



Deep Reinforcement Learning and Explainable AI (XAI)

IEIE (June 26th, 2019)

Prof. Joongheon Kim

Korea University, Seoul, Korea

<https://joongheon.github.io/>

Explainable AI (XAI)

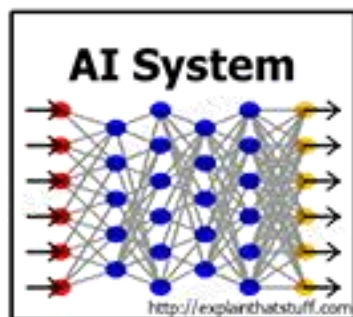
- David Gunning (DARPA),
IJCAI 2016 Workshop Presentation

Reinforcement Learning Review

Imitation Learning

Concluding Remarks





- We are entering a new age of AI applications
- Machine learning is the core technology
- Machine learning models are opaque, non-intuitive, and difficult for people to understand

DoD and non-DoD Applications

Transportation

Security

Medicine

Finance

Legal

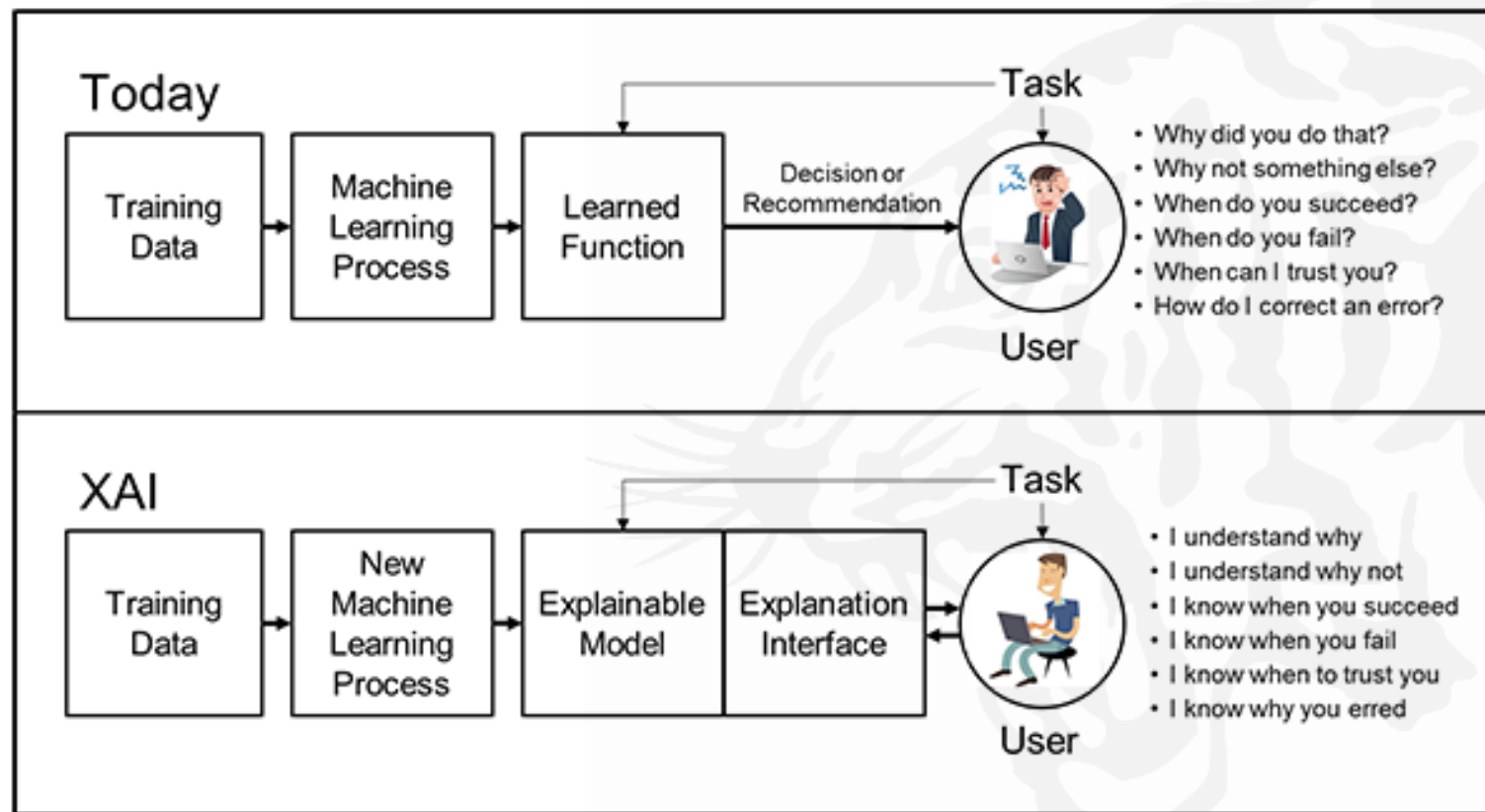
Military

User

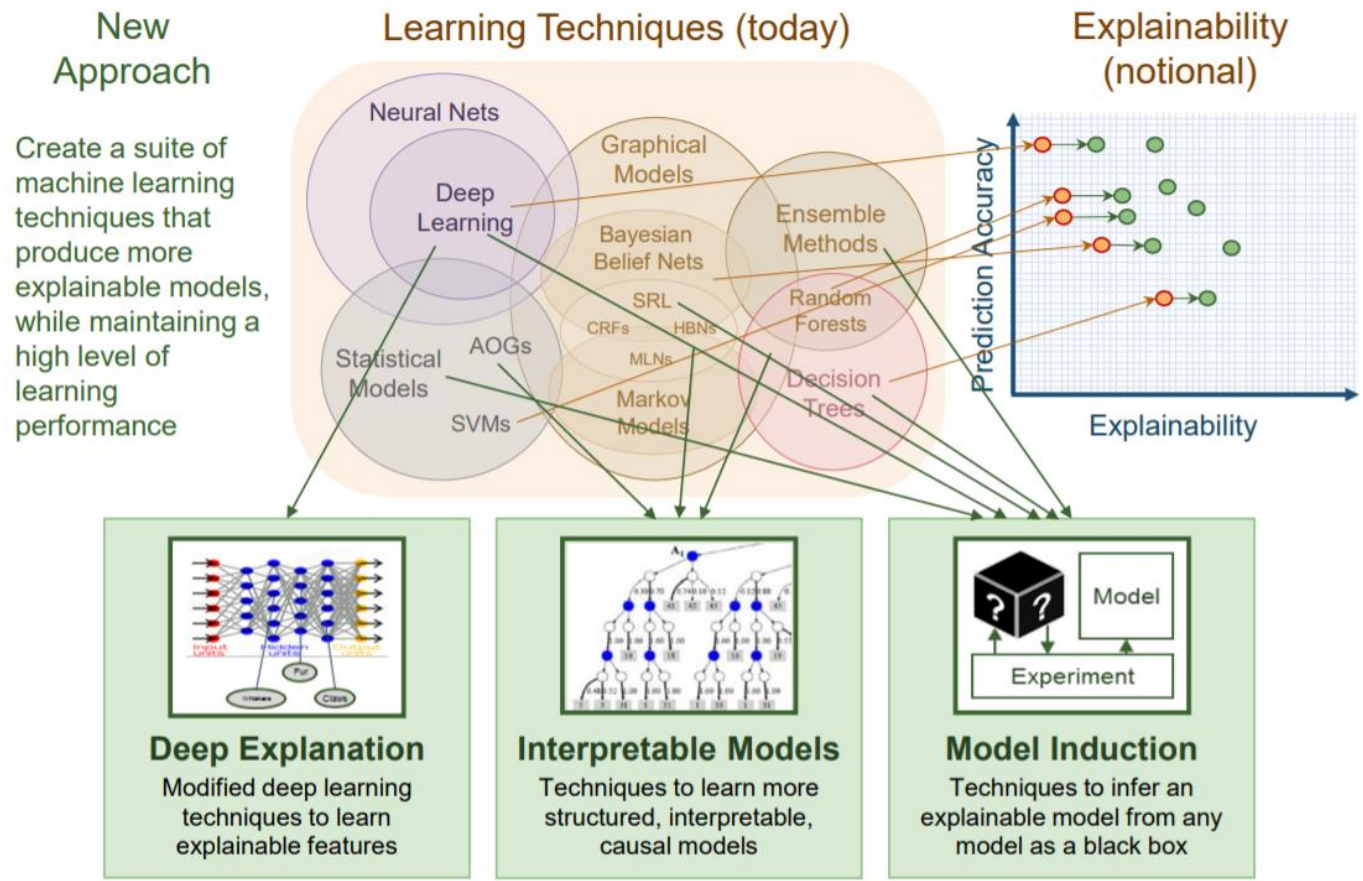


- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

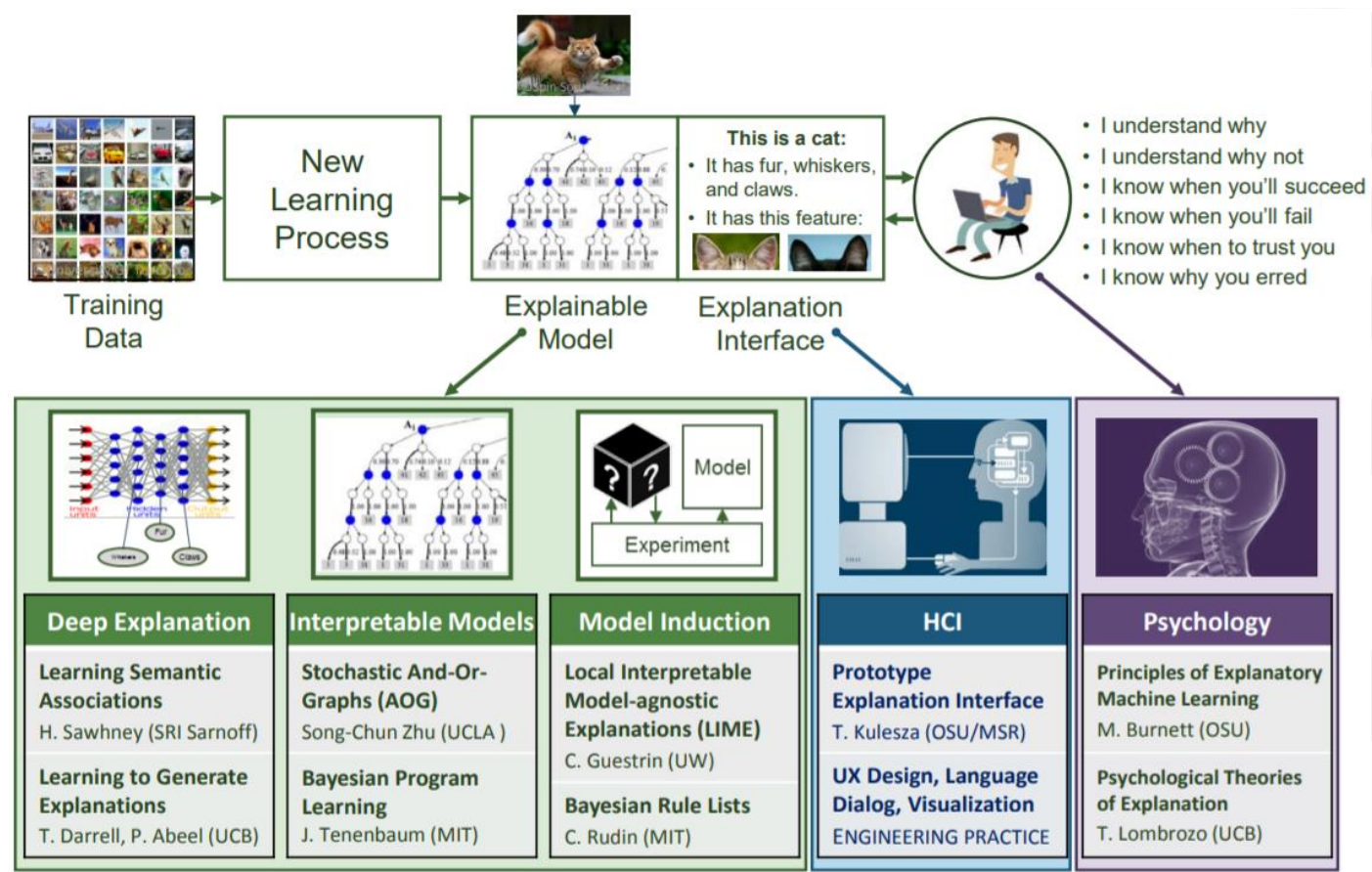
XAI Concept



XAI: Performance vs. Explainability



Why Do You Think It Will Be Successful



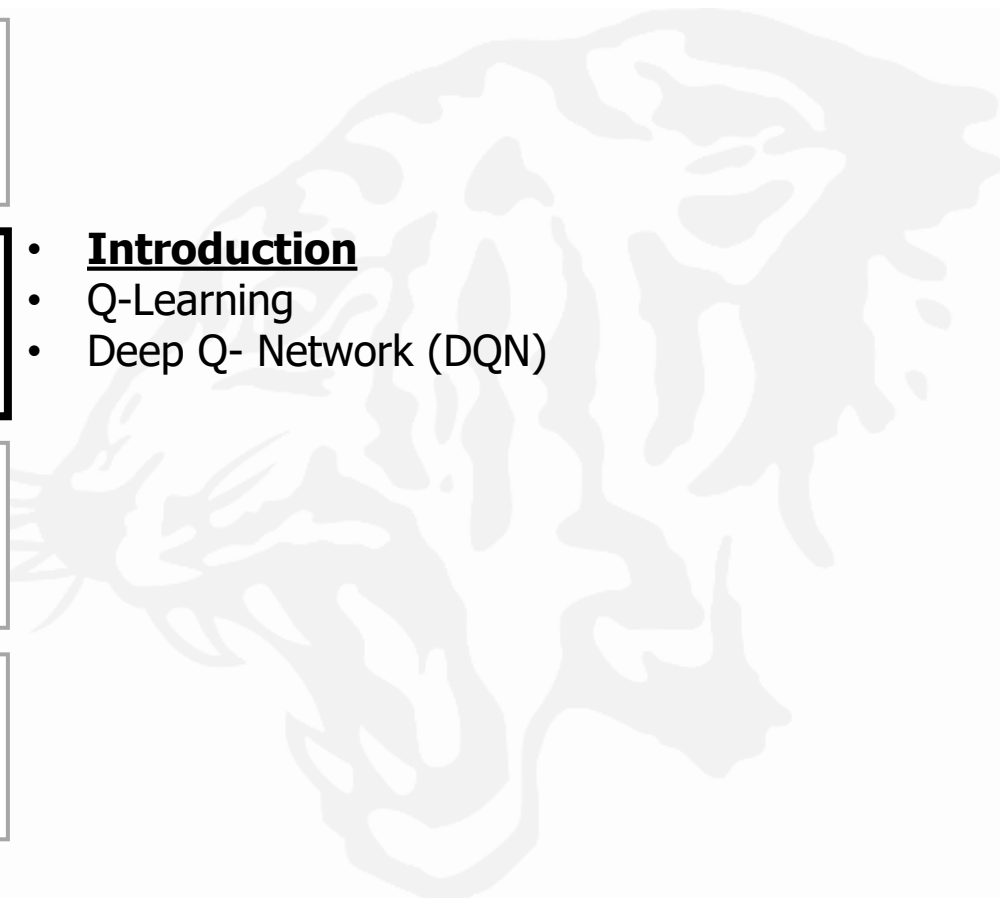
Explainable AI (XAI)

Reinforcement Learning Review

Imitation Learning and Automotive
Applications

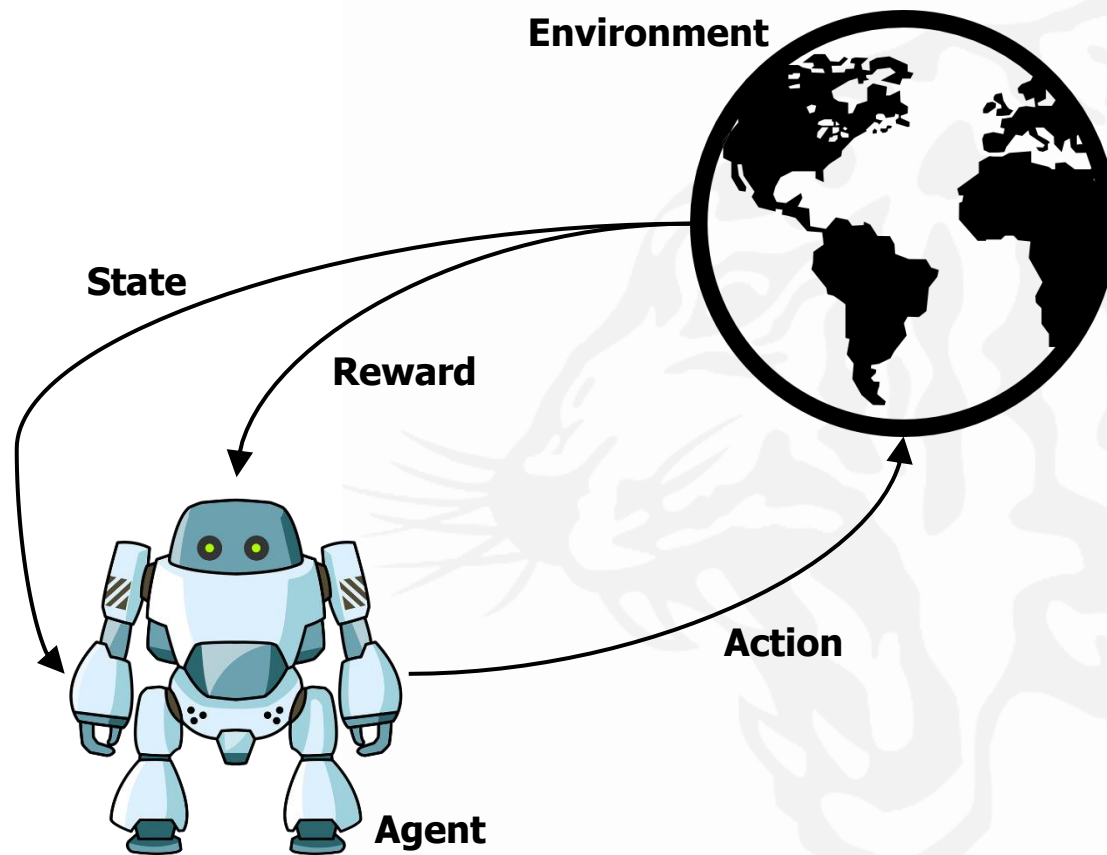
Concluding Remarks

- **Introduction**
- Q-Learning
- Deep Q- Network (DQN)



- Brief History and Successes
 - Minsky's PhD thesis (1954): Stochastic Neural-Analog **Reinforcement** Computer
 - Analogies with animal learning and psychology
 - Job-shop scheduling for NASA space missions (Zhang and Dietterich, 1997)
 - Robotic soccer (Stone and Veloso, 1998) – part of the world-champion approach
- When RL can be used?
 - Find the (approximated) **optimal action sequence** for **expected reward maximization** (not for single optimal solution)
 - Define **actions** and **rewards**. These are all we need to do.

Introduction to RL



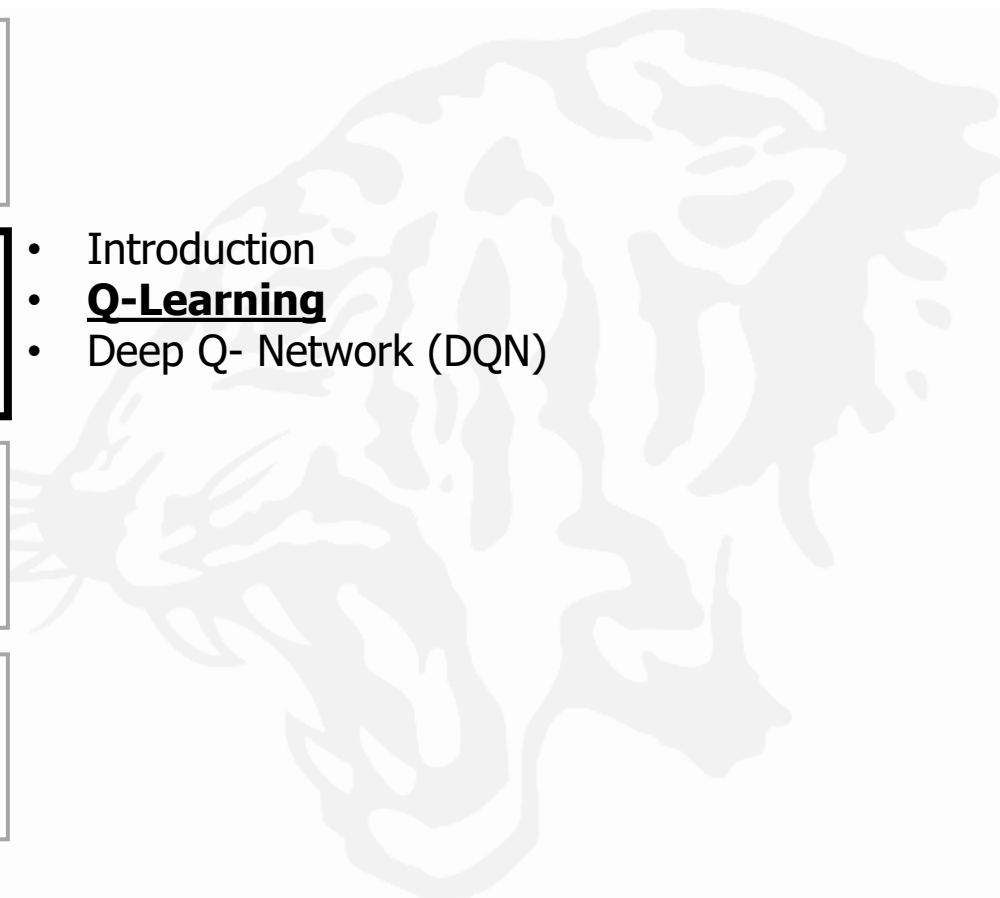
Explainable AI (XAI)

Reinforcement Learning Review

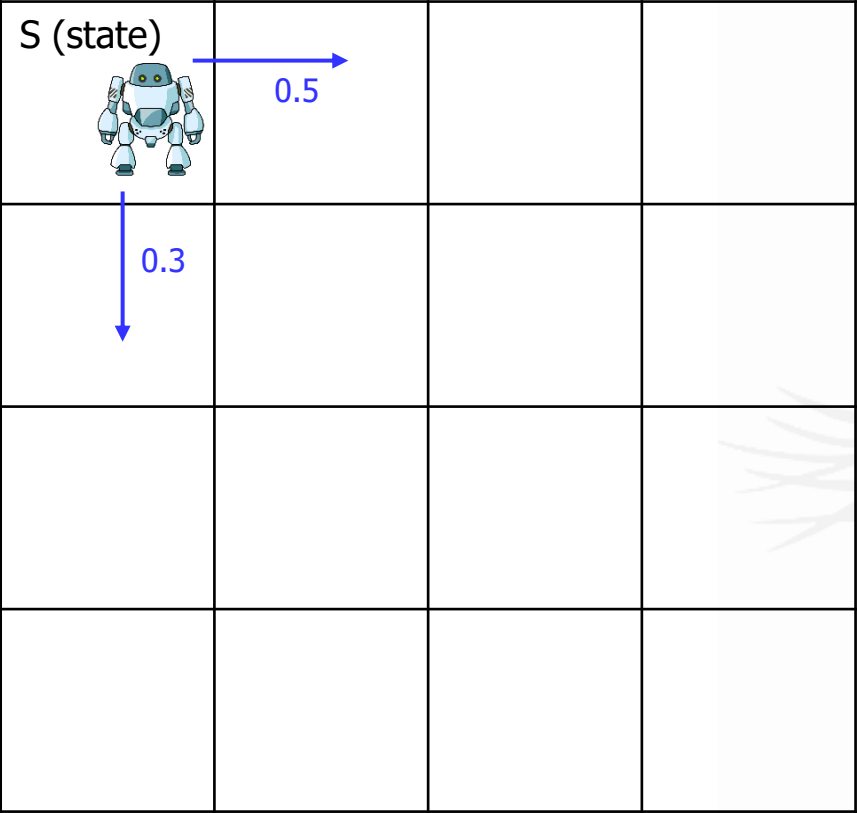
Imitation Learning and Automotive
Applications

Concluding Remarks

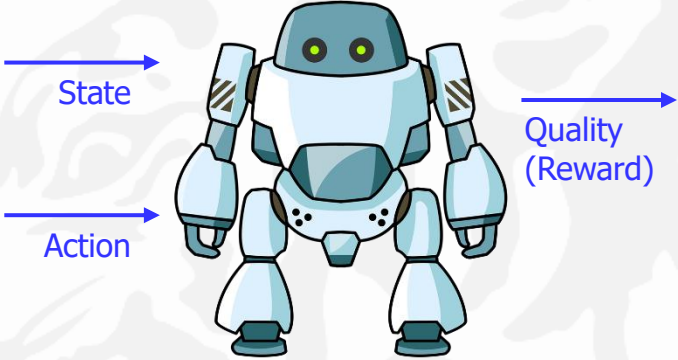
- Introduction
- **Q-Learning**
- Deep Q- Network (DQN)



Quick Example: Q-Learning



- Q-Function
 - State-action value function



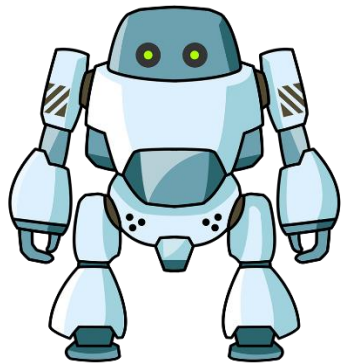
Q (state, action)

Q(s1, LEFT): 0.0
Q(s1, RIGHT): 0.5
Q(s1, UP): 0.0
Q(s1, DOWN): 0.3

Maximum

$$\text{RIGHT} \leftarrow \arg \max_{a \in A} Q(s_1, a)$$

Quick Example: Q-Learning



Q (state, action)

Q(s1, LEFT): 0.0

Q(s1, RIGHT): 0.5 —————→ Maximum

Q(s1, UP): 0.0

Q(s1, DOWN): 0.3

RIGHT ← $\arg \max_{a \in A} Q(s_1, a)$

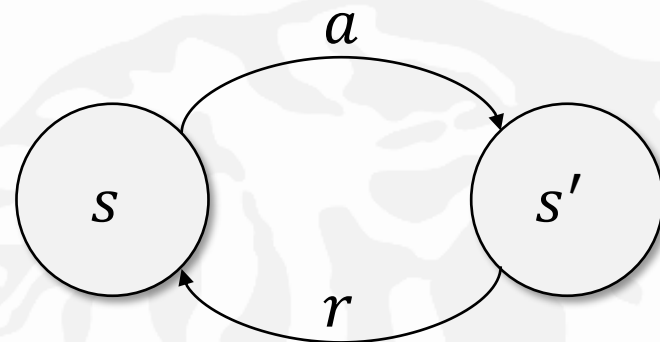
Optimal Policy π and Max Q

- $\text{Max } Q = \max_{a'} Q(s, a')$
- $\pi^*(s) = \arg \max_a Q(s, a)$

Quick Example: Q-Learning

- My condition
 - I am now in state s
 - When I do action a , I will go to s' .
 - When I do action a , I will get reward r
 - Q in s' , it means $Q(s', a')$ exists.
- How can we express $Q(s, a)$ using $Q(s', a')$?

$$Q(s, a) = r + \max_{a'} Q(s', a')$$



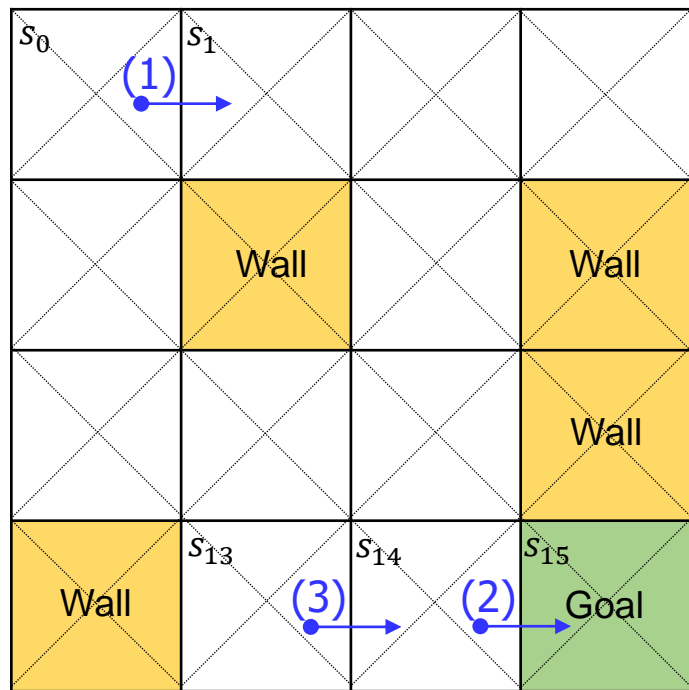
Recurrence (e.g., factorial)

```
F(x){  
    if (x != 1){ x * F(x-1) }  
    if (x == 1){ F(x) = 1 }  
}
```

$$\begin{aligned} 3! &= F(3) = 3 * F(2) \\ &= 3 * 2 * F(1) \\ &= 3 * 2 * 1 = 6 \end{aligned}$$

Quick Example: Q-Learning

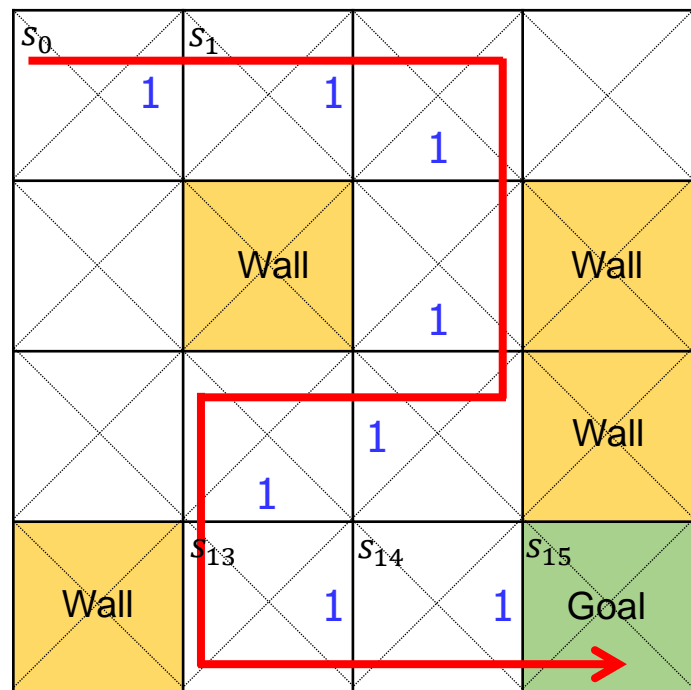
- 16 states and 4 actions (U, D, L, R)



- Initial Status
 - All 64 Q values are 0,
 - Reward are all zero except $r_{s_{15},L} = 1$
- For (1), from s_0 to s_1
 - $Q(s_0, a_R) = r + \max_a Q(s_1, a) = 0 + \max\{0,0,0,0\} = 0$
- For (2), from s_{14} to s_{15} (goal)
 - $Q(s_{14}, a_R) = r + \max_a Q(s_{15}, a) = 1 + \max\{0,0,0,0\} = 1$
- For (3), from s_{13} to s_{14}
 - $Q(s_{13}, a_R) = r + \max_a Q(s_{14}, a) = 0 + \max\{0,0,1,0\} = 1$

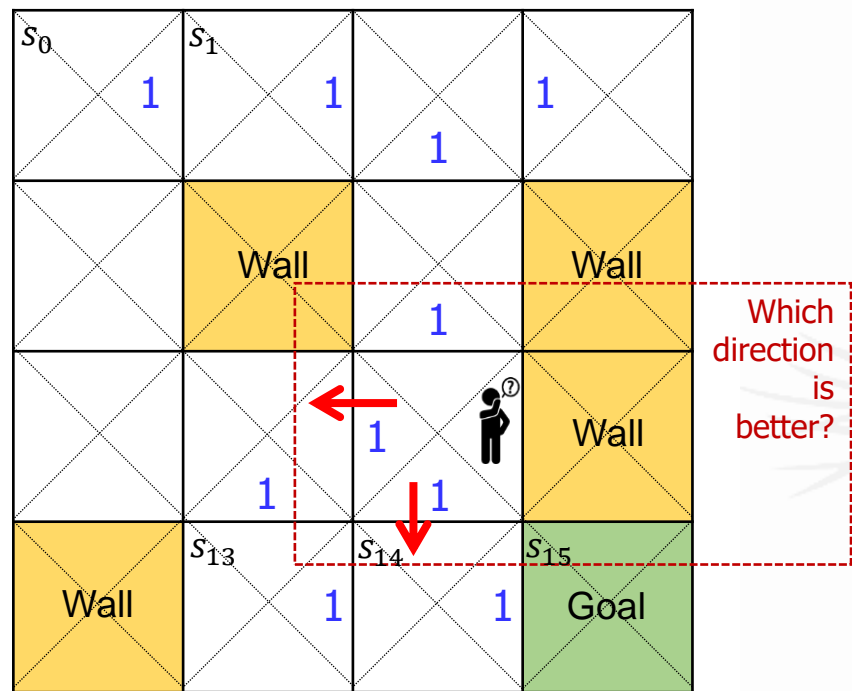
Quick Example: Q-Learning

- 16 states and 4 actions (U, D, L, R)



Quick Example: Q-Learning

- 16 states and 4 actions (U, D, L, R)



Learning $Q(s, a)$ with Discounted Reward

$$Q(s, a) = r + \gamma \cdot \max_a Q(s', a')$$

$$0 < \gamma \leq 1$$

Explainable AI (XAI)

Reinforcement Learning Review

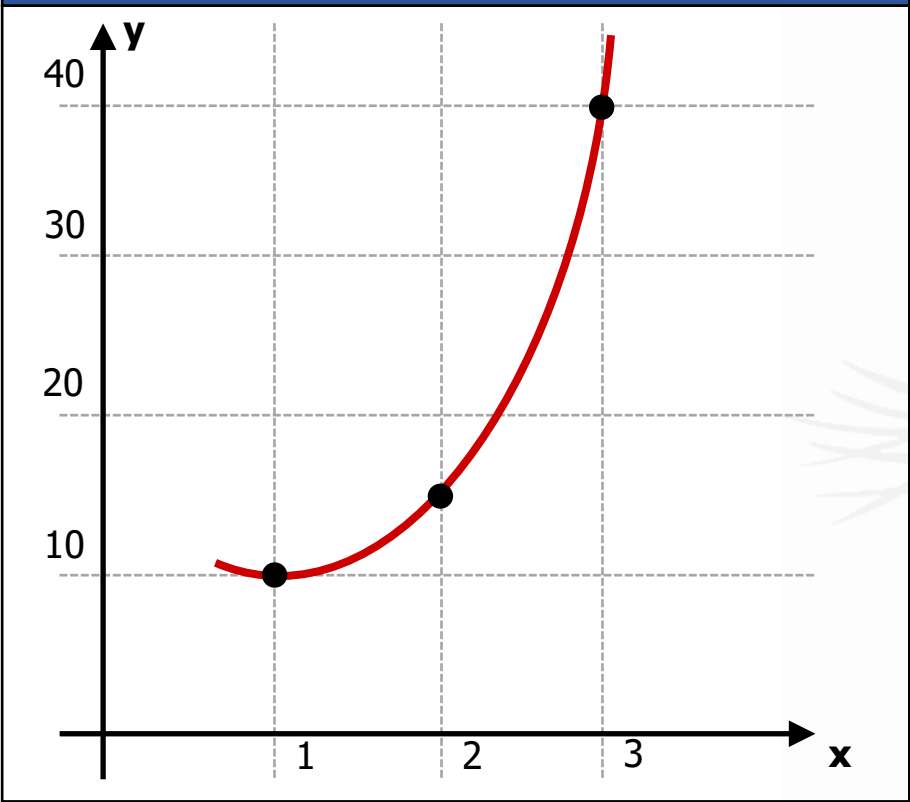
Imitation Learning

Concluding Remarks

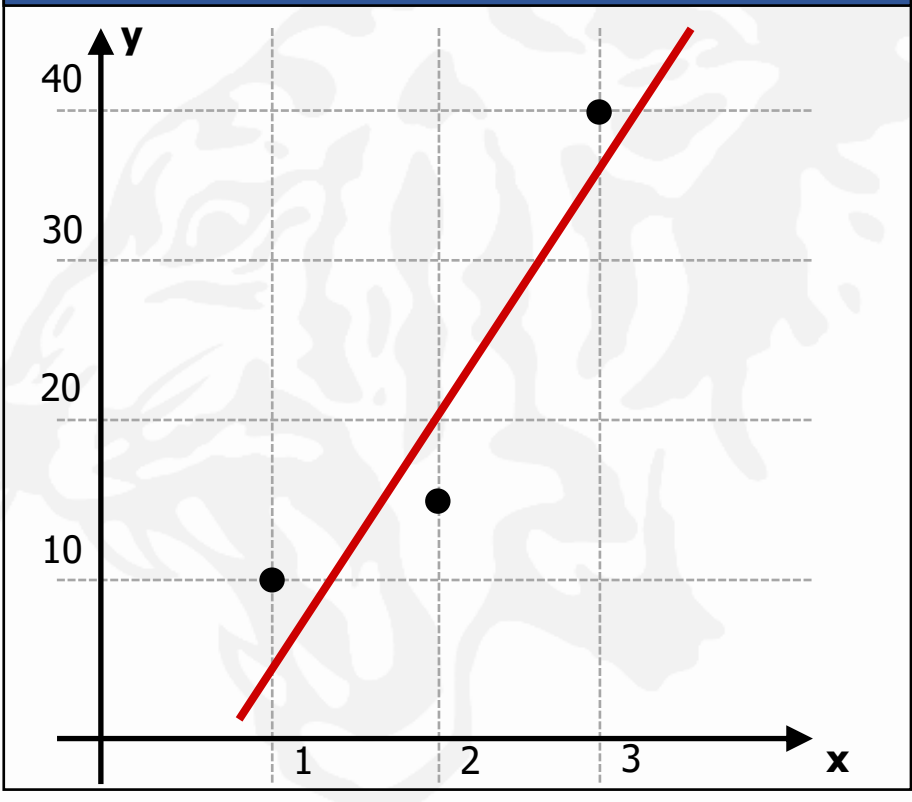
- Introduction
- Q-Learning
- **Deep Q- Network (DQN)**

Interpolation vs. Linear Regression

Interpolation

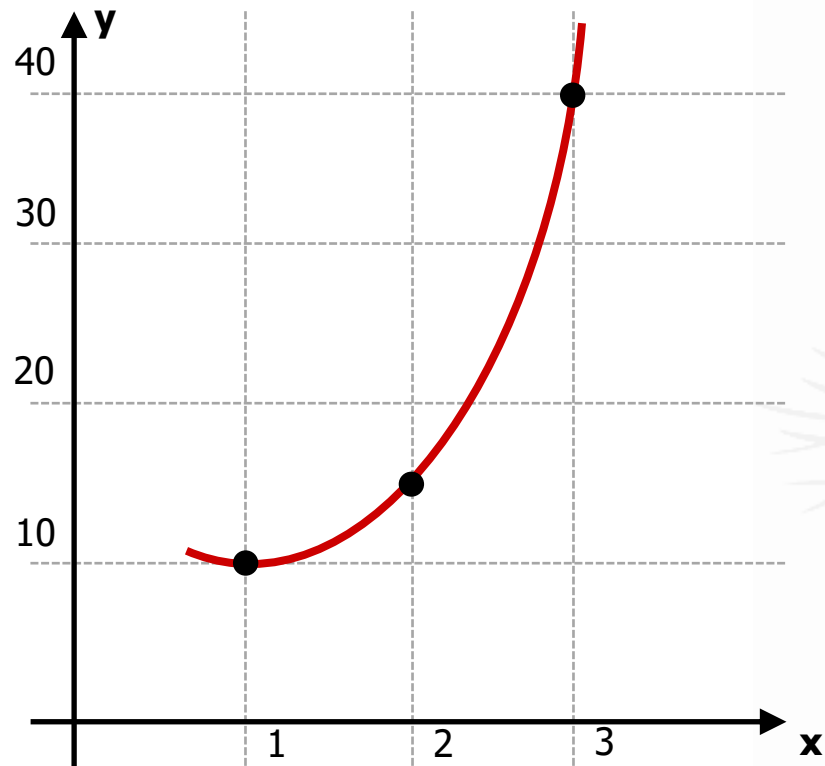


Linear Regression



Interpolation vs. Linear Regression

Interpolation



Interpolation with Polynomials

$$y = a_2x^2 + a_1x^1 + a_0$$

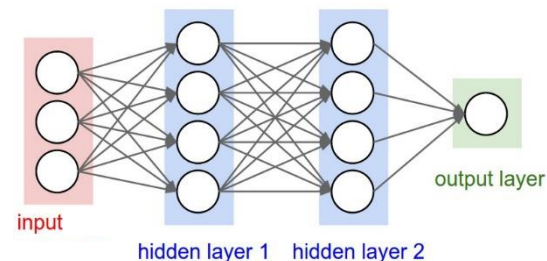
where three points are given.

→ Unique coefficients (a_0, a_1, a_2) can be calculated.



Is this related to
Neural Network Training?

Interpolation and Neural Network Training



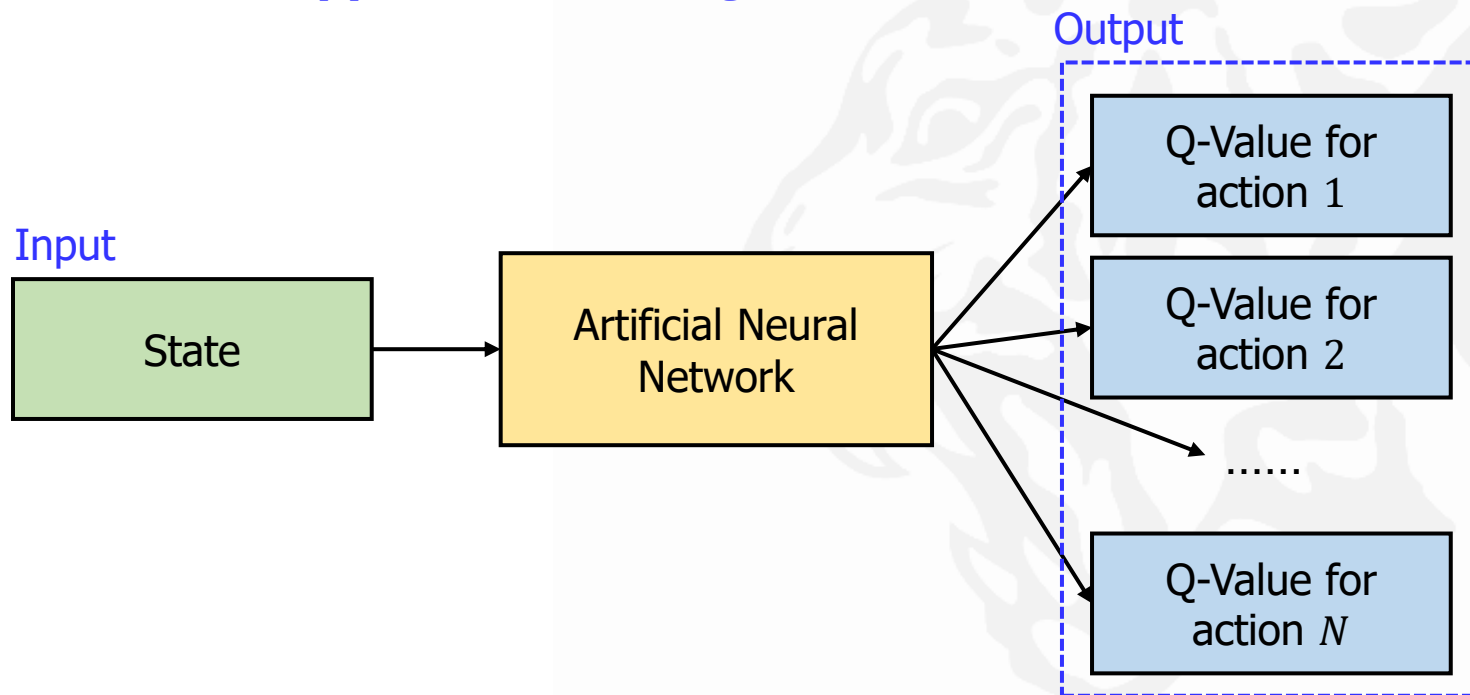
$$Y = a(a(a(X \cdot W_1 + b_1) \cdot W_2 + b_2) \cdot W_o + b_o)$$

where training data/labels (X : data, Y : labels) are given.

- Find $W_1, b_1, W_2, b_2, W_o, b_o$
- This is the mathematical meaning of neural network training.
- **Function Approximation**
- The most well-known function approximation with neural network:
Deep Reinforcement Learning

Deep Q-Network

- Large-Scale Q-Values
 - It is inefficient to make the Q-table for each state-action pair.
→ ANN is used to **approximate the Q-function**.



Outline

Explainable AI (XAI)

Reinforcement Learning Review

Imitation Learning

Concluding Remarks



- ICML 2018 Tutorial
 - <https://sites.google.com/view/icml2018-imitation-learning/>



Imitation Learning Tutorial ICML 2018

Introduction to Imitation Learning

- Gameplay

Pro-Gamer



Trained Agent



The goal of Imitation Learning is to train a policy to mimic
the expert's demonstrations

Introduction to Imitation Learning

- Problems of RL



1. Reward Shaping



2. Safe Learning



3. Exploration process

Imitation Learning **handles with** these problems
through the demonstration of the experts.

Introduction to Imitation Learning

- Starcraft2

States: s = minimap, screen

Action: a = select, drag

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$

States: s

Action: a

Policy: π_{θ}

- Policy maps states to actions : $\pi_{\theta}(s) \rightarrow a$
- Distributions over actions : $\pi_{\theta}(s) \rightarrow P(a)$

State Dynamics: $P(s'|s,a)$

- Typically not known to policy
- Essentially the simulator/environment

Rollout: sequentially execute $\pi_{\theta}(s_0)$ on initial state

- Produce trajectories τ

$P(\tau|\pi)$: distribution of trajectories induced by a policy

$P(s|\pi)$: distribution of states induced by a policy



Imitation Learning Applications: PPF/RFTN Injection Control in Medicine

- PPF/RFTN Injection Control in Medicine

States: $s = \text{BIS}, \text{BP}, \dots$

Action: $a = \text{PPF}, \text{RFTN}, \dots$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$



Autonomous Driving with Imitation Learning

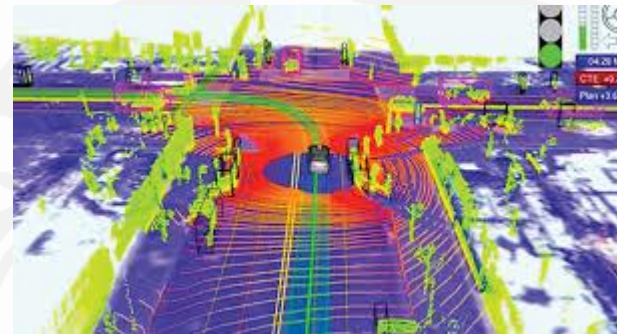
- Autonomous Driving Control

States: s = **sensors**

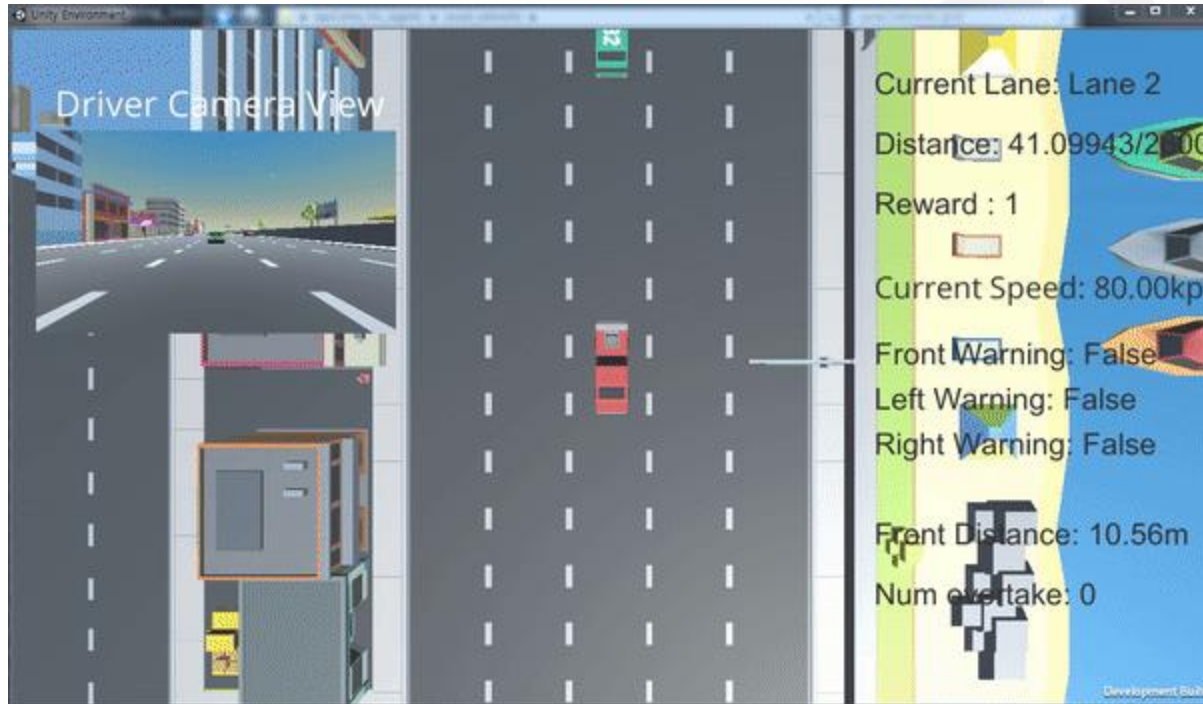
Action: a = **steering wheel, brake, ...**

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$



Autonomous Driving with Imitation Learning



Outline

Explainable AI (XAI)

Reinforcement Learning Review

Imitation Learning

Concluding Remarks



- **Explainable AI (XAI)**
- **Reinforcement Learning:** Q-Learning, DQN
- **Imitation Learning:** Reinforcement Learning, Imitation Learning
- Special Thanks to [MyungJae Shin \(CAU\)](#)
- More questions?
 - Joongheon@korea.ac.kr
- More details?
 - <https://joongheon.github.io/>

