

Deep Reinforcement Learning For Beginners

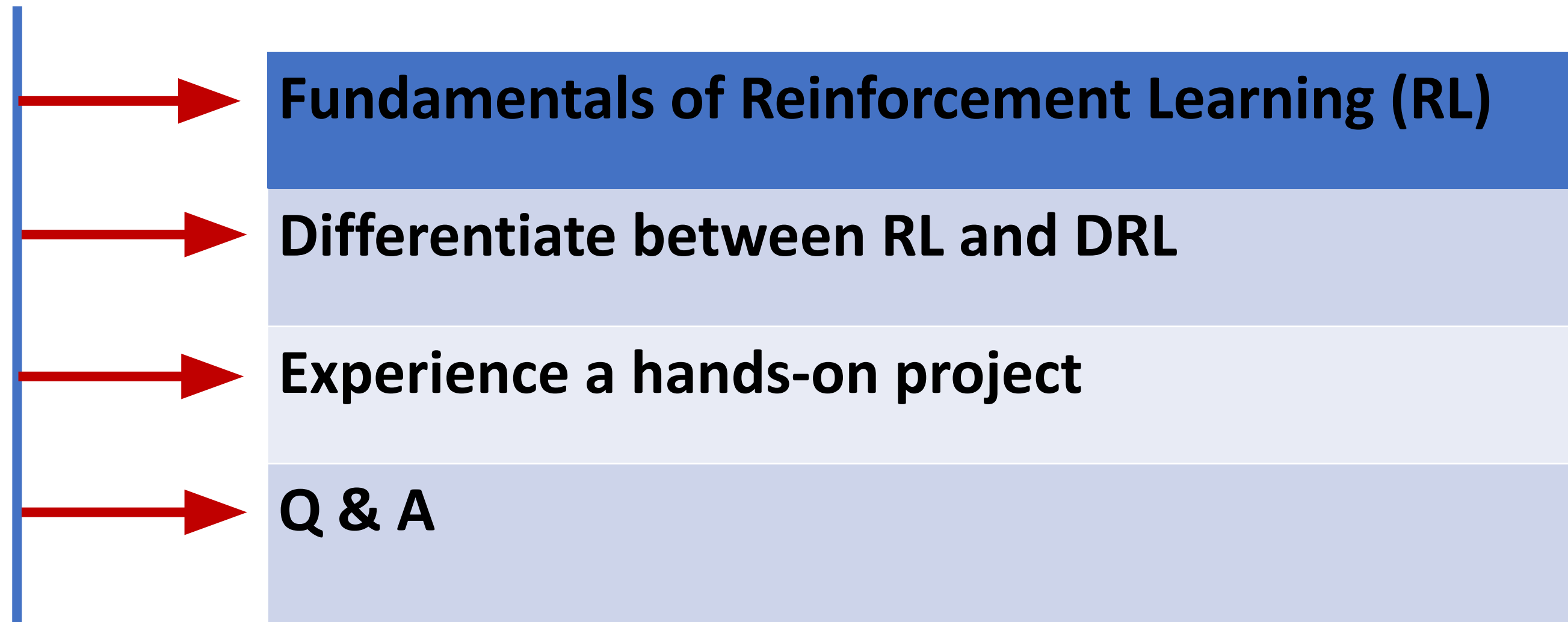
Ashab Uddin , PhD Student, ECE, University of Windsor

Ashab Uddin, PhD Candidates, Electrical and Computer Engineering

- **Research Interest:**

- ☐ **Connected Autonomous Vehicle Application(Edge computing, Computational offloading, Caching , etc.)**
- ☐ **Communication Resource Allocation for Vehicular Network**
- ☐ **Deep Reinforcement Learning and Game Theory Application**

Objectives of the session



Fundamentals of Reinforcement Learning (RL)

Introduction

- ❑ Before digging inside of the basics , let's visit the real world application:
- ❑ The first wave of urban robots is here | Challengers

Reinforcement Learning

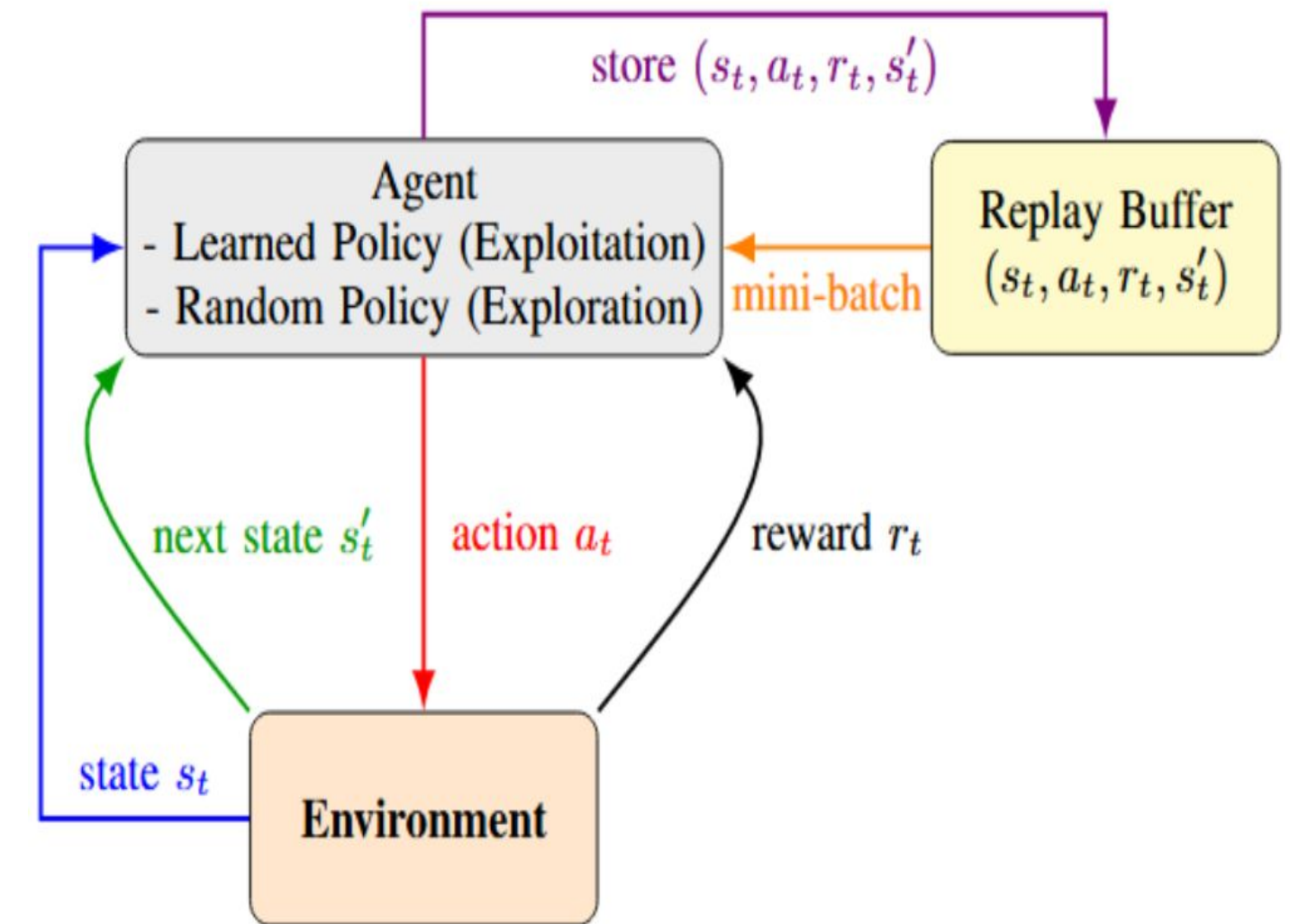
- ❖ Learns from interactions with the environment to achieve a long-term goal
- ❖ The goal is defined by a reward signal, which must be maximized
- ❖ The agent must be able to partially or fully sense the environment state
- ❖ The state is typically represented as a feature vector.
- ❖ State transitions must be related with actions

Exploration and Exploitation

- ❖ Agent learn through exploring action randomly
- ❖ Agent exploit the explored knowledge / learning during training in a controlled way that called exploitation.
- ❖ Example: Initially, a robot take steps (left, right, up, down) randomly with 50% probability and 50% from learned knowledge.

Reinforcement Learning Elements

- ❑ **Environment:** The whole system without the agent itself.
- ❑ **Agent:** A controller or decision maker.
- ❑ **State:** A snapshot of environment / Information Vector
- ❑ **Action:** Decision of Agent.
- ❑ **Policy:** A policy maps (a total map) states of the environment to actions



Policy Type

❖ Stochastic Policy :

- ❑ Policy outputs a probability distribution over actions (e.g., SoftMax).

- ❑ Defined as $\pi(a|s)$

- ❑ **Example:** A robot agent policy for current state, s is:

[Left, Right, Up, Down]==[0.3,0.20,0.35,0.15]==Up

❖ Deterministic Policy:

- ❑ Policy outputs a specific action value.

- ❑ Defined as $\pi(s) = a$

- ❑ **Example:** a car agent decide its velocity $a\%$ of max velocity.

Reinforcement Learning System Example

□ Environment: ?

□ Agent: ?


□ State: $s =$

$$[s_{00}, s_{01}, s_{02}, s_{10}, s_{11}, s_{12}, s_{20}, s_{21}, s_{22}]$$
$$s_0 = [1, 0, 0, 0, 0, 0, 0, 0, 0]$$

□ Action: $a = [P(\text{left}), P(\text{right}), p(\text{up}), p(\text{down})]$

□ Policy: will be developed by learning

	0	1	2
0	 -0.8	-0.8	-0.8
1	-0.8	-1.0	-0.8
2	-1.0	 -10	10 

	0	1	2
0	s_0	s_1	s_2
1	s_3	s_3	s_4
2	s_5	s_7	s_8 


Markov Process

□ Markov Property:

$$p(s'|s) = Pr\{S_{t+1} = s' | S_t = s\}$$

	0.5	0.25	0	0.25	0	0	0	0	0	0
	0.25	0.25	0.25	0.25	0	0	0	0	0	0
.....										
....										
.....										
....										

	0	1	2
0	 -0.8	-0.8	-0.8
1	-0.8	-1.0	-0.8
2	-1.0	 -10	10 

	0	1	2
0	s_0	s_1	s_2
1	s_3	s_3	s_4
2	s_5	s_7	s_8 

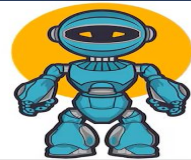

Non-Markovian

- Non-Markovian Example
- Depends on History:: How:

□ *Reward*

= *Cell Reward* – *$f(\text{historic energy drain so far})$*

□ Objective: Maximize reward with minimizing energy consumption

	0	1	2
0	 -0.8	-0.8	-0.8
1	-0.8	10	-0.8
2	-1.0	 -10	10 

Non-Markovian to Markovian

- ❑ Easily Solvable problem is Markovian

- ❑ Need to convert it to Markovian

- ❑ *Reward*

$$= \text{Cell Reward} - f(\text{energy drain in current step})$$

Markov Decision Process

Transition based on current state and action

$$\mathbf{p}(s', r | s, a) = \text{Pr}\{S_{t+1} = s', R_t = r | S_t = s, a_t = a\}$$

And MDP is a tuple presentation of (S, A, P, R, γ)

S = set of states

A = set of action

P = Probability Transition Matrix

R = set of rewards

γ = discounted factor

Probability Transition

Probability Transition Matrix : $\mathbf{p}(s', r | s, a)$

$$\sum_{s'} P(s', r | s, a) = 1$$

Left	1	0	0	0	0	0	0	0	0
	0	1	0	0	0	0	0	0	0
Up	1	0	0	0	0	0	0	0	0
Down	0	0	0	1	0	0	0	0	0

	0	1	2
0	s_0	s_1	s_2
1	s_3	s_3	s_4
2	s_5	s_7	s_8

Expected Return and Discounted Factor

Expected Return, $G(s_t) = r_{t+1} + \gamma * r_{t+2} + \gamma^2 * r_{t+3} + \dots$.

First Step Reward = r_{t+1}

Future Reward = $\gamma * r_{t+2} + \gamma^2 * r_{t+3} + \dots = \gamma * (r_{t+2} + \gamma^1 * r_{t+3} + \dots)$

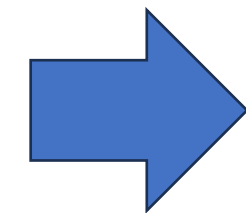
Put together, $G(s_t) = r_{t+1} + \gamma * G(s_{t+1})$

- The discount factor is like a “how fast it gains” maximization :
- Lower $\gamma \rightarrow$ **favors quick wins**(agent is short sighted)
- Higher $\gamma \rightarrow$ **favors total sum**, even if it takes longer(agent is far sighted)

State Value and Action Value Function

❖ State value function [1] :

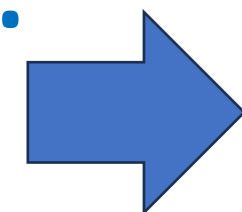
how good a state is,
following a policy π :



$$\begin{aligned}\text{❖ } V_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | s_t = s] = \mathbb{E}_{\pi}[r_t + \gamma G_{t+1} | s_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s'|s, a) (r + \gamma V_{\pi}(s')) \\ &= \sum_a \pi(a|s) Q_{\pi}(s, a) \dots (1)\end{aligned}$$

❖ Action value function [1]:

how good a (state , action)
pair is, following a policy π :



$$\begin{aligned}\text{❖ } Q_{\pi}(s, a) &= \sum_{s', r} p(s'|s, a) (r + \gamma V_{\pi}(s')) \\ &= \sum_{s', r} p(s'|s, a) [r + \gamma \sum_{a'} \pi(a'|s') Q_{\pi}(s', a')] \dots (2)\end{aligned}$$



Bellman Equation

❖ Bellman Optimality[1]:

$$V_{\pi^*}(s) = \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_{\pi^*}^*(s')] \dots\dots(3)$$

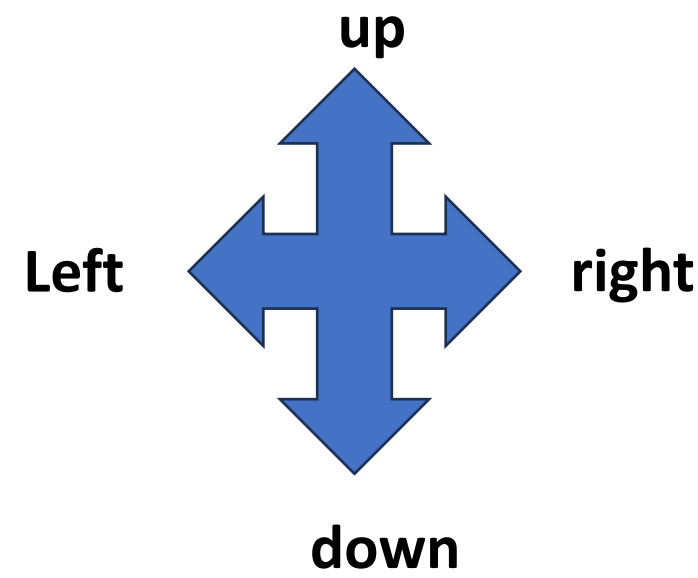
\max_a denotes that we are taking the maximum over all possible actions.


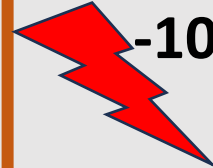

$$Q_{\pi^*}(s, a) = r(s, a) + \sum_{s'} p(s'|s, a) * \gamma * \max_{a'} Q_{\pi^*}^*(s', a') \dots\dots(4)$$

$\max_{a'}$ denotes that we are taking the maximum over all possible actions.

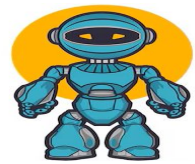
❖ π^* called the optimal policy among all policies π

How RL Converge ?



	0	1	2	3
0	 -0.8	-0.8	-0.8	-0.8
1	-0.8	-1.0	-0.8	-0.8
2	-1.0	 -10	-0.5	-0.8
3	-0.8	-1.0	 10	-0.8

-10



Agent,
Initial State




Agent,
Final State

- Objective is to reach final state with maximum reward.

- Optimal path: 

$(0,0) > (0,1) > (0,2) > (1,2) > (2,2) > (3,2)$: Return=7.1

- $Q((0,0), R) = (-0.8) + \max_a (Q((1,0), a) = (-0.8) + (-2.1 | s = (1,0), a = R))$

- Suboptimal Path:  Return=-3.1

- Suboptimal Path:  Return=-3.6

Differentiate between RL and DRL

RL:: Value Iteration Method

Bellman optimality update:

$$V_k(s) = \max_a \left[r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_{k-1}(s') \right]$$

Policy Extraction: when $V^*(s) = V_k(s)$

$$\pi^*(s) = \operatorname{argmax}_a \left[r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^*(s') \right]$$

Initialize $V(s) \leftarrow 0$ for all $s \in S$
theta \leftarrow small positive number
repeat

$\Delta \leftarrow 0$

$V_{\text{old}} \leftarrow V(s)$

for each $s \in S$ do

if s is terminal:

$V(s) \leftarrow 0$ // or the terminal value

continue

// compute $Q(s, a)$ for this state

$\text{best_q} \leftarrow -\infty$

for each $a \in A(s)$ do

$q \leftarrow \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma * V_{\text{old}}(s')]$

if $q > \text{best_q}$: $\text{best_q} \leftarrow q$

$\Delta \leftarrow \max(\Delta, |\text{best_q} - V(s)|)$

$V(s) \leftarrow \text{best_q}$

end for

until $\Delta \leq \text{theta}$

Next: Policy Evaluation

RL:: Policy Iteration Method

Bellman update with current policy:

$$V_k(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) (r(s, a, s') + \gamma V_{k-1}^\pi(s')) \rightarrow$$

Check Policy Improvement:

$$\pi^*(s) = \operatorname{argmax}_a \left[r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_k^*(s') \right] \rightarrow$$

Initialize: policy $\pi(s) \in A(s)$, $V(s) \leftarrow 0$ for all s , initialize θ

repeat

$\Delta \leftarrow 0$, $V_{\text{old}} \leftarrow V(s)$, $\theta \leftarrow \text{theta} \leftarrow$ small positive number

for each $s \in S$ do

if s is terminal:

$V(s) \leftarrow 0$

continue

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_{\text{old}}(s')]]$

$\Delta \leftarrow \max(\Delta, |V(s) - V_{\text{old}}(s)|)$

end for

until $\Delta < \theta$

policy_stable \leftarrow true

for each $s \in S$ do

if s is terminal

continue

$a_{\text{old}} \leftarrow \pi(s)$

$a^* \leftarrow \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]]$

$\pi(s) \leftarrow a^*$

if $a^* \neq a_{\text{old}}$: policy_stable \leftarrow false

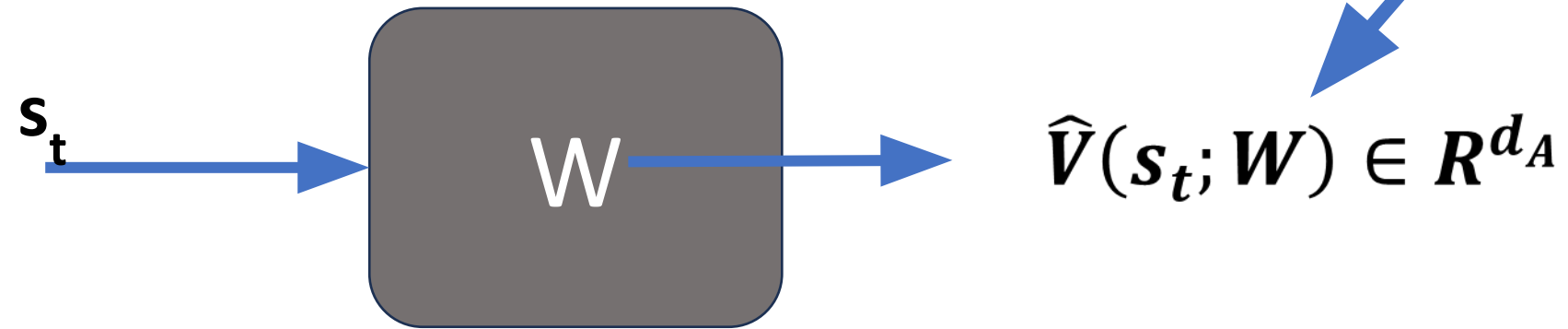
end for

until policy_stable

Output: π (optimal), $V (= V^\pi)$

DRL: Value Iteration by Value Function Approximation

- VFA used supervised learning building a parametrized model to map state to value function or action value function

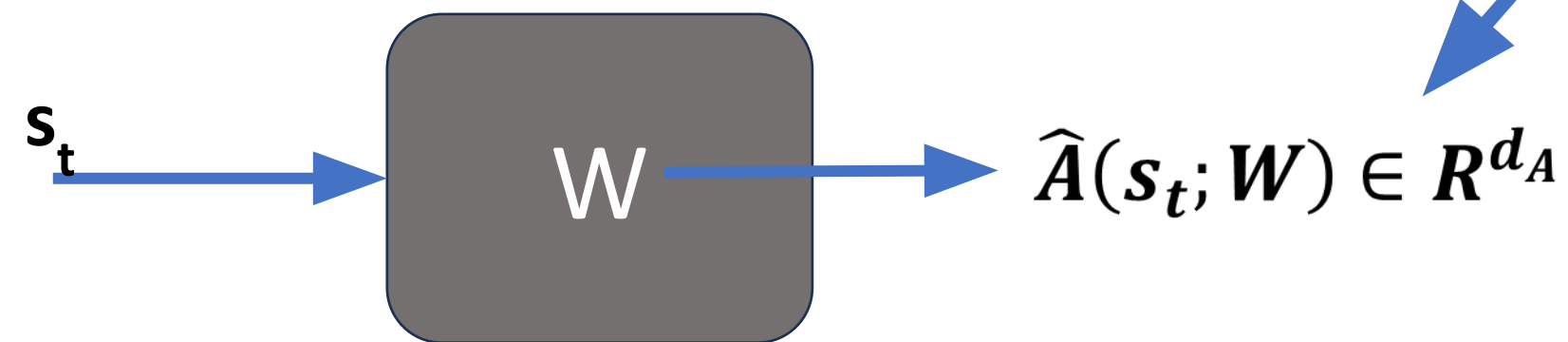


- **Value Iteration through VFA[2]:** minimize loss between estimated and expected value of V or Q without doing gradient on policy itself.

$$\square J(W) = E_{\pi} \left[\left(r + \gamma \hat{Q}_{\max, \pi}(s'_t, a'_t, W^-) - \hat{Q}_{\pi}(s_t, a_t, W) \right)^2 \right]$$

DRL: Policy Iteration by Value Function Approximation

- ❑ VFA used supervised learning building a parametrized model to map state to value function or action value function



- ❑ **Policy Gradient through VFA[2]:** minimize loss function with doing gradient on policy itself.

$$\text{❑ } J(W) = -E_{\pi}[\log \pi(a_t | s_t, W) * \hat{Q}_{\pi}(s_t, a_t, W')]$$

Pseudocode for Deep Q Learning Network (DQN)

```
1 Input:  $epoch\_no, \epsilon_{start}, \epsilon_{end}$ 
2 Output:  $loss, gradients$ 
3 //1. Initialization:
4 Initialize replay memory  $D$  with capacity  $N$ ;
5 Initialize action-value  $Q$  with random weights  $w$ ;
6 Initialize target action-value  $\hat{Q}$  weights  $w^- \leftarrow w$ ;
7 Initialize scores with window size;
8  $\epsilon \leftarrow \epsilon_{start}$ ;
9 for  $episode \leftarrow 1$  to  $M$  do
10   Initialize input raw data  $x_1$ ;
11   Preprocess initial state:  $S \leftarrow \phi(< x_1 >)$ ;
12   for  $time\ step: \tau \leftarrow 1$  to  $T_{max}$  do
13     // 2. Generate training data:
14     Select action  $A$  from state  $S$  using:
14      $\pi \leftarrow \epsilon - Greedy(\hat{Q}(S, A, w))$ ;
15     Take action  $A$ , Observe reward  $R$  and get next
16     input  $s_{\tau+1}$ ;
17     Preprocessing next state:  $S' \leftarrow \phi(s_{\tau+1})$ ;
18     Store experience tuple  $(S, A, R, S')$  in replay
19     memory  $D$ ;
20      $S' \leftarrow S$ ;
21     // 3. Learning:
22     Obtain random mini-batch of  $(s_j, a_j, r_j, s_{j+1})$ 
23     from  $D$ ;
24     if  $episode\ terminate\ at\ step\ j + 1$  then
25       | Set target  $\bar{F}_j \leftarrow r_j$ ;
26     else
27       | Set target  $\bar{F}_j \leftarrow r_j + \gamma \max_{a'} \hat{Q}(s_j, a, w^-)$ ;
28     Update:  $w \leftarrow w + \alpha \nabla_{w_j} L(w_j)$  with Adam;
29     Every  $\mathcal{N}$  steps, update:  $w^- \leftarrow w$ ;
30
31    $\epsilon \leftarrow \max(\epsilon_{end}, \epsilon * decay)$ ;
32   Store score for current episode;
```


Loss Function & Update for Model

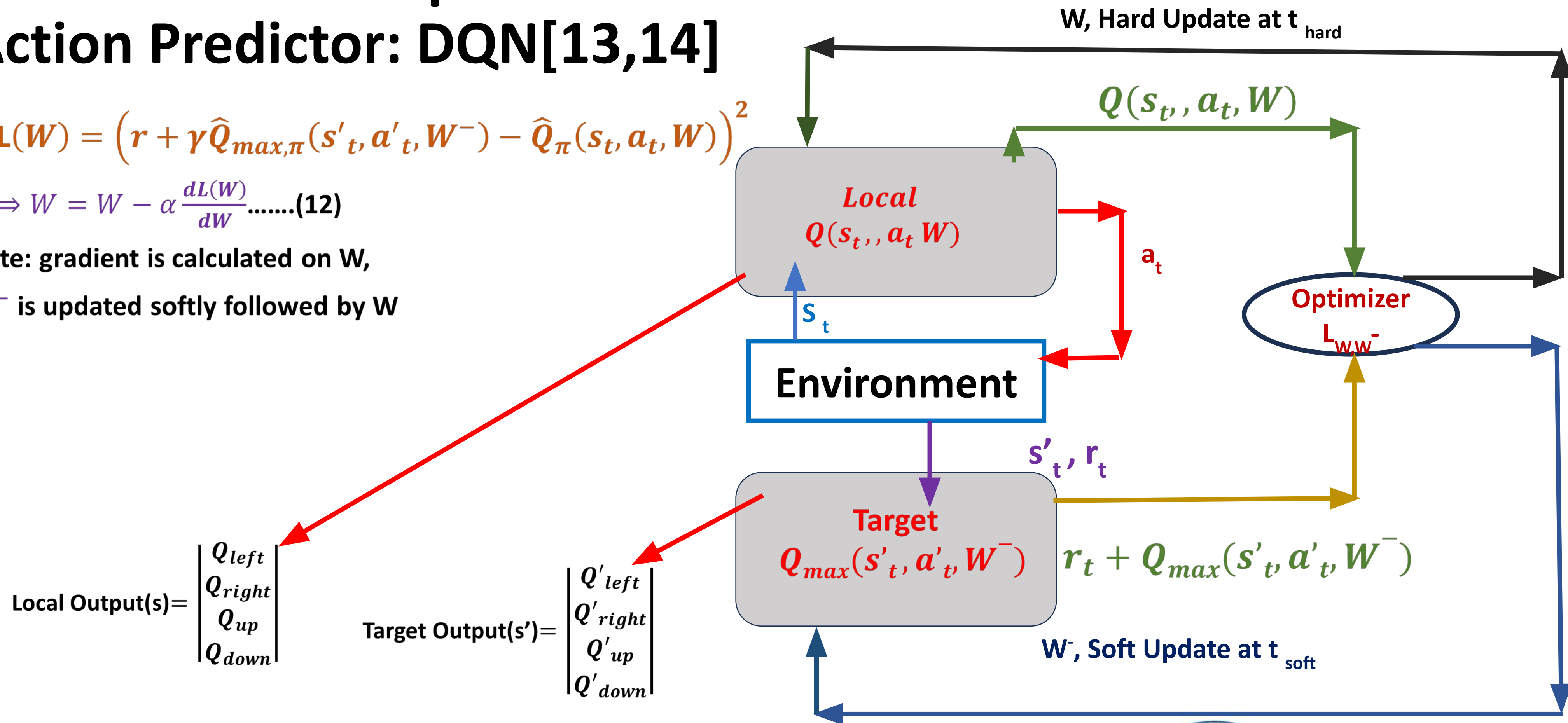
Action Predictor: DQN[13,14]

- $L(W) = \left(r + \gamma \hat{Q}_{max,\pi}(s'_t, a'_t, W^-) - \hat{Q}_{\pi}(s_t, a_t, W) \right)^2$

- $\Rightarrow W = W - \alpha \frac{dL(W)}{dW} \dots\dots(12)$

Note: gradient is calculated on W ,

W^- is updated softly followed by W



Why two Networks

- ❑ With one network
 - ❑ target estimation shifts every update
 - ❑ unstable feedback loop.

- ❑ With two networks
 - ❑ target estimation held constant for a while
 - ❑ smoother, more stable convergence of Q-values.

Tunable Parameters based on MDP

- Network architecture
- Learning Rate
- Batch size
- Target Network Update Parameters
- Discount Factor

Experience a hands-on project



How to setup

- 1. Go to Google Account Signup (<https://accounts.google.com/signup>)
- 2. Go to the Goggle Colab
<https://colab.research.google.com/>
- 3. On the Colab homepage, click File → New Notebook.
- 4. In a cell code : `print('Hello, DRL')`, it will create new Folder as Colab Notebook.
- 5. Upload the shared folder to that folder.
- 6. Open script as Colab Laboratory

Key Recap and Takeaways

- ❖ RL = learning from experience
- ❖ DRL = Machines learning from experience
- ❖ Basics: agent, state, action, reward → learning through trial & error
- ❖ Key of the Design: Markov Decision Process

References

- [11] R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction, 2nd ed. Cambridge, MA: MIT Press, 2018.
- [12] M. Riedmiller, "Neural Fitted Q Iteration – First Experiences with a Data Efficient Neural Reinforcement Learning Method," in Proc. 16th European Conf. Machine Learning (ECML 2005), pp. 317-328, 2005.

Thank You



Q & A



Reinforcement Learning

- ❖ learns from interactions with the environment to achieve a long-term goal
- ❖ The goal is defined by a reward signal, which must be maximized
- ❖ The agent must be able to partially or fully sense the environment state
- ❖ The state is typically represented as a feature vector.
- ❖ State transitions must be related with actions

