

# Lecture 12. Intro to Reinforcement Learning



Good afternoon, dear students. Today we are going to talk about reinforcement learning, it is a slightly separate topic that is not directly related to deep learning, but is included in machine learning of course. It is quite possible to read a semester course on it, but we will limit ourselves to today's one lecture.

## Lecture 12. Intro to Reinforcement Learning

### └ Five-minutes block

- Describe (or draw) the general architecture of the autoencoder
- What is the main advantage of the self-supervised learning method?
- Please, write down GAN loss functions

Ok, and as usual now we have five-minutes block. The last one in our course. So prepare yourself and send right answers to me in five minutes.

# Lecture 12. Intro to Reinforcement Learning

## └ Problem Statement

### Problem Statement

- Up to this point, we either restored the function from the training set  $(X, Y)$  (supervised learning) or searched for the structure in the set of objects  $X$  (unsupervised learning).
- But how does learning work in real world? Usually we do some action and get the result, gradually learning

Great, let's go on and let me remind you, that up to this point, we found the appropriate function  $f$ , based on the training set. And it is called supervised learning. Or also we found the inner structure in the set of objects, you know, texts or images, and it is called unsupervised learning, when we haven't answers  $y$ .



## Lecture 12. Intro to Reinforcement Learning

2022-11-26

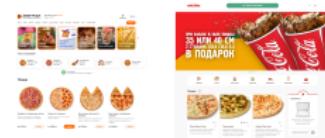
- └ Real world

Ok, so what we have in real world? Usually we do some action and get the result, and so we are learning gradually.

## Lecture 12. Intro to Reinforcement Learning

### └ One more motivation

One more motivation  
For example, we need to select the main page of a pizzeria website to attract the clients.



And another one important example. Suppose you need to choose the most attractive design for the main page of a pizzeria website, but you don't want to lose a lot of customers, time, money on experiments.

## Lecture 12. Intro to Reinforcement Learning

### └ What approaches are there?

- A/B testing — many users see potentially bad options
- multi-armed bandits — a special case of reinforcement learning

In A/B testing you divide your customers on two groups (A and B), show them different versions of main page and then measure the results, the conversions and make your conclusions. And what is multi-armed bandits?

## Lecture 12. Intro to Reinforcement Learning

Bernoulli's one-armed bandit



Probability of winning  $\theta = 0.05$

### └ Bernoulli's one-armed bandit

Let's begin from the one-armed bandit. So you pull the handle of a slot machine and you can win with some probability theta.

# Lecture 12. Intro to Reinforcement Learning

## └ Multi-Armed Bandits



$\theta_1 = 0.02$        $\theta_2 = 0.01(\min)$        $\theta_3 = 0.05$        $\theta_4 = 0.1(\max)$

We do not know the true probabilities, but we want to come up with a strategy that maximizes the payoff (reward)

In multi-armed bandits model you have several handles with some winning probabilities theta one, two and so on.

# Lecture 12. Intro to Reinforcement Learning

## └ Mathematical statement of the problem

### Mathematical statement of the problem

Given possible actions

$x_1, \dots, x_K$

At the next iteration  $t$ , for each action  $x_i^t$  performed, we get the answer

$$y_i^t \sim q(y|x_i^t, \Theta),$$

which brings us a reward

$$r_i^t = r(y_i^t)$$

There is an optimal action  $x_{i^*}$  (sometimes  $x_Q$ )

$$\forall i : E(r_Q^i) \geq E(r_i^t)$$

#### Question

How to evaluate different strategies?

Ok, now let's move on, to mathematical statement of the problem. ...

## Lecture 12. Intro to Reinforcement Learning

### └ Multi-Armed Bandits possible applications

## Areas:

- advertising banners
- recommendations (goods, music, movies etc.)
- slot machines in the casino

## Approaches:

- Thompson sampling
- Upper Confidence Bound (UCB)
- $\epsilon$ -greedy strategy

## Note

The main difference from more general reinforcement learning approach is that there is no environment state in multi-armed bandits

Ok, and what is the possible applications of multi-armed bandits? And there are models of implementing and optimizing the multi-armed bandits, but we won't discuss them in detail here.

# Lecture 12. Intro to Reinforcement Learning

## └ Multi-Armed Bandits

### Multi-Armed Bandits

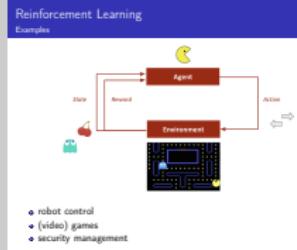
Disadvantage  
Do not take into account delayed effects  
For example, the effect of clickbait in advertising



So you may increase your metrics for some short period of time by clickbait, but in the long run of course you will lose clients and money consequently.

# Lecture 12. Intro to Reinforcement Learning

## └ Reinforcement Learning



Let's go on to the Reinforcement Learning model.

# Lecture 12. Intro to Reinforcement Learning

## └ Agent Model in a Changing Environment

**Definitions:**

- state  $s \in S$
- agent action  $a \in A$
- reward  $r \in R$
- state transition dynamics  
 $P(s_{t+1}|s_t, a_t, \dots, s_{t-j}, a_{t-j}, \dots, s_0, a_0)$
- win function

**Task:**

$$\pi(a|s) : \mathbb{E}_\pi[R] \rightarrow \max$$

- $r_t = r(s_t, a_t, \dots, s_0, a_0)$
- total reward  $R = \sum_t r_t$
- agent policy  $\pi(a|s)$

So, more formally you have states, actions and rewards. And your task is to find optimal policy to maximize the total reward.

# Lecture 12. Intro to Reinforcement Learning

## └ Cross-entropy method. Algorithm

Cross-entropy method. Algorithm

Trajectory  $\equiv [s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_T]$ Initialize strategy model  $\pi(a|s)$ 

Repeat:

- play  $N$  sessions
- choose the best  $K$  of them and take their trajectories
- adjust  $\pi(a|s)$  so that  $s$  is able to maximize the probabilities of actions from the best trajectories

The first simple algorithm to find appropriate policy is called cross-entropy method. Here you have some initialization, maybe random, as usual, then you play  $N$  sessions, choose best of them by your reward and update your strategy.



## Lecture 12. Intro to Reinforcement Learning

2022-11-26

└ There is a problem..

So the problem is, that we can not even store in memory the whole table of states.

# Lecture 12. Intro to Reinforcement Learning

## └ Approximate cross-entropy methods

### Possible solutions:

- split the state space into sections and treat them as states
- get probabilities from  $\pi_0(a|s)$  machine learning model: linear model, neural network, random forest
- often these probabilities need to be further specified later

So people invented Approximate cross-entropy methods.

## Lecture 12. Intro to Reinforcement Learning

### └ Disadvantages of the cross-entropy method

... So we need to correct these shortcomings.

- it is unstable for small samples
- in the case of a non-deterministic environment, it chooses lucky cases (random played in favor of the agent)
- it is focusing on behavior in simple states
- it ignores a lot of information
- there are tasks in which the end never comes (stock market game)

# Lecture 12. Intro to Reinforcement Learning

## └ Markov Decision Process

### Markov Decision Process

#### Definitions:

- state  $s \in S$
- agent action  $a \in A$
- reward  $r \in R$
- transition dynamics

$$P(a_{t+1}|s_t, a_t, \dots, s_{t-1}, a_{t-1}, \dots, s_0, a_0) = P(s_{t+1}|s_t, a_t)$$

(Markov's assumption)

- agent policy

$$\pi(a|s)$$

- value function

$$r_t = r(s_t, a_t)$$

- (Markov's assumption)

#### Task:

$$\pi(a|s) : \mathbb{E}_\pi [R] \rightarrow \max$$

Here we have Markov's assumption, which is actually not so strong, as we could expect. For example it is true in chess and Starcraft also. Only your current state determines everything, and not the history of the game.

# Lecture 12. Intro to Reinforcement Learning

## └ TD-training

We arbitrarily initialize the function  $V(s)$  and the strategy  $\pi$ . Then we repeat:

- initialize  $s$
- for each agent step:
  - select  $a$  by strategy  $\pi$
  - do action  $a$ , get result  $r$  and next state  $s'$
  - update function  $V(s)$  by formula

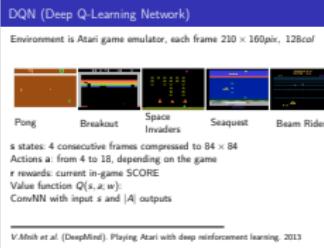
$$V(s) = V(s) + \alpha(r + \gamma V(s') - V(s))$$

- go to the next step by assigning  $s := s'$

The amazing fact is, that temporal difference training works.

# Lecture 12. Intro to Reinforcement Learning

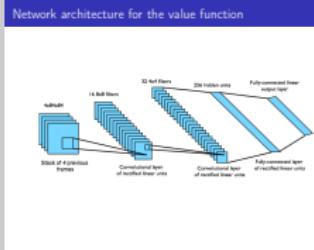
## └ DQN (Deep Q-Learning Network)



Ok, and the last reinforcement learning model in our lecture is the Deep Q-Learning Network.

## Lecture 12. Intro to Reinforcement Learning

### └ Network architecture for the value function



Here you see actual CNN structure for Atari game — quite simple, isn't it? And I think you know what happened next :)

## Temporary page!

$\text{\LaTeX}$  was unable to guess the total number of pages correctly. A  
was some unprocessed data that should have been added to the  
this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page  
away, because  $\text{\LaTeX}$  now knows how many pages to expect for the  
document.