

Off-policy deep reinforcement learning algorithms (seminar)

Oleksii Hrinchuk

Deep Q-Network (DQN)

1. Initialize Q-Network Q_θ , set initial exploration rate $\varepsilon = 1$. Run epsilon-greedy policy w.r.t. Q_θ (or random policy which is the same) for a number of steps to initialize replay buffer \mathcal{D} .
2. (Train loop)
 - Make a step in the environment with policy $\pi = \varepsilon\text{-greedy}(Q_\theta)$
 - Put transition $\mathbf{d}_t = \{\mathbf{s}_t, \mathbf{a}_t, \mathbf{r}_t, \mathbf{s}_{t+1}\}$ into replay buffer \mathcal{D}
 - If step % q_update_frequency == 0:
 - Sample mini-batch of transitions $\{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_B\} \sim \mathcal{D}$
 - Update Q-Network weights $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}_\theta$
 - If step % target_update_frequency == 0:
 - Update target network weights $\tilde{\theta} \leftarrow \theta$
 - Decay exploration rate $\varepsilon \leftarrow \varepsilon - \Delta\varepsilon$
3. Return greedy policy w.r.t. the learned Q-function Q_θ

Improvements to DQN

Double DQN (DDQN)

- Max operator in TD-targets leads to overoptimistic estimates of Q-function.

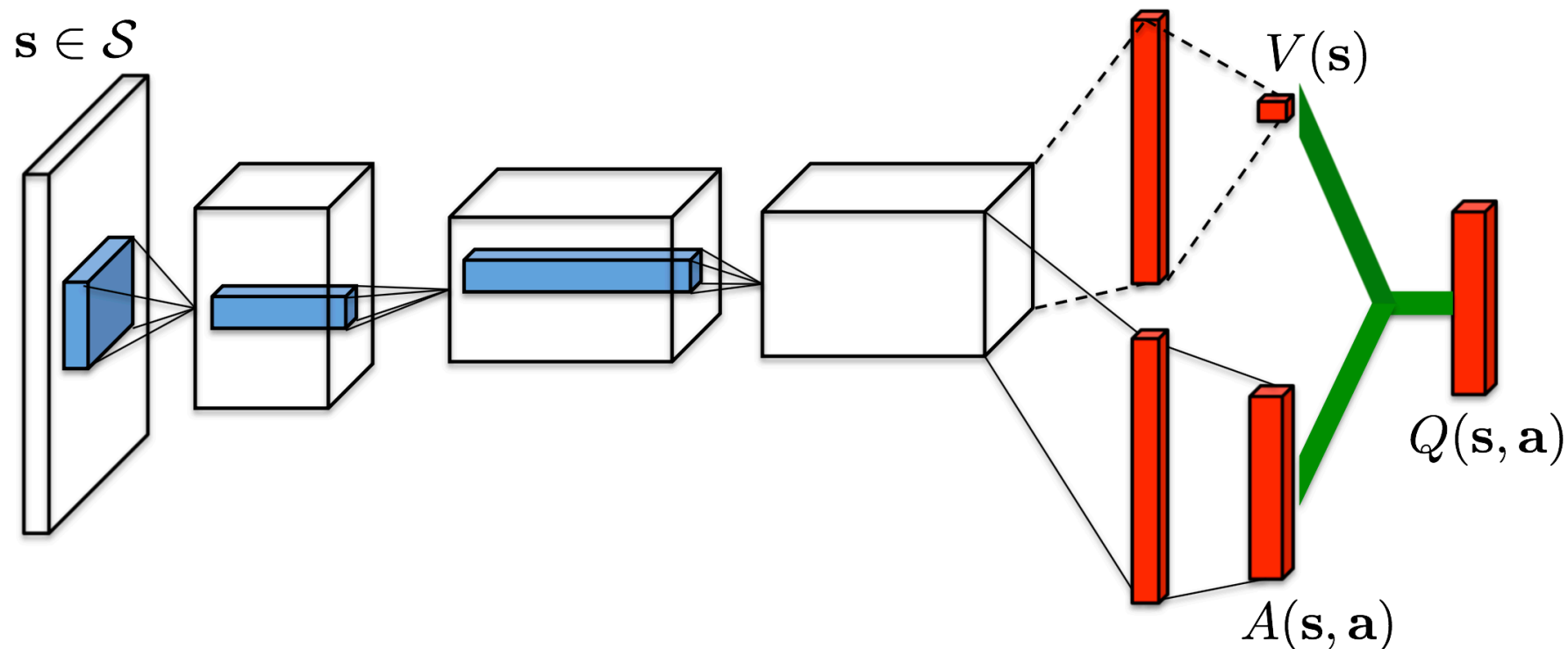
$$\begin{aligned} y_t &= \mathbf{r}_t + \gamma \max_{\mathbf{a}' \in \mathcal{A}} Q_{\tilde{\theta}}(\mathbf{s}_{t+1}, \mathbf{a}') \\ &= \mathbf{r}_t + \gamma Q_{\tilde{\theta}}(\mathbf{s}_{t+1}, \arg \max_{\mathbf{a}' \in \mathcal{A}} Q_{\tilde{\theta}}(\mathbf{s}_{t+1}, \mathbf{a}')) \end{aligned}$$

- **Q-Network** is used to **select** actions, target network is used to evaluate actions

$$y_t = \mathbf{r}_t + \gamma Q_{\tilde{\theta}}(\mathbf{s}_{t+1}, \arg \max_{\mathbf{a}' \in \mathcal{A}} Q_{\theta}(\mathbf{s}_{t+1}, \mathbf{a}'))$$

Dueling DDQN

- Split Q-function into two channels
 - Action independent value function $V(s)$
 - Action dependent advantage function $A(s, a)$



$$Q(s, a) = V(s) + A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} A(s, a)$$

Prioritized experience replay

- Transitions with high TD-error are sampled more often

$$\mathbf{d}_t = \{\mathbf{s}_t, \mathbf{a}_t, \mathbf{r}_t, \mathbf{s}_{t+1}\}$$

$$\delta_t = \mathbf{r}_t + \gamma \max_{\mathbf{a}' \in \mathcal{A}} Q(\mathbf{s}_{t+1}, \mathbf{a}') - Q(\mathbf{s}_t, \mathbf{a}_t)$$

- Define prioritization
 - proportional $p_t = |\delta_t| + \epsilon$
 - rank-based $p_t = 1/\text{rank}(t)$
- Sample transitions according to probability

$$P(\mathbf{d}_t) = \frac{p_t^\alpha}{\sum_{t'=1}^{|\mathcal{D}|} p_{t'}^\alpha}$$

Boltzmann exploration

- Epsilon-greedy exploration is bad in situations when picking random action can be disastrous
- Idea: pick exploratory actions based on Q-function:
 - good actions have high probability
 - bad actions have really low probability

$$\pi_{\mathcal{B}}(\mathbf{a}|\mathbf{s}) = \text{softmax} \left(\frac{Q(\mathbf{s}, \mathbf{a})}{\tau} \right) = \frac{\exp (Q(\mathbf{s}, \mathbf{a})/\tau)}{\sum_{\mathbf{a}' \in \mathcal{A}} \exp (Q(\mathbf{s}, \mathbf{a}')/\tau)}$$

$$\pi_{\mathcal{B}}(\mathbf{a}|\mathbf{s}) \rightarrow \arg \max_{\mathbf{a} \in \mathcal{A}} Q(\mathbf{s}, \mathbf{a}), \quad \tau \rightarrow 0$$

Additional off-policy references

- A Distributional Perspective on Reinforcement Learning. Bellemare et al., ICML (2017)
- Distributional Reinforcement Learning with Quantile Regression. Dabney et al., AAAI (2018)
- Rainbow: Combining Improvements in Deep Reinforcement Learning. Hessel et al., AAAI (2018)
- Self-Imitation Learning. Oh et al., ICML (2018)

Additional on-policy references

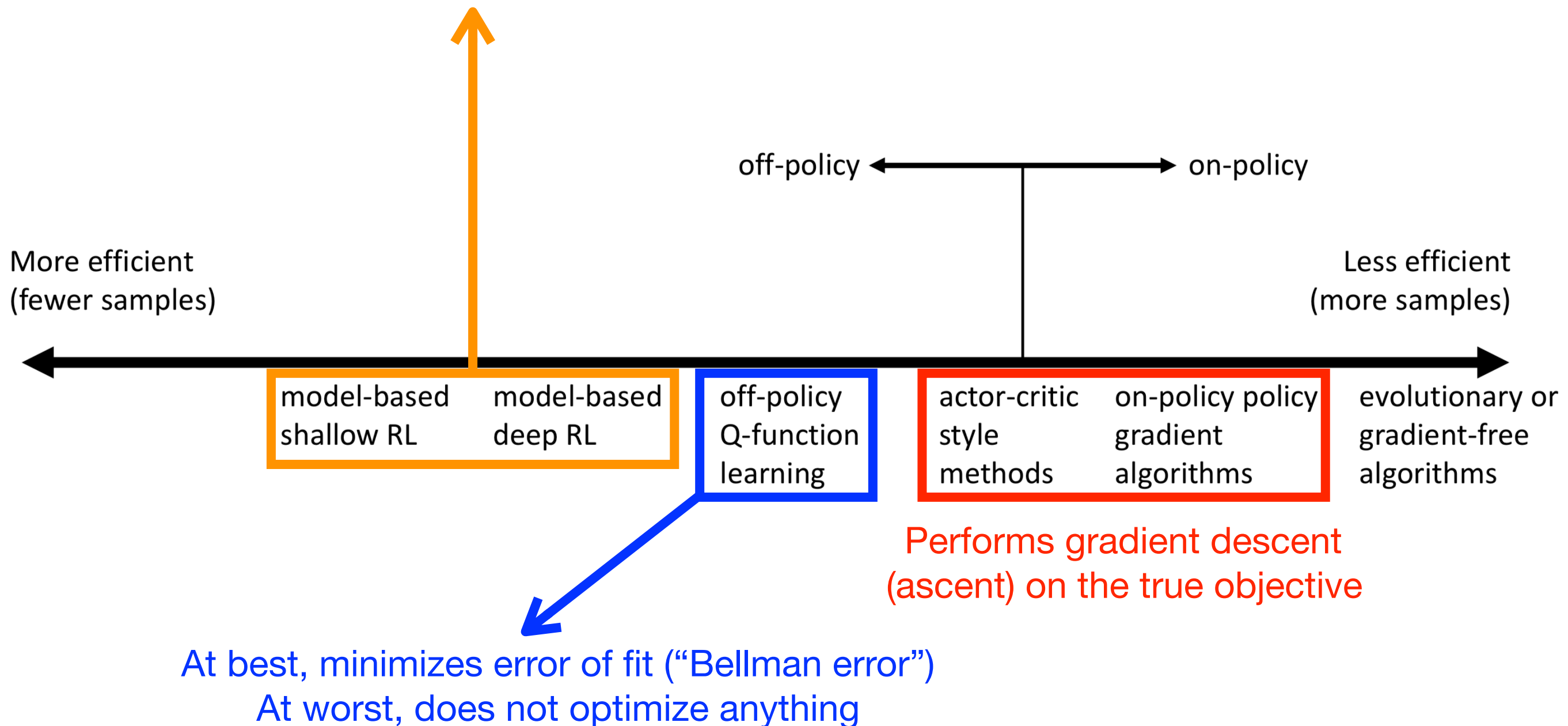
- Asynchronous Methods for Deep Reinforcement Learning. Mnih et al., ICML (2016)
- Trust Region Policy Optimization. Schulman et al., ICML (2015)
- Proximal Policy Optimization Algorithms. Schulman et al. (2017)

Why so many RL algorithms?

- Different tradeoffs
 - Sample efficiency
 - Stability and ease of use
- Different assumptions
 - Stochastic or deterministic?
 - Continuous or discrete?
 - Episodic or infinite horizon?
- Different things are easy or hard in different settings
 - Easier to represent the policy?
 - Easier to represent the model?

Sample efficiency

Good convergence properties
No guarantee that better model
leads to better policy



Task

- Train **the best agent** to play 10x10 Snake game
- You are allowed to adopt third-party open source code, but you **have to** cite the source
- You are not allowed to change the file [snake.py](#)
- You should submit
 - Plot (reward / number of training transitions)
 - Best agent score (mean and std of 1000 runs)
 - Source code which **reproduces** best agent results
 - Short report on what you try, what worked best (worst)
- Deadline: Sunday, Sep 30, 23:59