

RL: model은 보충해 MDP를 꼭지 algo

DRL: Deep neural net + RL.

Lecture H1. DRL - Value-based agent 1

Sim, Min Kyu, Ph.D., mksim@seoultech.ac.kr



서울과학기술대학교 데이터사이언스학과

1 I. Deep Reinforcement Learning (DRL)

2 II. Value-based agent

3 III. Deep Q-learning (DQN)

I. Deep Reinforcement Learning (DRL)

Agent's space to learn

- RL agent must explore its state and action space. The space can be

1. Discrete, finite and small set.
 - Skier's problem has only 8 states and 2 actions.
2. Discrete, finite, but large set.
 - Chess has ~~10^47~~ states. 10^{43}
 - Go has ~~10^{170}~~ states. 10^{170}
3. Discrete and infinite set.
4. Continuous set (thus infinite).
 - Fine-control problems such as pendulum problem
 - Require policy-based approach. ✓
 - May discretize it into discrete set.

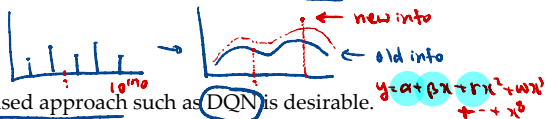


$\{0, 0.2, 0.4, 0.6, 0.8, 1\}$

Suitable approaches

$$q(s,a) \leftarrow q(s,a) + \alpha (R + \gamma V(s) - q(s,a))$$

- For 1, tabular approach is possible. For 2, 3, and 4, states are too many.
 - The storage space itself is burden, if the state has size of 10^{170} . 10^{170}
 - In tabular approach, (s, a, r, s') updates only a single state's value. Thus, processing time can be inhibitive.
 - For large or continuous space, functional approximation approach is necessary to search for all spaces.



- For large space,
 - Part II** • If the space is discrete, value-based approach such as DQN is desirable.
 - Part I.** • If the space is continuous, then policy-based approach is required.

- Generalization plays an important role in the acquisition of intelligence by humans and machines
 - Remind that the extension from an action to a policy is a form of generalization.
 - Tabular form to functional form has also an implication of generalization.

II. Value-based agent

Connecting deep learning with value functions.

- Suppose the a function is parameterized with θ . Then, an approximation function's (e.g. neural net) loss function is typically written as

$$L(\theta) = (\underline{y_{true}} - \underline{y_{\theta}})^2,$$

↑ 02 parameterized 된 값이 주어져.

where y_{θ} is an estimated response using the function parameterized by θ .

- If the function approximates state-value function $V(s)$, then

$$L(\theta) = (\underline{V_{true}(s)} - \underline{V_{\theta}(s)})^2 \quad w. \theta.$$

- If the function approximates state-value function $q(s, a)$, then

$$L(\theta) = (q_{true}(s, a) - q_{\theta}(s, a))^2$$

min $L(\theta)$
 θ

- Find something missing here? How to aggregate the argument s and a so that $L(\theta)$ is a single value?

- Possible approach is to aggregate over all s .
 - Is this possible? Not realistic in current simulated approach because we have to create a stochastic path that goes through all states and only once for each state.
 - Is this desirable? Not really, we should want to focus more on frequent states or frequent trajectory.
- If you agree with the above two points, then see the following alternative.

$$L(\theta) = \mathbb{E}_{\pi} \left[(V_{true}(s) - V_{\theta}(s))^2 \right]$$

or

$$L(\theta) = \mathbb{E}_{\pi} \left[(q_{true}(s, a) - q_{\theta}(s, a))^2 \right]$$

- In this approach, the update goes through a stochastic path through π .
- During the agent's learning process, the current π is the best at the moment during the learning. Thus, the visited states following π is most important concern. This learning process is natural in a sense of the Einstein's quote, "Learn from yesterday, live for today."

- ↑
02 parametrized g g_{θ}
DQN
"g-net"

learning rate

III. Deep Q-learning (DQN)

Main algorithm

- $q(s, a)$
 - $V(s)$ is value, and $\pi(s)$ is policy. Then, what is really $q(s, a)$?
 - Remind that having an accurate $q(s, a)$ leads to π^* . In this sense, $q(s, a)$ can be regarded as *implicit policy*.
 - Thus, $q(s, a)$ is value and also implicit policy.
- In `pol_eval_Q()`, Q-learning updates $q(s, a)$ in the following way.

$$✓ \quad q(s, a) \leftarrow q(s, a) + \alpha(r_t + \gamma \max_{a' \in \mathcal{A}} q(s', a') - q(s, a)), \quad \forall s, a$$

Handwritten: $Q - tq_t$ (with a bracket under the term being added)

- In `pol_eval_DQN()`, DQN updates $q(s, a)$ in the following way.

$$✓ \quad q_\theta(s, a) \leftarrow q_\theta(s, a) + \alpha(r_t + \gamma \max_{a' \in \mathcal{A}} q_\theta(s', a') - q_\theta(s, a)), \quad \forall s, a$$

Handwritten: Deep Q - tqt (with a bracket under the term being added)

Experienced replay

~~0~~~~0~~000000000000

- In “Atari 13” paper, the authors introduced the notion of *replay buffer*.
 - This is related to what samples should be used for deep learning updates or so-called model.fit(). *sample*.
 - *Replay buffer* stores the recent records of N episodes, shuffle them, and select only some small portion of data to update Q-network. Data structure type of deque is popularly used.
 - This reduces correlation of the data. Is this desirable?
 - Is experienced replay desirable for stationary or non-stationary environment?
 - Experienced replay works for *off-policy algorithm*, because the study material data is not only stemmed from current policy.
- ✓ • The implementation is not complex.

- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M. (2013). Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602.

Target network

- Default setting of DQN has a single neural net for $q(s, a)$. This network is used to setting a Q-target (q_tgt) and the value of Q-target (q_tgt) is used to update the deep Q-network ($q_net()$).
- This may deteriorate the stability, resulting in a high variance during the learning.
- The idea is to prepare two separate networks for $q(s, a)$. The two networks are repetitively and interchangeably used for prediction and updating.
- The implementation is not complex.

"Learn from yesterday, live for today, hope for tomorrow. The important thing is not to stop questioning. - Albert Einstein"