# RL seminar #5: DDPG & Bandits

...

MIPT, Deep learning lab & iPavlov.ai

2 Dec 2017

# Outline

# Table of Contents

# Assignments

## Coding (Programming assignment #2)
Deadline: ... Dec 2017

## Quiz
Quiz N4 to be issued.

## Questions
...

## Course
...

# Table of Contents

# DPG: Overview

### Task
Continuous actions in high dimensions

### Motivation

- Deterministic policy gradient can be estimated much more efficiently than the usual stochastic policy gradient.
- Deterministic policy gradient is simple and natural

from David Silver's lecture on Deterministic Policy Gradient Algorithms

# Nontation

- Continuous state space $s \in \mathbb{R}^l$
- Continuous action space $a \in A^m$
- Parameter vector $\theta \in \mathbb{R}^n$
- Deterministic policy $a = \mu_\theta(s)$
- Stochastic policy $\pi_\theta(s, a) = p(a|s; \theta)$

Find parameters $\theta$ optimizing performance of policy

from David Silver's lecture on Deterministic Policy Gradient Algorithms

# Policy Gradient theorem

Stochastic case:

$$J(\theta) = \mathbb{E}_{s \sim p^\pi(s)} \Big[ \int_a \pi_\theta(s, a) R(s, a) da \Big]$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{s,a} \Big[ \nabla_\theta \log \pi_\theta(s, a) Q^\pi(s, a) \Big]$$

Deterministic case:

$$\nabla_\theta J(\theta) = \mathbb{E}_s \Big[ \nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a)|_{a=\mu_\theta(s)} \Big]$$

Update policy in the direction that most improves Q

**Algorithm 1** DDPG algorithm

1: Randomly initialize critic n-k $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights $\theta^Q$ and $\theta^\mu$.
2: Initialize target network $Q'$ and $\mu'$ with weights $\theta^{Q'} \leftarrow \theta^Q$, $\theta^{\mu'} \leftarrow \theta^\mu$
3: Initialize replay buffer $R$
4: **for** episode $= 1$, M **do**
5:     Initialize a random process $\mathcal{N}$ for action exploration
6:     Receive initial observation state $s_1$
7:     **for** t $= 1$, T **do**
8:         Select action $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$ by current policy and exploration noise
9:         Execute action $a_t$ and observe reward $r_t$ and observe new state $s_{t+1}$
10:         Store transition $(s_t, a_t, r_t, s_{t+1})$ in $R$
11:         Sample a random minibatch of $N$ transitions $(s_i, a_i, r_i, s_{i+1})$ from $R$
12:         Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$
13:         Update critic by minimizing the loss: $L = \frac{1}{N}\sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$
14:         Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N}\sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

15:         Update the target networks:

$$\theta^{Q'} \leftarrow \tau\theta^Q + (1-\tau)\theta^{Q'}, \theta^{\mu'} \leftarrow \tau\theta^\mu + (1-\tau)\theta^{\mu'}$$

16:     **end for**
17: **end for**

# Table of Contents

# DDPG code

- https://github.com/fgvbrt/nips_rl
- https://github.com/Scitator/Run-Skeleton-Run
- https://github.com/vy007vikas/PyTorch-ActorCriticRL

# Table of Contents

# MA bandits: overview

### Task definition
...

# Upper confidence bound

Optimal exploration

...

# Applications

...

# Table of Contents

# Questions

1. L12 – 'Closer look at backward pass'.
2. L12 – 'Benefits of soft optimality'.
3. L12 – 'More efficient sample-based updates / Importance sampling'.

# Table of Contents

# Next steps

## Plan for the week

- Quiz N4 (?)
- Home assignment N3 (?)

## Reading

Lectures 15-17 of CS294.
Please, post your questions about lectures in google doc:
`https://goo.gl/qN6jmJ`

## Rating

ref: `goo.gl/yxqhBg`
HW scores are coming (week 4-10 Dec)