

### Lecture 12: Deep Reinforcement Learning

Deep Learning @ UvA

UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES

### Reinforcement Learning

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP REINFORCEMENT LEARNING - 2



## What is Reinforcement Learning?

- General purpose framework for learning Artificial Intelligence models
- RL assumes that *the agent* (our model) can take *actions*
- These actions affect *the environment* where *the agent* operates
   more specifically *the state* of the environment and *the state* of the agent
- Given the state of the environment and the agent, an action taken from the agent causes a reward
  - can be positive or negative
- Goal: the goal of an RL agent is to learn how to take actions that maximize future rewards

### Some examples of RL

### Some examples of RL

- Controlling physical systems
  - Robot walking, jumping, driving
- Logistics
  - Scheduling, bandwidth allocation
- o Games
  - Atari, Go, Chess, Pacman
- Learning sequential algorithms
  - Attention, memory

### Reinforcement Learning: An abstraction



• Experience is a series of observations, actions and rewards  $o_1, r_1, a_1, o_2, r_2, a_2, \dots, o_t, r_t$ 

• The state is the summary of experience so far  $s_t = f(o_1, r_1, a_1, o_2, r_2, a_2, \dots, o_t, r_t)$ 

• If we have fully observable environments, then  $s_t = f(o_t)$ 

• Policy is the agent's behavior function

 $\circ$  The policy function maps the state input  $s_t$  to an action output  $a_t$ 

- Deterministic policy:  $a_t = f(s_t)$
- Stochastic policy:  $\pi(a_t|s_t) = \mathbb{P}(a_t|s_t)$

• A value function is the prediction of the future reward • Given the state  $s_t$  what will my reward be if I do action  $a_t$ 

• The Q-value function gives the expected future reward

• Given state  $s_t$ , action  $a_t$ , a policy  $\pi$  the Q-value function is  $Q^{\pi}(s_t, a_t)$ 

### How do we decide about actions, states, rewards?

• We model the policy and the value function as machine learning functions that can be optimized by the data

• The *policy function*  $a_t = \pi(s_t)$  selects an action given the current state

• The value function  $Q^{\pi}(s_t, a_t)$  is the expected total reward that we will receive if we take action  $a_t$  given state  $s_t$ 

 $\,\circ\,$  What should our goal then be?

### Goal: Maximize future rewards!

• Learn the policy and value functions such that the action taken at the t-th time step  $a_t$  maximizes the expected sum of future rewards

$$Q^{\pi}(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t, a_t)$$

 $\circ \gamma$  is a discount factor. Why do we need it?

• Learn the policy and value functions such that the action taken at the t-th time step  $a_t$  maximizes the expected sum of future rewards

$$Q^{\pi}(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t, a_t)$$

- $\circ \gamma$  is a discount factor. Why do we need it?
  - The further into the future we look t + 1, ..., t + T, the less certain we can be about our expected rewards  $r_{t+1}, ..., r_{t+T}$

### **Bellman equation**

#### • How can we rewrite the value function in more compact form $Q^{\pi}(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t, a_t) =?$

• How can we rewrite the value function in more compact form  $Q^{\pi}(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t, a_t)$   $= \mathbb{E}_{s',a'}(r + \gamma Q^{\pi}(s', a') | s_t, a_t)$ 

• This is the *Bellman equation* 

• How can we rewrite the value function in more compact form  $Q^{\pi}(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t, a_t)$   $= \mathbb{E}_{s',a'}(r + \gamma Q^{\pi}(s', a') | s_t, a_t)$ 

• This is the *Bellman equation* 

• How can we rewrite the optimal value function  $Q^*(s_t, a_t)$ ?

• Optimal value function  $Q^*(s, a)$  is attained with the optimal policy  $\pi^*$  $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$ 

• After we have found the optimal policy 
$$\pi^*$$
 we do the optimal action  $\pi^* = \operatorname*{argmax}_a Q^*(s, a)$ 

• By expanding the optimal value function

$$Q^{*}(s,a) = r_{t+1} + \gamma \max_{a_{t+1}} Q^{*}(s_{t+1}, a_{t+1})$$
$$Q^{*}(s,a) = \mathbb{E}_{s'}\left(r + \gamma \max_{a'} Q^{*}(s',a') \middle| s,a\right)$$

- The model is learnt from experience
- The model acts as a replacement for the environment
- When planning, the agent can interact with the model
- For instance look ahead search to estimate the future states given actions



## Approaches to Reinforcement Learning

- Policy-based
  - $^{
    m o}$  Learn directly the optimal policy  $\pi^*$
  - $\circ$  The policy  $\pi^*$  obtains the maximum future reward
- o Value-based
  - Learn the optimal value function  $Q^*(s, a)$
  - This value function applies for any policy
- Model-based
  - Build a model for the environment
  - Plan and decide using that model

### How to make RL deep?

### How to make RL deep?

- Use Deep Networks for the
  - Value function
  - Policy
  - Model

#### o Optimize final loss with SGD

### How to make RL deep?



## Deep Reinforcement Learning

- Non-linear function approximator: Deep Networks
- o Input is as raw as possible, e.g. image frame
  - Or perhaps several frames (When needed?)
- Output is the best possible action out of a set of actions for maximizing future reward
- **Important:** no strict need anymore to compute the actual value of the action-value function and take the maximum:  $\arg \max_{\alpha} Q_{\theta}(s, \alpha)$
- Instead, one can teach the network to return directly the optimal action

### Value-based Deep RL

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP REINFORCEMENT LEARNING - 23



### Q-Learning

### • Optimize for Q value function $Q^{\pi}(s_t, a_t) = \mathbb{E}_{s'}(r + \gamma Q^{\pi}(s', a')|s_t, a_t)$

• What does this imply in terms of learning? What do we need?

• Optimize for Q value function  $Q^{\pi}(s_t, a_t) = \mathbb{E}_{s'}(r + \gamma Q^{\pi}(s', a')|s_t, a_t)$ 

- What does this imply in terms of learning? What do we need?
- Target/"Ground truth" *Q* values
- We set  $r + \gamma \max_{a'} Q_t(s', a')$  to be the learning target
- Then we minimize the loss

$$\min\left(r + \gamma \max_{a'} Q_t(s', a') - Q_t(s, a)\right)^2$$

- In the beginning of learning the function Q(s, a) is incorrect
- $\circ$  But the hypothesis is that as training progresses, so does our Q predictions

### Q-Learning

• Value iteration algorithms solve the Bellman equation  $Q_{t+1}(s,a) = \mathbb{E}_{s'}\left(r + \gamma \max_{a'} Q_t(s',a') \middle| s,a\right)$ 

 $\circ$  In the simplest case  $Q_t$  is a table

 $\,\,\circ\,$  To the limit iterative algorithms converge to  $Q^*$ 

• However, a table representation for  $Q_t$  is not always enough. Why/when?

### Q-Learning

• Value iteration algorithms solve the Bellman equation  $Q_{t+1}(s,a) = \mathbb{E}_{s'}\left(r + \gamma \max_{a'} Q_t(s',a') \middle| s,a\right)$ 

- $\circ$  In the simplest case  $Q_t$  is a table
  - $\,\,\circ\,$  To the limit iterative algorithms converge to  $Q^*$
- However, a table representation for  $Q_t$  is not always enough. Why/when?
- When the state space is enormous or near infinite
  - Imagine making a table with all possible images in the rows

• The objective is the mean squared-error in Q-values  $\mathcal{L}(\theta) = \mathbb{E}[\left(\hat{Q} - Q(s, a, \theta)\right)^2]$ 

where  $\hat{Q} = r + \gamma \max_{a'} Q(s', a', \theta)$  is a pre-computed scalar target value.

• So, what does this means for the gradient? How will it look like?

• The objective is the mean squared-error in Q-values  $\mathcal{L}(\theta) = \mathbb{E}[\left(\hat{Q} - Q(s, a, \theta)\right)^2]$ 

where  $\hat{Q} = r + \gamma \max_{a'} Q(s', a', \theta)$  is a pre-computed scalar target value.

- $\circ$  So, what does this means for the gradient? How will it look like?
- The Q-Learning gradient then becomes  $\frac{\partial \mathcal{L}}{\partial \theta} = \mathbb{E}[-2 \cdot \left(\hat{Q} - Q(s, a, \theta)\right) \frac{\partial Q(s, a, \theta)}{\partial \theta}]$
- Backprop to get the gradient
- Optimize end-to-end with SGD

### A system perspective



https://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/

## Algorithmically

- 1. Do a feedforward pass for the current state *s* to predict Q-values for all actions
- 2. Select an action
- a. either by ε-greedy policy
- b. or the one with the best Q value  $a = \max_{a'} Q(s', a', \theta)$
- 3. Perform the action, get a new state s' with reward r. Store the <s, a, r, s'>
- 4. Set Q-value target to  $r + \gamma \max_{a'} Q(s', a', \theta)$ 
  - use the max calculated in step 2
  - For all other action classes, set the Q-value target to the same as originally returned from step 1
  - Makes the error 0 for those action classe  $\rightarrow$  Zero gradient
- 5. Update the weights using backpropagation (in fact, between step 3 and 4 you probably should do experience replay)

## Deep Q Networks on Atari

- End-to-end learning from raw pixels
- o Input: last 4 frames
- Output: 18 joystick positions
- Reward: change of score



### Deep Q Networks on Atari



### Stability in Deep Reinforcement Learning





## Stability problems

# Naively, Q-Learning oscillates or diverges with neural networks Why?

- Naively, Q-Learning oscillates or diverges with neural networks
- o Why?
- Sequential data breaks IID assumption
  - Highly correlated samples break SGD
- However, this is not specific to RL, as we have seen earlier

## Stability problems

# Naively, Q-Learning oscillates or diverges with neural networks Why?

• The learning objective is

$$\mathcal{L}(\theta) = \mathbb{E}\left[\left(r + \gamma \max_{a'} Q(s', a', \theta) - Q(s, a, \theta)\right)^2\right]$$

- The target depends on the Q function also. This means that if we update the current Q function with backprop, the target will also change
- Non-stationarity (BIG problem)
- Plus, we know neural networks are highly non-convex
- $\circ$  Policy changes will change fast even with slight changes in the Q function
  - Policy might oscillate
  - Distribution of data might move from one extreme to another

## Stability problems

# Naively, Q-Learning oscillates or diverges with neural networks Why?

- Not easy to control the scale of the Q values  $\rightarrow$  gradients are unstable Q
- $\circ$  Remember, the Q function is the output of a neural network
- There is no guarantee that the outputs will lie in a certain range
   Unless care is taken
- $_{\odot}$  Naïve Q gradients can be too large, or too small  $\rightarrow$  generally unstable and unreliable
- Where else did we observe a similar behavior?

## Improving stability: Experience replay

- Replay memory/Experience replay
- Store memories  $\langle s, a, r, s' \rangle$
- Train using random stored memories instead of the latest memory transition
- Breaks the temporal dependencies SGD works well if samples are roughly independent
- Learn from all past policies

- Take action  $a_t$  according to  $\varepsilon$ -greedy policy
- Store transition  $(s_t, a_t, r_{t+1}, s_{t+1})$  in replay memory D
- Sample random mini-batch of transitions (s, a, r, s') from D
- Optimize mean squared error using the mini-batch  $\mathcal{L}(\theta) = \mathbb{E}_{(s,a,r,s')\sim D} \left[ \left( r + \gamma \max_{a'} Q(s',a',\theta) - Q(s,a,\theta) \right)^2 \right]$
- Effectively, update your network using random past inputs (experience), not the ones the agent currently sees

## Improving stability: Freeze target Q network

- Instead of having "moving" targets, have two networks
  - One Q-Learning and one Q-Target networks
- $\circ$  Copy the Q network parameters to the target network every K iterations
  - Otherwise, keep the old parameters between iterations
  - The targets come from another (Q-Target) network with slightly older parameters
- $\,\circ\,$  Optimize the mean squared error as before, only now the targets are defined by the "older" Q function

$$\mathcal{L}(\theta) = \mathbb{E}\left[\left(r + \gamma \max_{a'} Q(s', a', \theta_{old}) - Q(s, a, \theta)\right)^2\right]$$

• Avoids oscillations

### Improving stability: Take care of rewards

- Clip rewards to be in the range [-1, +1]
- Or normalize them to lie in a certain, stable range
- Can't tell the difference between large and small rewards

	Q-learning	Q-learning	Q-learning	Q-learning
			+ Replay	+ Replay
		+ Target Q		+ Target Q
Breakout	3	10	241	317
Enduro	29	142	831	1006
River Raid	1453	2868	4103	7447
Seaquest	276	1003	823	2894
Space Invaders	302	373	826	1089

- Skipping frames
  - Saves time and computation
  - Anyways, from one frame to the other there is often very little difference
- $\circ \epsilon$ -greedy behavioral policy with annealed temperature during training
  - $\circ$  Select random action (instead of optimal) with probability  $\varepsilon$
  - In the beginning of training our model is bad, no reason to trust the "optimal" action
- Alternatively: Exploration vs exploitation
  - $\circ$  early stages  $\rightarrow$  strong exploration
  - $\circ$  late stages ightarrow strong exploitation

### Policy-based Deep RL

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP REINFORCEMENT LEARNING - 47



- $\circ$  Problems with modelling the Q-value function
  - Often too expensive → must take into account all possible states, actions →Imagine when having continuous or high-dimensional action spaces
  - $\circ$  Not always good convergence  $\leftarrow$  Oscillations
- Often learning directly a policy  $\pi_{\theta}(a|s)$  that gives the best action without knowing what its expected future reward is easier
- o Also, allows for stochastic policies ← no exploration/exploitation dilemma
- Model optimal action value with a non-linear function approximator  $Q^*(s, a) \approx Q(s, a; w)$



UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES



- Train learning agent for the optimal policy  $\pi_w(a|s)$  given states s and possible actions a
- The policy class can be either deterministic or stochastic

Slides inspired by P. Abbeel

• Use a deep networks as non-linear approximator that finds optimal policy by maximizing  $Q(s, a; \theta)$ 

$$\begin{aligned} \mathcal{L}(w) &= Q(s, a; w) \\ &= \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots | \pi_w(s_t, a_t)] \end{aligned}$$

• If policy is deterministic

$$\frac{\partial \mathcal{L}}{\partial w} = \mathbb{E}\left[\frac{\partial \log \pi(a|s,w)}{\partial w}Q^{\pi}(s,a)\right]$$

• If policy is stochastic 
$$a = \pi(s)$$
  
$$\frac{\partial \mathcal{L}}{\partial w} = \mathbb{E}\left[\frac{\partial Q^{\pi}(s,a)}{\partial a}\frac{\partial a}{\partial w}\right]$$

• To compute gradients use the log-derivative trick (REINFORCE algorithm (Williams, 1992))  $\nabla_{\theta} \log p(x; \theta) = \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)}$ 

• Use a deep networks as non-linear approximator that finds optimal policy by maximizing  $Q(s, a; \theta)$ 

$$\mathcal{L}(w) = Q(s, a; w)$$
  
=  $\mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | \pi_w(s_t, a_t)]$ 

• We can rewrite the score function as  $\mathcal{L}(w) = \mathbb{E}_{\pi_w}[R(\tau)]$ 

Expected future reward given a policy

- Note:  $\mathcal{L}(w)$  is the score, not loss  $\rightarrow$  we want to maximize it  $\frac{\partial \mathcal{L}}{\partial w} = \nabla_{w} \mathbb{E}_{\pi_{w}}[R(\tau)] = \sum_{t} \nabla_{w} \pi_{w}(t) R(t)$
- We use the log-derivative trick

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

• Given stochastic policy  $\pi(a|s,w)$  $\frac{\partial \mathcal{L}(w)}{\partial w} = \mathbb{E}\left[\frac{\partial \log \pi(a|s,w)}{\partial w}Q^{\pi}(s,a)\right]$ 

• Since we have a stochastic/random quantity/variable, i.e., the policy  $(\pi(a|s,w))$ , this means that to compute our loss we must take expectations w.r.t. that RV

• Given deterministic policy 
$$a = \pi(s)$$
  
$$\frac{\partial \mathcal{L}(w)}{\partial w} = \mathbb{E}\left[\frac{\partial Q^{\pi}(s,a)}{\partial a}\frac{\partial a}{\partial w}\right]$$

for continuous a and differentiable Q

• Since the policy is deterministic, we simply need to take gradients of the final objective  $(Q^{\pi})$  w.r.t. the optimized variables (w), taking into account all computational paths (via  $a = \pi(s)$ )

## Asynchronous Advantage Actor-Critic (A3C)

- o Learn both
  - Policy function and
  - Value function
- Multiple agents trained at the same time
- o Global Network consists of
  - ConvNet to model spatial correlations
  - LSTM to model temporal correlations



## A3C Training



## A3C Training

• Estimate Value function

$$V(s,v) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \cdots |s]$$

• Estimate the Q value after *n* steps

$$q_{t} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^{n} V(s_{t+n}, v)$$

# • Update actor by $\frac{\partial \mathcal{L}_{actor}}{\partial w} = \frac{\partial \log \pi(a_t | s_t, w)}{\partial w} (q_t - V(s_t, v))$

### Advantage Estimates

• Discounted future rewards  $R = \gamma(r)$ 

• The model must learn how good the actions taken are

- With Advantage Estimates A = R V(s)
  - Now just how good the actions taken are
  - Also, how much better the actions where than expected

• Model will focus on the areas of the parameter space that it was lacking

## A3C in labyrinth

o End-to-end learning of softmax policy from pixels

- Observations are the raw pixels
- The state is implemented as an LSTM
- Outputs value V(s) and softmax over actions  $\pi(a|s)$
- o Task
  - Collect apples (+1)
  - escape (+10)

#### o <u>Demo</u>



### Model-based Deep RL

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP REINFORCEMENT LEARNING - 61



### Learning models of the environment

• Often quite challenging because of cumulative errors

• Errors in transition models accumulate over trajectories

• Planning trajectories are different from executed trajectories

• At the end of a long trajectory final rewards are wrong

• Can be better if we know the rules

- At least 10<sup>10<sup>48</sup></sup> possible game states
   Chess has 10<sup>120</sup>
- Monte Carlo Tree Search used mostly
  - Start with random moves and evaluate how often they lead to victory
  - Learn the value function to predict the quality of a move
  - Exploration-exploitation trade-off



Tic-Tac-Toe possible game states

- AlphaGo relies on a tree procedure for search
- AlphaGo relies on ConvNets to guide the tree search
- A ConvNet trained to predict human moves achieved 57% accuracy
  - Humans make intuitive moves instead of thinking too far ahead
- For Deep RL we don't want to predict human moves
  - Instead, we want the agent to learn the optimal moves
- Two policy networks (one per side) + One value network
- Value network trained on 30 million positions while policy networks play



### AlphaGo

- Both humans and Deep RL agents play better end games
  - Maybe a fundamental cause?
- In the end the value of a state is computed equally from Monte Carlo simulation and the value network output
  - Combining intuitive play and thinking ahead
- Where is the catch?



- Both humans and Deep RL agents play better end games
  - Maybe a fundamental cause?
- In the end the value of a state is computed equally from Monte Carlo simulation and the value network output
  - Combining intuitive play and thinking ahead
- Where is the catch?
- State is not the pixels but positions
- Also, the game states and actions are highly discrete



- What is allowed
  - 1 A4 page with whatever you want on it (handwritten or printed)
- What is not allowed
  - Everything else (internet, phones, messaging, etc)
- And if you are interested in the coming months in research in any of the topics discussed in the course
  - especially in temporal sequences, deep dynamics, video, causality, generative models and/or oncology (together with the Netherlands Cancer Insitute)

Drop me a line (<u>egavves@uva.nl</u>)

### Summary

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP REINFORCEMENT LEARNING - 68

- Reinforcement Learning
- o Q-Learning
- Deep Q-Learning
- Policy-based Deep RL
- Model-based Deep RL
- Making Deep RL stable