

# Reinforcement Learning for playing Connect Four

Simon Hölck, Florian Cimander, Tim Löh

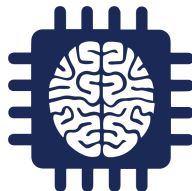
Mathematical Data Science

August 23, 2020



# Table of Content

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



Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhner

Reinforcement  
Learning

Q Learning

Deep Q Learning  
Training

Best Practices

Demo

Discussion

# Reinforcement Learning

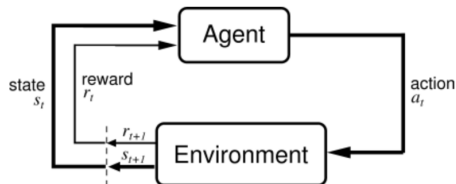
# Reinforcement Learning - Terminologies

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhner

## Markov Decision Process (MDP)

- ▶ Agent
- ▶ Environment
- ▶ Action
- ▶ State
- ▶ Reward



Reinforcement  
Learning

Q Learning

Deep Q Learning  
Training

Best Practices

Demo

Discussion

# Overview of Reinforcement Learning Algorithms

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löh

Reinforcement  
Learning

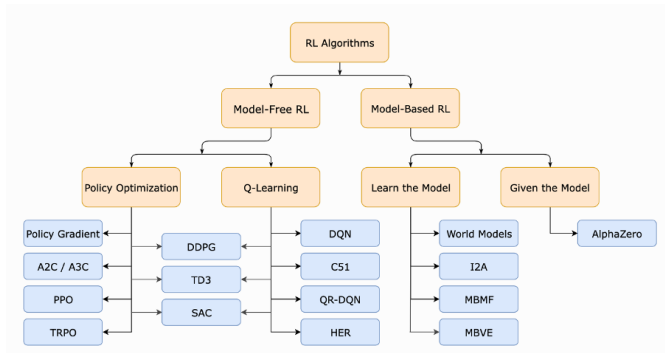
Q Learning

Deep Q Learning  
Training

Best Practices

Demo

Discussion

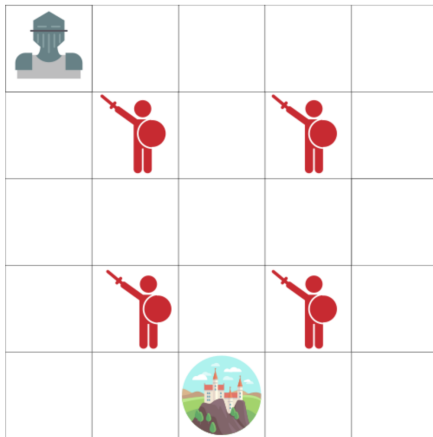


# Q Learning

# Q Learning

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löh



Reinforcement  
Learning

Q Learning

Deep Q Learning  
Training

Best Practices

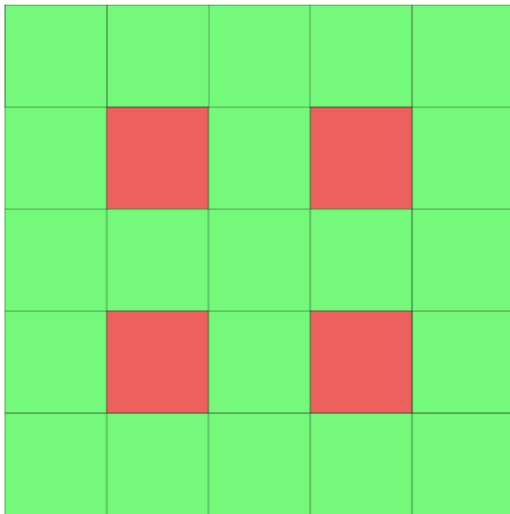
Demo

Discussion

# Q Learning

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhrl



Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

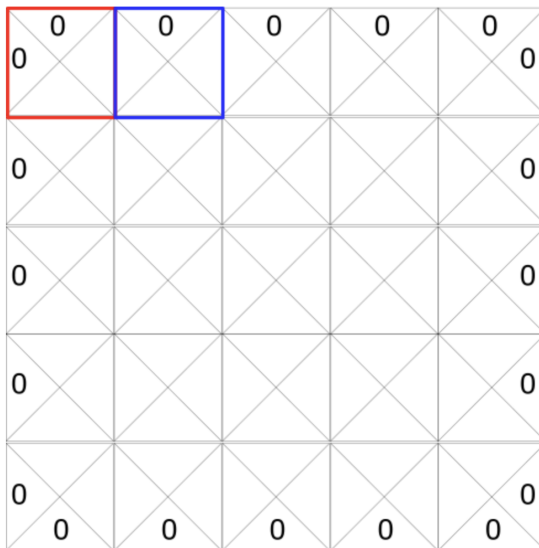
Discussion



# Q Learning

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhrl



Reinforcement  
Learning

Q Learning

Deep Q Learning  
Training

Best Practices

Demo

Discussion

# Q Learning

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhrl

Q Learning is a value-based RL algorithm

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

Q-Values for the state  
given a particular state

Expected discounted  
cumulative reward

Given the state and action

Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

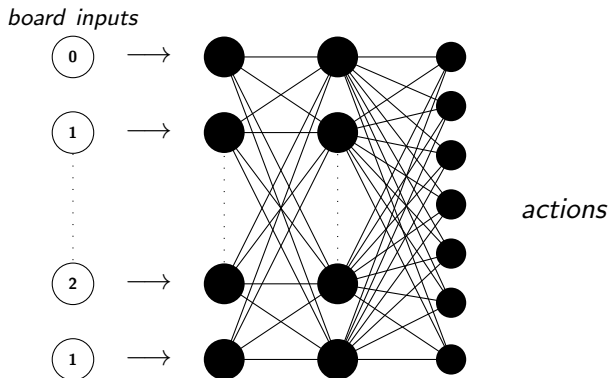
Demo

Discussion

# Deep Q Learning

## Deep Q Learning

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



# Implementation of a Deep Q Learning Agent

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhner

## 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
- ▶ Exploitation: Agent takes actions given by the network (random)

→ We use *epsilon decay* to find the right balance between the two

Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

# Implementation of a Deep Q Learning Agent

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimdander,  
Tim Löhner

Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

## 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 Löhner

Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

## 3. Learning Process

- ▶ for every action during a game the following quintuple is saved and stored to the agents memory  
(observation, action, reward, new observation, done)
- ▶ rewards are determined when game is finished
- ▶ the agent learns the saved transitions in batches that are chosen randomly from its memory

→ solves the problem of correlated inputs/outputs

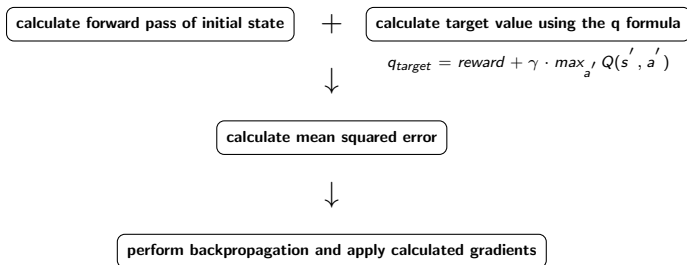
# Implementation of a Deep Q Learning Agent

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimeter,  
Tim Löh

## 3. Learning Process

- ▶ for each transition in the batch we perform the following steps



Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

→ problem of non-stationarity



# Implementation of a Deep Q Learning Agent

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimeter,  
Tim Löh

Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

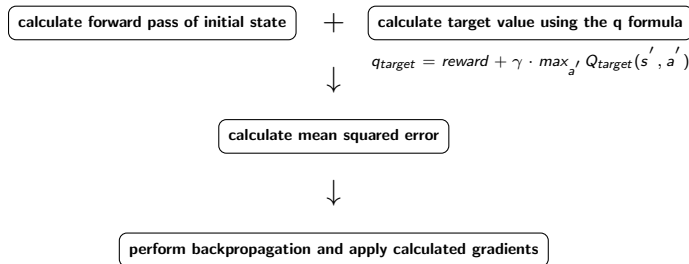
Best Practices

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



→ targets stay constant and problem becomes more stationary

# Training

Reinforcement  
Learning

Q Learning

Deep Q Learning

**Training**

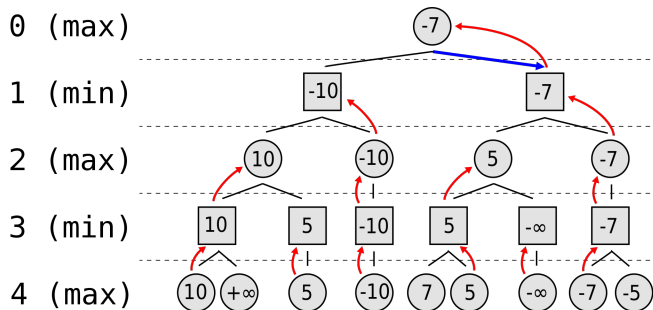
Best Practices

Demo

Discussion

# Opponent: Minimax

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





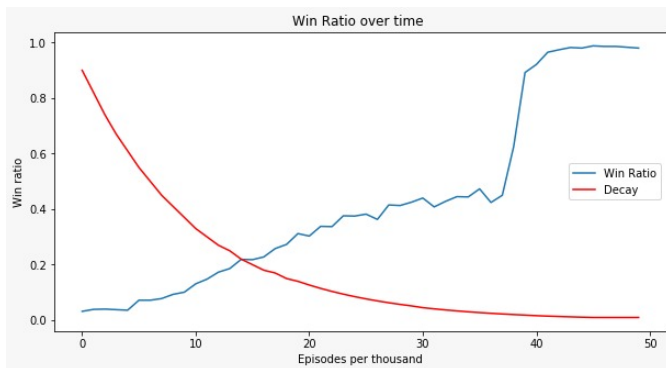


# Learning over time

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhrl

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



Reinforcement  
Learning

Q Learning

Deep Q Learning

Training

Best Practices

Demo

Discussion

# Save the trained model and start all-over again

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 Löhner

Reinforcement  
Learning

Q Learning

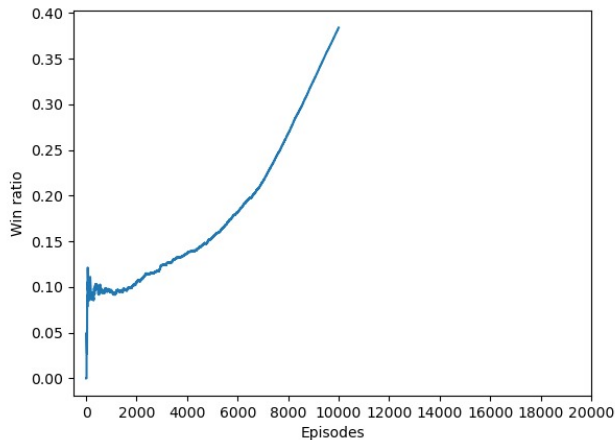
Deep Q Learning

**Training**

Best Practices

Demo

Discussion



# Best Practices



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

Best set of Hyperparameters so far:

- ▶ Randomness 100% with quadratic fall down to fixed 5%
- ▶ Batch Size = 128
- ▶ Memory size = 50000
- ▶ Learning Rate for the Neural Network = 0.01
- ▶ Episodes 50000 on depth = 1: *20 hours*
- ▶ Episodes 10000 on depth = 2: *10 hours*

# Demo

## Demo

## Reinforcement Learning for playing Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhre

## It is time to show a Demo in the GUI



## Demo

# Discussion

# Discussion

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhrl

- ▶ 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

Q Learning

Deep Q Learning  
Training

Best Practices

Demo

Discussion

Thanks for listening to us!

# Sources

- ▶ Figure 1: <https://www.pngwave.com/png-clip-art-jdcjo>
- ▶ Figure 2: Reinforcement Learning: An Introduction (Sutton, Barto)
- ▶ Figure 4: <https://www.afcea.org/content/artificial-intelligence-will-change-human-values>

Reinforcement  
Learning for  
playing  
Connect Four

Simon Hölck,  
Florian Cimander,  
Tim Löhner

Reinforcement  
Learning

Q Learning

Deep Q Learning  
Training

Best Practices

Demo

Discussion