Mobile
Robotics,
Reinforce-
ment
Learning and
Deep Rein-
forcement
Learning for
Mobile
Robots

# Mobile Robotics, Reinforcement Learning and Deep Reinforcement Learning for Mobile Robots

Material based on Reinforcement Learning: an Introduction, 2nd Edition [sect. 5.1-5.3, 6.1-6.3, 6.5], Reinforcement Learning course offered by Prof. Pascal Poupart at Univ. of Waterloo

# Summary

- Introduction to Reinforcement Learning
- Deep Reinforcement Learning
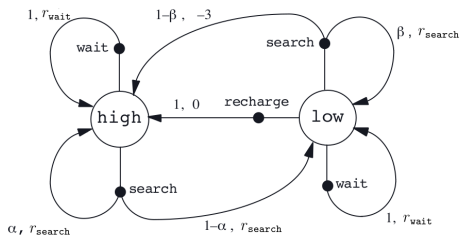- Deep Q Network
- DQN for Mapless Navigation

Guide an MDP without knowing the dynamics

- do not know which states are good/bad (no $R(s, a)$)
- do not know where actions will lead us (no $T(s, a, s')$)
- hence we must try out actions/states and collect the reward

Planning

Learning

- <u>Model-Based</u> methods try to **learn a model**
    - \+ avoid repeating bad states/actions
    - \+ fewer execution steps
    - \+ efficient use of data
- <u>Model-Free</u> methods try to **learn Q-function and policy** directly
    - \+ simplicity, no need to build and use a model
    - \+ no bias in model design

# Q-Learning: pseudo-code

---

**Algorithm 1** Tabular Q-Learning

---

1:  Initialize $Q(s, a)$ arbitrarily
2:  Initialize $s$ {observe current state}
3:  **loop**
4:      Select and execute action $a$
5:      Observe new state $s'$ receive immediate reward $r$
6:      $Q(s, a) = Q(s, a) + \alpha[r + \gamma \max'_a Q(s', a') - Q(s, a)]$
7:      update state $s \leftarrow s'$
8:  **end loop**

---

$\diamondsuit$  $\epsilon$-greedy: choose best action most of the time, but every once in a while (with probability $\epsilon$) choose randomly amongst all action (with equal probabiliy)

$\diamond$ For many real world domains we can not explicitly represent key functions for RL ($\pi(s)$, $V(s)$, $Q(s, a)$)

$\diamond$ We can try to approximate them

- Linear approximation
- Neural Network approximation
    - Deep RL

$\diamond$ Deep Q Network approximates $Q(s, a)$ with a DNN

# Gradient Q-Learning

$\diamond$ approximate $Q(s, a)$ with a parametrized function $Q_{\mathbf{w}}(s, a)$

$\diamond$ Minimize squared error between estimate and target

- Estimate $Q_{\mathbf{w}}(s, a)$
- Target: $r(s, a, s') + \gamma \max_{a'} Q_{\overline{\mathbf{w}}}(s', a')$

$\diamond$ squared error:

$$Err(\mathbf{w}) = (Q_{\mathbf{w}}(s, a) - r(s, a, s') - \gamma \max_{a'} Q_{\overline{\mathbf{w}}}(s', a'))^2$$

$\diamond$ gradient:

$$\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = 2(Q_{\mathbf{w}}(s, a) - r(s, a, s') - \gamma \max_{a'} Q_{\overline{\mathbf{w}}}(s', a'))\frac{\partial Q_{\mathbf{w}}(s,a)}{\partial \mathbf{w}}$$

(Scalar 2 is a constant factor and not important for update)

# Gradient Q-Learning Algorithm

---

**Algorithm 2** Gradient Q-Learning

---

1: Initialize weights $\mathbf{w}$ randomly in $[-1, 1]$

2: Initialize $s$ {observe current state}

3: **loop**

4:     Select and execute action $a$

5:     Observe new state $s'$ receive immediate reward $r$

6:     $\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = (Q_{\mathbf{w}}(s, a) - r - \gamma \max_{a'} Q_{\mathbf{w}}(s', a')) \frac{\partial Q_{\mathbf{w}}(s,a)}{\partial \mathbf{w}}$

7:     update weights $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}}$

8:     update state $s \leftarrow s'$

9: **end loop**

---

Mobile
Robotics,
Reinforce-
ment
Learning and
Deep Rein-
forcement
Learning for
Mobile
Robots

$\diamondsuit$ Non-linear approximation of Q(s,a), $Q(s, a) \approx g(\mathbf{x}; \mathbf{w})$

$\diamondsuit$ **gradient Q-Learning may not converge**

$\diamondsuit$ Issue:

- we update the weights to reduce error for a specific experience (i.e., a specific $(s, a)$) **but** by changing the weights we may end up changing the $Q(s, a)$ potentially everywhere.

◇ Two main approaches to **mitigate** divergence:

1 experience replay
2 use two different networks
   - Q-network
   - Target network

# Experience replay

Mobile Robotics, Reinforcement Learning and Deep Reinforcement Learning for Mobile Robots

$\Diamond$ Store previous experiences (i.e., $(s, a, s', r)$) and use them at each step

- Store previous $(s, a, s', r)$ in a dedicated memory buffer
- At each step sample a mini-batch from this buffer and use the mini-batch to update the weights

$\Diamond$ Benefits

1. reduces correlation between successive samples (increase stability)
2. reduces number of interaction with the environment (increase data efficiency)

# Target Network

$\diamondsuit$ Maintain a separate **target** network and update this network periodically (not with every experience)

- Q-network $Q_{\boldsymbol{w}}(s, a)$
- Target network $Q_{\overline{\boldsymbol{w}}}(s, a)$

$\diamondsuit$ repeat for every $(s, a, s', r)$ in the mini-batch update the Q-network

- $\boldsymbol{w} \leftarrow \boldsymbol{w} - \alpha_t (Q_{\boldsymbol{w}}(s, a) - r - \gamma \max_{a'} Q_{\overline{\boldsymbol{w}}}(s', a')) \frac{\partial Q_{\boldsymbol{w}}(s, a)}{\partial \boldsymbol{w}}$

$\diamondsuit$ update the target network

- $\overline{\boldsymbol{w}} \leftarrow \boldsymbol{w}$

# Deep Q Network

Mobile
Robotics,
Reinforce-
ment
Learning and
Deep Rein-
forcement
Learning for
Mobile
Robots

◇ Human-level control through deep reinforcement learning (V. Mnih et al., Nature 2015)

◇ Gradient Q-Learning

- Deep neural networks to approximate Q(s,a)
- Experience Replay and Target network



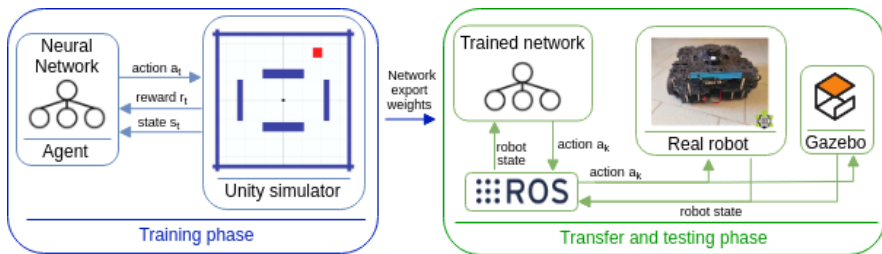◇ above human-level performance in many Atari video games

**Algorithm 3** DQN

1: Initialize weights $\mathbf{w}$ and $\overline{\mathbf{w}}$ randomly in $[-1, 1]$
2: Initialize $s$ {observe current state}
3: **loop**
4:      Select and execute action $a$
5:      Observe new state $s'$ receive immediate reward $r$
6:      Add $(s, a, s', r)$ to experience buffer
7:      Sample mini-batch $MB$ of experiences from buffer
8:      **for** $(\hat{s}, \hat{a}, \hat{s}', \hat{r}) \in MB$ **do**
9:          $\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = (Q_{\mathbf{w}}(\hat{s}, \hat{a}) - \hat{r} - \gamma \max_{\hat{a}'} Q_{\overline{\mathbf{w}}}(\hat{s}', \hat{a}')) \frac{\partial Q_{\mathbf{w}}(\hat{s}, \hat{a})}{\partial \mathbf{w}}$
10:          update weights $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}}$
11:      **end for**
12:      update state $s \leftarrow s'$
13:      every c steps, update target: $\overline{\mathbf{w}} \leftarrow \mathbf{w}$
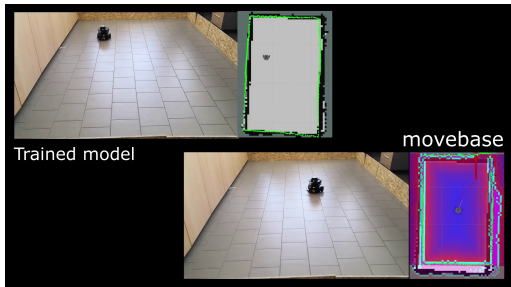14: **end loop**

Acting in the real environment is difficult/dangerous

- train in a synthetic environment

◇ Most DRL approaches considers continuous action space
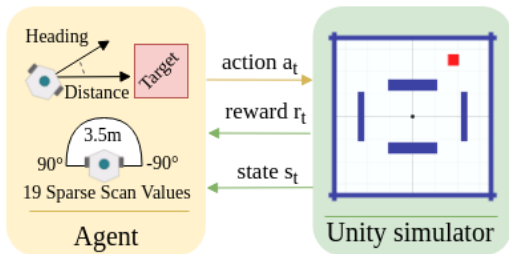- Proximal Policy Optimization (PPO)
- Policy Gradient (DDPG)

◇ Continuous approaches are time-consuming compared to discrete action space (e.g., DQN)

◇ **Mapless Navigation**: Navigate in the environment, avoiding obstacles, without a map

- laser scans and target heading as input
- angular velocities as output
- dense reward (i.e., distance from target)



Trained model

movebase

Proper discretization of the action space and the use of DQN can result in a significantly shorter training time (from 20 hours to 1 hour), maintaining comparable performance.

Mobile
Robotics,
Reinforce-
ment
Learning and
Deep Rein-
forcement
Learning for
Mobile
Robots

# Deep RL: current trends

◇ Formal verification of DRL models
- ensure the learned model respect safety properties

◇ Transfer of Learning/Curricula learning
- ToL: train the model in an environment and deploy in another one
- Curricula Learning: learn a difficult task by training on a series of simpler tasks

◇ DRL for robotics
- Adaptation to environment is critical but interacting with the environment is difficult, expensive and potentially dangerous.

◇ MADRL
- A set of agents/robots that learn at the same time in the same environment

Mobile
Robotics,
Reinforce-
ment
Learning and
Deep Rein-
forcement
Learning for
Mobile
Robots

- **PPO**: Schulman et al. (2017) Proximal Policy Optimization
- **DDPG**: Lillicrap et al. (2015) Continuous control with deep reinforcement learning
- **Mapless**: Tai et al. (2017) "Virtual-to-real deep reinforcement learning: Continuous control of mobile robots for mapless navigation"