# Reinforcement Learning for playing Connect Four

Simon Hölck, Florian Cimander, Tim Löhr
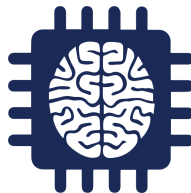
Mathematical Data Science

August 24, 2020

FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

NATURWISSENSCHAFTLICHE
FAKULTÄT

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Table of Content

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

- ▶ Reinforcement Learning
- ▶ Q Learning
- ▶ Deep Q Learning
- ▶ Training
- ▶ Best Practices
- ▶ Demo
- ▶ Discussion

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Reinforcement Learning

# Reinforcement Learning - Terminologies

**Markov Decision Process (MDP)**

► Agent

► Environment

► Action

► State

► Reward

# Overview of Reinforcement Learning Algorithms

Simon Hölck,
Florian Cimander,
Tim Löhr

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Q Learning

# Q Learning

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Q Learning

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Q Learning

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Q Learning

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

Q Learning is a value-based RL algorithm

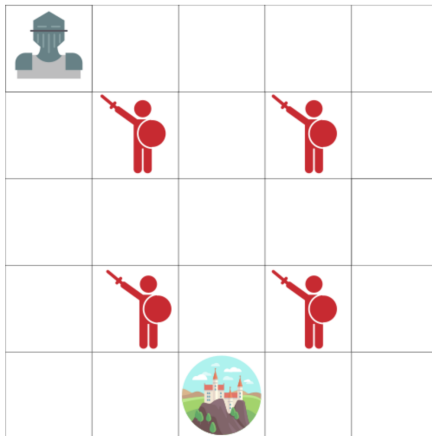$$Q^{\pi}(s_t, a_t) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + ... | s_t, a_t]$$

Q-Values for the state
given a particular state

Expected discounted
cumulative reward

Given the state and action

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Deep Q Learning

# Deep Q Learning

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

► How can we use the basics of Q Learning without having to store a q-table?
  $\rightarrow$ universal approximation theorem
  $\rightarrow$ use neural network as approximation to q function

*board inputs*



*actions*

# Implementation of a Deep Q Learning Agent

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
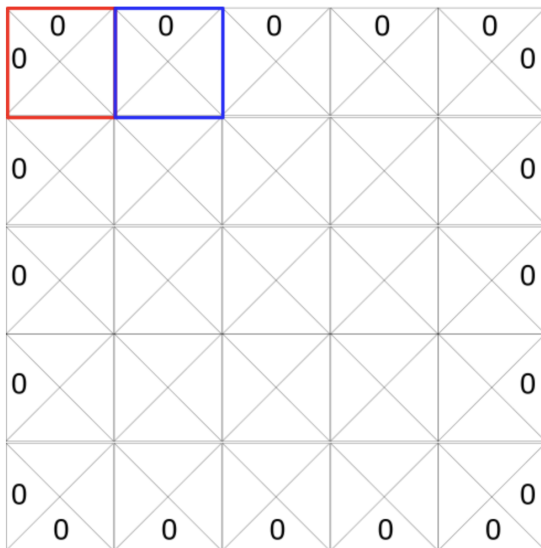Tim Löhr

Reinforcement
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

1. **Neural Network**
   - ▶ used eight fully connected layers to have enough depth for good approximation
   - ▶ added three dropout layers to prevent overfitting
   - ▶ used relu as activation function in forward pass

2. **Exploration vs. Exploitation**
   - ▶ Exploration: Agent makes out of character decisions that are not given by the network (random)
   - ▶ Exploitation: Agent takes actions given by the network
   - → We use *epsilon decay* to find the right balance between the two

# Implementation of a Deep Q Learning Agent

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

3. **Learning Process**
   - ▶ Two problems
     - a) correlated inputs/outputs
     - b) non-stationarity

# Implementation of a Deep Q Learning Agent

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

### 3. **Learning Process**

- ▶ for every action during a game the following quintuple is stored to the agents *experience replay memory*

  (state, action, reward, new state, done)

- ▶ rewards are determined when game is finished
- ▶ the agent learns the saved transitions in batches that are chosen randomly from its memory

$\rightarrow$ solves the problem of correlated inputs/outputs

# Implementation of a Deep Q Learning Agent

3. **Learning Process**
   ▶ for each transition in the batch we perform the following steps

| calculate forward pass of initial state | $+$ | calculate target value using the q formula |

$$q_{target} = reward + \gamma \cdot \max_{a'} Q(s', a')$$

$\downarrow$

calculate mean squared error

$\downarrow$

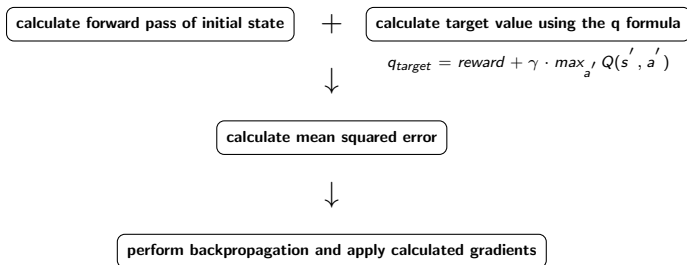perform backpropagation and apply calculated gradients

$\rightarrow$ problem of non-stationarity

# Implementation of a Deep Q Learning Agent

Reinforcement Learning for playing Connect Four

Simon Hölck, Florian Cimander, Tim Löhr

Reinforcement Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

3. **Learning Process**
    - ▶ we deal with this problem by using a *target network*
    - ▶ the target network is an old version of the underlying DQN network
    - ▶ it gets updated very infrequently

$$\boxed{\text{calculate forward pass of initial state}} \quad + \quad \boxed{\text{calculate target value using the q formula}}$$

$$q_{target} = reward + \gamma \cdot max_{a'} \, Q_{target}(s', a')$$

$$\downarrow$$

$$\boxed{\text{calculate mean squared error}}$$

$$\downarrow$$

$$\boxed{\text{perform backpropagation and apply calculated gradients}}$$

$\rightarrow$ targets stay constant and problem becomes more stationary

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Training

# Opponent: Minimax

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

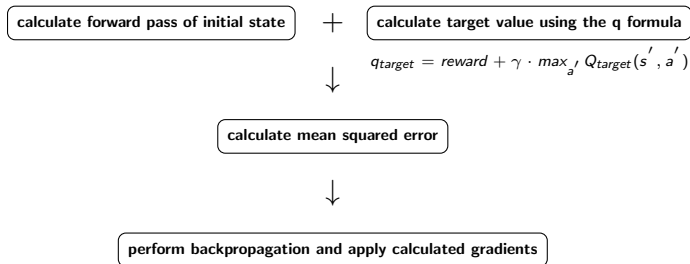Reinforcement
Learning

Q Learning

Deep Q Learning

Training
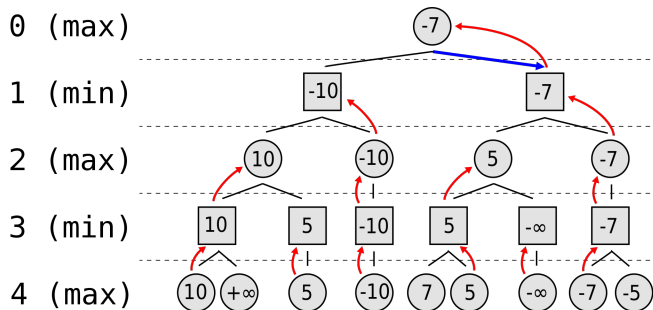
Best Practices

Demo

Discussion

Self-play RL like Alpha Zero is too time consuming.
Solution: Let the Neural Network practice against Minimax

# Minimax depth = 1 vs Neural Network

- **Yellow**: Minimax
- **Red**: Neural Network

*7000 Episodes and Neural Network win-ratio of 11%*

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Minimax depth = 1 vs Neural Network

- **Yellow**: Neural Network
- **Red**: Minimax

*50000 Episodes and Neural Network win-ratio of 99%*

Reinforcement
Learning for
playing
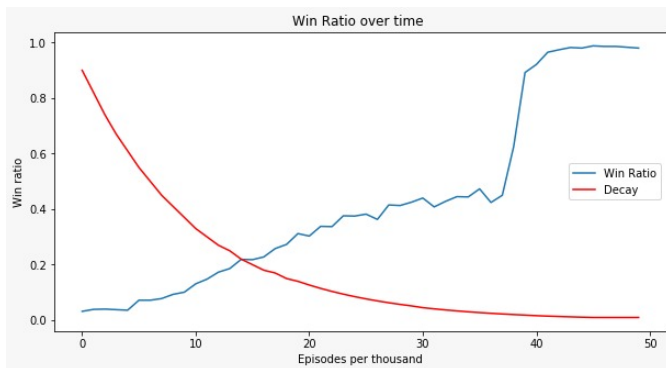Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Learning over time

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

We started by letting the Neural Network train **50000** Episodes against Minimax with depth $= 1$

# Save the trained model and start all-over again

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

Now this 99% ratio Neural Network plays and practices
against the Minimax with depth $= 2$

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Best Practices

# Best Practices

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
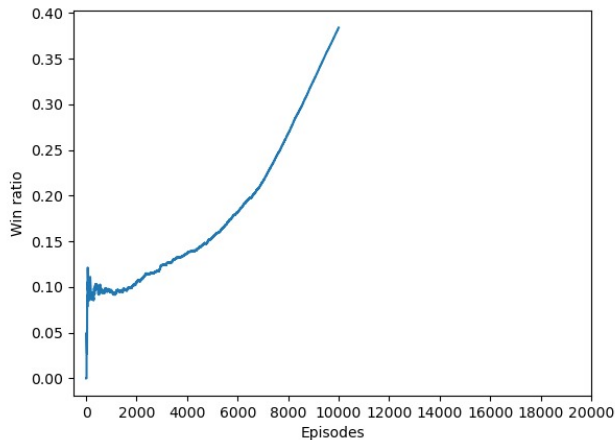Florian Cimander,
Tim Löhr

Reinforcement
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

*Trial and Error, Error ... and more Errors*

Best set of Hyperparameters so far:

- ▶ Randomness 100% with exponential fall down to fixed 1%
- ▶ Batch Size $= 128$
- ▶ Memory size $= 50000$
- ▶ Learning Rate for the Neural Network $= 0.01$
- ▶ Episodes 50000 on depth $= 1$: *15 hours*
- ▶ Episodes 10000 on depth $= 2$: *10 hours*

Simon Hölck,
Florian Cimander,
Tim Löhr

# Demo

# Demo

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

It is time to show a Demo in the GUI

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

# Discussion

# Discussion

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

▶ Does a 100% win ratio against Minimax of depth 5 or 6
plays good against human opponents?

▶ Can the Neural Network generalize better if we would
have used the Monte Carlo Tree Search (MCTS) as
opponent instead of the Minimax?

▶ Could a ResNet train a better AI against human players
with everything else kept the same than our very basic
Neural Network architecture?

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

Thanks for listening to us!

# Sources

Reinforcement
Learning for
playing
Connect Four

Simon Hölck,
Florian Cimander,
Tim Löhr

▶ Figure 1: https://www.pngwave.com/png-clip-art-jdcjo
▶ Figure 2: Reinforcement Learning: An Introduction (Sutton, Barto)
▶ Figure 3: https://www.afcea.org/content/artificial-intelligence-will-change-human-values
▶ Figures Knight & Princess Game: https://www.freecodecamp.org/news/diving-deeper-into-reinforcement-learning-with-q-learning-c18d0db58efe/