.1 Reguli alternative de sah

Antisah (Antichess)

Această variantă a jocului are următoarele reguli:

- Jucătorul care își pierde toate piesele (inclusiv regele) sau rămâne fără mutări legale câștigă partida.
- Capturarea pieselor adversarului este obligatorie. Dacă un jucător poate captura o piesă adversă, este obligat să o facă.
- Dacă există mai multe opțiuni de captură, jucătorul poate alege care piesă să captureze, dar trebuie să captureze.

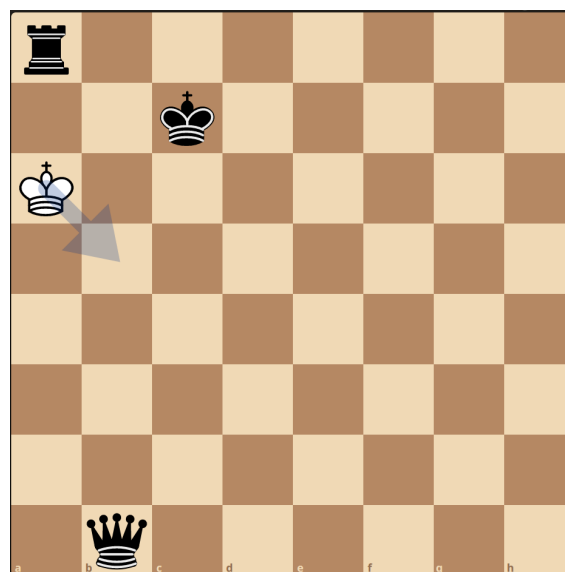


Figure 1: Imagine creată cu ajutorul lichess

În imaginea de mai sus, jucătorul cu piesele albe are o poziție câștigătoare: când muta regele pe b5, jucătorul cu piesele negre este obligat să îi captureze regele.

Secvența de mutări câștigătoare : 1.Kb5 Qxb5#

King Of The Hill

Această variantă adaugă o singură regulă nouă : un jucător poate câștiga partida prin două metode: fie prin regulile standard de șah (șah-mat, abandon, sau remiză), fie prin aducerea regelui său în centrul tablei: jocul este câștigat imediat dacă regele unui jucător ajunge în unul dintre cele patru pătrate centrale (d4, d5, e4, e5) și nu se află sub atac (nu este în șah).

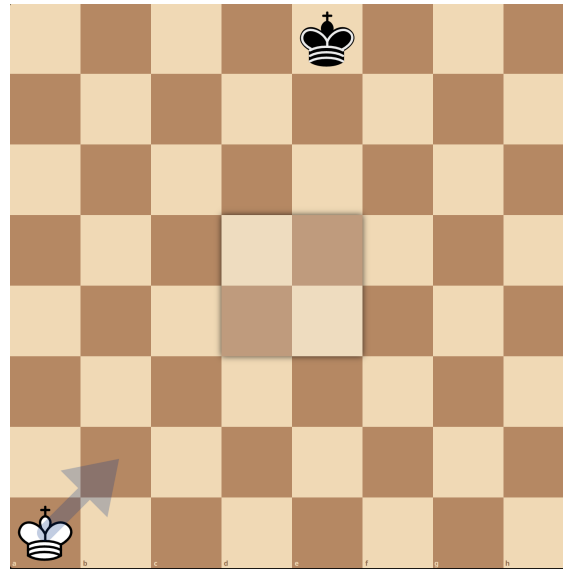


Figure 2: Imagine creată cu ajutorul lichess

În imaginea de mai sus, pentru că muta primul, jucătorul cu piesele albe câștigă, pentru că regele său ajunge primul pe unul dintre cele 4 pătrate din centru.

Secvența de mutări câștigătoare: 1.Kb2 Ke7 2.Kc3 Ke6 3.Kd4#

Three-Check Chess

În 3-Check, jucătorii pot câștiga partida în două moduri:

- Prin regulile tradiționale de șah (șah-mat, abandon, remiză).
- Prin darea de trei șahuri regelui adversarului. Primul jucător care reușește să dea trei șahuri câștigă automat partida.

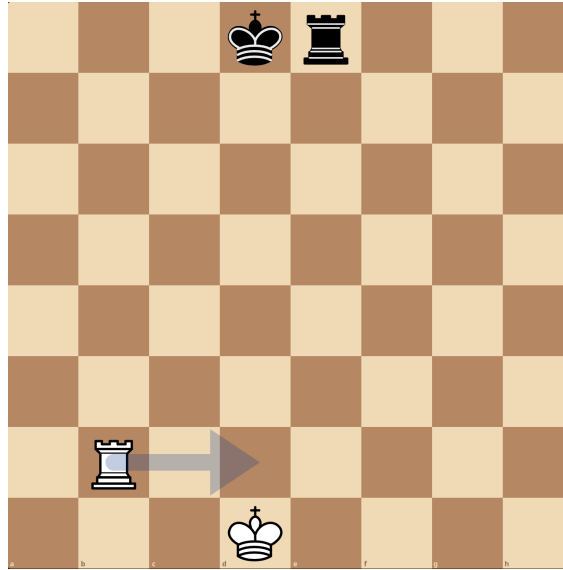


Figure 3: Imagine creata cu ajutorul lichess

În imaginea de mai sus, se observă o diferență dintre șahul standard și varianta 3-check. Jucând cu regulile standard, această poziție este o remiză, deoarece niciun jucător nu poate progresa. Însă, în această variantă, jucătorul alb câștigă, deoarece jucătorul cu piesele negre nu își poate apăra regele de șah pentru că tura sa este pe același rând cu regele.

Secvența de mutări câștigătoare: 1.Rd2+ Kc7 2.Rc2+ Kb8 3.Rb2# (aici există mai multe posibilități).

1.2 Motivatie

Domeniul inteligenței artificiale pentru jocul de șah a înregistrat progrese semnificative de-a lungul anilor, ducând la dezvoltarea unor engine-uri de șah puternice și a unor algoritmi avansați. Folosind metode actuale de a estima abilitatea unui jucător de șah - sistemul ELO - Stockfish sau orice alt engine modern are o probabilitate mai mare de 99% de a bate campionul mondial de șah. În realitate, inteligența artificială domină meciurile împotriva jucătorilor umani din anul 1997 (3).

Utilizarea reinforcement learning în șah este un subiect interesant deoarece nu a fost considerată o metodă tradițională în acest domeniu. Pentru o perioadă lungă de timp, un engine de șah avea următoarea structură: algoritmul de bază era minimax cu alpha-beta pruning, unde, fiindcă spațiul de căutare era prea mare, funcția de evaluare era scrisă manual și se baza pe diferite euristici. La această structură se adăuga o bază de date cu diferite deschideri, care erau bazate pe meciurile celor mai buni jucători de șah din lume. Această abordare, îmbunătățită de-a lungul timpului, a adus rezultate semnificative pentru o perioadă destul de lungă de timp, dovadă că cel mai performant engine de șah, Stockfish, a utilizat rețele neuronale pentru prima dată abia în anul 2020.

Reinforcement learning a apărut pentru prima dată în inteligența artificială pentru șah în anul 2018. Dintre cele cinci cele mai populare engine-uri de șah, doar AlphaZero și Leela Chess Zero utilizează reinforcement learning, în timp ce Stockfish, Komodo și Houdini se bazează exclusiv pe metode clasice, cum ar fi MCTS (Monte Carlo Tree Search), alpha-beta pruning și funcții de evaluare bazate pe euristici. Merita menționat că, dintre toate engine-urile numite, AlphaZero și LeelaChessZero sunt semnificativ mai noi (ambele în 2018), ceea ce sugerează că subdomeniul reinforcement learning în inteligența artificială pentru șah are potențial de viitor.

Majoritatea abordărilor curente - la fel ca toți algoritmi de inteligența artificială - se bazează pe regulile standard de șah. Câteva abordări existente sunt descrise în capitolul 2, însă aceste abordări nu sunt complete, în cazul CrazyAre, unde unele variante lipsesc cu totul, sau sunt bazate pe cod proprietar, care nu este disponibil public, cum este AlphaZero. Folosirea unor diferite variante de șah oferă mai multe avantaje, cum ar fi:

- Imbunatatirea strategiilor existente: regulile diferite si situatiile care apar datorita acestora pot crea o multime mai diversa de strategii pentru un engine.
- Imbunatatirea unui engine de sah: antrenarea unui engine cu reguli diferite poate fi vazuta ca o metoda de a imbunatati inteligenta artificiala: spre exemplu, unele situatii pot aparea mult mai des intr-o forma de sah decat in alta; aceste situatii diferite pot fi vazute ca un set de date de antrenare mai mare.
- Imbunatatirea abilitatilor jucatorilor umani: de-a lungul timpului, rezultatele obtinute de inteligenta artificiala in sah au reprezentat atat o sursa de inspiratie cat si o modalitate de a-si testa teoriile pentru jucatorii umani. Unele reguli pot pune accentul pe abilitatile unui jucator in partea de inceput a meciului (dechiderea), iar altele se axeaza pe sfarsitul jocului (endgame), deci cunostintele dobandite cu ajutorul inteligentei artificiale va crea jucatori de sah mai buni.

1.3 Abordari anterioare

CrazyAra

Crazy Ara este un engine de şah bazat pe reţele neuronale, conceput special pentru a juca varianta de şah numită Crazyhouse, varianta de sah care se bazeaza pe urmatoarea regula: atunci când un jucător capturează o piesă adversă, acea piesă este transformată într-o piesă de aceeaşi culoare ca jucătorul care a capturat-o. Motorul utilizează algoritmi de reinforcement learning pentru a evalua poziţiile pe tablă şi a prezice cele mai bune mutări. Combină o reţea neuronală antrenată pe un set de date de 500.000 de jocuri jucate de oameni cu un algoritm de căutare Monte Carlo Tree Search (MCTS) pentru a selecta mutările în timpul jocului. Această abordare a fost inspirată de succesul lui AlphaZero de la Google, care utilizează, de asemenea, reinforcement learning şi MCTS pentru a obtine rezultate importante. Dezvoltarea Crazy Ara include versiuni scrise atât în Python, cât şi în C++, versiunea mai nouă în C++ oferind o performanţă îmbunătăţită. Este compatibil cu diverse interfeţe GUI precum CuteChess şi WinBoard pentru o utilizare convenabilă. Stilul de joc al lui Crazy Ara este adesea descris ca fiind mai intuitiv şi dinamic în comparaţie cu motoarele tradiţionale de şah, care se bazează foarte mult pe algoritmi de căutare extensivi. Crazy Ara tinde să prioritizeze poziţiile mai bune pe tablă în detrimentul valorii imediate a pieselor individuale, ceea ce duce la un joc mai agresiv şi orientat spre iniţiativă.

În acest articol (4) din 2020, autorii compara performanta CrazyAra cu alte engine-uri de sah ce pot juca varianta Crazyhouse. În toate meciurile jucate, CrazyAra a obţinut rezultate pozitive împotriva tuturor engine-urilor în afara de Stockfish.

Un dezavantaj semnificativ adus de CrazyAra este faptul ca se limiteaza doar la varianta Crazyhouse, deci nu poate juca multe alte reguli populare cum ar fi : Chess960 (Fisher Random Chess, o varianta unde jucatorii incep cu o pozitie aleatoare a pieselor, King of the Hill, Atomic Chess etc.

AlphaZero

AlphaZero este o inteligenţă artificială creată de DeepMind, capabilă să înveţe şi să joace la nivel înalt jocuri complexe precum şah, shogi şi Go. AlphaZero a reprezentat prima dezvoltare serioasă a unor algoritmi de reinforcement learning pentru inteligenta artificiala în jocuri. Foloseşte o combinaţie de reţele neuronale, MCTS şi reinforcement learning:

- reţele neuronale profunde pentru a evalua poziţiile de pe tablă şi a selecta mişcărilor. Reţeaua neuronală ia starea tabloului ca intrare şi produce două informaţii cheie: politică- o distribuţie a probabilităţilor asupra tuturor mişcărilor posibile, indicând probabilitatea ca fiecare mişcare să fie cea mai bună si valoare- o valoare scalară care estimează probabilitatea de a câştiga din poziţia curentă.

- reinforcement learning: AlphaZero joacă un număr mare de jocuri împotriva sa, generând date despre pozițiile jocului, mișcări și rezultate. Datele generate din jocurile cu sine însuși sunt stocate într-un buffer, rețeaua neuronală fiind antrenată pe aceste date, învățând să prezică politica și valoarea.
- MCTS: Combină predicțiile rețelei neuronale cu explorarea bazată pe căutare. Începând de la nodul rădăcină (starea curentă a tabloului), algoritmul selectează noduri copil bazate pe un echilibru între exploatare (folosind politica rețelei neuronale) și explorare (încercând noduri mai puțin vizitate). Când se ajunge la un nod frunză (un nod nevizitat anterior), acesta este extins prin adăugarea unui nou nod corespunzător unei mișcări posibile. Rețeaua neuronală evaluează noul nod, furnizând o estimare a politicii și a valorii. Estimarea valorii este propagată înapoi în arbore, actualizând statisticile pentru fiecare nod de-a lungul traseului.

În aceasta lucrare (13) din 2020, AlphaZero este antrenat pe 9 variante de sah diferite pentru a trage concluzii despre efectul regulilor asupra :

- tendinței unui meci de sah de a ajunge la remiza
- avantajul pe care îl conferă prima mutare
- diversitatea deschiderilor
- cât de decisive sunt meciurile (în acest caz, un rezultat decisiv înseamnă un meci fără remiza)
- dinamica jocului (numărul de mutări și cât de "încet" este un meci)
- utilizarea mutărilor speciale introduse de acea variantă.

Aceste rezultate, cât și alte rezultate obținute când AlphaZero a fost antrenat pe regulile standard de sah, demonstrează că AlphaZero este cea mai promitoare abordare de reinforcement learning pentru sah. Când vine vorba de performanță, AlphaZero a fost testat împotriva altor engine-uri de sah în mod direct doar cu reguli standard, unde a obținut rezultate bune, având un scor ELO asemănător cu Stockfish.

Singurul dezavantaj pe care îl aduce AlphaZero este faptul că nu este open-source. Spre deosebire de Stockfish, rezultatele unui engine nu pot fi comparate direct cu AlphaZero. Bazat pe o putere de calcul mare, el este antrenat doar intern de către DeepMind, compania care l-a creat. Spre exemplu, algoritmul dezvoltat în această lucrare de licență, nu poate fi comparat cu AlphaZero. Totuși, merita menționat că DeepMind a publicat diferite articole în care detaliază modul în care funcționează AlphaZero și algoritmii care stau în spatele inteligenței artificiale, astfel că pot exista mai multe implementări care seamănă cu AlphaZero.

Fairy Stockfish

Fairy Stockfish este o extensie a engine-ului de șah Stockfish, proiectată pentru a suporta o gamă largă de variante de șah. Aceasta se bazează pe arhitectura robustă și optimizată a Stockfish, dar include funcții suplimentare pentru a gestiona regulile unice ale variantelor diferite de sah.

Fairy Stockfish moștenește arhitectura de bază a Stockfish, care include:

- Alpha-Beta Pruning: Un algoritm de căutare optimizat care evaluează mișcările posibile prin eliminarea ramurilor care nu influențează decizia finală.
- Reprezentarea Bitboard: O reprezentare eficientă a tablei de șah și a pieselor, permițând calculuri și manipulări rapide.
- Funcția de Evaluare: O funcție sofisticată care evaluează valoarea unei poziții pe tablă, ținând cont de echilibrul material, activitatea pieselor, siguranța regelui și alți factori.

Fairy Stockfish introduce mai multe modificări și extensii: include un cadru flexibil pentru definirea regulilor diferitelor variante de șah. Acesta îi permite să ajusteze modelele de mișcare ale pieselor, condițiile de câștig și alte reguli specifice jocului. Definițiile Pieselor: Fairy Stockfish poate defini

noi tipuri de piese și mișcările lor. De exemplu, piesele din variante precum Crazyhouse sau Atomic Chess au reguli de mișcare și interacțiune unice. De asemenea, funcția de evaluare din Fairy Stockfish este ajustată pentru a ține cont de aspectele unice ale diferitelor variante de șah: funcția de evaluare include euristici adaptate caracteristicilor specifice fiecărei variante. De exemplu, în Crazyhouse, se ia în considerare valoarea pieselor aflate în mână pentru reintroducere. Se fac ajustări pentru a evalua corect mobilitatea și influența pieselor specifice variantelor. De asemenea, Fairy Stockfish evaluează controlul asupra zonelor critice ale tablei, care poate varia semnificativ între variante.

Fiind derivat din Stockfish, Fairy Stockfish este cel mai performant engine pentru variante non-standard de șah. Deoarece este performant și suporta un număr mare de seturi de reguli, este de obicei folosit pentru a testa alte engine-uri și algoritmi (va fi folosit pentru testare și în lucrarea de față). În schimb, Fairy Stockfish nu folosește reinforcement learning. Acest lucru nu este neapărat un dezavantaj, însă nu ne poate oferi informații în plus despre cum poate fi folosit reinforcement learning pentru jocul de șah.

Leela Chess Zero (Lc0)

Leela Zero Chess este un proiect open-source care vizează regulile standard de șah; însă arhitectura sa flexibilă îi permite să fie folosit și pentru alte seturi de reguli. Lc0 folosește reinforcement learning, anume o combinație de rețele neuronale și MCTS, metoda fiind inspirată de cea folosită de AlphaZero. Deși nu a fost folosit în mod oficial în acest sens, Lc0 poate fi reantrenat pe seturi de reguli diferite pentru a putea juca Antichess, Crazyhouse sau alte variante.

2 Descrierea solutiei

Algoritmul dezvoltat se va baza pe Deep Q-Learning, deci vom folosi retele neuronale adanci pentru a aproxima functia Q-valoare.

2.1 Tehnologii folosite

1. TensorFlow

Pentru a crea retea neuronală adancă care sta la baza algoritmului de Deep Q-learning, vom folosi TensorFlow. Pentru a crea retea neuronală, putem alege mai multe moduri de a reprezenta tabla de șah, acesta de mai jos este doar unul dintre ele.

```
1 def create_chess_model():
2     model = Sequential()
3     model.add(Conv2D(64, (3, 3), activation='relu', input_shape=(8, 8, 12)))
4     model.add(Conv2D(128, (3, 3), activation='relu'))
5     model.add(Conv2D(256, (3, 3), activation='relu'))
6     model.add(Flatten())
7     model.add(Dense(512, activation='relu'))
8     model.add(Dense(1, activation='sigmoid'))
9
10    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['
11    accuracy'])
12    return model
```

Listing 1: Exemplu utilizare TensorFlow

În acest exemplu, am folosit un tensor de forma (8,8,12), unde tabla este reprezentată ca o matrice pătratică de dimensiune 8 cu 12 nivele, unde fiecare nivel reprezintă prezența unei piese.

2. Gym, Stable Baseline 3

Gym este o bibliotecă ce aparține ecosistemului OpenAI și este folosită pentru dezvoltarea algoritmilor de reinforcement learning. Este folosită atât pentru că este simplă de utilizat cât și pentru că are un număr mare de "medii" (environment) ce pot fi folosite. În acest caz, putem folosi environmentul gym-chess.

Stable Baseline 3 este o bibliotecă care implementează algoritmi de reinforcement learning în Python. Este un fork al bibliotecii Stable Baseline și este scris utilizând bibliotecă PyTorch. Conține mai mulți algoritmi de bază de reinforcement learning, cum ar fi: Deep Q-Learning (DQN), Advantage Actor-Critic (A2C), Proximal Policy Optimization (PPO), Soft Actor-Critic (SAC) etc. Pentru simplitate, în exemplul de mai jos am folosit algoritmul DQN din bibliotecă SB3. Acest lucru nu este obligatoriu; acest algoritm poate fi și implementat manual.

```
1
2 env = gym.make('Chess-v0')
3
4
5 env = make_vec_env(lambda: env, n_envs=1)
6
7
8 class CustomCNNPolicy(DQN):
9     def __init__(self, *args, **kwargs):
10         super(CustomCNNPolicy, self).__init__(*args, **kwargs)
11         self.q_net = create_chess_model()
12
13     def forward(self, obs):
14         return self.q_net(obs)
15
16
17 model = DQN(CustomCNNPolicy, env, verbose=1)
18
19
```

```

20 checkpoint_callback = CheckpointCallback(save_freq=10000, save_path='./
    checkpoints/', name_prefix='dqn_chess')
21 model.learn(total_timesteps=1000000, callback=checkpoint_callback)
22
23
24 model.save("dqn_chess_model")
25
26 #Testare
27 obs = env.reset()
28 for i in range(1000):
29     action, _states = model.predict(obs, deterministic=True)
30     obs, rewards, dones, info = env.step(action)
31     env.render()
32     if dones:
33         obs = env.reset()

```

Listing 2: Exemplu utilizare Gym

3. Fairy Stockfish

Dupa cum am descris in mod anterior, acesta este cel mai performant engine pentru variante non standard de sah. Il vom folosi pentru a testa performanta algoritmului dezvoltat. Presupunand ca avem algoritmul dezvoltat, vom testa performanta acestuia prin jocuri impotriva lui Fairy Stockfish, cu un script care arata asa:

```

1 def test_algorithm_vs_stockfish(stockfish_path, algorithm, num_games=100):
2     env = AntichessEnv(stockfish_path)
3     results = {"wins": 0, "losses": 0, "draws": 0}
4
5     for game in range(num_games):
6         state = env.reset()
7         while not env.board.is_game_over():
8             algorithm_move = algorithm.get_move(state)
9             if algorithm_move not in [move.uci() for move in env.board.
legal_moves]:
10                 raise ValueError(f"Invalid move {algorithm_move} by the
algorithm.")
11
12             state, done = env.step(algorithm_move)
13             if done:
14                 break
15
16             state, done = env.step(env._parse_best_move(env._get_response()))
17             if done:
18                 break
19
20             if env.board.is_checkmate():
21                 if env.board.turn:
22                     results["wins"] += 1
23                 else:
24                     results["losses"] += 1
25             else:
26                 results["draws"] += 1
27
28     env.close()
29     return results
30

```

Listing 3: Exemplu Script Testare

4. Python Chess, pygame

Pentru a dezvolta interfata grafica, am folosit librariile Python Chess si pygame. Python chess ofera suport pentru toate cele 3 variante de sah pe care le vom testa, cat si suport pentru diferite functionalitati necesare pentru interfata grafica. Spre exemplu, verificarea corectitudinii mutarilor, verificarea daca jocul s-a incheiat sau nu, lista cu toate mutarile facute etc. nu au mai trebuit implementate manual.

Pygame este una dintre cele mai folosite librarii pentru implementarea unei interfete grafice si este potrivit si in acest caz, pentru ca ofera metode simple de a implementa functionalitatile necesare.

2.2 Alte resurse

În această secțiune, vom menționa implementări asemănătoare sau care pot oferi informații despre cum ar trebui structurat algoritmul, sau cum putem să îmbunătățim performanța algoritmului.

1. Implementare bazată pe AlphaZero

Aceasta este o implementare de algoritmi de reinforcement learning care se bazează pe arhitectura AlphaZero. După cum descrie autorul, algoritmul face o căutare MCTS modificată, unde selecția nu este aleatorie ci bazată pe rețeaua neuronală, iar simularea este înlocuită de valoarea primită de rețeaua neuronală.

2. Implementare bazată pe Q-Learning

Aceasta este o implementare bazată pe Deep Q-Learning. Este un algoritm fără model, deci nu are o modalitate de a vedea mișcările înainte. Din acest motiv, nu este o implementare la fel de performantă.

3 Bibliografie

1. Sistemul de rating ELO
2. Fairy Stockfish
3. Meciul dintre Garry Kasparov si Deep Blue
4. Chess Variants with Neural Networks
5. Variante de sah cu reinforcement learning
6. Gym
7. Stable Baseline 3
8. Python Chess
9. CrazyAra
10. CrazyAra, github
11. AlpahZero description
12. AlphaZero Variants Article
13. AlphaZero Variants Source
14. AlphaZero Based implementation
15. Antichess rules
16. Three-check rules
17. King of the hill rules
18. Crazyhouse rules
19. Editor pentru crearea pozitiilor
20. Houdini Chess Engine
21. Stockfish Chess Engine
22. Komodo Chess Engine
23. Leela Chess Zero Chess Engine
24. Monte Carlo Tree Search
25. Articol despre AlphaZero si Reinforcement Learning
26. AlphaZero si MCTS
27. Tutorial implementare AlphaZero
28. AlphaZero explicat