

## lecture 5

Previously...

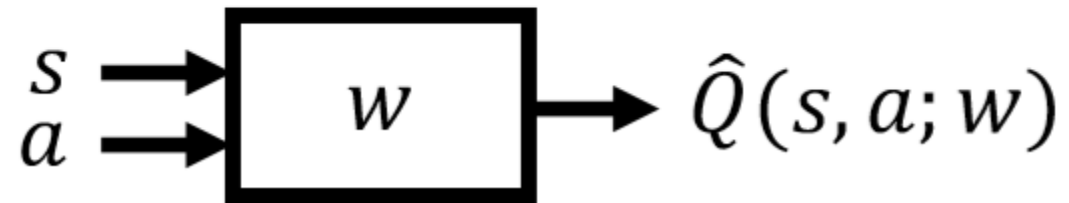
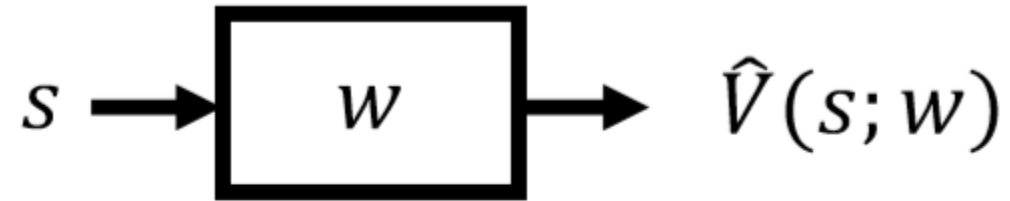
Q. Difference between tabular representation  
vs functional approximation?

Previously...

Q. Difference between tabular representation vs functional approximation?

A. Tabular representation is a table that holds probabilities/likelihood of every possible states as a result of current state + action, whilst functional approximation gives a more compact representation using parameters to represent the tabular representations.

## Value Functional Approximation



↖

Tabular Representation (lec2)

$$\begin{pmatrix} V(s_1) \\ \vdots \\ V(s_N) \end{pmatrix} = \begin{pmatrix} R(s_1) \\ \vdots \\ R(s_N) \end{pmatrix} + \gamma \begin{pmatrix} P(s_1|s_1) & \cdots & P(s_N|s_1) \\ P(s_1|s_2) & \cdots & P(s_N|s_2) \\ \vdots & \ddots & \vdots \\ P(s_1|s_N) & \cdots & P(s_N|s_N) \end{pmatrix} \begin{pmatrix} V(s_1) \\ \vdots \\ V(s_N) \end{pmatrix}$$

Value Functional Approximation

=

Generalization

1. Reduce memory required

2. Reduce computation

3. Reduce experience...

Q1. I understand that it can be also described as reducing data required.. But how does it reduce experience..?

A1. 근사치를 구하니까 필요한 데이터도 줄어 수 있다.

## lecture 5

From just  
policy  
evaluation



To Value Functional  
Approximation Prediction

From: Having a look up table of value estimates and then updating the value estimates each episode or steps

To: reapproximating function when every time new data is given (every step/run)

Feature vectors

$$\mathbf{x}(s) = \begin{pmatrix} x_1(s) \\ x_2(s) \\ \dots \\ x_n(s) \end{pmatrix}$$

$$\mathbf{x}(s, a) = \begin{pmatrix} x_1(s, a) \\ x_2(s, a) \\ \dots \\ x_n(s, a) \end{pmatrix}$$

## lecture 5

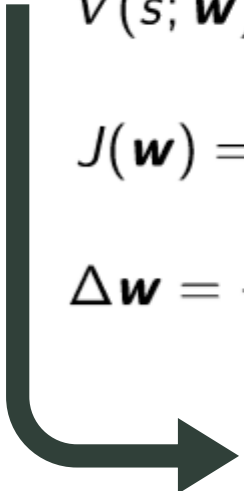
### Update Linear VFA for Prediction with...

Oracle

$$\hat{V}(s; \mathbf{w}) = \sum_{j=1}^n x_j(s) w_j = \mathbf{x}(s)^T \mathbf{w}$$

$$J(\mathbf{w}) = \mathbb{E}_{\pi}[(V^{\pi}(s) - \hat{V}(s; \mathbf{w}))^2]$$

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$


$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \left( 2 \left( V^{\pi}(s) - \hat{V}(s; \mathbf{w}) \right) \right) \mathbf{x}(s)$$

The equation is annotated with dashed boxes and arrows. A dashed box around  $\frac{1}{2}$  has an arrow pointing to the text 'step-size'. A dashed box around  $2(V^{\pi}(s) - \hat{V}(s; \mathbf{w}))$  has an arrow pointing to the text 'prediction error'. A dashed box around  $\mathbf{x}(s)$  has an arrow pointing to the text 'feature value'.

Update = step-size x prediction error x feature value



## lecture 5

### Update Linear VFA for Prediction with...

Oracle  $\Delta \mathbf{w} = -\frac{1}{2} \alpha \left( 2 \left( V^\pi(s) - \hat{V}(s; \mathbf{w}) \right) \right) \mathbf{x}(s)$

$$\Delta \mathbf{w} = \alpha \left( V^\pi(s) - \hat{V}(s; \mathbf{w}) \right) \mathbf{x}(s)$$

Monte  
Carlo

$$\begin{aligned} \Delta \mathbf{w} &= \alpha (G_t - \hat{V}(s_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{V}(s_t; \mathbf{w}) \\ &= \alpha (G_t - \hat{V}(s_t; \mathbf{w})) \mathbf{x}(s_t) \\ &= \alpha (\boxed{G_t} - \mathbf{x}(s_t)^T \mathbf{w}) \mathbf{x}(s_t) \end{aligned}$$

During algorithm for MC linear  
VFA policy evaluation

$$G_t(s) = \sum_{j=t}^{L_k} r_{k,j}$$

Gamma is set a 1, and this is no  
problem because MC itself is  
episodic = bound to terminate =  
return is bounded

## lecture 5

Update Linear VFA for Prediction with...

Oracle  $\Delta \mathbf{w} = -\frac{1}{2} \alpha \left( 2 \left( V^\pi(s) - \hat{V}(s; \mathbf{w}) \right) \right) \mathbf{x}(s)$

$$\Delta \mathbf{w} = \alpha \left( V^\pi(s) - \hat{V}(s; \mathbf{w}) \right) \mathbf{x}(s)$$

Monte  
Carlo

$$\Delta \mathbf{w} = \alpha (G_t - \mathbf{x}(s_t)^T \mathbf{w}) \mathbf{x}(s_t)$$

Tempo  
ral  
Differe  
nce

$$\begin{aligned} \Delta \mathbf{w} &= \alpha \left( \overset{\text{TD target}}{\boxed{r + \gamma \hat{V}^\pi(s'; \mathbf{w})}} - \hat{V}^\pi(s; \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{V}^\pi(s; \mathbf{w}) \\ &= \alpha (r + \gamma \hat{V}^\pi(s'; \mathbf{w}) - \hat{V}^\pi(s; \mathbf{w})) \mathbf{x}(s) \\ &= \alpha (r + \gamma \boxed{\mathbf{x}(s')^T \mathbf{w}} - \mathbf{x}(s)^T \mathbf{w}) \mathbf{x}(s) \end{aligned}$$

Q2. Can I point at this and say that bootstrapping is used?

## lecture 5

### Convergence Guarantees

"The Markov Chain defined by a MDP with a particular policy will eventually converge to a probability distribution over states  $d(s)$ "

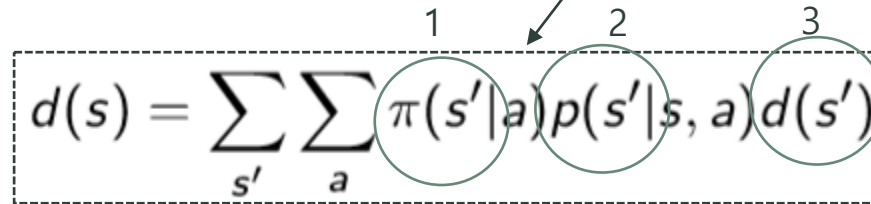


Q3. 옹? 무슨 뜻?

## Convergence Guarantees

"The Markov Chain defined by a MDP with a particular policy will eventually converge to a probability distribution over states  $d(s)$ "

Q4. 엥? 무슨 뜻?

$$d(s) = \sum_{s'} \sum_a \pi(s'|a) p(s'|s, a) d(s')$$


모든 action과 state을 거쳐서 합한 게  $d(s)$ 이고 이는 1이다. 다 1로 합해지는 것은 당연. 혹시 다른 의미가 있을까

# lecture 5

## Convergence Guarantees

Stationary distribution



Q5. what does it mean by "stationary" in stationary distribution?



$$MSVE(\mathbf{w}) = \sum_{s \in S} d(s) (V^\pi(s) - \hat{V}^\pi(s; \mathbf{w}))^2$$

$$MSVE(\mathbf{w}_{MC}) = \min_{\mathbf{w}} \sum_{s \in S} d(s) (V^\pi(s) - \hat{V}^\pi(s; \mathbf{w}))^2$$

same

$$MSVE(\mathbf{w}_{TD}) \leq \frac{1}{1-\gamma} \min_{\mathbf{w}} \sum_{s \in S} d(s) (V^\pi(s) - \hat{V}^\pi(s; \mathbf{w}))^2$$

Error from bootstrapping

## lecture 5

### Control using VFA

Interleave

1. Policy evaluation
2. E-greedy policy improvement

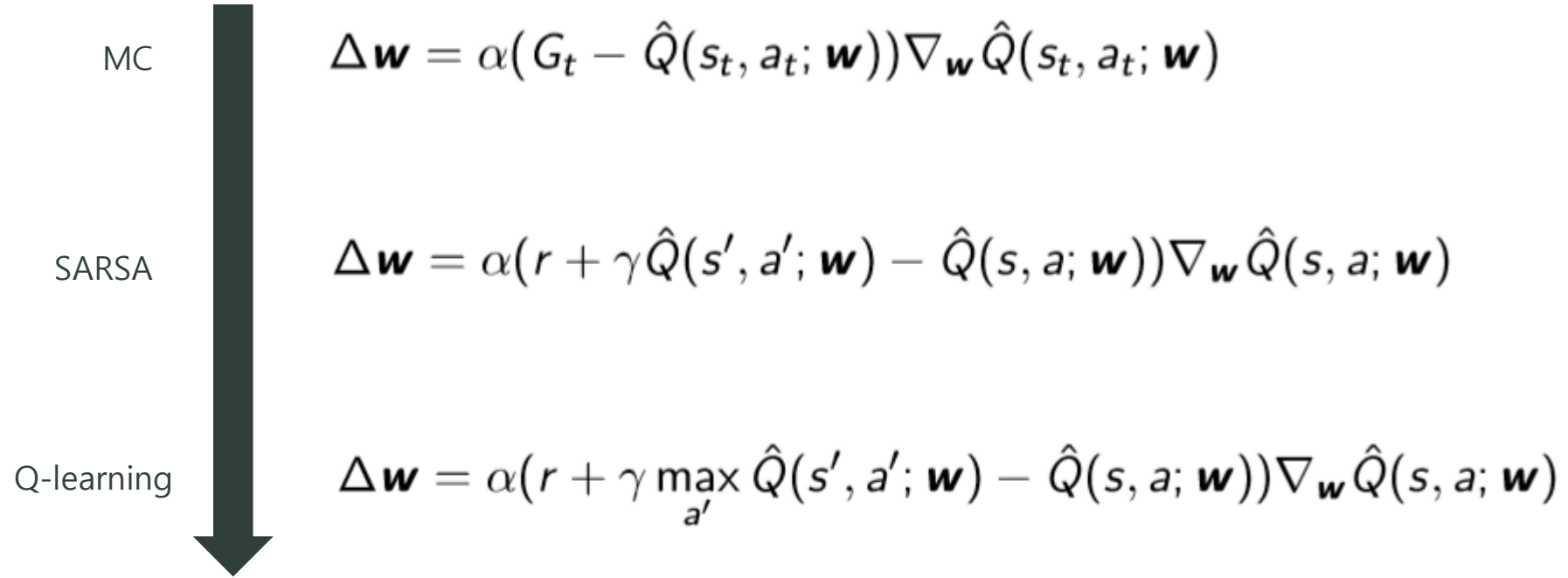
But unstable      Deadly Triad

- Functional approximation
- Bootstrapping
- Off-policy learning

## lecture 5


Control using VFA

Incremental model-free approach



Control using VFA

Q6. What is the difference between linear and nonlinear VFA?



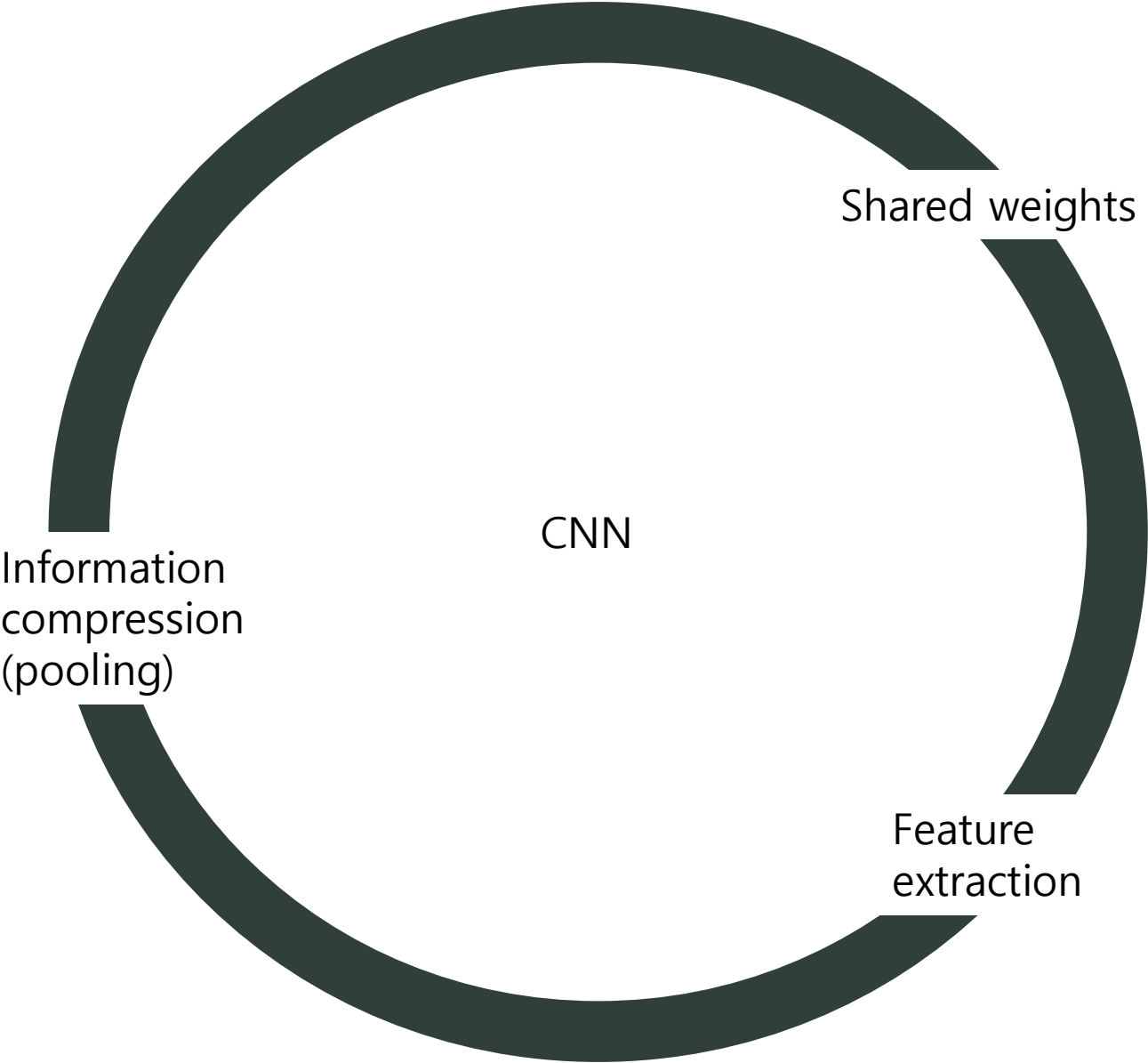
Algorithm	Tabular	Linear VFA	Nonlinear VFA
Monte-Carlo Control	converges	Converges but might have some oscillation	nope
SARSA	converges	Converges but might have some oscillation	Nope
Q-learning	converges	Nope	Nope



## lecture 6

DNN

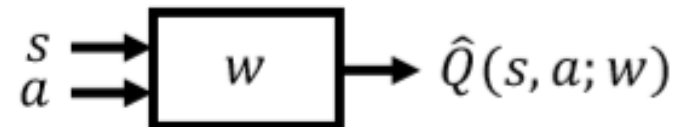
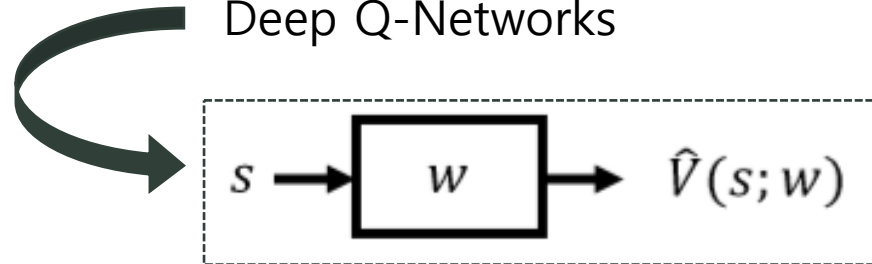
CNN



Deep Q-Learning

Deep Reinforcement Learning

Deep Q-Networks



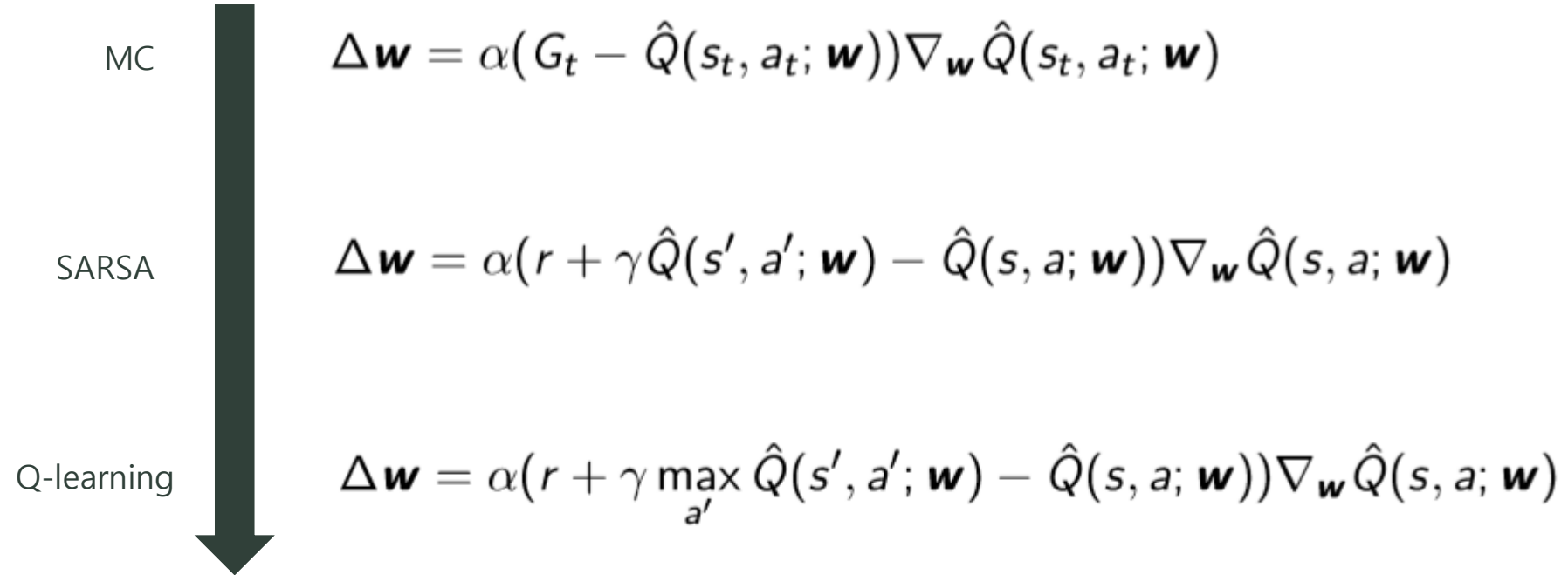
Q7. Don't need the Value function in the DQN's only Q function (parameterized) required. Isn't it?



From lecture 5

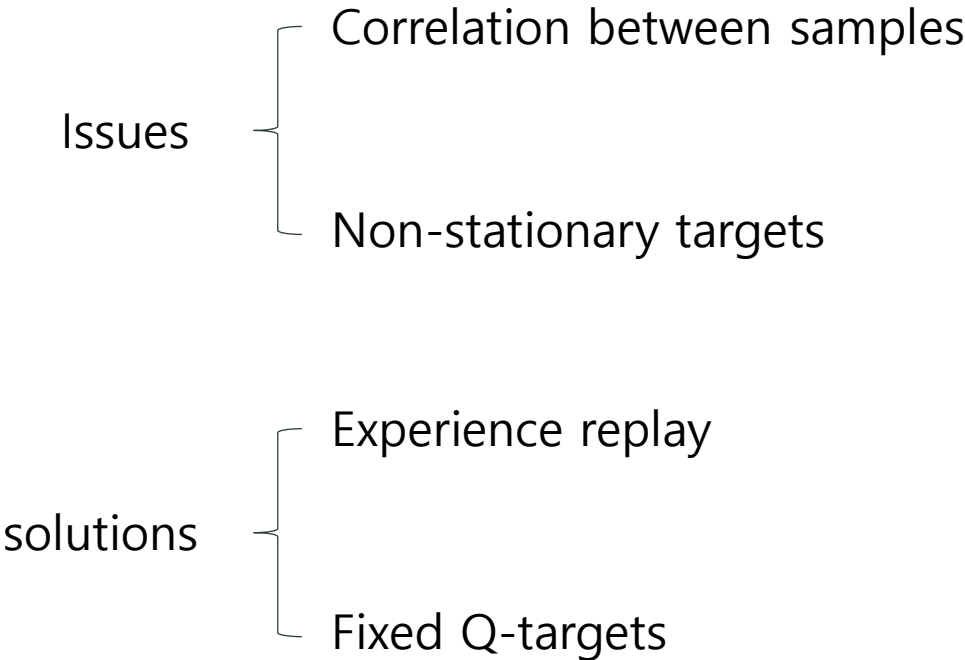
Control using VFA

Incremental model-free approach



lecture 6

DQN



$$\text{Double DQN} \quad \Delta \mathbf{w} = \alpha \left( r + \gamma \underbrace{\hat{Q}(\arg \max_{a'} \hat{Q}(s', a'; \mathbf{w}^-); \mathbf{w}^-)}_{\text{Action selection: } \mathbf{w}} - \hat{Q}(s, a; \mathbf{w}) \right)$$

Action evaluation:  $\mathbf{w}^-$

$$\text{Prioritized order replay} \quad p_i = \left| r + \gamma \max_{a'} Q(s_{i+1}, a'; \mathbf{w}^-) - Q(s_i, a_i; \mathbf{w}) \right|$$

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

$$\text{Dueling DQN} \quad A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$



$$\text{Double DQN} \quad \Delta \mathbf{w} = \alpha \left( r + \gamma \underbrace{\hat{Q}(\arg \max_{a'} \hat{Q}(s', a'; \mathbf{w}); \mathbf{w}^-)}_{\text{Action evaluation: } \mathbf{w}^-} - \underbrace{\hat{Q}(s, a; \mathbf{w})}_{\text{Action selection: } \mathbf{w}} \right)$$

Q7. Difference between fixed Q-target.

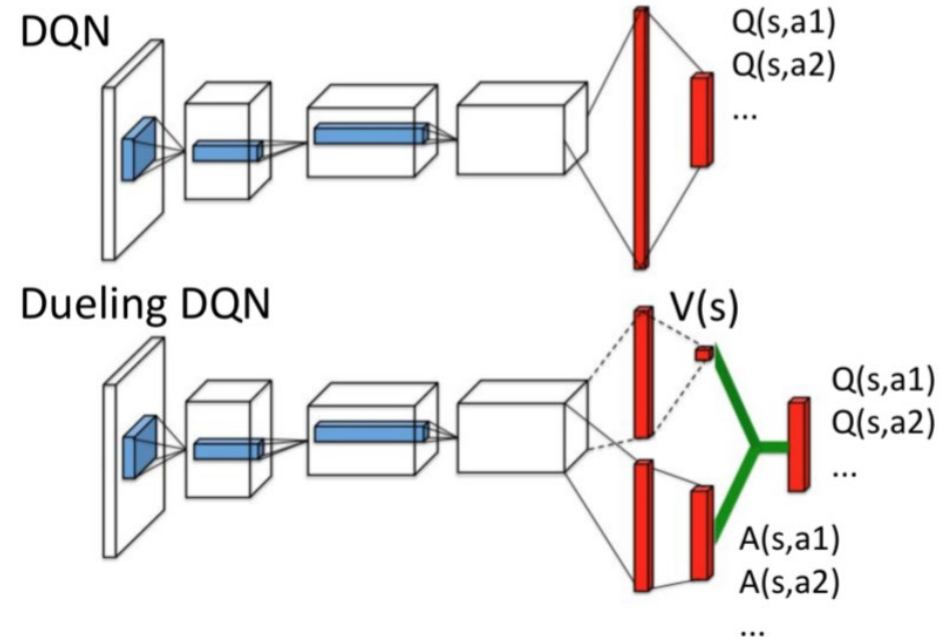
TD Error - Priority of a tuple is proportional to DQN error

Prioritized  
order replay

$$p_i = \left[ r + \gamma \max_{a'} Q(s_{i+1}, a'; \mathbf{w}^-) - Q(s_i, a_i; \mathbf{w}) \right]$$

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

## lecture 6



How much better or worse  
taking a particular action versus  
following the current policy

Identifiability : Whether there exists a  
unique Q for A and V  
given  $\pi$  (policy)

Dueling DQN 
$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

## lecture 7

## lecture 7

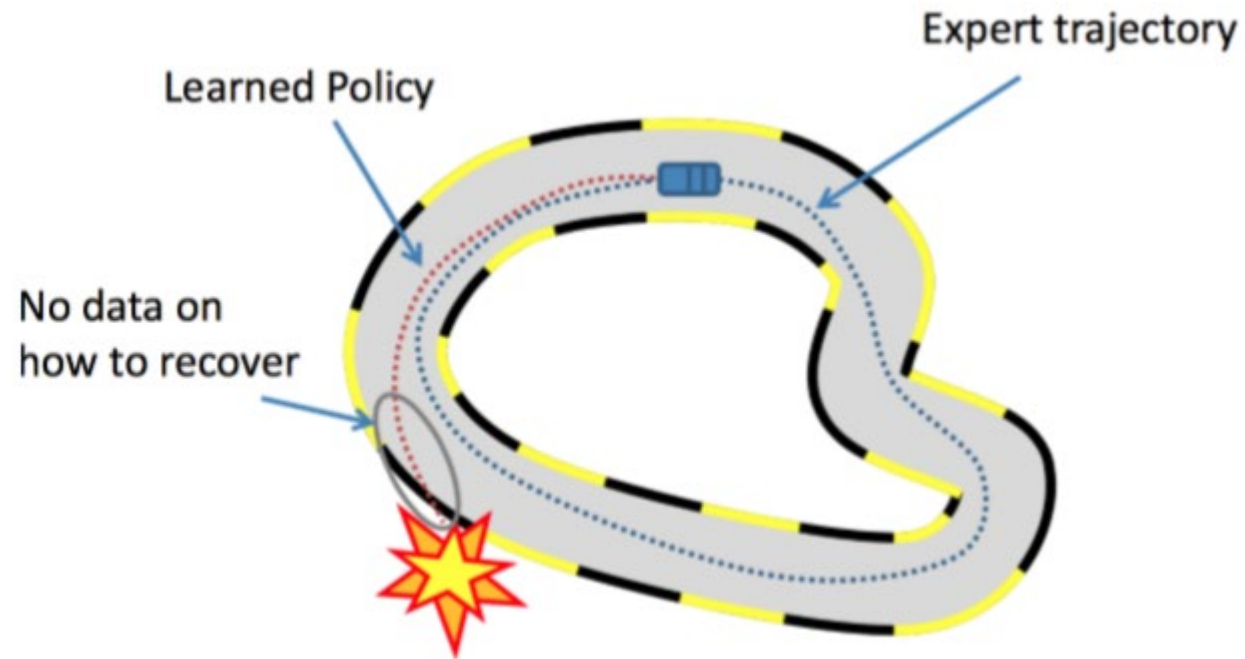
DQN may require too large number of samples to learn a good policy

Or even a large number of samples used to learn might not promise a good policy improvement

→ Imitation learning

## lecture 7

- Imitation learning
- Behavioral Cloning : Estimate policy from training examples
- Problem : Compounding Error



- Solution : DAGGER ( or is it? )

## lecture 7

→ Imitation learning

→ Apprenticeship learning via Inverse RL

Inverse RL

1. Identify weight vector given demonstrations (from teacher)

2. Use that to calculate function of policy  $\pi$  

Distribution  
of states

3.  $\mu$  : discounted sum of state features under policy  $\pi$

= Q8. Can I  
understand it as  
Indirectly gaining  
information on  
the policy?

Q9. So are we calculating  
anything at the 'just'  
Inverse RL stage? Or is it a  
method to combine with  
apprenticeship learning?

## lecture 7

→ Imitation learning

→ Apprenticeship learning via Inverse RL

Inverse RL

1. Identify weight vector given demonstrations (from teacher)
2. Use that to calculate function of policy  $\pi$
3.  $\mu$  : discounted sum of state features under policy  $\pi$

Possibly wrong

Maybe room for some discussion?



## lecture 7

→ Imitation learning

→ Apprenticeship learning via Inverse RL

Find  $w$  such that...

$$w^{*T} \mu(\pi^*) \geq w^{*T} \mu(\pi), \forall \pi \neq \pi^*$$

$$\arg \max_w \max_{\gamma} s.t. w^T \mu(\pi^*) \geq w^T \mu(\pi) + \gamma \quad \forall \pi \in \{\pi_0, \pi_1, \dots, \pi_{i-1}\}$$

Q10. So... how do we obtain  $\mu$ ?  ~~$\mu$  is the demonstration itself?~~ Only calculation of weight  $w$  is sufficient in finding out the reward function?

## Ambiguity

- There is an infinite number of reward functions with the same optimal policy.
- There are infinitely many stochastic policies that can match feature counts
- Which one should be chosen?

Maybe room for some discussion?