



Hochschule Karlsruhe
Technik und Wirtschaft
UNIVERSITY OF APPLIED SCIENCES



inovex

AI Labor - Sommersemester 2020

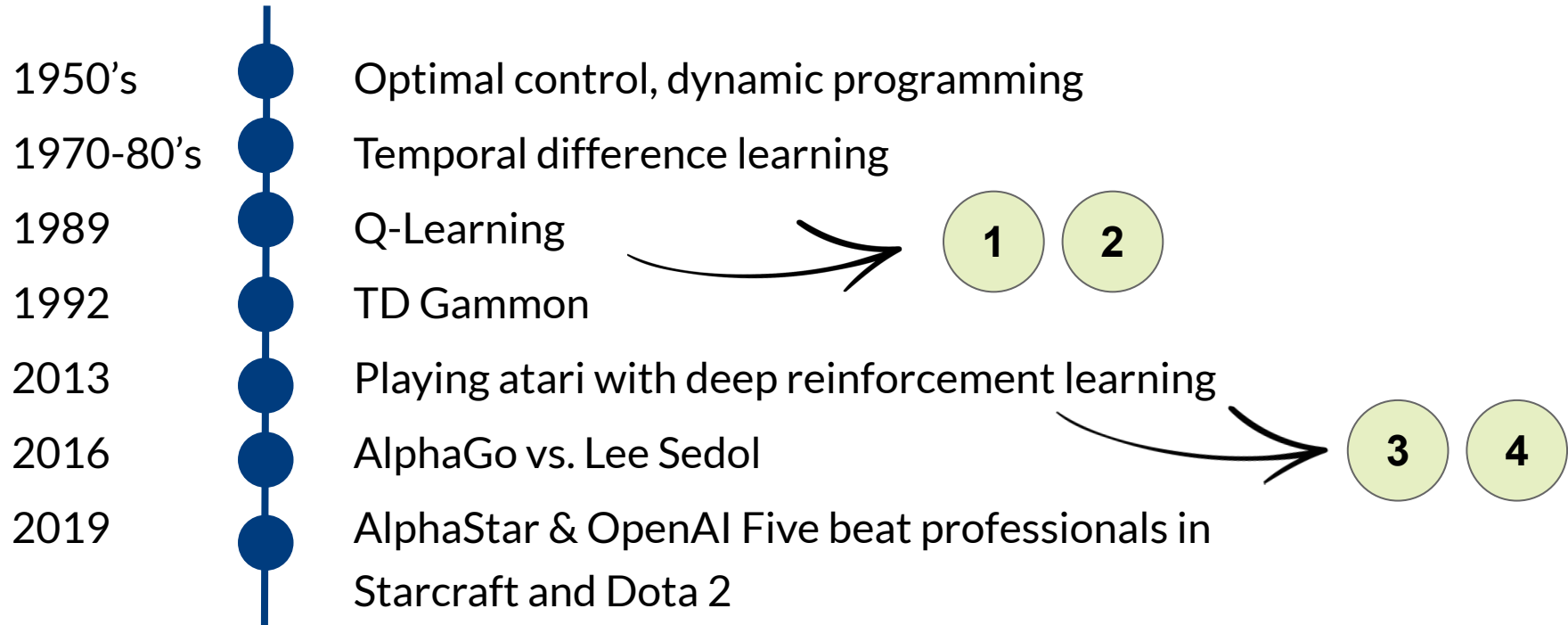
Reinforcement **L**earning
Sprintwechsel

Frederik Martin, Sebastian Blank

Karlsruhe, 12. Juni 2020

Reinforcement Learning

Meilensteine im Reinforcement Learning



Reinforcement Learning

› **Theorie**

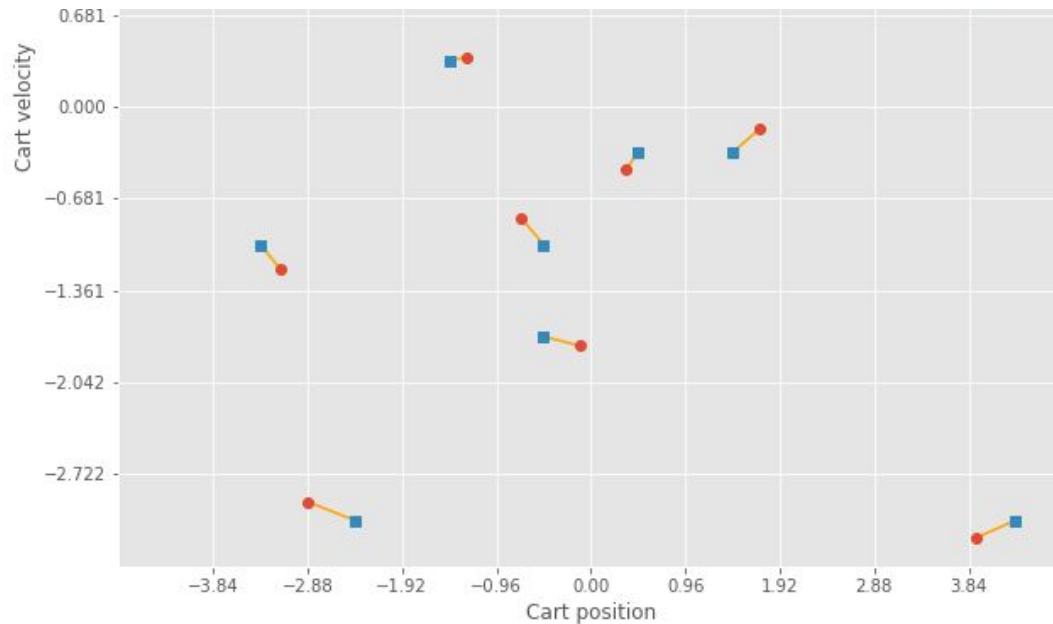
- Deep Q-Network
- Experience Replay
- Target Model
- Vorverarbeitung für Pixel-basierte Atari Games (Framestacking, etc.)

› **Praxis**

- CartPole Gym mit Deep Q-Learning (Aufgabe 3)
- Pong (Pixel-basiert) mit Deep Q-Learning (Aufgabe 4)

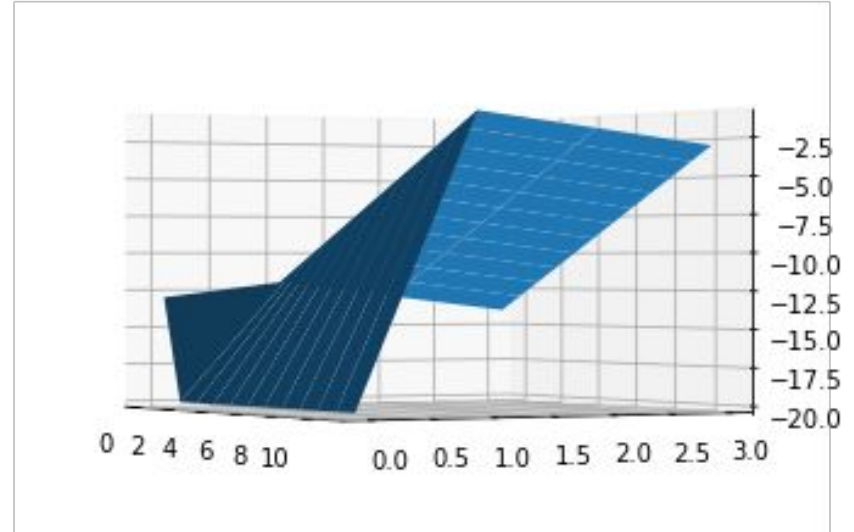
Was bisher geschah ...

Diskretisierung



Funktionsapproximation

Beispiel Cliffwalking

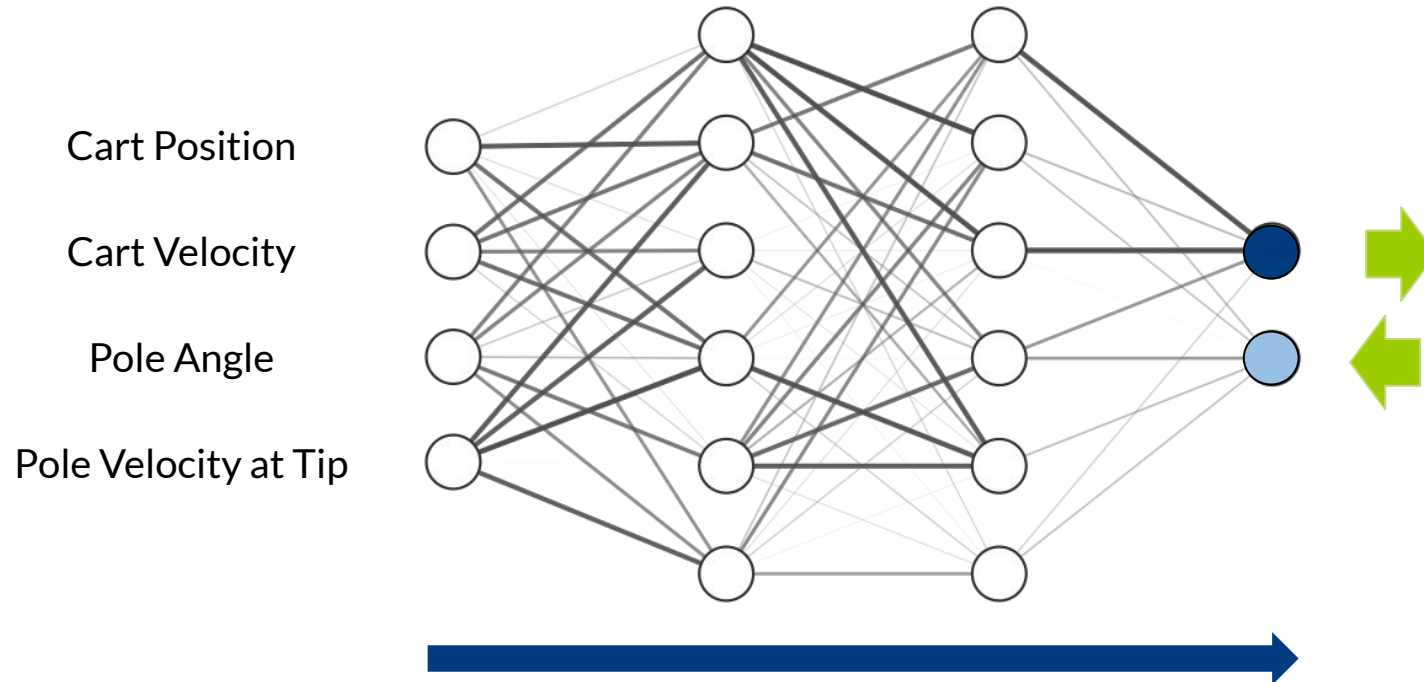


$$\hat{v}(s, \theta) \approx v_{\pi}(s), \theta \in \mathbb{R}^d$$

\nearrow **Approximation** \uparrow **Value-Function** \uparrow **Gewichtsvektor**

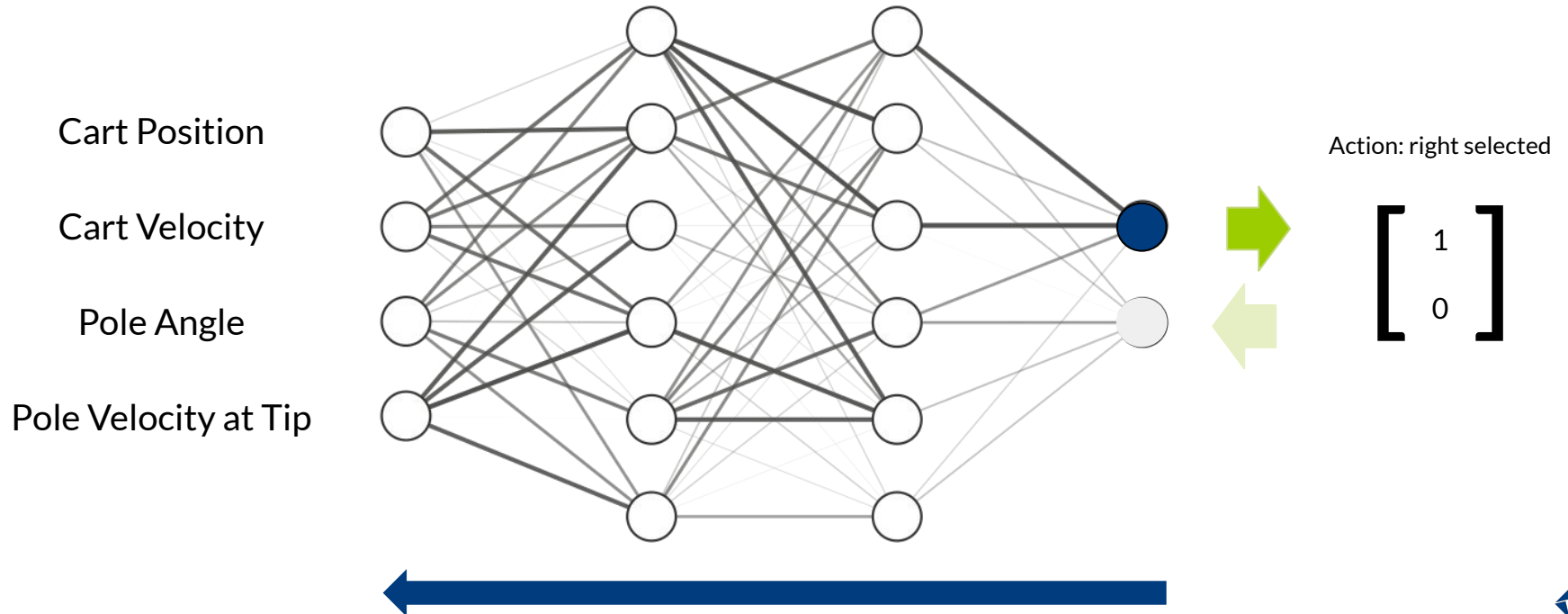
Funktionsapproximation

Beispiel CartPole: Forward-Pass



Funktionsapproximation

Beispiel CartPole: Backward-Pass



Funktionsapproximation

Deep Q-Network

- Approximation der Q-Funktion mit NN
- Optimierung mit SGD
 - Minimierung des Abstands zwischen Schätzer und Target

$$\underline{L_i(\theta_i)} = \mathbb{E}_{\underline{(s,a,r,s') \sim U(D)}} [\underline{(y_i - \hat{q}(s, a, \theta_i))^2}]$$

Diagram illustrating the components of the loss function $L_i(\theta_i)$ in the context of a Deep Q-Network (DQN):

- Loss in Iteration i**: Points to L_i .
- Weights**: Points to θ_i .
- Experience Replay**: Points to the expectation operator $\mathbb{E}_{(s,a,r,s') \sim U(D)}$.
- Target Q-Value**: Points to y_i .
- approximierter Q-Value**: Points to $\hat{q}(s, a, \theta_i)$.

Experience Replay

- Problem
 - Starke Korrelation der States erschwert das Lernen
- Lösung
 - Letzte N Experiences werden in **Replay Memory** gespeichert
 - Random Uniform Sampling

$e_{t+N} = (s_{t+N}, a_{t+N}, s_{t+N+1}, r_{t+N+1})$
...
$e_t = (s_t, a_t, s_{t+1}, r_{t+1})$

Replay Memory

Target Network

- Kopie der eigentlichen Architektur mit fixen Gewichten
- Update alle c Iterationen

$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim U(D)} [\underbrace{(y_i - \hat{q}(s, a, \theta_i))^2}_{\text{Q-Network}}]$$

Q-Network

$$r + \gamma \max_{a'} Q(s', a', \theta_i^-)$$

Target Network

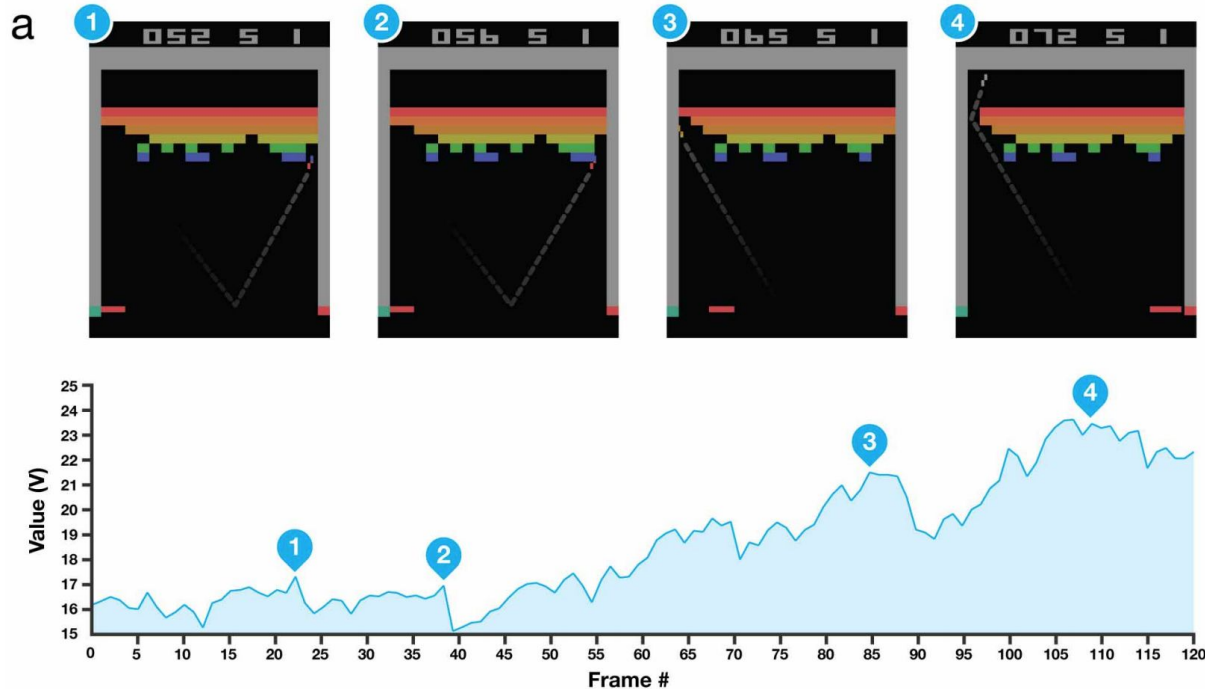
Aufgabe 3: CartPole Gym mit Deep Q-Learning

› Jupyter Lab Notebook



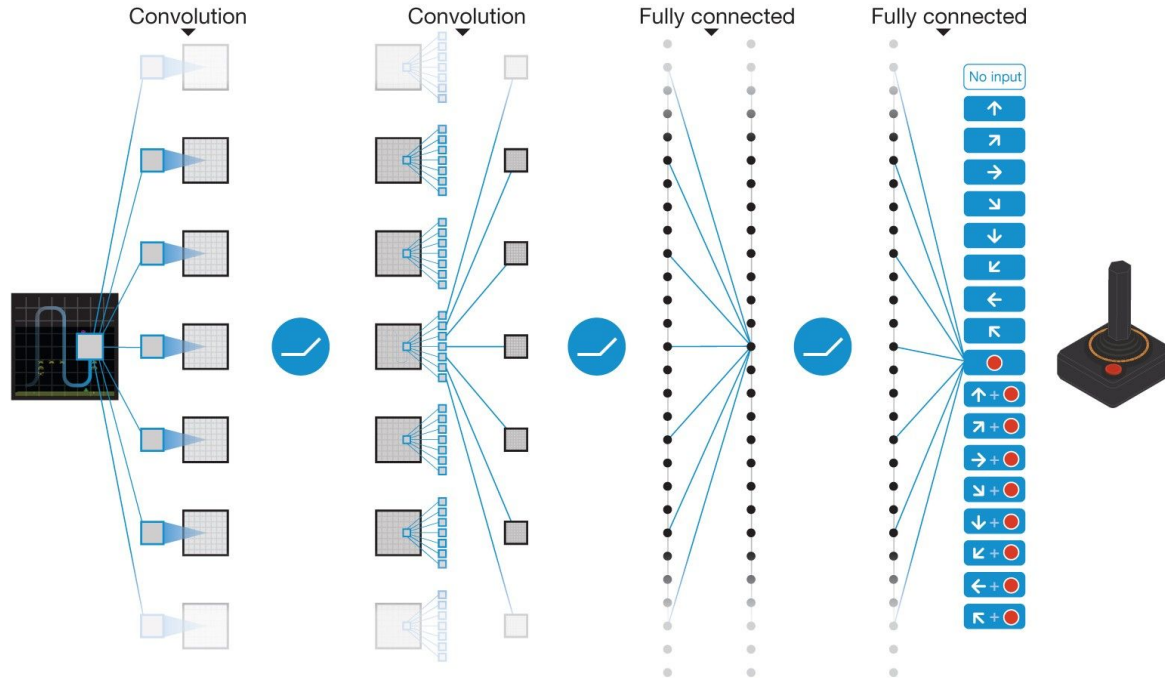
Deep Q-Network

Pixel-basierte Atari Games



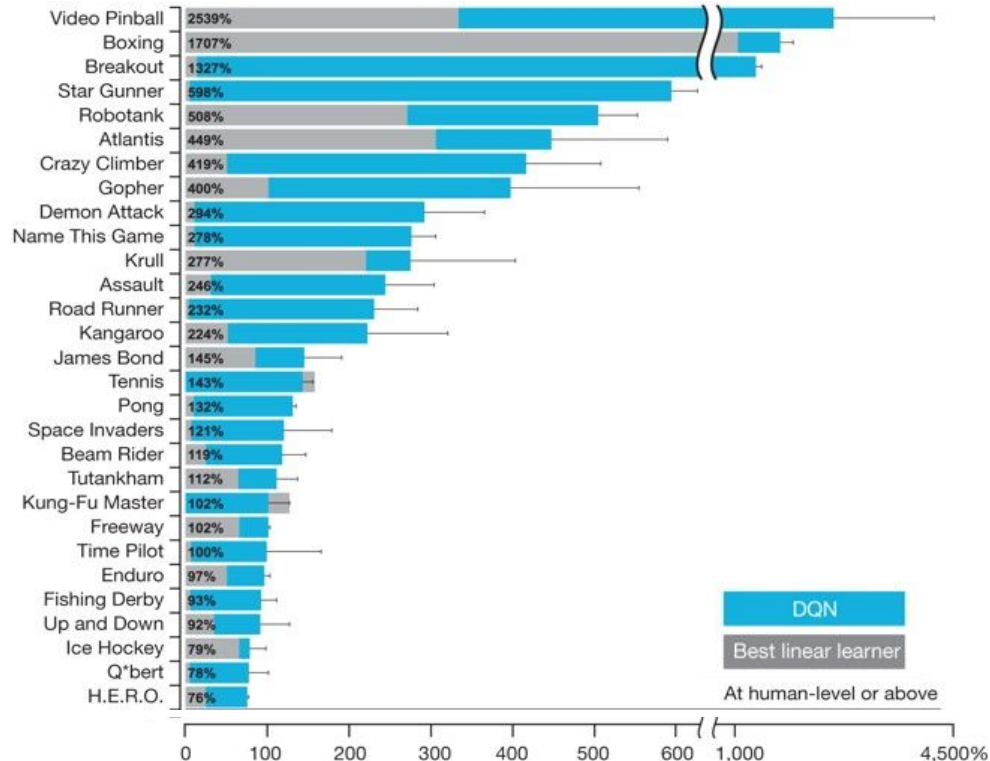
Deep Q-Network

Pixel-basierte Atari Games

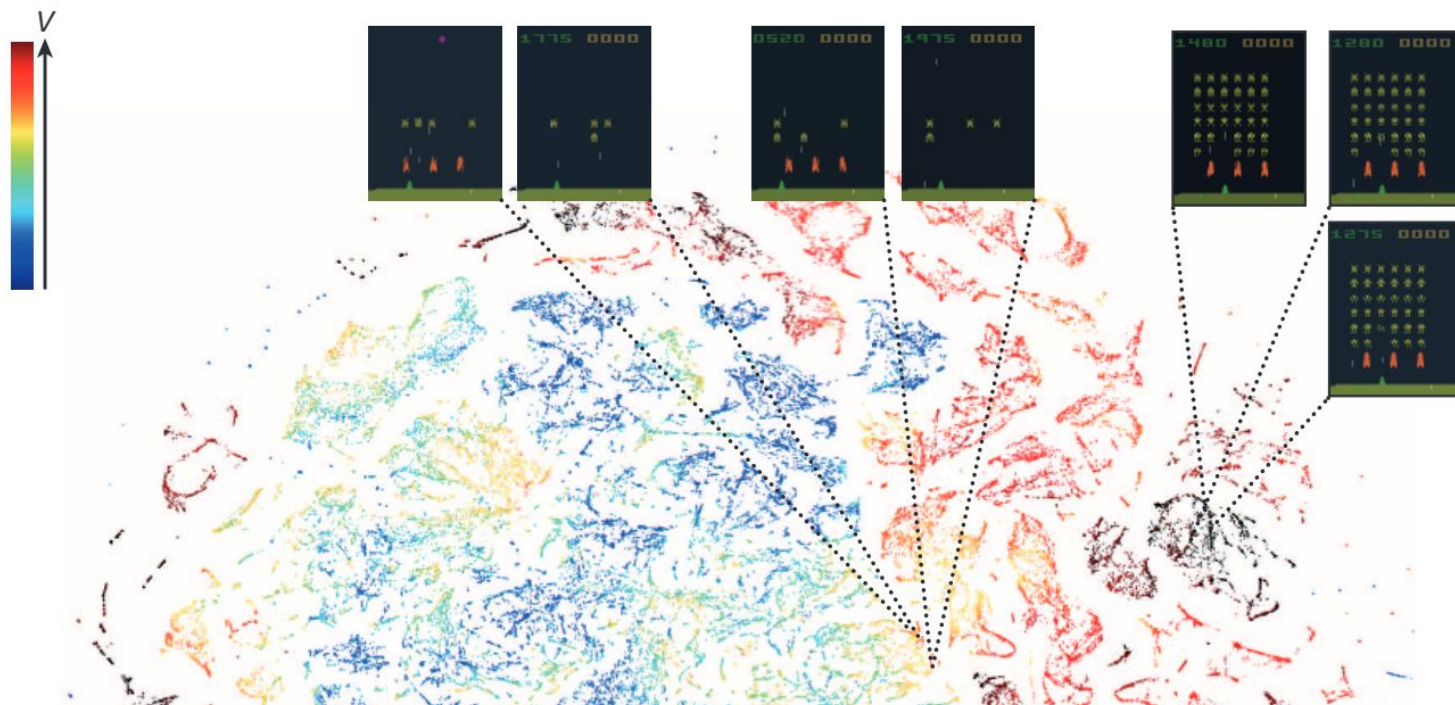


Deep Q-Network

Pixel-basierte Atari Games



Deep Q-Network



Deep Q-Network

Preprocessing Pixel-basierter Inputs

- **Warp Frame**
 - Konvertierung der Frames in Graufstufen
 - Downsampling auf 84x84 Pixel
- **Framestack**
 - Kombination von vier aufeinanderfolgenden Frames um Bewegungen nachvollziehen zu können
 - Auswahl der maximalen Helligkeitswerten
 - Ausführung der gewählten Aktion für Framestack

Deep Q-Network

Besonderheiten Pong / Atari

- **No-Ops after Reset**
 - Skippen der ersten 30 Steps
- **Fire Reset**
 - Automatisches Drücken der FIRE-Taste als erste Aktion
- **Episodic Life**
 - Verlust eines Lebens bedeutet das Ende der Episode
 - Game Over setzt die Environment komplett zurück
- **Reward Clipping**
 - Sieg: +1 Niederlage: -1 Sonst: 0

Deep Q-Learning

Algorithmus

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Erstes Paper
(2013) ohne
Target Network

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right]. \quad (3)$$

Aufgabe 4: Pong mit Deep Q-Learning

› Jupyter Lab Notebook



Literatur

- › Kostenlose "Standard"-Lektüre für den Einstieg in RL:
Reinforcement Learning: An Introduction (Sutton and Barto), siehe
<http://incompleteideas.net/book/RLbook2018.pdf>
- › Ausführlich und gut erklärter Einstieg in RL (Video-Lektionen):
UCL Course on RL (David Silver, Google DeepMind), siehe
<https://www.davidsilver.uk/teaching/>
- › *Algorithms in Reinforcement Learning* von Csaba Szepesvári, siehe
<https://sites.ualberta.ca/~szepesva/papers/RLAlgsInMDPs.pdf>
- › Blog mit Videos zum Einstieg in RL und Q-Learning, DQN und vieles mehr:
Reinforcement Learning – Introducing Goal Oriented Intelligence, siehe
<https://deeplizard.com/learn/video/nyjbcRQ-uQ8>

Feedback



<https://forms.gle/fPisrqSjqHJCCV8Q6>

Vielen Dank

Frederik Martin

fmartin@inovex.de

Sebastian Blank

sblank@inovex.de

