# MCTS - Connect four

Markov Decision Processes and Reinforcement Learning

Samuele Angheben

Trento University

February 9, 2025
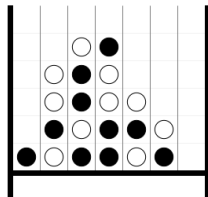
# Overview

The objective of this project is to explore different reinforcement learning algorithm on the game **Connect Four**, starting from very simple mcts implementation to alphazero.

**Algorithms:**

- **Naive MCTS:** A basic Monte Carlo Tree Search approach without learning.
- **REINFORCE with Baseline:** A policy gradient method incorporating a baseline.
- **AlphaZero:** A self-play RL algorithm combining MCTS with deep learning.

**Connect Four** is a two-player game played on a vertical 7-column, 6-row grid where players take turns dropping colored discs into the slots, aiming to connect four discs of their color horizontally, vertically, or diagonally.

- **Deterministic:** Outcome fully determined by player actions. $S_{t+1}, R_t = p(S_t, a_t)$
- **Episodic:** Sequence of actions leading to a terminal state. Episode: $\{(S_1, a_1), \ldots, (S_T, a_T)\}$
- **Two-Player:** Players 1 and 2 alternate moves, strategies $\pi_1$ and $\pi_2$.
- **Zero-Sum:** Rewards sum to zero. $R_1 = -R_2$. Terminal rewards: Win $(+1/-1)$, Draw $(0/0)$
- **Perfect Information:** Both players have complete knowledge of the current state.

- **Trajectory Length:** Max 42 moves (7 columns x 6 rows).
- **State Space:**
    - Naive: $3^{42} \approx 1.1 \times 10^{20}$
    - Improved: $\approx 4.5 \times 10^{12}$ states
    - Storage: $4.5 \times 10^{12}$ states $\times 8$ bytes/state $= 36$ TB (64 bit state representation)
- **Tabular Methods:** Infeasible due to large state space.
- **Solvability:** Solved. Player 1 can always win with optimal play (Wikipedia). No readily available solvers for testing.
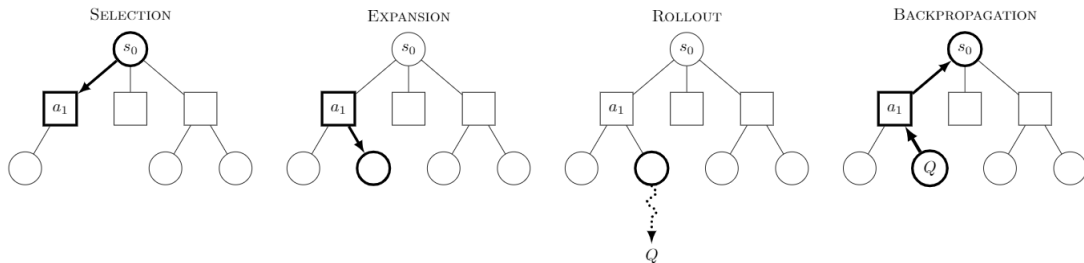
**Framework:** JAX provides a high-performance framework for computational efficiency and flexibility, allowing for scalable and efficient development of machine learning algorithm.

**JAX Features:**

- JIT Compilation
- GPU/TPU Acceleration
- Efficient Vectorization (vmap)
- Automatic Parallelization (pmap)
- Functional Programming (Pure Functions)

- `mctx`: JAX-native MCTS implementations.
- `pgx`: Game environment simulation, diverse games, fast parallelization.
- `Haiku`: Neural network building and training in JAX.
- `optax`: Gradient processing and optimization.

Monte Carlo Tree Search is a decision-making algorithm that explores a decision space by constructing a search tree through repeated simulations. Given a computation budget it balances exploration and exploitation to refine the root node policy.

To perform a MCTS search using the MCTX library, two main functions must be provided:

**RootFnOutput(state):**
- Specifies the representation of the root state.
- Returns the prior logits and the estimated value of the root state.

**recurrent_fn(state, action):**
- Encapsulates the environment dynamics.
- Returns the reward, discount factor, and for the new state, the prior logits and value.

The search returns **action_weight**, representing the updated root node policy.

In alternating-turn games, a negative discount factor ($\gamma = -1$) is used to invert value estimates between players:

$$V(S_t) = r_t + \gamma \cdot V(S_{t+1}) = r_t + (-1) \cdot V(S_{t+1})$$

- $S_t$: Current state (Player 1's turn).
- $S_{t+1}$: Next state (Player 2's turn), where the value $V(S_{t+1})$ represents the opponent's optimal value.

This inversion maintains the zero-sum property by ensuring that a high value for one player corresponds to a low value for the opponent.

The neural network is a custom **ResNet**, a convolutional network with residual connections and batch normalization. It has two output heads:

$$l, v' = \text{ResNet}_\theta(S_t)$$

**1. Policy Head:** Produces a probability distribution over possible actions:

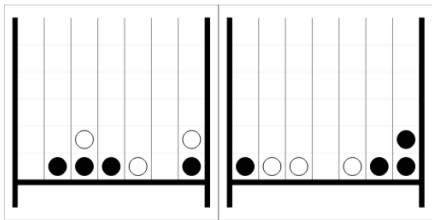$$\pi(a \mid S_t) \approx \hat{\pi}(a \mid S_t) = \text{softmax}(l)$$

**2. Value Head:** Outputs a scalar value approximating the value function:

$$V(S_t) \approx \hat{V}(S_t) = v'$$

Naive implementation of Monte Carlo Tree Search:

- **no learning**
- MCTS
  - policy: uniform
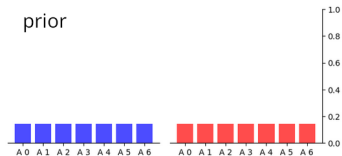  - value: single random rollout

```
61
62 def RecurrentFn(s, a)
63     s', reward <- env_step(s, a)
64     V, pi <- random_rollout(s'), uniform_prior
65     discount <- -1.0
66     return s', pi, V, reward, discount
67
68 def RunMCTS(s)
69     root = (V, pi(a|s)) <- random_rollout(s), uniform_prior
70     policy_output <- MCTS(root, RecurrentFn, num simulation)
71     return arg max_a policy_output
72
```
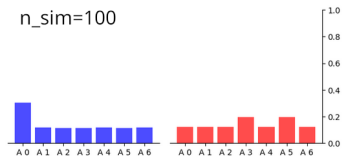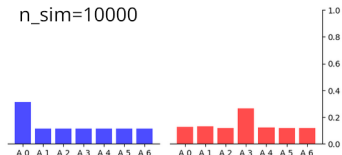
Value = [-1, -1]



Policy output

prior

n_sim=100

n_sim=10000

Value of the state: [-1. -1.]

Reinforce baseline algorithm:

- Policy based algorithm
- no MCTS
- on policy, batch
- neural network function approximation

```
1   Initialize: Neural Network, Optimizer, Replay Buffer        32
2                                                               33
3   Main Loop:                                                  34
4       for iteration in 1..MaxIterations:                      35
5           # Self-play Phase                                   36          # Training Phase
6           Collect Games:                                      37          Train Network:
7               for N games:                                    38              batch ~ Replay Buffer # sample (s_t, a_t, G_t)
8                   def step_fn(s):                             39              compute_loss:
9                       root = V_nn, pi_nn(a|s) <- nn_forward(s)40                  V_nn, pi_nn(a|s_t) <- nn_forward(s_t)
10                      run MCTS(root, num_simulation):         41                  A(s_t, a_t) <- G_t - V_nn(s_t) # compute advantage
11                          s', reward <- env                   42                  policy_loss <- -mean[log(pi_nn(a_t|s_t) * A(s_t, a_t)]
12                          V', pi <- nn_forward                43                  value_loss  <- MSE(V_nn(s_t), G_t)
13                                                              44                  entropy_loss <- entropy(pi_nn)
14                      a <- MCTS.policy_out.action             45
15                                                              46                  loss <- policy_loss + value_loss + 0.1 entropy_loss
16                      s', r <- env_step(s, a)                 47
17                                                              48              Compute gradient (loss)
18                      return (s, policy_out, s', r, discount=-1),49              Update model parameters
19                              s'                              50
20                                                              51      # Saving Phase (periodic)
21                  s_init <- env.reset()                       52      if iteration % saving_interval == 0:
22                  data <- jax.scan(step_fn, s_init)           53          Save checkpoint
23                                                              54
24                  samples <- compute_loss_input(data) # cumulative rewards55
25                  Save samples to Replay Buffer # (s_t, pi_t, G_t:value target)56
26                                                              57
27                                                              58
28          Shuffle samples and make minibatches (Replay Buffer)59
29                                                              60
```

**Total Loss:**

$$L = \mathcal{L}_{\text{policy}} + \mathcal{L}_{\text{value}} + 0.1 \times \mathcal{L}_{\text{entropy}}$$
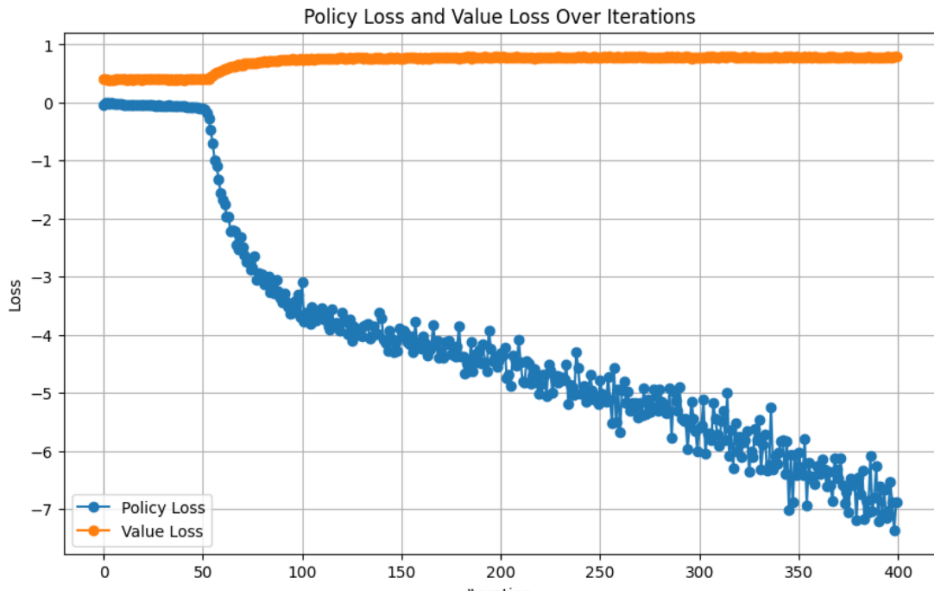
**Where:**

$$\mathcal{L}_{\text{policy}} = -\mathbb{E}\left[\log \hat{\pi}(a_t \mid S_t) \cdot A(S_t, a_t)\right]$$

$$A(S_t, a_t) = G_t - \hat{V}(S_t)$$

$$\mathcal{L}_{\text{value}} = \mathbb{E}\left[\left(\hat{V}(S_t) - G_t\right)^2\right]$$

$$\mathcal{L}_{\text{entropy}} = -\mathbb{E}\left[\sum_a \pi(a \mid S_t) \log \pi(a \mid S_t)\right]$$

**Gradient Clipping**: ensures that the gradients do not exceed a specified threshold, preventing exploding gradients.

Policy Loss and Value Loss Over Iterations

Even with loss improvements, divergence still occurs.

**The Deadly Triad:**

    Bootstrapping                   ×

    Function Approximation   ✓

    Off-Policy (No Batch)     ×

Alphazero algorithm

- MCTS
  - policy: neural network policy head
  - value: neural network value head
- neural network function approximation

```
1  Initialize: Neural Network, Optimizer, Replay Buffer
2
3  Main Loop:
4      for iteration in 1..MaxIterations:
5          # Self-play Phase
6          Collect Games:
7              for N games:
8                  def step_fn(s):
9                      root = V_nn, pi_nn(a|s) <- nn_forward(s)
10                     run MCTS(root, num_simulation):
11                         s', reward <- env
12                         V', pi <- nn_forward
13
14                     a <- MCTS.policy_out.action
15
16                     s', r <- env_step(s, a)
17
18                     return (s, policy_out, s', r, discount=-1),
19                         s'
20
21                 s_init <- env.reset()
22                 data <- jax.scan(step_fn, s_init)
23
24                 samples <- compute_loss_input(data) # cumulative rewards
25                 Save samples to Replay Buffer # (s_t, pi_t, G_t:value target)
26
27
28         Shuffle samples and make minibatches (Replay Buffer)
29
```

```
28
29
30
31
32     # Training Phase
33     Train Network:
34         batch ~ Replay Buffer # sample (s_t, policy_out_t, G_t)
35         compute_loss:
36             V_nn, pi_nn(a|s_t) <- nn_forward(s_t)
37             policy_loss <- cross_entropy(pi_nn(a|s_t), policy_out)
38             value_loss  <- MSE(V_nn(s_t, G_t))
39
40             loss <- policy_loss + value_loss + entropy_loss
41
42         Compute gradient (loss)
43         Update model parameters
44
45     # Saving Phase (periodic)
46     if iteration % saving_interval == 0:
47         Save checkpoint
48
49
50
51
52
53
54
55
56
```
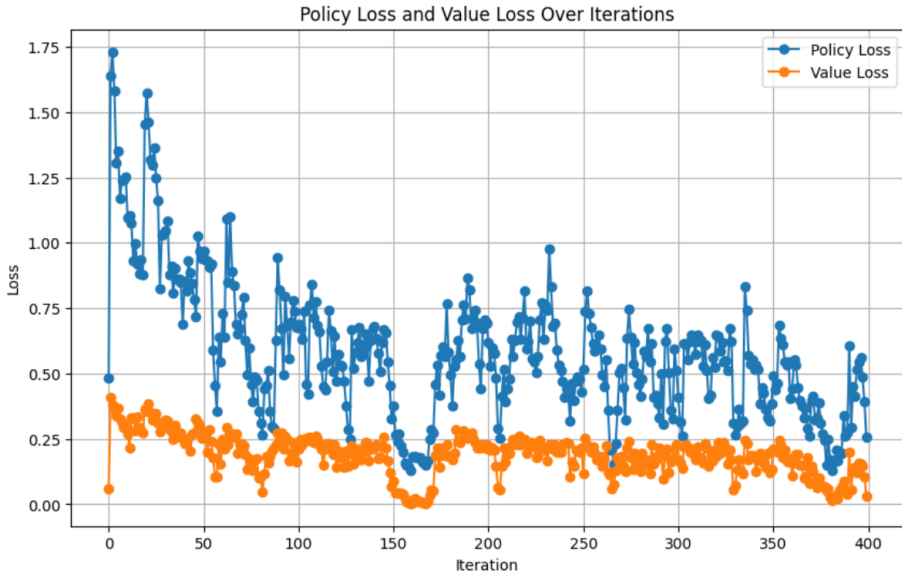
**Total Loss:**

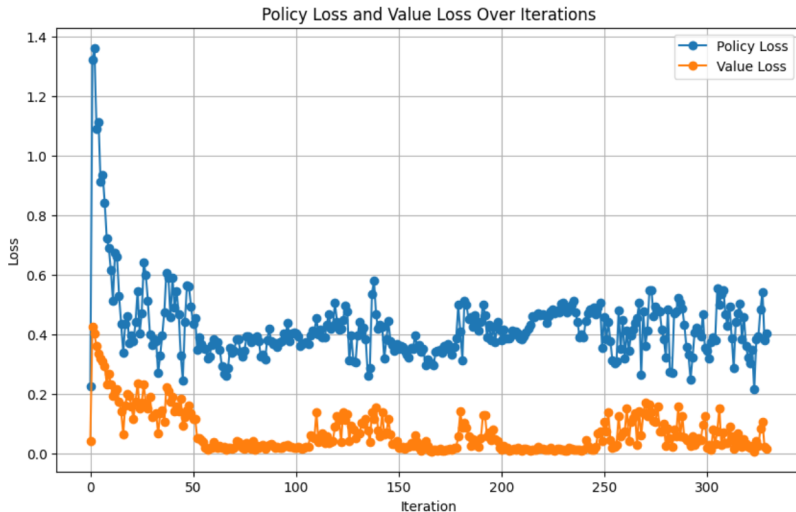$$L = \mathcal{L}_{\text{policy}} + \mathcal{L}_{\text{value}}$$

**Where:**

$$\mathcal{L}_{\text{policy}} = -\mathbb{E}_{a \sim policy\_out} \left[ \log \hat{\pi}(a) \right]$$

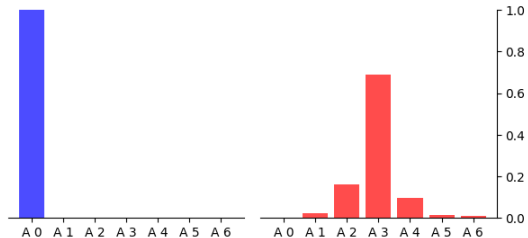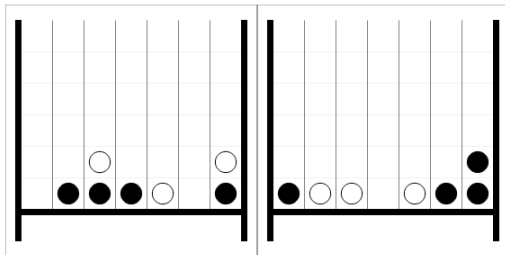$$\mathcal{L}_{\text{value}} = \mathbb{E} \left[ \left( \hat{V}(S_t) - G_t \right)^2 \right]$$

Policy Loss and Value Loss Over Iterations

With better training: increase num. MCTS simulation, batch size, neural network



Policy Loss and Value Loss Over Iterations

```
61
62 def RecurrentFn(s, a)
63     s', reward <- env_step(s, a)
64     V, pi <- nn_forward(s')
65     discount <- -1.0
66     return s', pi, V, reward, discount
67
68 def RunMCTS(s)
69     root = (V, pi(a|s)) <- nn_forward(s)
70     policy_output <- MCTS(root, RecurrentFn, num simulation)
71     return arg max_a policy_output
72
```

```
74
75 def RunOneShot(s)
76     pi_nn(a|s) <- nn_forward(s)
77     return arg max_a pi_nn(a|s)
```

Array([-0.93581444,  0.85516936], dtype=float32)

# **Notebook games and demo play**

- hyper-parameter tuning for better training
- longer training
- network architecture
  - increase capacity
  - transformer architecture
- MCTS: test different selection algorithm

# References

📄 Danihelka, I., Guez, A., Schrittwieser, J., and Silver, D. (2022).
Policy improvement by planning with gumbel.
https://www.openreview.net/forum?id=bERaNdoegnO.
Accessed: 2025-02-07.

📄 Koyamada, S., Okano, S., Nishimori, S., Murata, Y., Habara, K., Kita, H., and Ishii, S. (2024).
Pgx: Hardware-accelerated parallel game simulators for reinforcement learning.

📄 Silver, D., Huang, A., Maddison, C. J., and et al. (2016).
Mastering the game of go with deep neural networks and tree search.
*Nature*, 529(7587):484–489.

📄 Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., Lillicrap, T., Simonyan, K., and Hassabis, D. (2017a).
Mastering chess and shogi by self-play with a general reinforcement learning algorithm.

📄 Silver, D., Schrittwieser, J., Simonyan, K., and et al. (2017b).
Mastering the game of go without human knowledge.
*Nature*, 550:354–359.

# Thank you for your attention!