

RWTH AACHEN  
UNIVERSITY

# Advanced Machine Learning Lecture 20

## Deep Reinforcement Learning II

02.02.2017

Bastian Leibe  
RWTH Aachen  
<http://www.vision.rwth-aachen.de/>  
leibe@vision.rwth-aachen.de

Advanced Machine Learning Winter'16

RWTH AACHEN  
UNIVERSITY

## This Lecture: *Advanced Machine Learning*

- Regression Approaches
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
  - Gaussian Processes
- Approximate Inference
  - Sampling Approaches
  - MCMC
- Deep Learning
  - Linear Discriminants
  - Neural Networks
  - Backpropagation & Optimization
  - CNNs, ResNets, RNNs, **Deep RL**, etc.

B. Leibe

Advanced Machine Learning Winter'16

RWTH AACHEN  
UNIVERSITY

## Topics of This Lecture

- Recap: Reinforcement Learning
  - Key Concepts
  - Temporal Difference Learning
- Deep Reinforcement Learning
  - Value based Deep RL
  - Policy based Deep RL
  - Model based Deep RL
- Applications

B. Leibe

4

Advanced Machine Learning Winter'16

RWTH AACHEN  
UNIVERSITY

## Recap: Reinforcement Learning

- Motivation
  - General purpose framework for decision making.
  - Basis: **Agent** with the capability to **interact** with its **environment**
  - Each **action** influences the agent's future **state**.
  - Success is measured by a scalar **reward** signal.
  - Goal: **select actions to maximize future rewards**.

```

    graph TD
      Agent -- action --> Environment
      Environment -- "observation, reward" --> Agent
  
```

- Formalized as a partially observable Markov decision process (POMDP)

Slide adapted from: David Silver, Sergey Levine

5

Advanced Machine Learning Winter'16

RWTH AACHEN  
UNIVERSITY

## Recap: Reward vs. Return

- Objective of learning
  - We seek to maximize the **expected return**  $G_t$  as some function of the reward sequence  $R_{t+1}, R_{t+2}, R_{t+3}, \dots$
  - Standard choice: **expected discounted return**

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

where  $0 \leq \gamma \leq 1$  is called the **discount rate**.

- Difficulty
  - We don't know which past actions caused the reward.
  - ⇒ Temporal credit assignment problem

B. Leibe

6

Advanced Machine Learning Winter'16

RWTH AACHEN  
UNIVERSITY

## Recap: Policy

- Definition
  - A policy determines the agent's behavior
  - Map from state to action  $\pi: \mathcal{S} \rightarrow \mathcal{A}$
- Two types of policies
  - Deterministic policy:  $a = \pi(s)$
  - Stochastic policy:  $\pi(a|s) = \Pr\{A_t = a | S_t = s\}$
- Note
  - $\pi(a|s)$  denotes the probability of taking action  $a$  when in state  $s$ .

B. Leibe

7

Advanced Machine Learning Winter'16

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

## Recap: Value Function

- **Idea**
  - Value function is a prediction of future reward
  - Used to evaluate the goodness/badness of states
  - And thus to select between actions
- **Definition**
  - The **value of a state**  $s$  under a policy  $\pi$ , denoted  $v_\pi(s)$ , is the expected return when starting in  $s$  and following  $\pi$  thereafter.
 
$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]$$
  - The **value of taking action**  $a$  in state  $s$  under a policy  $\pi$ , denoted  $q_\pi(s, a)$ , is the expected return starting from  $s$ , taking action  $a$ , and following  $\pi$  thereafter.
 
$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a]$$

B. Leibe 8

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

## Recap: Optimal Value Functions

- **Bellman optimality equations**
  - For the **optimal state-value function**  $v_*$ :
 
$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a)$$

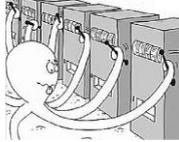
$$= \max_{a \in \mathcal{A}(s)} \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$
  - $v_*$  is the unique solution to this system of nonlinear equations.
  - For the **optimal action-value function**  $q_*$ :
 
$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]$$
  - $q_*$  is the unique solution to this system of nonlinear equations.
- ➔ If the dynamics of the environment  $p(s', r | s, a)$  are known, then in principle one can solve those equation systems.

B. Leibe 9

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

## Recap: Exploration-Exploitation Trade-off

- **Example: N-armed bandit problem**
  - Suppose we have the choice between  $N$  actions  $a_1, \dots, a_N$ .
  - If we knew their value functions  $q_*(s, a_i)$ , it would be trivial to choose the best.
  - However, we only have estimates based on our previous actions and their returns.
- **We can now**
  - **Exploit** our current knowledge
    - And choose the **greedy** action that has the highest value based on our current estimate.
  - **Explore** to gain additional knowledge
    - And choose a **non-greedy** action to improve our estimate of that action's value.

B. Leibe 10  
Image source: research.microsoft.com

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

## Recap: TD-Learning

- **Policy evaluation (the prediction problem)**
  - For a given policy  $\pi$ , compute the state-value function  $v_\pi$ .
- **One option: Monte-Carlo methods**
  - Play through a sequence of actions until a reward is reached, then backpropagate it to the states on the path.
 
$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

Target: the actual return after time  $t$
- **Temporal Difference Learning - TD( $\lambda$ )**
  - Directly perform an update using the estimate  $V(S_{t+\lambda+1})$ .
 
$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Target: an estimate of the return (here: TD(0))

B. Leibe 11

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

## Recap: SARSA - On-Policy TD Control

- **Idea**
  - Turn the TD idea into a control method by always updating the policy to be greedy w.r.t. the current estimate
- **Procedure**
  - Estimate  $q_\pi(s, a)$  for the current policy  $\pi$  and for all states  $s$  and actions  $a$ .
  - TD(0) update equation
 
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$
  - This rule is applied after every transition from a nonterminal state  $S_t$ .
  - It uses every element of the quintuple  $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$ .
  - ➔ Hence, the name SARSA.

B. Leibe 12  
Image source: Sutton & Barto

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

## Recap: Q-Learning - Off-Policy TD Control

- **Idea**
  - Directly approximate the optimal action-value function  $q_*$ , independent of the policy being followed.
- **Procedure**
  - TD(0) update equation
 
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$
  - Dramatically simplifies the analysis of the algorithm.
  - All that is required for correct convergence is that all pairs continue to be updated.

B. Leibe 13  
Image source: Sutton & Barto

RWTH AACHEN UNIVERSITY

## Approaches Towards RL

- Value-based RL
  - Estimate the **optimal value function**  $q_*(s, a)$
  - This is the maximum value achievable under any policy
- Policy-based RL
  - Search directly for the **optimal policy**  $\pi_*$
  - This is the policy achieving maximum future reward
- Model-based RL
  - Build a model of the environment
  - Plan (e.g. by lookahead) using model

Advanced Machine Learning Winter'16

Slide credit: David Silver

B. Leibe

14

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Recap: Reinforcement Learning
  - Key Concepts
  - Temporal Difference Learning
- Deep Reinforcement Learning
  - Value based Deep RL
  - Policy based Deep RL
  - Model based Deep RL
- Applications

Advanced Machine Learning Winter'16

B. Leibe

15

RWTH AACHEN UNIVERSITY

## Deep Reinforcement Learning

- RL using deep neural networks to approximate functions
  - Value functions
    - Measure goodness of states or state-action pairs
  - Policies
    - Select next action
  - Dynamics Models
    - Predict next states and rewards

Advanced Machine Learning Winter'16

Slide credit: Sergey Levine

B. Leibe

16

RWTH AACHEN UNIVERSITY

## Deep Reinforcement Learning

- Use deep neural networks to represent
  - Value function
  - Policy
  - Model
- Optimize loss function by stochastic gradient descent

Advanced Machine Learning Winter'16

Slide credit: David Silver

B. Leibe

17

RWTH AACHEN UNIVERSITY

## Q-Networks

- Represent value function by **Q-Network** with weights  $w$

$$Q(s, a, w) = q_*(s, a)$$

Advanced Machine Learning Winter'16

Slide credit: David Silver

B. Leibe

18

RWTH AACHEN UNIVERSITY

## Deep Q-Learning

- Idea
  - Optimal Q-values should obey Bellman equation
 
$$Q_*(s, a) = \mathbb{E} \left[ r + \gamma \max_{a'} Q(s', a') \mid s, a \right]$$
  - Treat the right-hand side  $r + \gamma \max_{a'} Q(s', a', w)$  as a target
  - Minimize MSE loss by stochastic gradient descent
 
$$L(w) = \left( r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right)^2$$
  - This converges to  $Q_*$  using a lookup table representation.
  - Unfortunately, it **diverges** using neural networks due to
    - Correlations between samples
    - Non-stationary targets

Advanced Machine Learning Winter'16

Slide adapted from David Silver

B. Leibe

19

RWTH AACHEN UNIVERSITY

## Deep Q-Networks (DQN): Experience Replay

- Adaptations
  - To remove correlations, build a dataset from agent's own experience
 

$s_1, a_1, r_1, s_2$
$s_2, a_2, r_2, s_3$
$s_3, a_3, r_3, s_4$
...
$s_t, a_t, r_t, s_{t+1}$

 $\rightarrow s, a, r, s'$
- Perform minibatch updates to samples of experience drawn at random from the pool of stored samples
  - $(s, a, r, s') \sim U(D)$  where  $D = \{(s_t, a_t, r_{t+1}, s_{t+1})\}$  is the dataset
- Advantages
  - Each experience sample is used in many updates (more efficient)
  - Avoids correlation effects when learning from consecutive samples
  - Avoids feedback loops from on-policy learning

Slide adapted from David Silver B. Leibe 20

RWTH AACHEN UNIVERSITY

## Deep Q-Networks (DQN): Experience Replay

- Adaptations
  - To remove correlations, build a dataset from agent's own experience
 

$s_1, a_1, r_1, s_2$
$s_2, a_2, r_2, s_3$
$s_3, a_3, r_3, s_4$
...
$s_t, a_t, r_t, s_{t+1}$

 $\rightarrow s, a, r, s'$
  - Sample from the dataset and apply an update
 
$$L(w) = (r + \gamma \max_{a'} Q(s', a', w^-) - Q(s, a, w))^2$$
  - To deal with non-stationary parameters  $w^-$ , are held fixed.
    - Only update the target network parameters every  $C$  steps.
    - I.e., clone the network  $Q$  to generate a target network  $\hat{Q}$ .
    - $\Rightarrow$  Again, this reduces oscillations to make learning more stable.

Slide adapted from David Silver B. Leibe 21

RWTH AACHEN UNIVERSITY

## Application: Deep RL in Atari

- Goal: Learning to play Atari games

V. Mnih et al., *Human-level control through deep reinforcement learning*, Nature Vol. 518, pp. 529-533, 2015  
B. Leibe Image source: Vladimír Mnih et al. 22

RWTH AACHEN UNIVERSITY

## Idea Behind the Model

- Interpretation
  - Assume finite number of actions
  - Each number here is a real-valued quantity that represents the **Q function** in Reinforcement Learning
- Collect experience dataset:
  - Set of tuples  $\{(s, a, s', r), \dots\}$
  - (State, Action taken, New state, Reward received)
- L2 Regression Loss
 
$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right]$$

*Current reward + estimate of future reward, discounted by  $\gamma$*

Slide credit: Andrei Karpaty B. Leibe 23

RWTH AACHEN UNIVERSITY

## Results: Breakout

B. Leibe Video source: Vladimír Mnih et al. 24

RWTH AACHEN UNIVERSITY

## Results: Space Invaders

B. Leibe Video source: Vladimír Mnih et al. 25



Advanced Machine Learning Winter'16

## Policy Gradients

- How to make high-value actions more likely
  - The gradient of the stochastic policy  $\pi(s, \mathbf{u})$  is given by
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial}{\partial \mathbf{u}} \mathbb{E}[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | \pi(\cdot, \mathbf{u})]$$

$$= \dots ?$$
- Wait - how do we calculate that?
  - Any ideas?

Slide adapted from David Silver B. Leibe 32

Advanced Machine Learning Winter'16

## Policy Gradients

- Deriving the gradient of an expectation
  - General case
 
$$\nabla_{\theta} \mathbb{E}_{p(x; \theta)}[f(x)] = \nabla_{\theta} \sum_x p(x; \theta) f(x)$$

$$= \sum_x \nabla_{\theta} p(x; \theta) f(x)$$

$$= \sum_x p(x; \theta) \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)} f(x)$$

$$= \sum_x p(x; \theta) \nabla_{\theta} \log p(x; \theta) f(x)$$

$$= \mathbb{E}_{p(x; \theta)}[\nabla_{\theta} \log p(x; \theta) f(x)]$$

B. Leibe 33

Advanced Machine Learning Winter'16

## Policy Gradients

- How to make high-value actions more likely
  - The gradient of a stochastic policy  $\pi(s, \mathbf{u})$  is given by
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial}{\partial \mathbf{u}} \mathbb{E}_{\pi} [r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | \pi(\cdot, \mathbf{u})]$$

$$= \mathbb{E}_{\pi} \left[ \frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q_{\pi}(s, a) \right]$$
  - The gradient of a deterministic policy  $a = \pi(s)$  is given by
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbb{E}_{\pi} \left[ \frac{\partial Q_{\pi}(s, a)}{\partial a} \frac{\partial a}{\partial \mathbf{u}} \right]$$

if  $a$  is continuous and  $Q$  is differentiable.

Slide adapted from David Silver B. Leibe 34

Advanced Machine Learning Winter'16

## Actor-Critic Algorithm

- Procedure
  - Estimate value function  $Q(s, a, \mathbf{w}) \approx Q_{\pi}(s, a)$
  - Update policy parameters  $\mathbf{u}$  by stochastic gradient ascent
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q(s, a, \mathbf{w})$$
 **stochastic policy**
  - or
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial Q(s, a, \mathbf{w})}{\partial a} \frac{\partial a}{\partial \mathbf{u}}$$
 **deterministic policy**

Slide adapted from David Silver B. Leibe 35

Advanced Machine Learning Winter'16

## Asynchronous Advantage Actor-Critic (A3C)

- Further improvement
  - Estimate state-value function
 
$$V(s) \approx \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \dots | s]$$
  - Q-value estimated by an  $n$ -step sample
 
$$q_t = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n V(s_{t+n}, \mathbf{v})$$
  - Actor is updated towards target
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial \log \pi(a_t | s_t, \mathbf{u})}{\partial \mathbf{u}} (q_t - V(s_t, \mathbf{v}))$$
  - Critic is updated to minimize MSE w.r.t. target
 
$$L_v = (q_t - V(s_t, \mathbf{v}))^2$$

⇒ Combined effect: 4x mean Atari score vs. Nature DQN

Slide credit: David Silver B. Leibe 36

Advanced Machine Learning Winter'16

## Deep Policy Gradients (DPG)

- DPG is the continuous analogue of DQN
  - Experience replay: build data-set from agent's experience
  - Critic estimates value of current policy by DQN
 
$$L_w(\mathbf{w}) = (r + \gamma Q(s', \pi(s', \mathbf{u}^-), \mathbf{w}^-) - Q(s, a, \mathbf{w}))^2$$
  - To deal with non-stationarity, targets  $\mathbf{u}^-, \mathbf{w}^-$  are held fixed
  - Actor updates policy in direction that improves Q
 
$$\frac{\partial L_u(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial Q(s, a, \mathbf{w})}{\partial a} \frac{\partial a}{\partial \mathbf{u}}$$
  - In other words critic provides loss function for actor.

Slide credit: David Silver B. Leibe 37

RWTH AACHEN  
UNIVERSITY

## Summary

- The future looks bright!
  - Soon, you won't have to play video games anymore...
  - Your computer can do it for you (and beat you at it)
- Reinforcement Learning is a very promising field
  - Currently limited by the need for data
  - At the moment, mainly restricted to simulation settings

38

Advanced Machine Learning Winter'16

B. Leibe

RWTH AACHEN  
UNIVERSITY

## Topics of This Lecture

- Recap: Reinforcement Learning
  - Key Concepts
  - Temporal Difference Learning
- Deep Reinforcement Learning
  - Value based Deep RL
  - Policy based Deep RL
  - Model based Deep RL
- Applications

39

Advanced Machine Learning Winter'16

B. Leibe

RWTH AACHEN  
UNIVERSITY

## Often Used in Games, E.g. Alpha Go



40

Advanced Machine Learning Winter'16

B. Leibe

RWTH AACHEN  
UNIVERSITY

## References and Further Reading

- More information on Reinforcement Learning can be found in the following book



Richard S. Sutton, Andrew G. Barto  
Reinforcement Learning: An Introduction  
MIT Press, 1998

- The complete text is also freely available online  
<https://webdocs.cs.ualberta.ca/~sutton/book/ebook/the-book.html>

43

Advanced Machine Learning Winter'16

B. Leibe

RWTH AACHEN  
UNIVERSITY

## References and Further Reading

- DQN paper
  - [www.nature.com/articles/nature14236](http://www.nature.com/articles/nature14236)
- AlphaGo paper
  - [www.nature.com/articles/nature16961](http://www.nature.com/articles/nature16961)



44

Advanced Machine Learning Winter'16

B. Leibe