

# **Apprenticeship learning with DNN (with experiments)**

Jeong Gwan Lee

KAIST

(Korea Advanced Institute of Science and Technology)

# What I implemented,

---

1. LSPI
2. basis function : sigmoid, RBF, ANNCP, DAN(h2, output)
3. Batch, off-policy and model-free Apprenticeship learning(LSTD-mu)
4. Deep Action Network(DAN) + DQN

# Batch, Off-Policy and Model-Free Apprenticeship Learning(BOMAL)

Klein, Edouard, Matthieu Geist, and Olivier Pietquin. "Batch, off-policy and model-free apprenticeship learning." *European Workshop on Reinforcement Learning*. Springer, Berlin, Heidelberg, 2011.

# Apprentice Learning,

---

is a learning approach, which assumes that **the reward function is unknown**, but instead, *expert knowledge* is available.

This *expert knowledge* is in the form of sequences of states and actions, trajectories, where the goal state is achieved.

The main idea is to find a optimal policy close to the *expert policy*  $\pi^E$ , to use these trajectories.

# Inverse Reinforcement Learning(IRL)

Reward function is parameterized as **linear combination** of features,

$$R(s) = \mathbf{w}^T \phi(s) \quad s.t. \quad \phi : S \rightarrow [0, 1]^k,$$

w is the weight vector to be learned,  
 $\phi$  is the features vector(basis function),  
k is the number of features.

Feature expectations, denoted as  $\mu$ ,

$$\mu(\pi) = E\left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | \pi\right] \in \mathbb{R}^k.$$

Estimated feature expectations for expert policy,

$$\hat{\mu}_E = \frac{1}{m} \sum_{i=0}^m \sum_{t=0}^{\infty} \gamma^t \phi(s_t^{(i)}),$$

Learning process would be repeated until,  $\|\mu(\pi) - \hat{\mu}_E\| < \epsilon$

# MDP\mathcal{R} set-up

---

Why "reward( $R$ )", "parameter( $\theta$ )" and "basis function( $\phi(s)$ )" should be bounded?

$$\mathcal{H}_\phi = \{\theta^T \phi(s), \theta \in \mathbb{R}^p\}, \quad R = \theta^T \phi(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \cdots + \theta_p \phi_p(s)$$

- (*paper*)  $R \in \mathbb{R}^S$  the reward function (assumed to be absolutely bounded by 1)
 
$$|R| \leq 1$$
- (*paper*) we assume the basis functions to be bounded by 1:
 
$$|\phi_i(s)| \leq 1 \quad 1 \leq i \leq p$$
- (*paper*) In order to ensure that rewards are bounded, we impose that  $\|\theta\|_2 \leq 1$ 

$$\|\theta\|_2 \leq 1$$

$$\mu^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | s_0 = s, \pi\right] \quad |V^{\pi_E}(s_0) - V^{\tilde{\pi}}(s_0)| = |\theta^T(\mu^{\pi_E}(s_0) - \mu^{\tilde{\pi}}(s_0))| \leq \|\mu^{\pi_E}(s_0) - \mu^{\tilde{\pi}}(s_0)\|_2$$

# MDP\mathcal{R} set-up

Why "reward( $R$ )", "parameter( $\theta$ )" and "basis function( $\phi(s)$ )" should be bounded?

$$\mathcal{H}_\phi = \{\theta^T \phi(s), \theta \in \mathbb{R}^p\}, \quad R = \theta^T \phi(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \cdots + \theta_p \phi_p(s)$$

- (*paper*)  $R \in \mathbb{R}^S$  the reward function (assumed to be absolutely bounded by 1)

$$|R| \leq 1$$

- (*paper*) we assume the basis functions to be bounded by 1:

$$|\phi_i(s)| \leq 1 \quad 1 \leq i \leq p$$

- (*paper*) In order to ensure that rewards are bounded, we impose that  $\|\theta\|_2 \leq 1$

$$\|\theta\|_2 \leq 1$$

$$\mu^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | s_0 = s, \pi\right] \quad |V^{\pi_E}(s_0) - V^{\tilde{\pi}}(s_0)| = |\theta^T(\mu^{\pi_E}(s_0) - \mu^{\tilde{\pi}}(s_0))| \leq \|\mu^{\pi_E}(s_0) - \mu^{\tilde{\pi}}(s_0)\|_2$$

# MDP\R set-up

Why "reward( $R$ )", "parameter( $\theta$ )" and "basis function( $\phi(s)$ )" should be bounded?

$$\mathcal{H}_\phi = \{\theta^T \phi(s), \theta \in \mathbb{R}^p\}, \quad R = \theta^T \phi(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \cdots + \theta_p \phi_p(s)$$

- (*paper*)  $R \in \mathbb{R}^S$  the reward function (assumed to be absolutely bounded by 1)

$$|R| \leq 1$$

- (*paper*) we assume the basis functions to be bounded by 1:

$$|\phi_i(s)| \leq 1 \quad 1 \leq i \leq p$$

- (*paper*) In order to ensure that rewards are bounded, we impose that  $\|\theta\|_2 \leq 1$

$$\|\theta\|_2 \leq 1$$

$$\theta = \begin{bmatrix} 0.7 \\ 0.7 \end{bmatrix} \quad \|\theta\|_2 = \sqrt{0.7^2 + 0.7^2} = \sqrt{0.98} \leq 1 \quad \phi(s) = \begin{bmatrix} \phi_1(s) \\ \phi_2(s) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad R = \theta^T \phi(s) = 0.7 + 0.7 = 1.4 > 1$$

# MDP\R set-up

Why "reward( $R$ )", "parameter( $\theta$ )" and "basis function( $\phi(s)$ )" should be bounded?

$$\mathcal{H}_\phi = \{\theta^T \phi(s), \theta \in \mathbb{R}^p\}, \quad R = \theta^T \phi(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \dots + \theta_p \phi_p(s)$$

- (*paper*)  $R \in \mathbb{R}^S$  the reward function (assumed to be absolutely bounded by 1)

$$|R| \leq 1$$

- (*paper*) we assume the basis functions to be bounded by 1:

$$|\phi_i(s)| \leq 1 \quad 1 \leq i \leq p$$

- (*paper*) In order to ensure that rewards are bounded, we impose that  $\|\theta\|_2 \leq 1$

$$\|\theta\|_1 \leq 1 \longrightarrow |\theta_1| + |\theta_2| + \dots + |\theta_p| \leq 1$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} \quad \phi(s) = \begin{bmatrix} \phi_1(s) \\ \phi_2(s) \\ \vdots \\ \phi_p(s) \end{bmatrix}$$

$ \theta_1  \phi_1(s)  +  \theta_2  \phi_2(s)  + \dots +  \theta_p  \phi_p(s)  \leq 1$	$(\because  \phi_i(s)  \leq 1)$
$ \theta_1\phi_1(s)  +  \theta_2\phi_2(s)  + \dots +  \theta_p\phi_p(s)  \leq 1$	$(\because  A  B  =  AB )$
$ \theta_1\phi_1(s) + \theta_2\phi_2(s) + \dots + \theta_p\phi_p(s)  \leq 1$	$(\because  A + B  \leq  A  +  B )$
$R = \theta^T \phi(s) = \theta_1\phi_1(s) + \theta_2\phi_2(s) + \dots + \theta_p\phi_p(s) \leq 1$	

# MDP\mathcal{R} set-up

Why "reward( $R$ )", "parameter( $\theta$ )" and "basis function( $\phi(s)$ )" should be bounded?

$$\mathcal{H}_\phi = \{\theta^T \phi(s), \theta \in \mathbb{R}^p\}, \quad R = \theta^T \phi(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \cdots + \theta_p \phi_p(s)$$

- (paper)  $R \in \mathbb{R}^S$  the reward function (assumed to be absolutely bounded by 1)

$$|R| \leq 1$$

- (paper) we assume the basis functions to be bounded by 1:

$$|\phi_i(s)| \leq 1 \quad 1 \leq i \leq p$$

- (paper) In order to ensure that rewards are bounded, we impose that  $\|\theta\|_2 \leq 1$

$$\|\theta\|_2 \leq 1$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \dots \\ \theta_p \end{bmatrix} \quad \phi(s) = \begin{bmatrix} \phi_1(s) \\ \phi_2(s) \\ \dots \\ \phi_p(s) \end{bmatrix} \quad R = \theta^T \phi(s) < C \quad (C : \text{constant}, p \text{ is bounded.})$$

# MDP\mathcal{R} set-up

Why "reward( $R$ )", "parameter( $\theta$ )" and "basis function( $\phi(s)$ )" should be bounded?

$$\mathcal{H}_\phi = \{\theta^T \phi(s), \theta \in \mathbb{R}^p\}, \quad R = \theta^T \phi(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \cdots + \theta_p \phi_p(s)$$

- (*paper*)  $R \in \mathbb{R}^S$  the reward function (assumed to be absolutely bounded by 1)

$$|R| \leq 1$$

- (*paper*) we assume the basis functions to be bounded by 1:

$$|\phi_i(s)| \leq 1 \quad 1 \leq i \leq p$$

- (*paper*) In order to ensure that rewards are bounded, we impose that  $\|\theta\|_2 \leq 1$

$$\|\theta\|_2 \leq 1$$

$$\|\theta\|_2 > 1$$

$$\mu^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | s_0 = s, \pi\right] \quad |V^{\pi_E}(s_0) - V^{\tilde{\pi}}(s_0)| = |\theta^T(\mu^{\pi_E}(s_0) - \mu^{\tilde{\pi}}(s_0))|$$

$$\|\mu^{\pi_E}(s_0) - \mu^{\tilde{\pi}}(s_0)\|_2 \leq \epsilon \text{ for some (small) } \epsilon > 0$$

# IRL Algorithm

initial policy  $\rightarrow \mu(s) \rightarrow R = \theta^T \phi(s) \rightarrow \text{policy(LSPI)} \rightarrow \mu(s) \rightarrow R = \theta^T \phi(s) = \rightarrow \text{policy(LSPI)} \dots$

1. Starts with some initial policy  $\pi^{(0)}$  and compute  $\mu^{\pi^{(0)}}(s_0)$ . Set  $j = 1$ ;
2. Compute  $t^{(j)} = \max_{\theta: \|\theta\|_2 \leq 1} \min_{k \in \{0, j-1\}} \theta^T (\mu^{\pi_E}(s_0) - \mu^{\pi^{(k)}}(s_0))$  and let  $\theta^{(j)}$  be the value attaining this maximum. At this step, one searches for the reward function which maximizes the distance between the value of the expert at  $s_0$  and the value of *any* policy computed so far (still at  $s_0$ ). This optimization problem can be solved using a quadratic programming approach or a projection algorithm [1];

[1] Abbeel, Pieter, and Andrew Y. Ng. "Apprenticeship learning via inverse reinforcement learning." *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004.

# IRL Algorithm – Projection method

---

### 3.1. A simpler algorithm

The algorithm described above requires access to a QP (or SVM) solver. It is also possible to change the algorithm so that no QP solver is needed. We will call the previous, QP-based, algorithm the **max-margin** method, and the new algorithm the **projection** method. Briefly, the projection method replaces step 2 of the algorithm with the following:

- Set  $\bar{\mu}^{(i-1)} = \bar{\mu}^{(i-2)} + \frac{(\mu^{(i-1)} - \bar{\mu}^{(i-2)})^T (\mu_E - \bar{\mu}^{(i-2)})}{(\mu^{(i-1)} - \bar{\mu}^{(i-2)})^T (\mu^{(i-1)} - \bar{\mu}^{(i-2)})} (\mu^{(i-1)} - \bar{\mu}^{(i-2)})$   
(This computes the orthogonal projection of  $\mu_E$  onto the line through  $\bar{\mu}^{(i-2)}$  and  $\mu^{(i-1)}$ .)
- Set  $w^{(i)} = \mu_E - \bar{\mu}^{(i-1)}$
- Set  $t^{(i)} = \|\mu_E - \bar{\mu}^{(i-1)}\|_2$

# IRL Algorithm

initial policy  $\rightarrow \mu(s) \rightarrow R = \theta^T \phi(s) \rightarrow \text{policy(LSPI)} \rightarrow \mu(s) \rightarrow R = \theta^T \phi(s) = \rightarrow \text{policy(LSPI)} \dots$

3. if  $t^{(j)} \leq \epsilon$ , terminate. The algorithm outputs a set of policies  $\{\pi^{(0)}, \dots, \pi^{(j-1)}\}$  among which the user chooses manually or automatically the closest to the expert (see [1] for details on how to choose this policy). Notice that the last policy is not necessarily the best (as illustrated in Section 4);

## LSPI(Least Square Policy Iteration)<sup>[2]</sup>: LSTD-Q + Policy Evaluation

4. solve the MDP with the reward function  $R^{(j)}(s) = (\theta^{(j)})^T \phi(s)$  and denote  $\pi^{(j)}$  the associated optimal policy. Compute  $\mu^{\pi^{(j)}}(s_0)$ ;

**state-action LSTD- $\mu$**

[2] Lagoudakis, Michail G., and Ronald Parr. "Least-squares policy iteration." *Journal of machine learning research* 4.Dec (2003): 1107-1149.

# IRL Algorithm

1. Initialize  $\mu^E(s_0)$   $\pi^{(0)}$   $\mu^{(0)}(s_0)$

2. Projection method

$$\bar{\mu}^{(0)} = \mu^{(0)}(s_0) \quad \theta^{(1)} = \mu^E(s_0) - \bar{\mu}^{(0)} \quad t^{(1)} = \|\mu^E(s_0) - \bar{\mu}^{(0)}\|_2$$

3. if  $t^{(j)} \leq \epsilon$ , terminate.

4. Solve MDP using LSPI, get  $\mu^{\pi^{(1)}}(s_0)$  using LSTD- $\mu$ .

$$R^{(1)}(s) = (\theta^{(1)})^T \phi(s) \xrightarrow{\text{LSPI}} \pi^{(1)} \xrightarrow{\text{LSTD- } \mu} \mu^{\pi^{(1)}}(s_0)$$


---

2. Projection method

$$\bar{\mu}^{(1)} = \frac{\text{Set } \bar{\mu}^{(i-1)} = \bar{\mu}^{(i-2)} + \frac{(\mu^{(i-1)} - \bar{\mu}^{(i-2)})^T (\mu_E - \bar{\mu}^{(i-2)})}{(\mu^{(i-1)} - \bar{\mu}^{(i-2)})^T (\mu^{(i-1)} - \bar{\mu}^{(i-2)})} (\mu^{(i-1)} - \bar{\mu}^{(i-2)})}{\theta^{(2)} = \mu^E(s_0) - \bar{\mu}^{(1)}} \quad t^{(2)} = \|\mu^E(s_0) - \bar{\mu}^{(1)}\|_2$$

3. if  $t^{(j)} \leq \epsilon$ , terminate.

4. Solve MDP using LSPI, get  $\mu^{\pi^{(2)}}(s_0)$  using LSTD-mu.

$$R^{(2)}(s) = (\theta^{(2)})^T \phi(s) \xrightarrow{\text{LSPI}} \pi^{(2)} \xrightarrow{\text{LSTD- } \mu} \mu^{\pi^{(2)}}(s_0)$$

# IRL Algorithm – LSPI(Least Square Policy Iteration)

**LSTDQ** ( $D, k, \phi, \gamma, \pi$ ) // Learns  $\hat{Q}^\pi$  from samples

//  $D$  : Source of samples  $(s, a, r, s')$   
 //  $k$  : Number of basis functions  
 //  $\phi$  : Basis functions  
 //  $\gamma$  : Discount factor  
 //  $\pi$  : Policy whose value function is sought

$\tilde{\mathbf{A}} \leftarrow \mathbf{0}$  //  $(k \times k)$  matrix  
 $\tilde{\mathbf{b}} \leftarrow \mathbf{0}$  //  $(k \times 1)$  vector

for each  $(s, a, r, s') \in D$  random action "a"  
 $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top$   
 $\tilde{\mathbf{b}} \leftarrow \tilde{\mathbf{b}} + \phi(s, a)r$

$\tilde{w}^\pi \leftarrow \tilde{\mathbf{A}}^{-1} \tilde{\mathbf{b}}$

return  $\tilde{w}^\pi$

$$\hat{Q}^\pi = \Phi w^\pi$$

**LSPI** ( $D, k, \phi, \gamma, \epsilon, \pi_0$ ) // Learns a policy from samples

//  $D$  : Source of samples  $(s, a, r, s')$   
 //  $k$  : Number of basis functions  
 //  $\phi$  : Basis functions  
 //  $\gamma$  : Discount factor  
 //  $\epsilon$  : Stopping criterion  
 //  $\pi_0$  : Initial policy, given as  $w_0$  (default:  $w_0 = 0$ )

$\pi' \leftarrow \pi_0$  //  $w' \leftarrow w_0$

repeat  
 $\pi \leftarrow \pi'$  //  $w \leftarrow w'$   
 $\pi' \leftarrow \text{LSTDQ } (D, k, \phi, \gamma, \pi)$  //  $w' \leftarrow \text{LSTDQ } (D, k, \phi, \gamma, w)$   
until  $(\pi \approx \pi')$  // until  $(\|w - w'\| < \epsilon)$

return  $\pi$  // return  $w$

# IRL Algorithm – LSPI(Least Square Policy Iteration)

Approximate Q(prediction)

$$\hat{Q}^\pi = \Phi w^\pi \quad \Phi = \begin{pmatrix} \phi(s_1, a_1)^\top \\ \vdots \\ \phi(s, a)^\top \\ \vdots \\ \phi(s_{|\mathcal{S}|}, a_{|\mathcal{A}|})^\top \end{pmatrix}$$

Find  $w^\pi$

$$w^\pi = \left( \Phi^\top \Delta_\mu (\Phi - \gamma \mathbf{P} \Pi_\pi \Phi) \right)^{-1} \Phi^\top \Delta_\mu \mathcal{R}$$

$$w^\pi = A^{-1}b$$

$$\mathbf{A} = \Phi^\top \Delta_\mu (\Phi - \gamma \mathbf{P} \Pi_\pi \Phi) \quad \text{and} \quad b = \Phi^\top \Delta_\mu \mathcal{R}$$

# IRL Algorithm – LSPI(Least Square Policy Iteration)

$$\begin{aligned}
 \mathbf{A} &= \Phi^\top \Delta_\mu (\Phi - \gamma \mathbf{P} \Pi_\pi \Phi) \\
 &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \phi(s, a) \mu(s, a) \left( \phi(s, a) - \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \phi(s', \pi(s')) \right)^\top \\
 &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mu(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \left[ \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top \right]
 \end{aligned}$$

$$\begin{aligned}
 b &= \Phi^\top \Delta_\mu \mathcal{R} \\
 &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \phi(s, a) \mu(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') R(s, a, s') \\
 &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mu(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \left[ \phi(s, a) R(s, a, s') \right] .
 \end{aligned}$$

$$D = \left\{ (s_i, a_i, r_i, s'_i) \mid i = 1, 2, \dots, L \right\}$$

$$\widetilde{\mathbf{A}} = \frac{1}{L} \sum_{i=1}^L \left[ \phi(s_i, a_i) \left( \phi(s_i, a_i) - \gamma \phi(s'_i, \pi(s'_i)) \right)^\top \right] \quad \lim_{L \rightarrow \infty} \widetilde{\mathbf{A}} = \Phi^\top \Delta_{\mu_D} (\Phi - \gamma \mathbf{P} \Pi_\pi \Phi) \quad \text{and} \quad \lim_{L \rightarrow \infty} \widetilde{b} = \Phi^\top \Delta_{\mu_D} \mathcal{R}$$

$$\widetilde{b} = \frac{1}{L} \sum_{i=1}^L \left[ \phi(s_i, a_i) r_i \right] ,$$

# state-action LSTD- $\mu$

Compute  $\mu^{\pi^{(j)}}(s_0)$  LSTD- $\mu$  : to estimate feature expectation of intermediate policies

$$\mu_i^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t \phi_i(s_t) | s_0 = s, \pi\right].$$

$$\mu^\pi(s) = \xi^T \psi(s) \quad \xi \in \mathbb{R}^q$$

$$V^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) | s_0 = s, \pi\right]$$

$$R = \theta^T \phi(s)$$

$$w^\pi = \left(\Phi^\top \Delta_\mu (\Phi - \gamma P \Pi_\pi \Phi)\right)^{-1} \Phi^\top \Delta_\mu \mathcal{R}$$

$$\xi_i^* = \left( \sum_{t=1}^n \psi(s_t, a_t) \left( \psi(s_t, a_t) - \gamma \psi(s'_t, \pi(s'_t)) \right)^T \right)^{-1} \sum_{t=1}^n \psi(s_t, a_t) \phi(s_t, a_t) \quad w_i^* = \left( \sum_{t=1}^n \phi(s_t, a_t) \left( \phi(s_t, a_t) - \gamma \phi(s'_t, \pi(s'_t)) \right)^T \right)^{-1} \sum_{t=1}^n \phi(s_t, a_t) r(s_t, a_t)$$

a set of transitions  $\{(s_t, r_t, s_{t+1})_{1 \leq t \leq n}\}$  sampled according to the policy  $\pi$

$$(\hat{\mu}^\pi(s_0))^T = \psi(s_0)^T (\Psi^T \Delta \Psi)^{-1} \Psi^T \Phi$$

$$\hat{\psi}(s_0) = \frac{\sum_{i=1, a \sim \text{random}}^N \psi(s_0, a)}{N}$$

# IRL Algorithm – LSPI(Least Square Policy Iteration)

**LSTDQ** ( $D, k, \phi, \gamma, \pi$ ) // Learns  $\hat{Q}^\pi$  from samples

```
// D : Source of samples ( $s, a, r, s'$ )
// k : Number of basis functions
//  $\phi$  : Basis functions
//  $\gamma$  : Discount factor
//  $\pi$  : Policy whose value function is sought
```

 $\tilde{\mathbf{A}} \leftarrow \mathbf{0}$  // ( $k \times k$ ) matrix
 $\tilde{\mathbf{b}} \leftarrow \mathbf{0}$  // ( $k \times 1$ ) vector

```
for each  $(s, a, r, s') \in D$  random action "a"
     $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top$ 
     $\tilde{\mathbf{b}} \leftarrow \tilde{\mathbf{b}} + \phi(s, a)r$ 
```

 $\tilde{w}^\pi \leftarrow \tilde{\mathbf{A}}^{-1} \tilde{\mathbf{b}}$ 
**return**  $\tilde{w}^\pi$ 

$$\hat{Q}^\pi = \Phi w^\pi$$

**LSPI** ( $D, k, \phi, \gamma, \epsilon, \pi_0$ ) // Learns a policy from samples

```
// D : Source of samples ( $s, a, r, s'$ )
// k : Number of basis functions
//  $\phi$  : Basis functions
//  $\gamma$  : Discount factor
//  $\epsilon$  : Stopping criterion
//  $\pi_0$  : Initial policy, given as  $w_0$  (default:  $w_0 = 0$ )
```

 $\pi' \leftarrow \pi_0$  //  $w' \leftarrow w_0$ 
**repeat**
 $\pi \leftarrow \pi'$  //  $w \leftarrow w'$ 
 $\pi' \leftarrow \text{LSTDQ } (D, k, \phi, \gamma, \pi)$  //  $w' \leftarrow \text{LSTDQ } (D, k, \phi, \gamma, w)$ 
**until** ( $\pi \approx \pi'$ ) // until ( $\|w - w'\| < \epsilon$ )

**return**  $\pi$  // **return**  $w$

# LSTD vs. LSTD-Q

LSTD for estimating  $\hat{v} \approx v_\pi$  ( $O(d^2)$  version)

Input: feature representation  $\mathbf{x}(s) \in \mathbb{R}^d$ , for all  $s \in \mathcal{S}$ ,  $\mathbf{x}(\text{terminal}) \doteq \mathbf{0}$

$$\widehat{\mathbf{A}}^{-1} \leftarrow \varepsilon^{-1} \mathbf{I}$$

$$\widehat{\mathbf{b}} \leftarrow \mathbf{0}$$

Repeat (for each episode):

Initialize  $S$ ; obtain corresponding  $\mathbf{x}$

Repeat (for each step of episode):

Choose  $A \sim \pi(\cdot | S)$  **behavior from target policy**

Take action  $A$ , observe  $R, S'$ ; obtain corresponding  $\mathbf{x}'$

$$\begin{aligned} \mathbf{v} &\leftarrow \widehat{\mathbf{A}}^{-1}^\top (\mathbf{x} - \gamma \mathbf{x}') \\ \widehat{\mathbf{A}}^{-1} &\leftarrow \widehat{\mathbf{A}}^{-1} - (\widehat{\mathbf{A}}^{-1} \mathbf{x}) \mathbf{v}^\top / (1 + \mathbf{v}^\top \mathbf{x}) \end{aligned}$$

$$\widehat{\mathbf{b}} \leftarrow \widehat{\mathbf{b}} + R \mathbf{x}$$

$$\theta \leftarrow \widehat{\mathbf{A}}^{-1} \widehat{\mathbf{b}}$$

$$S \leftarrow S'; \mathbf{x} \leftarrow \mathbf{x}'$$

until  $S'$  is terminal

on-policy

target = behavior

**LSTDQ** ( $D, k, \phi, \gamma, \pi$ )

// Learns  $\widehat{Q}^\pi$  from samples

//  $D$  : Source of samples  $(s, a, r, s')$

//  $k$  : Number of basis functions

//  $\phi$  : Basis functions

//  $\gamma$  : Discount factor

//  $\pi$  : Policy whose value function is sought

$$\widetilde{\mathbf{A}} \leftarrow \mathbf{0} \quad // (k \times k) \text{ matrix}$$

$$\widetilde{\mathbf{b}} \leftarrow \mathbf{0} \quad // (k \times 1) \text{ vector}$$

for each  $(s, a, r, s') \in D$  **random generation**

$$\widetilde{\mathbf{A}} \leftarrow \widetilde{\mathbf{A}} + \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top$$

$$\widetilde{\mathbf{b}} \leftarrow \widetilde{\mathbf{b}} + \phi(s, a) r$$

$$\widetilde{\mathbf{w}}^\pi \leftarrow \widetilde{\mathbf{A}}^{-1} \widetilde{\mathbf{b}}$$

return  $\widetilde{\mathbf{w}}^\pi$

off-policy

$$\widehat{Q}^\pi = \Phi w^\pi$$

## LSTD

Learns the **state** value function  $V^\pi$

Basis functions of **state**

Samples **cannot** be reused **on-policy**

Training samples collected **by**  $\pi$  **on-policy**

Sampling distribution **cannot** be controlled

Bias by the **stationary distribution of**  $\pi$

## LSTDQ

Learns the **state-action** value function  $Q^\pi$

Basis functions of **state and action**

Samples **can** be reused **off-policy**

Training samples collected **arbitrarily**

Sampling distribution **can** be controlled

Bias by the **sampling distribution**

}

# Backup slide : on-policy vs. off-policy

---

When the model of environment (such as transition probability) is unknown, the alternative way is depending on “sample experience”.

For episodic task, experience = “multiple episodes”

Obviously, a policy can make episodes(sample experience).

On-policy :

policy being evaluated and improved = policy generating data  
(therefore, need  $\varepsilon$ -greedy for exploration)

Off-policy :

policy being evaluated and improved(target policy)  $\neq$  policy generating data(behavior policy)

# Why LSTD- $\mu$ is off-policy method?

$$\begin{aligned}\mathbf{A} &= \Phi^\top \Delta_\mu (\Phi - \gamma \mathbf{P} \Pi_\pi \Phi) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \phi(s, a) \mu(s, a) \left( \phi(s, a) - \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \phi(s', \pi(s')) \right)^\top \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mu(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \left[ \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top \right]\end{aligned}$$

**for each**  $(s, a, r, s') \in D$

$$\begin{aligned}\tilde{\mathbf{A}} &\leftarrow \tilde{\mathbf{A}} + \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top \\ \tilde{b} &\leftarrow \tilde{b} + \phi(s, a) r\end{aligned}$$

$$\begin{aligned}b &= \Phi^\top \Delta_\mu \mathcal{R} \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \phi(s, a) \mu(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') R(s, a, s') \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mu(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \left[ \phi(s, a) R(s, a, s') \right] .\end{aligned}$$

TD(0) or one-step TD update rule

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$$

Additional degree of freedom ( $a_0 = a$ ) allows off-policy learning.(LSTD-Q)  
 LSTD-Q (Q-function)  $\rightarrow$  LSTD- $\mu$  (state-action feature expectation)

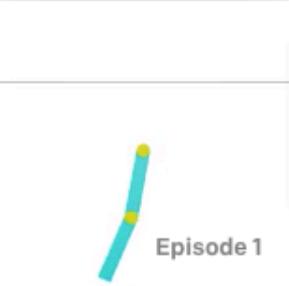
TD의 경우 behavior policy가 필요하다. 그래야  $a$ 를 해서  $s'$ 를 얻을 수 있으니까  
 그래서 평가되고 향상되는 policy와 sampling (behavior) policy 가 같으면 on-policy  
 서로 다르거나, 랜덤하게 제시해주는 경우 off-policy 이다!

# Implementation & Experiment

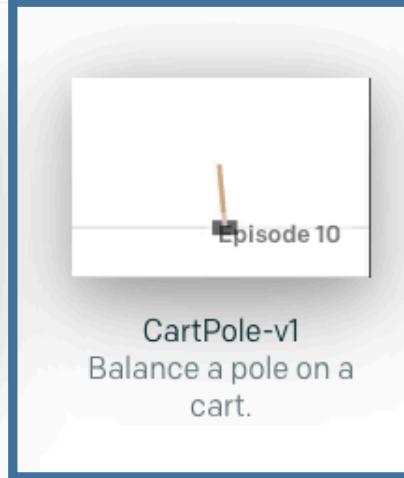
# Experiment Setting

## Classic control

Control theory problems from the classic RL literature.



Acrobot-v1  
 Swing up a two-link robot.



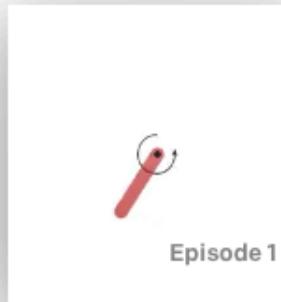
CartPole-v1  
 Balance a pole on a cart.



MountainCar-v0  
 Drive up a big hill.



MountainCarContinuous-v0  
 Drive up a big hill with continuous control.



Pendulum-v0  
 Swing up a pendulum.

## CartPole

state dimension : 4  
# of actions : 2

# My LSPI's Critical Problem

**LSPI** ( $D, k, \phi, \gamma, \epsilon, \pi_0$ ) // Learns a policy from samples

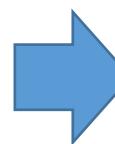
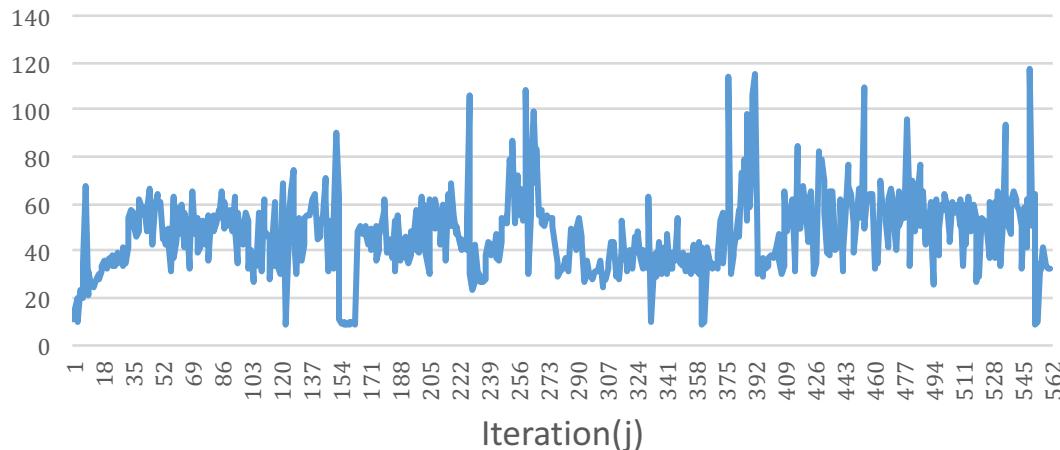
```
// D : Source of samples ( $s, a, r, s'$ )
// k : Number of basis functions
//  $\phi$  : Basis functions
//  $\gamma$  : Discount factor
//  $\epsilon$  : Stopping criterion
//  $\pi_0$  : Initial policy, given as  $w_0$  (default:  $w_0 = 0$ )  
  

 $\pi' \leftarrow \pi_0$  //  $w' \leftarrow w_0$   
  

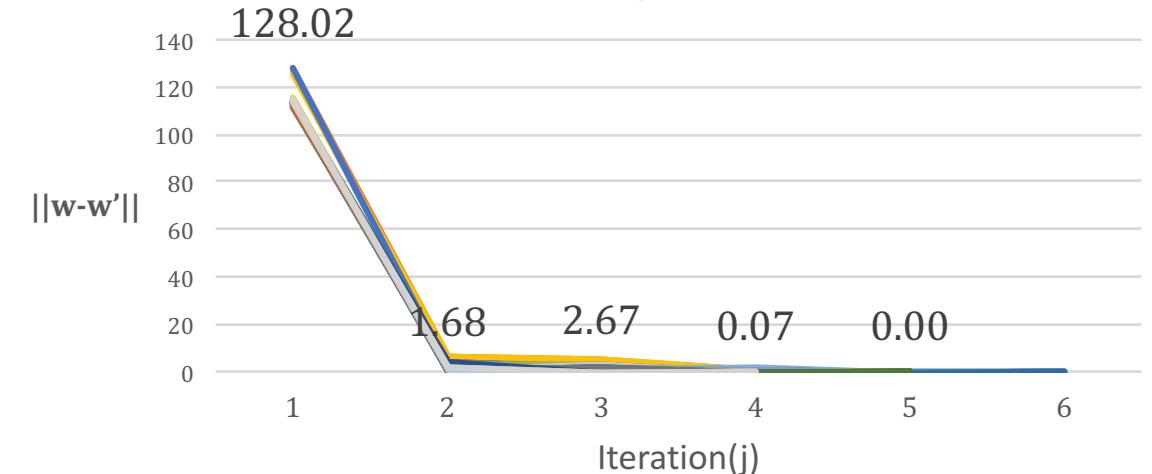
repeat
     $\pi \leftarrow \pi'$  //  $w \leftarrow w'$ 
     $\pi' \leftarrow \text{LSTDQ} (D, k, \phi, \gamma, \pi)$  //  $w' \leftarrow \text{LSTDQ} (D, k, \phi, \gamma, w)$ 
until ( $\pi \approx \pi'$ ) // until ( $\|w - w'\| < \epsilon$ )  
  

return  $\pi$  // return  $w$ 
```

Previous LSPI



LSPI Convergence



loop\_budget : 20

# My LSPI's stabilization

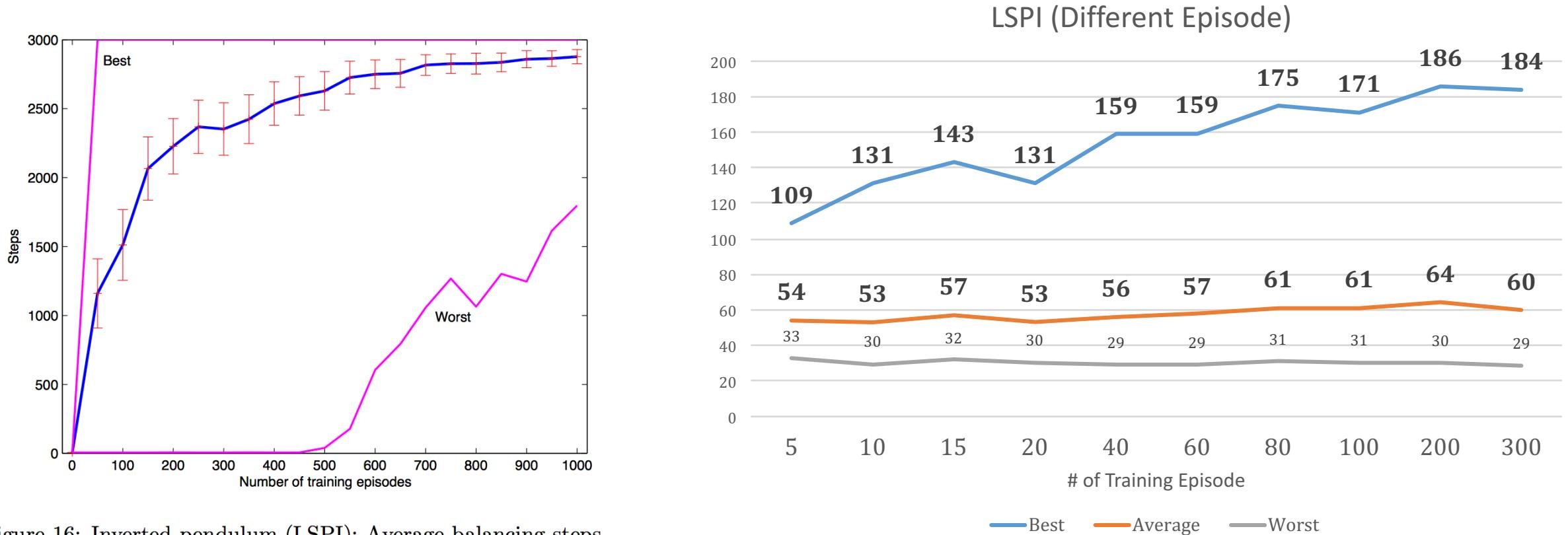


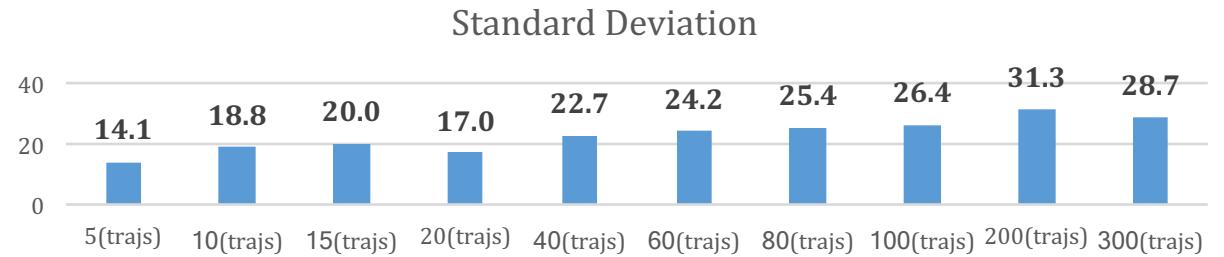
Figure 16: Inverted pendulum (LSPI): Average balancing steps.

## Experiment Setting

# Evaluation for each policy( $w^\pi$ ) : 1000

# Trials(Experiments) : 100

basis function : RBF



# Backup slide : Radial Basis Functions

RBFs are the natural generalization of coarse coding(binary features {0, 1}) to continuous-valued features([0, 1]).

A typical RBF feature has a Gaussian response dependent on the distance between the state,  $s$ , and the feature's center state,  $c_i$ , and the feature's width,  $\sigma_i$ .

$$x_i(s) \doteq \exp\left(-\frac{\|s - c_i\|^2}{2\sigma_i^2}\right)$$

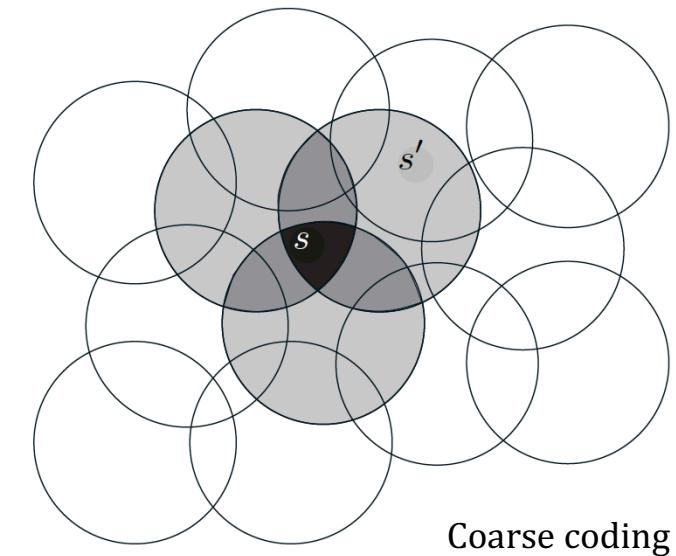
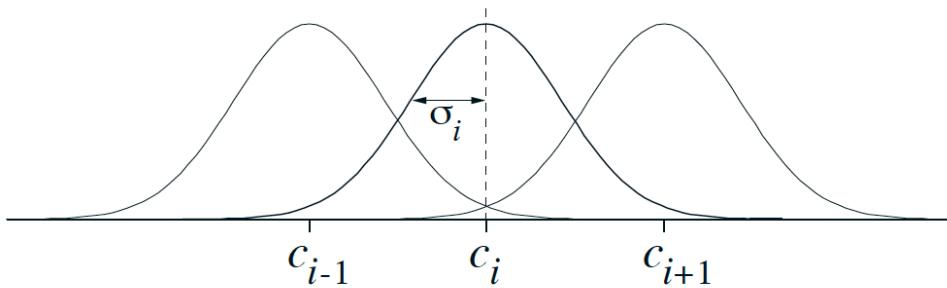


Figure 9.13: One-dimensional radial basis functions.

Advantage : vary smoothly, and differentiable.

# Basis function (RBF; gaussian)

N : dimension of states

A : # of actions

P : # of basis features

$$s = (s_1, s_2, \dots, s_N)^T$$

$$\mu_j = [\mu_{j,1}, \mu_{j,2}, \dots, \mu_{j,N}]^T, j \in \{1, 2, \dots, P\}$$

$\mu_{j,i}$  is a random float number in (-1, 1)

Radial Basis Functions;

given discrete true action  $\hat{a}$ ,

$$\phi_j(s) = \exp(-\|s - \mu_j\|^2) = \exp\left(-\frac{\|s - \mu_j\|^2}{2\sigma^2}\right) (\sigma = 1, \text{skip } \frac{1}{2})$$

$$\phi(s) = [1, \phi_1(s), \phi_2(s), \dots, \phi_P(s)] \in R^P$$

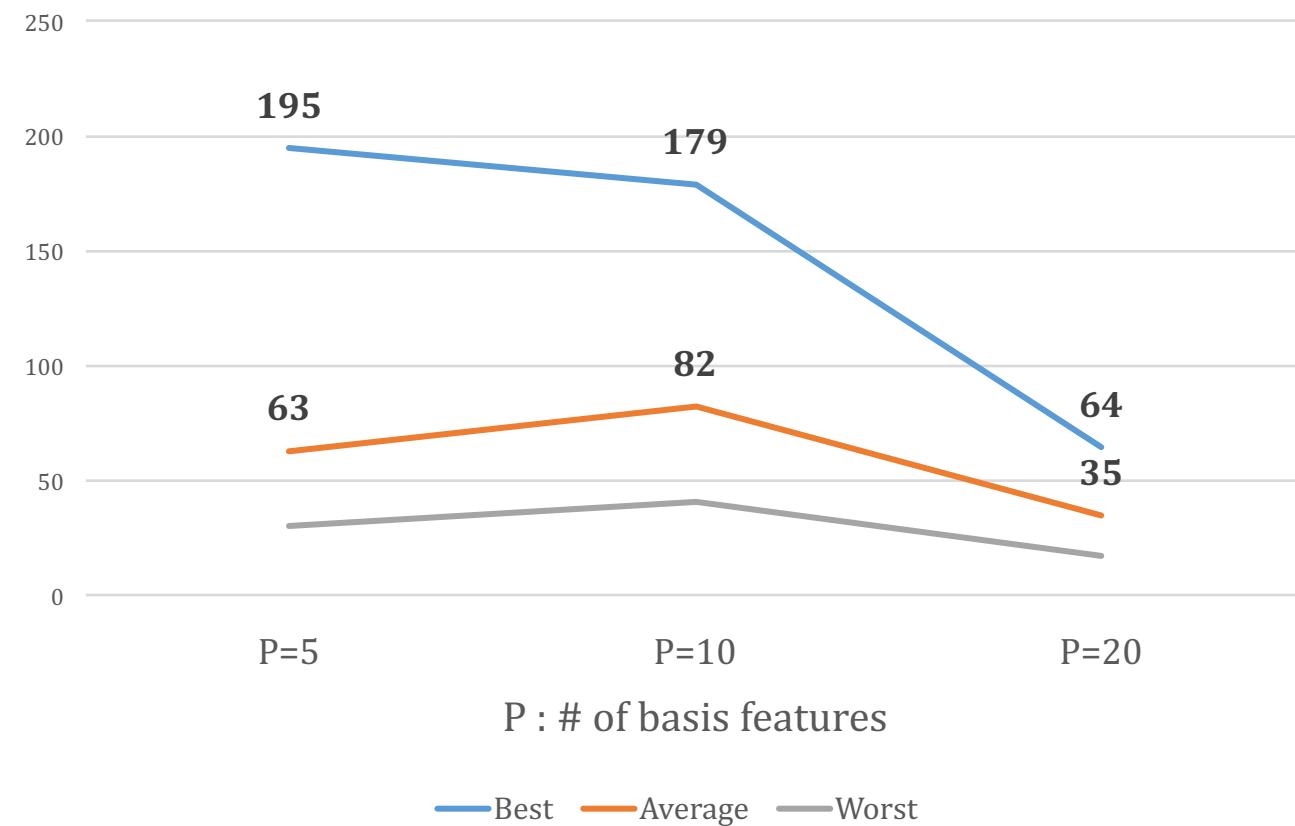
$$\Phi(s; a = \hat{a}) = [\underbrace{[0, 0, \dots, 0]}_{a \neq 0}, \dots, \underbrace{\phi(s)}_{a = \hat{a}}, \dots, \underbrace{[0, 0, \dots, 0]}_{a \neq A-1}] \in R^{A \times P}$$

# Episode for training  $w^\pi$ (LSTD-Q) : 100 (2000 samples)

# Evaluation for each policy( $w^\pi$ ) : 1000

# Trials(Experiments) : 100

Gaussian, different # of basis function



Basis Function Dimension : (A(=2) X P(=5,10,20)) → 10,20,40

# RBF

**P : # of basis features**

$$t = \|w^{\pi_{j-1}} - w^{\pi_j}\|_2$$

```

LSPI ( $D, k, \phi, \gamma, \epsilon, \pi_0$ ) // Learns a policy from samples
    //  $D$  : Source of samples  $(s, a, r, s')$ 
    //  $k$  : Number of basis functions
    //  $\phi$  : Basis functions
    //  $\gamma$  : Discount factor
    //  $\epsilon$  : Stopping criterion
    //  $\pi_0$  : Initial policy, given as  $w_0$  (default:  $w_0 = 0$ )
     $\pi' \leftarrow \pi_0$  //  $w' \leftarrow w_0$ 
    repeat
         $\pi \leftarrow \pi'$  //  $w \leftarrow w'$ 
         $\pi' \leftarrow \text{LSTDQ} (D, k, \phi, \gamma, \pi)$  //  $w' \leftarrow \text{LSTDQ} (D, k, \phi, \gamma, w)$ 
    until  $(\pi \approx \pi')$  // until  $(\|w - w'\| < \epsilon)$ 
    return  $\pi$  // return  $w$ 

```

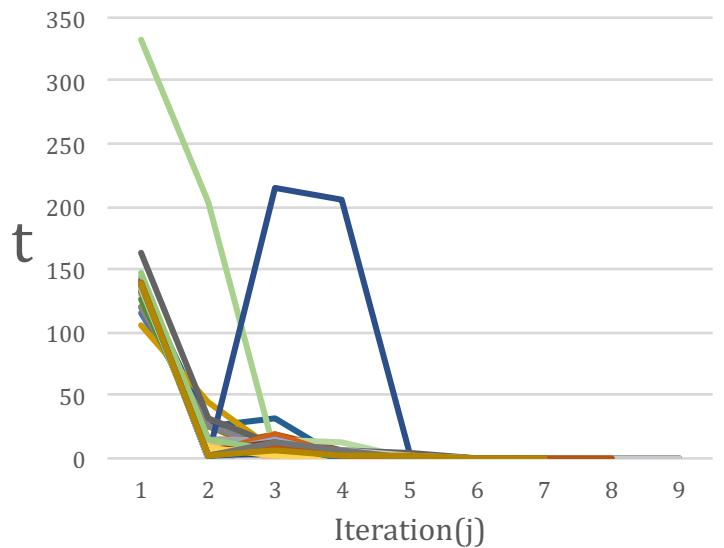
  

```

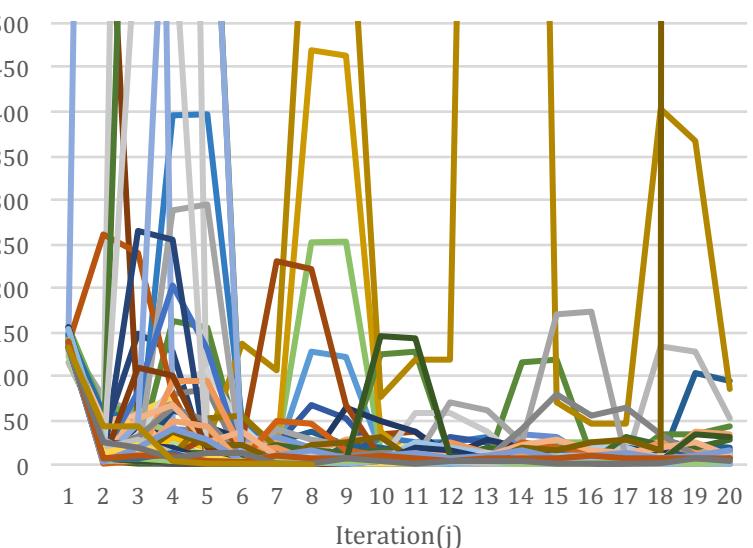
LSTDQ ( $D, k, \phi, \gamma, \pi$ ) // Learns  $\hat{Q}^\pi$  from samples
    //  $D$  : Source of samples  $(s, a, r, s')$ 
    //  $k$  : Number of basis functions
    //  $\phi$  : Basis functions
    //  $\gamma$  : Discount factor
    //  $\pi$  : Policy whose value function is sought
     $\tilde{\mathbf{A}} \leftarrow \mathbf{0}$  //  $(k \times k)$  matrix
     $\tilde{\mathbf{b}} \leftarrow \mathbf{0}$  //  $(k \times 1)$  vector
    for each  $(s, a, r, s') \in D$ 
         $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s, a) (\phi(s, a) - \gamma \phi(s', \pi(s')))^\top$ 
         $\tilde{\mathbf{b}} \leftarrow \tilde{\mathbf{b}} + \phi(s, a)r$ 
     $\tilde{w}^\pi \leftarrow \tilde{\mathbf{A}}^{-1}\tilde{\mathbf{b}}$ 
    return  $\tilde{w}^\pi$ 

```

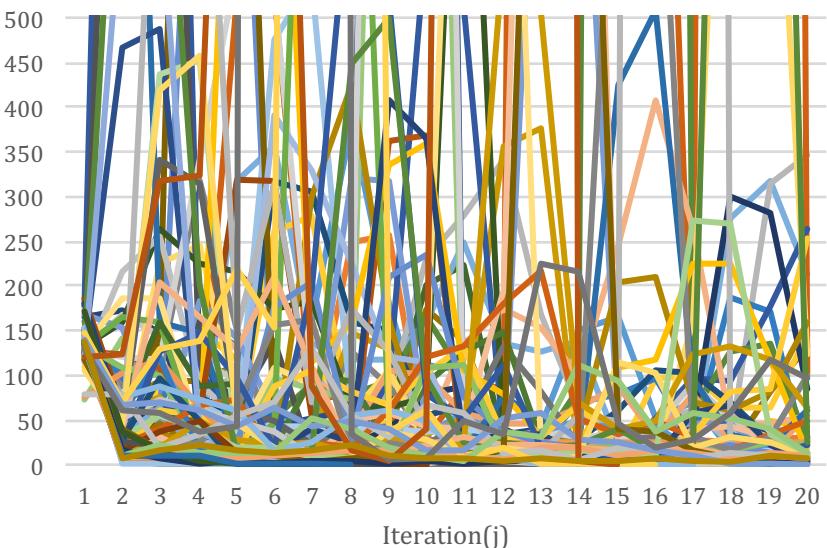
P=5 Convergence



P=10 Convergence



P=20 convergence



# Basis function (Sigmoid)

N : dimension of states

A : # of actions

P : # of basis features

$$s = (s_1, s_2, \dots, s_N)^T$$

$$\mu_j = [\mu_{j,1}, \mu_{j,2}, \dots, \mu_{j,N}]^T, j \in \{1, 2, \dots, P\}$$

$\mu_{j,i}$  is a random float number in (-1, 1)

sigmoid;

given discrete true action  $\hat{a}$ ,

$$\sigma(a) = \frac{1}{1+exp(-a)}$$

$$\phi_j(s) = \sigma(\|s - \mu_j\|_2)$$

$$\phi(s) = [1, \phi_1(s), \phi_2(s), \dots, \phi_P(s)] \in R^P$$

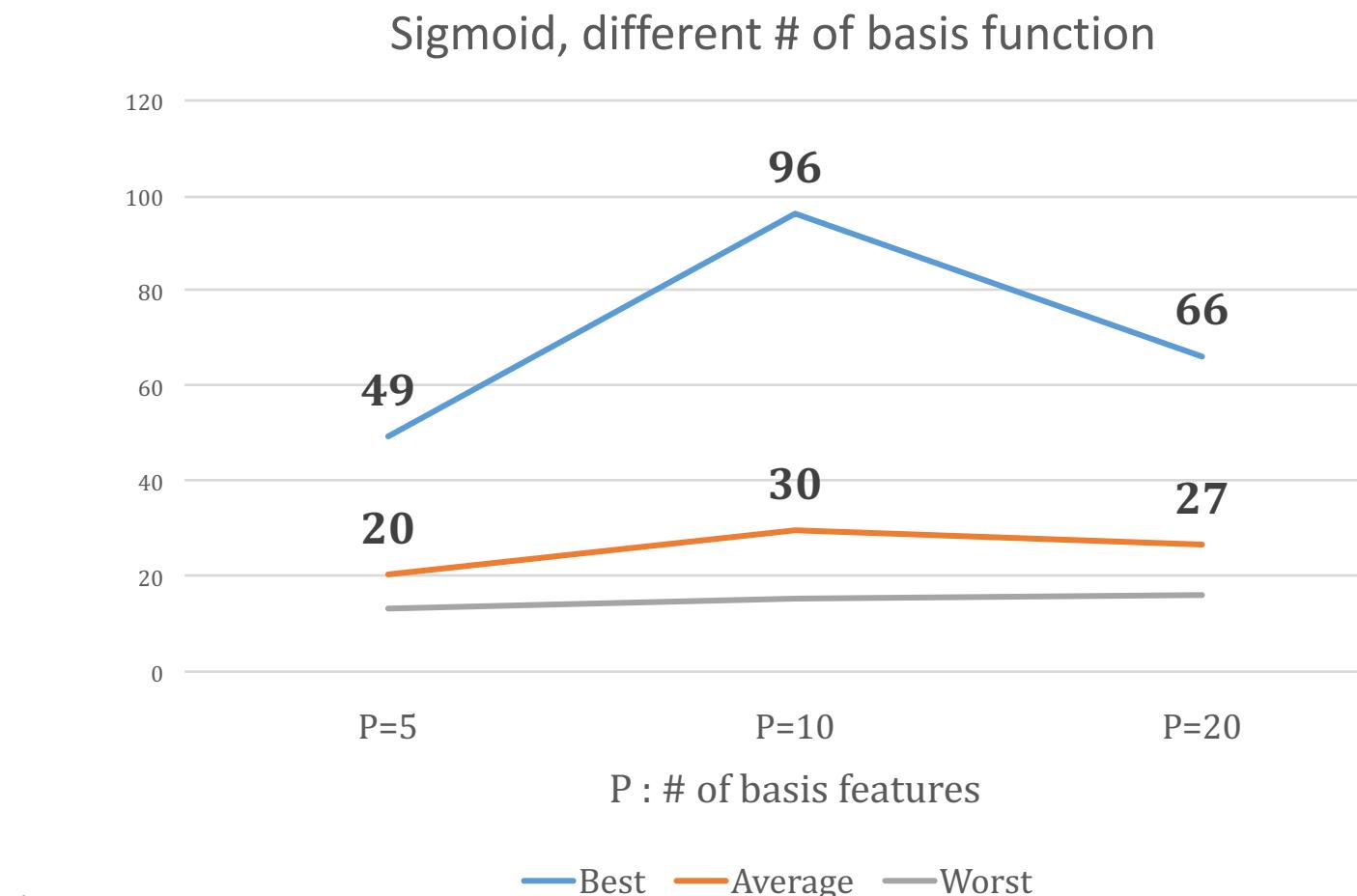
$$\Phi(s; a = \hat{a}) = [\underbrace{[0, 0, \dots, 0]}_{a \neq 0}, \dots, \underbrace{\phi(s)}_{a = \hat{a}}, \dots, \underbrace{[0, 0, \dots, 0]}_{a \neq A-1}] \in R^{A \times P}$$

# Episode for training  $w^\pi$ (LSTD-Q) : 100 (2000 samples)

# Evaluation for each policy( $w^\pi$ ) : 1000

# Trials(Experiments) : 100

# basis feature(P) = {5, 10, 20}

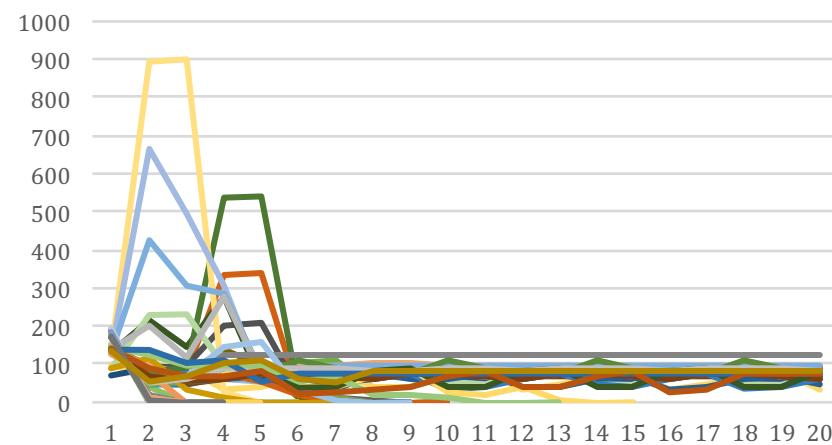


Basis Function Dimension : (A(=2) X P(=5,10,20)) → 10,20,40

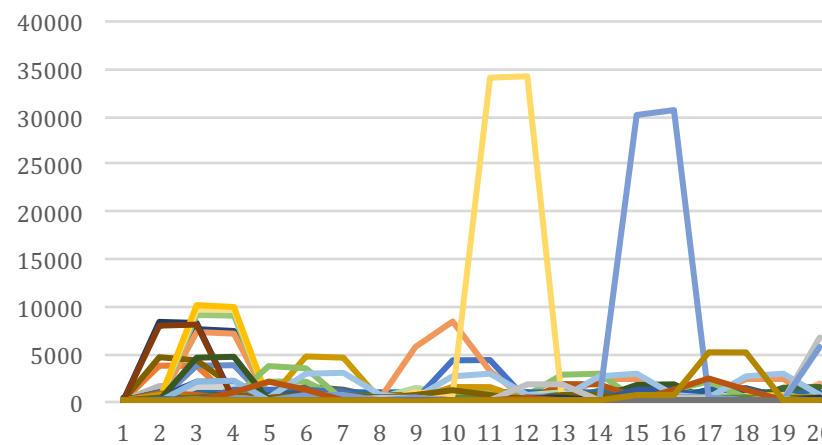
# Basis function (Sigmoid)

$$t = \|w^{\pi_{j-1}} - w^{\pi_j}\|_2$$

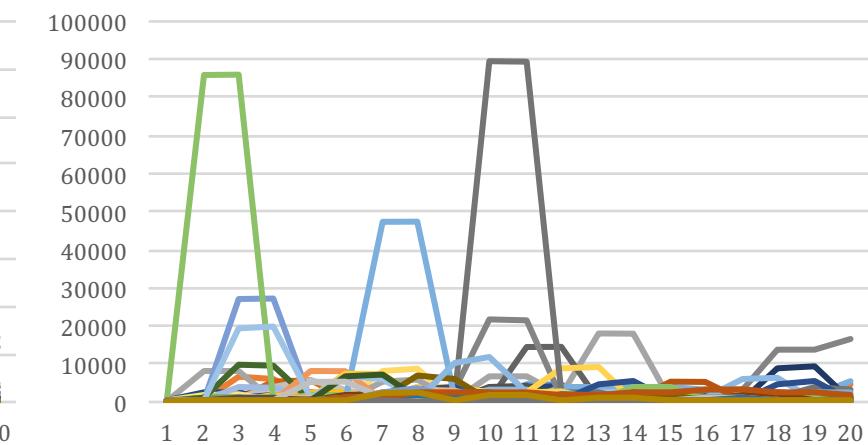
P=5



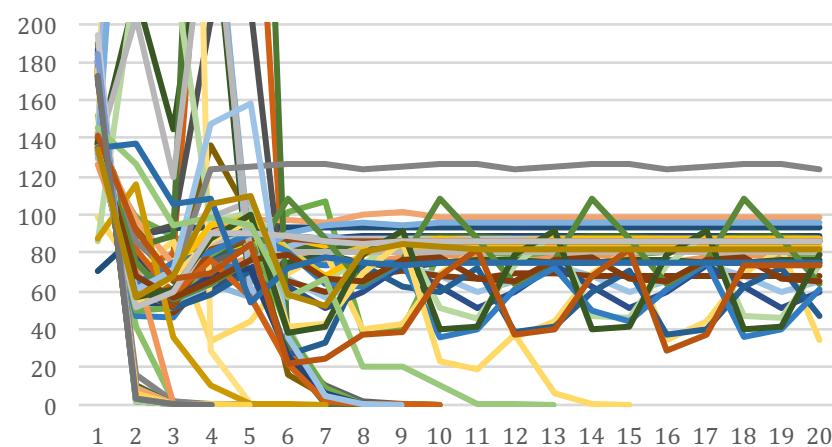
P=10



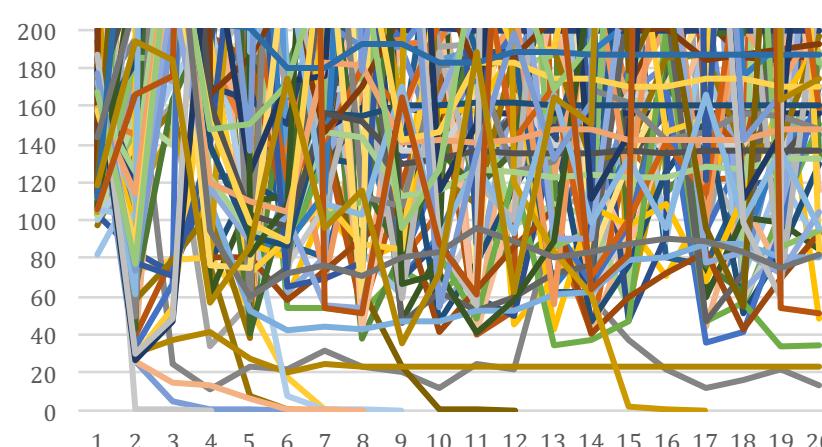
P=20



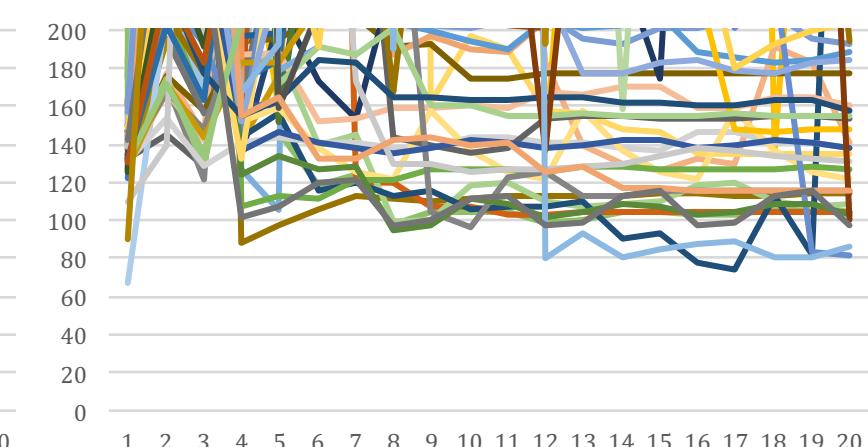
P=5 (0~200)



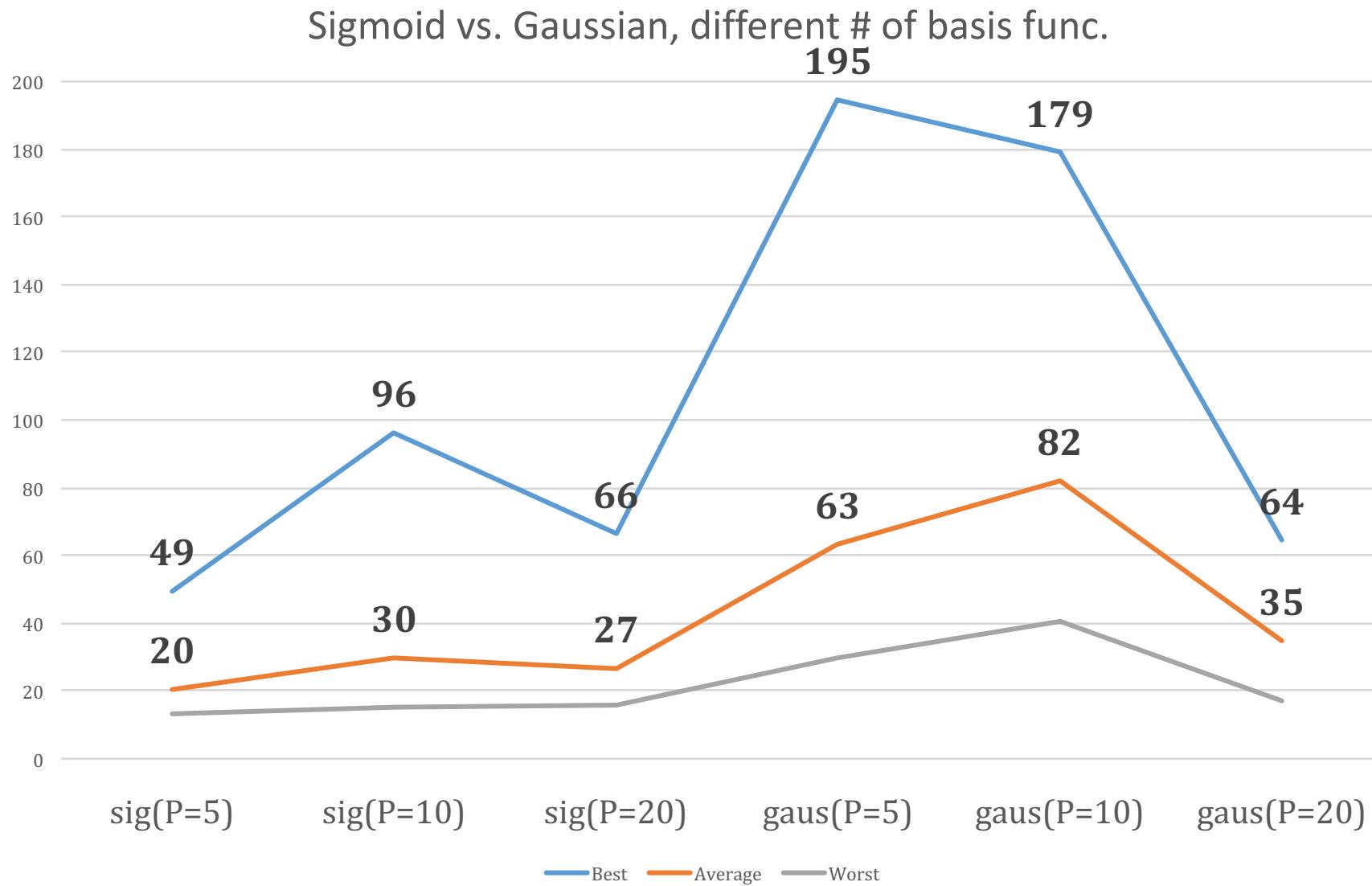
P=10 (0~200)



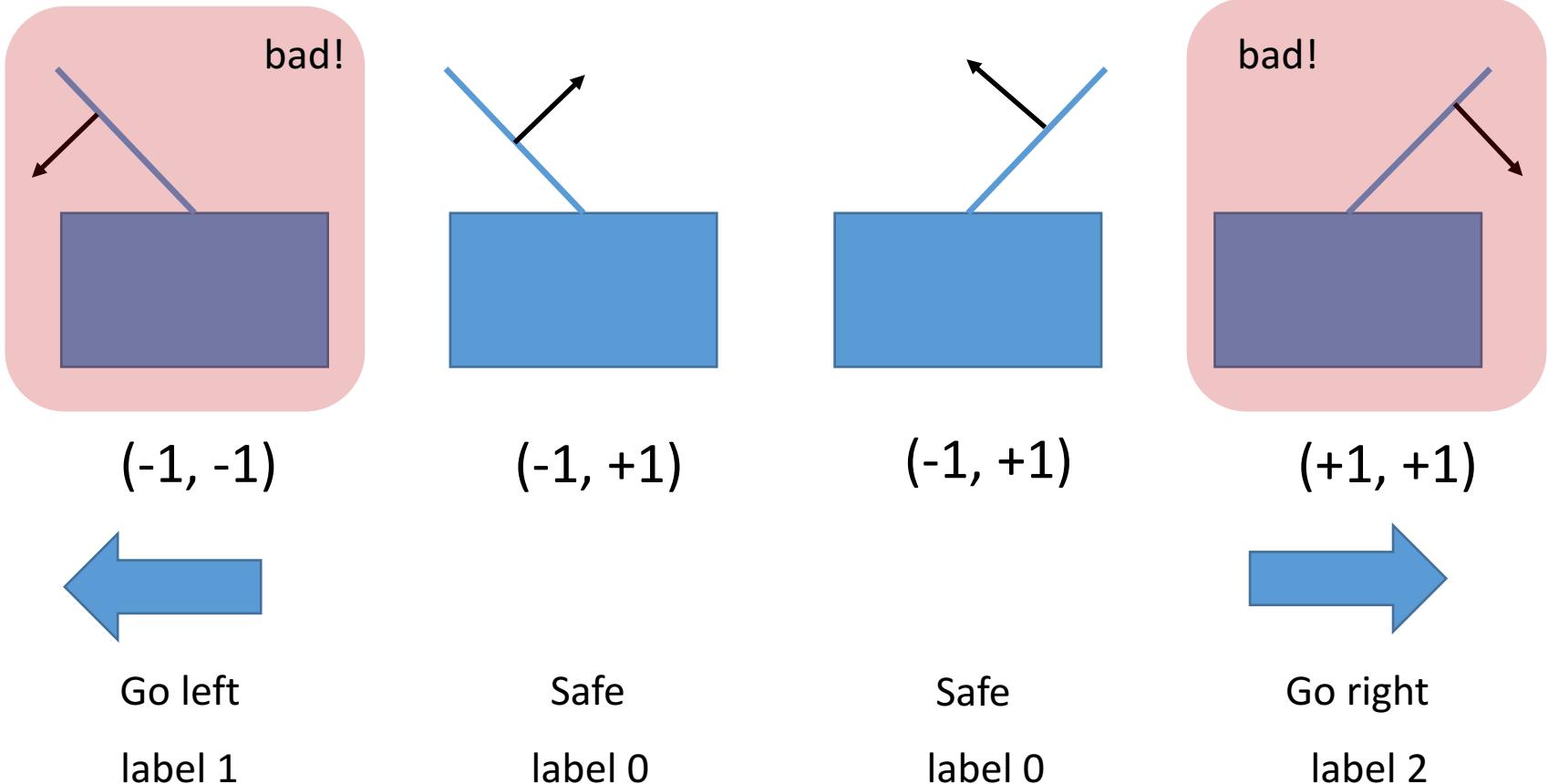
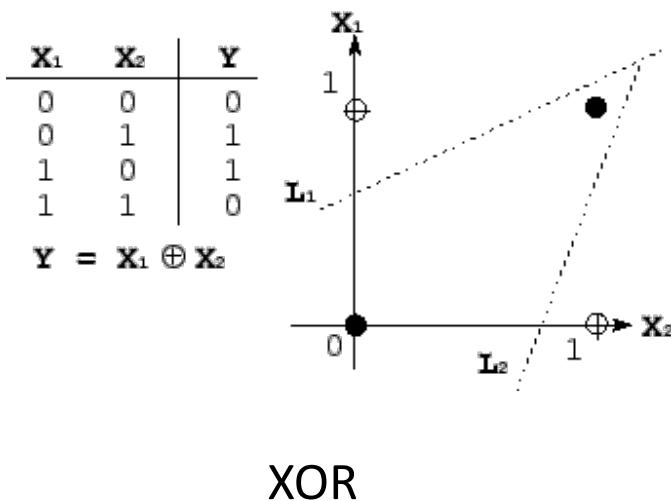
P=20 (0~200)



# Basis function (Gaussian vs. Sigmoid)



# CartPole Problem Analysis (manual policy)

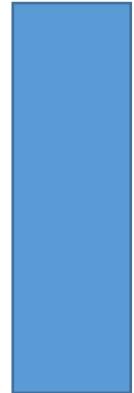


We need ANN with activation function!

# ANN for CartPole

## ANN for CartPole

Input : 4



$W1 : 4 \times 2$   
 $b1 : 2$



$h1 : 2$



$W2 : 2 \times 3$   
 $b2 : 3$

Softmax

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

output : 3

$$\begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

*softmax\_cross\_entropy*

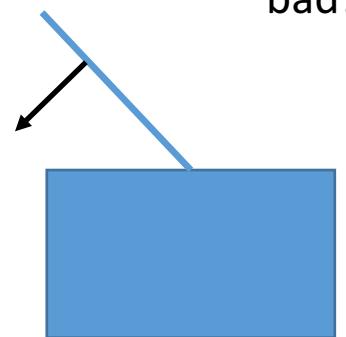
$$\text{loss} = - \sum_i^N a_t \log(\hat{a}_t)$$

label 0(safe)

label 1(go left)

label 2(go right)

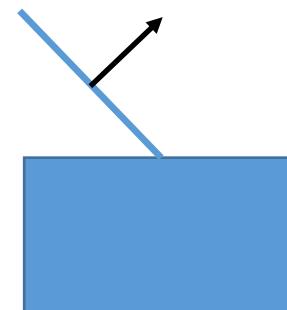
bad!



P

0.01  
0.98  
0.01

left



P

0.98  
0.01  
0.01

Safe

(-1, +1)

bad!



P

0.01  
0.01  
0.98

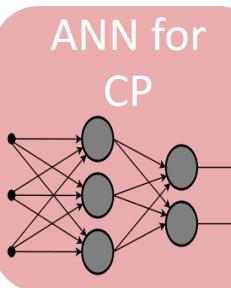
right

(+1, +1)

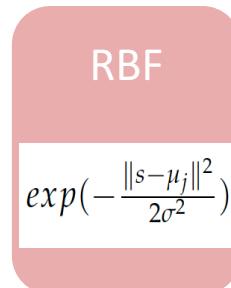
# Basis function (ANN for CartPole + RBF)

ANN for CartPole + RBF

$$S = \begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ S_4 \end{bmatrix}$$

ANN for CP


$$P = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

RBF


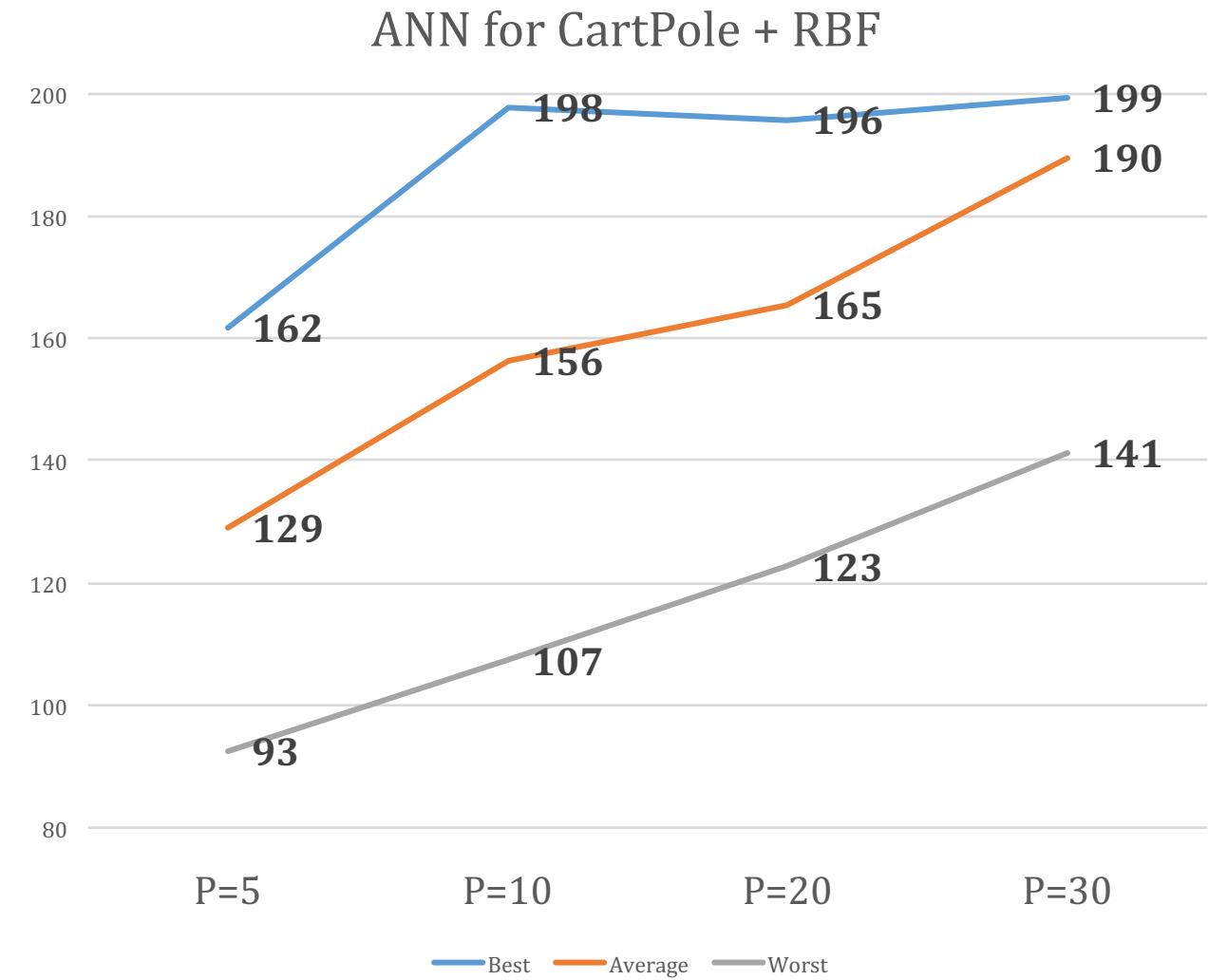
$$\phi_j(P)$$

# Episode for training  $w^\pi$ (LSTD-Q) : 100 (2000 samples)

# Evaluation for each policy( $w^\pi$ ) : 1000

# Trials(Experiments) : 100

# basis feature(P) = {5, 10, 20, 30}

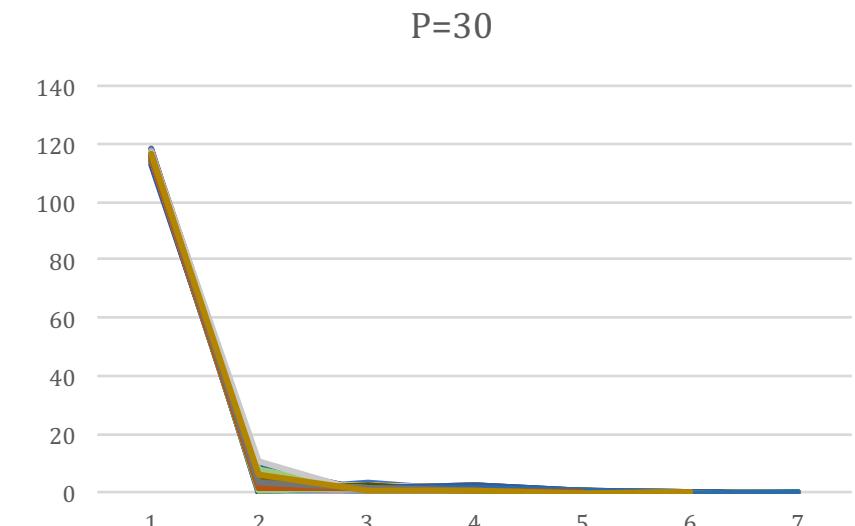
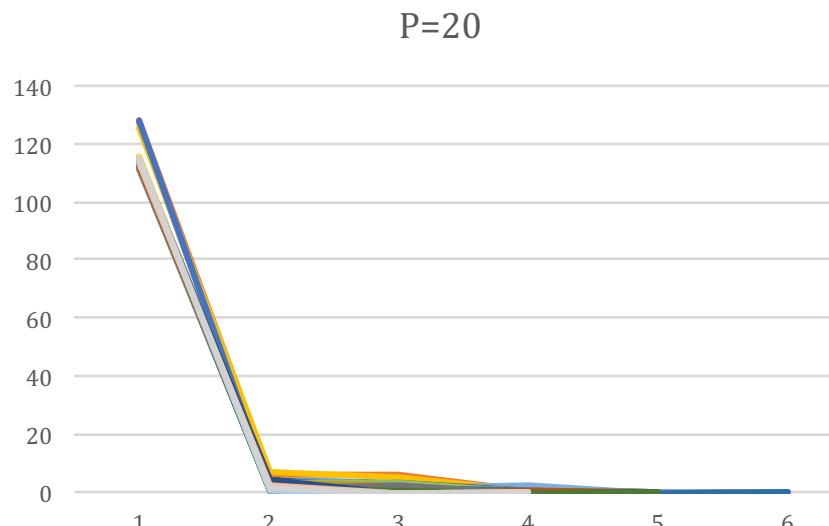
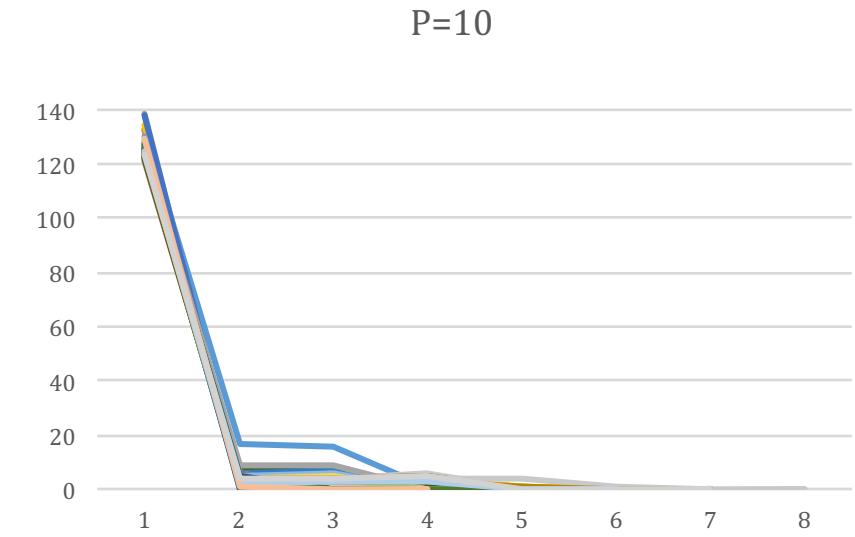
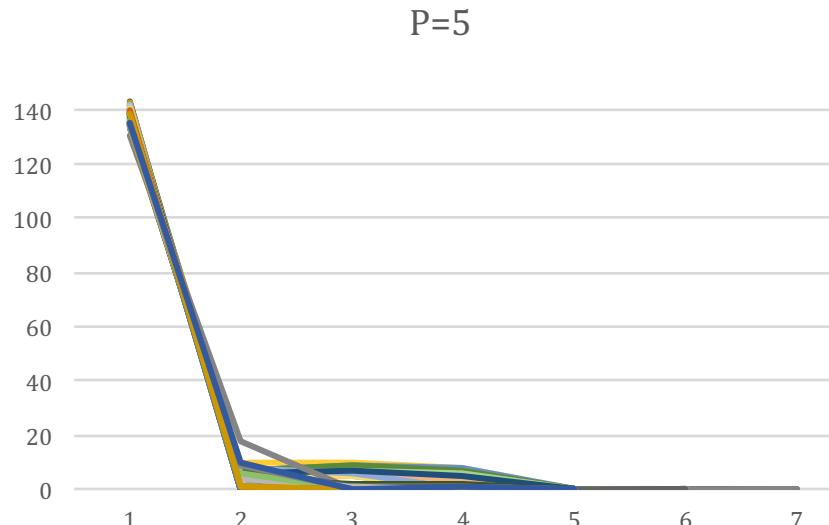


# Basis function (ANN for CartPole + RBF)

P : # of basis features

# Trials(Experiments) : 100

$$t = \|w^{\pi_{j-1}} - w^{\pi_j}\|_2$$



# Experiment Setting

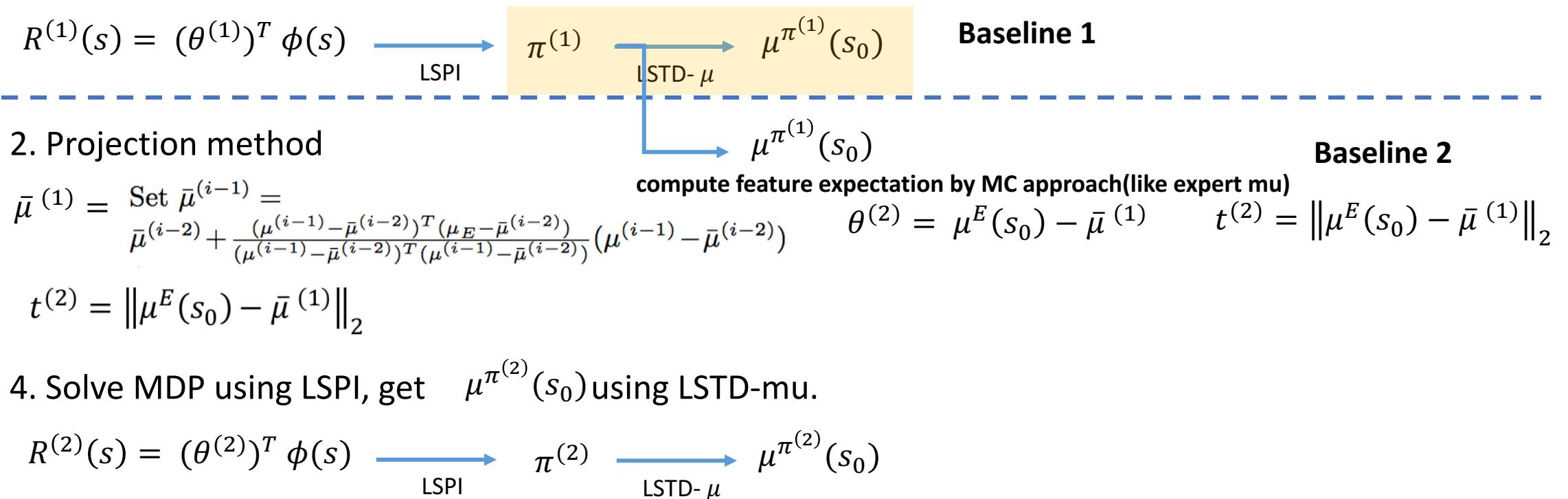
1. Initialize  $\mu^E(s_0)$   $\pi^{(0)}$   $\mu^{(0)}(s_0)$

2. Projection method

$$\bar{\mu}^{(0)} = \mu^{(0)}(s_0) \quad \theta^{(1)} = \mu^E(s_0) - \bar{\mu}^{(0)} \quad t^{(1)} = \|\mu^E(s_0) - \bar{\mu}^{(0)}\|_2$$

$$t^{(2)} = \|\mu^E(s_0) - \bar{\mu}^{(1)}\|_2$$

4. Solve MDP using LSPI, get  $\mu^{\pi^{(1)}}(s_0)$  using LSTD- $\mu$ .



# IRL + MC approach

# of Expert Episode : 100 ( $100 \times 200 = 20000$  samples( $s, a, s', r$ ))

# of Episode for training  $w^\pi$ (LSTD-Q) : 100 ( $100 \times 20 = 2000$  samples)

# of Episode for calculating  $\mu^\pi$  : 100 ( $100 \times ?$  samples)

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

Reward basis function : ANN cartpole + RBF

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



# IRL + MC approach

# of Expert Episode : **1 (= 200 samples ( $s, a, s', r$ ))**

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

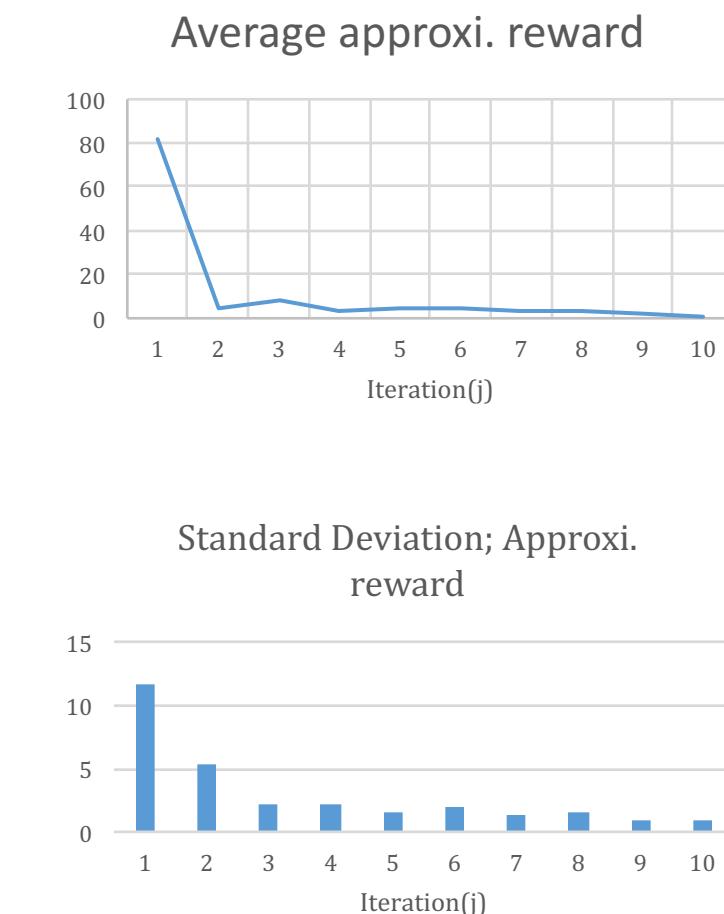
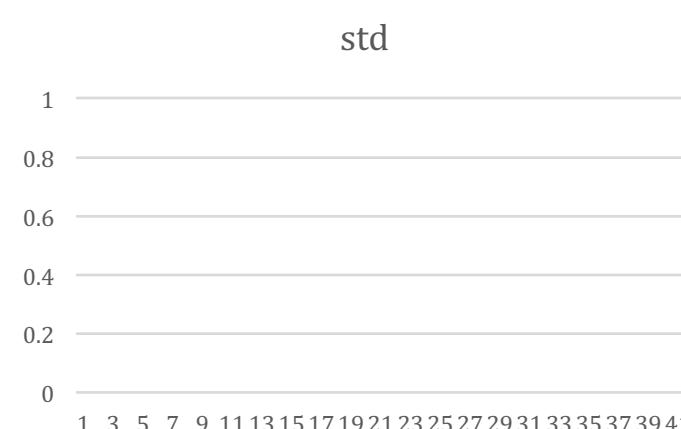
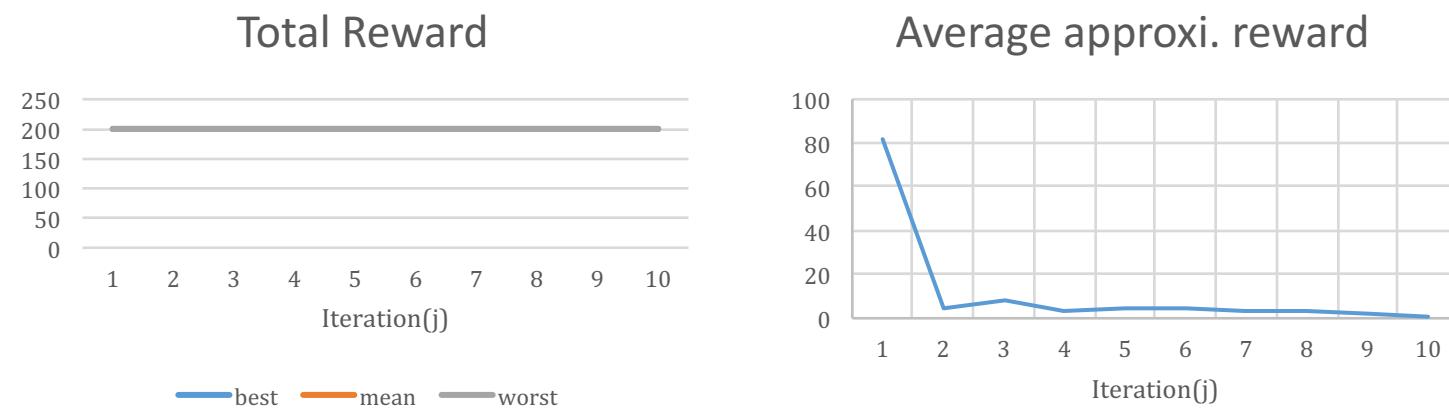
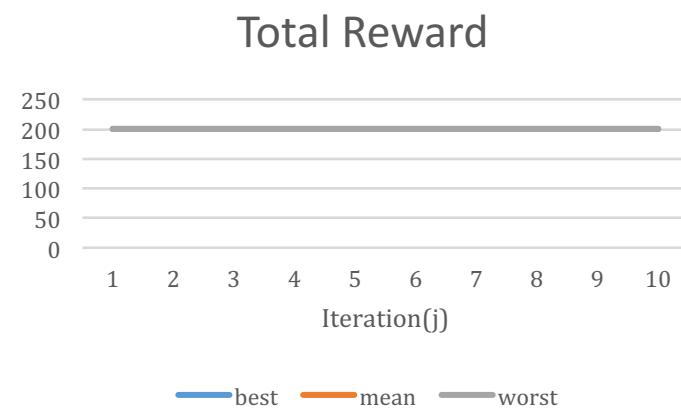
Reward basis function : ANN cartpole + RBF

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



→ 200 expert trajectories are enough!

# IRL + MC approach

for each  $(s, a, r, s') \in D$

$$\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s, a) (\phi(s, a) - \gamma \phi(s', \pi(s')))^T$$
$$\tilde{b} \leftarrow \tilde{b} + \phi(s, a)r$$

# of Expert Episode : 100

# of Episode for training  $w^\pi$  (LSTD-Q) : 10 (10 x 20 = 200 samples, |D|=200 )

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^T \phi_{reward}(s)$$

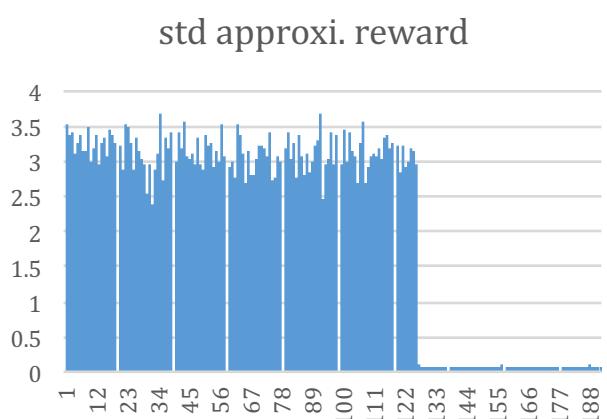
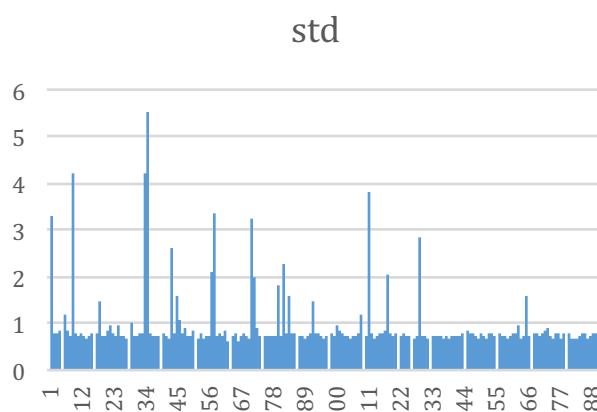
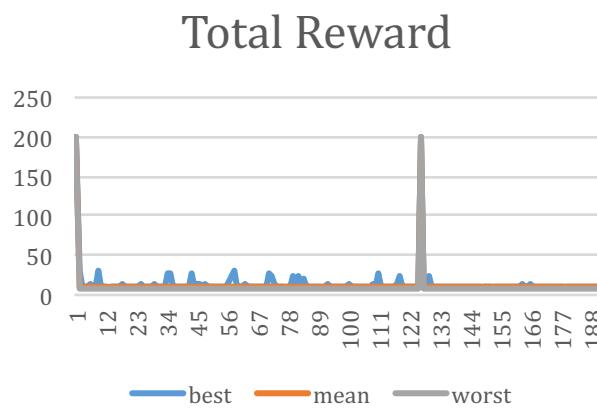
Reward basis function : ANN cartpole + RBF

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



# IRL + MC approach

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 1 (1 x 20 = 20 samples, |D|=20)

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

Reward basis function : ANN cartpole + RBF

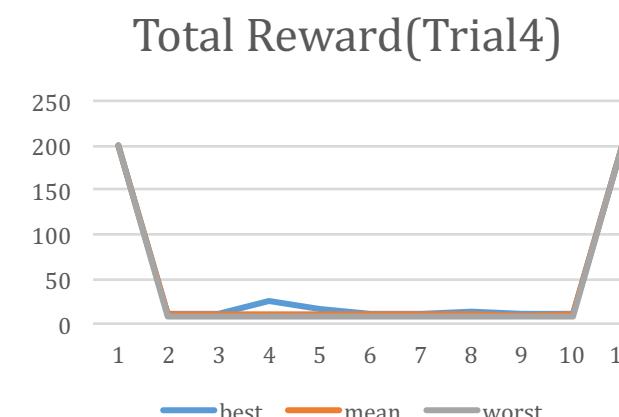
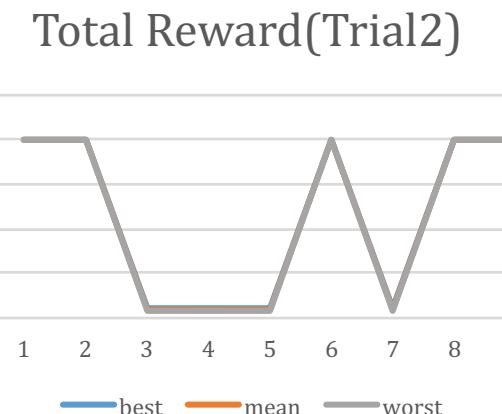
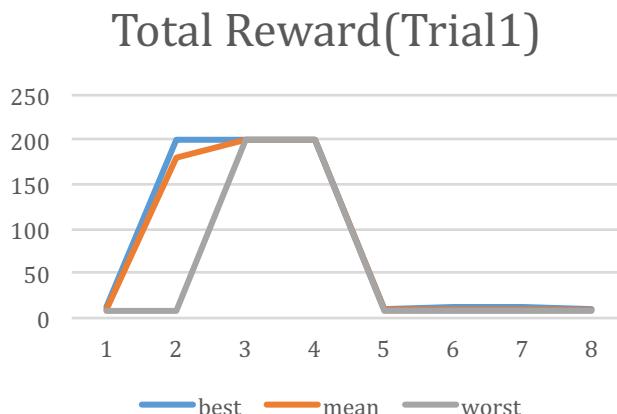
Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5

```
for each (s, a, r, s') ∈ D
     $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s, a) (\phi(s, a) - \gamma \phi(s', \pi(s')) )^\top$ 
     $\tilde{b} \leftarrow \tilde{b} + \phi(s, a)r$ 
```



→ 20 trajectories for training  $w^\pi$  are enough!

# IRL + MC approach

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s) \quad exp\left(-\frac{\|s - \mu_j\|^2}{2\sigma^2}\right)$$

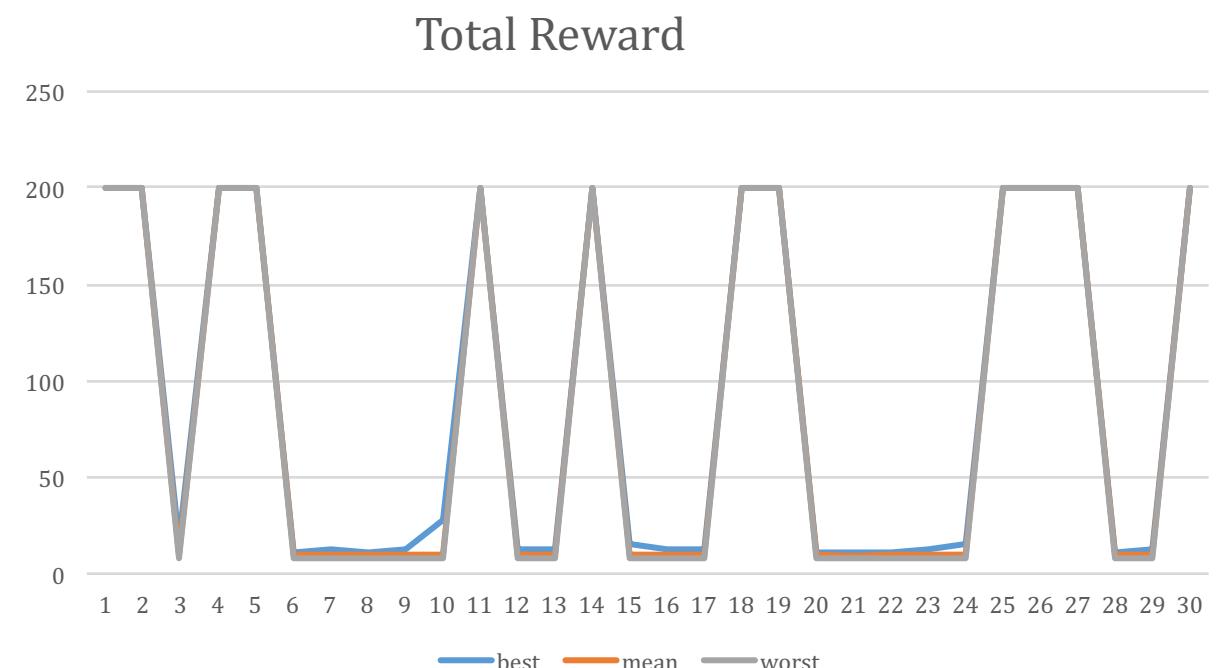
Reward basis function : RBF (ANN cartpole + RBF)

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



# IRL + MC approach

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

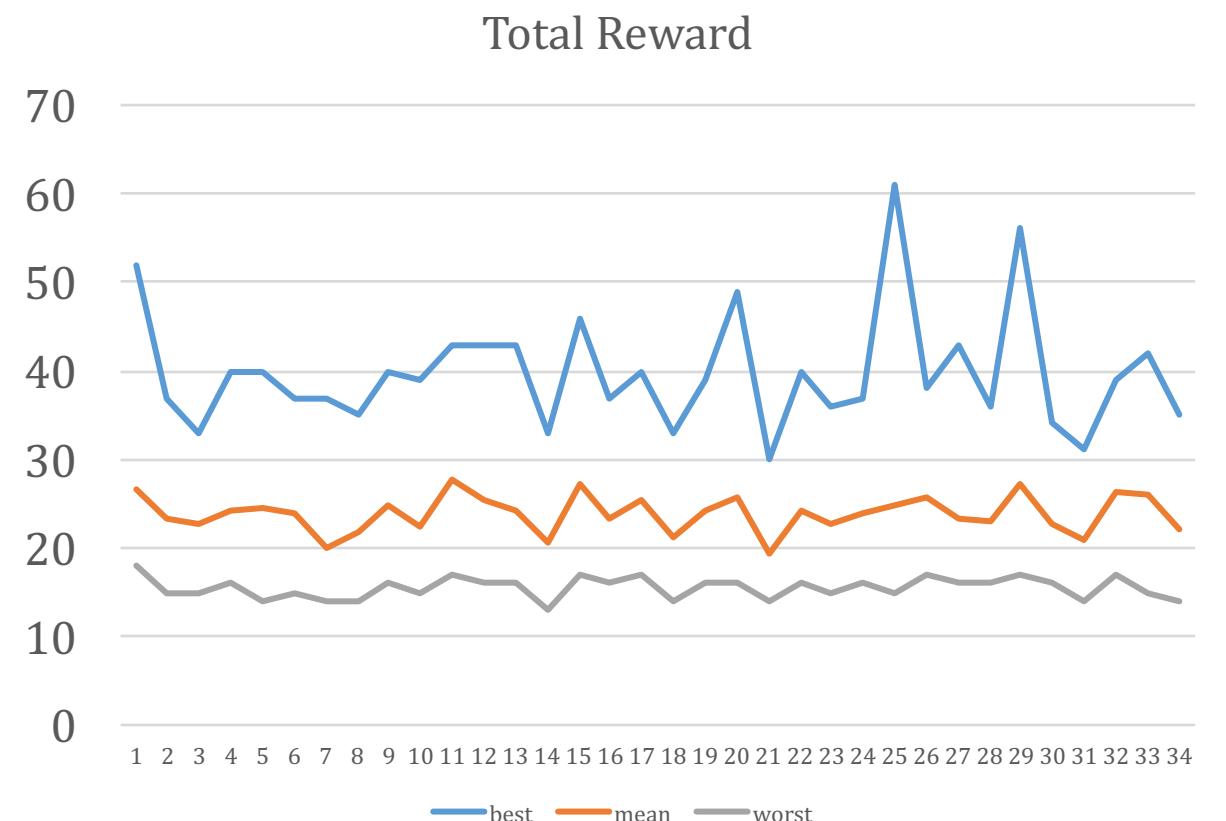
Reward basis function : ANN cartpole + RBF

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ~~RBF (ANN cartpole + RBF)~~

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



# IRL + MC approach

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

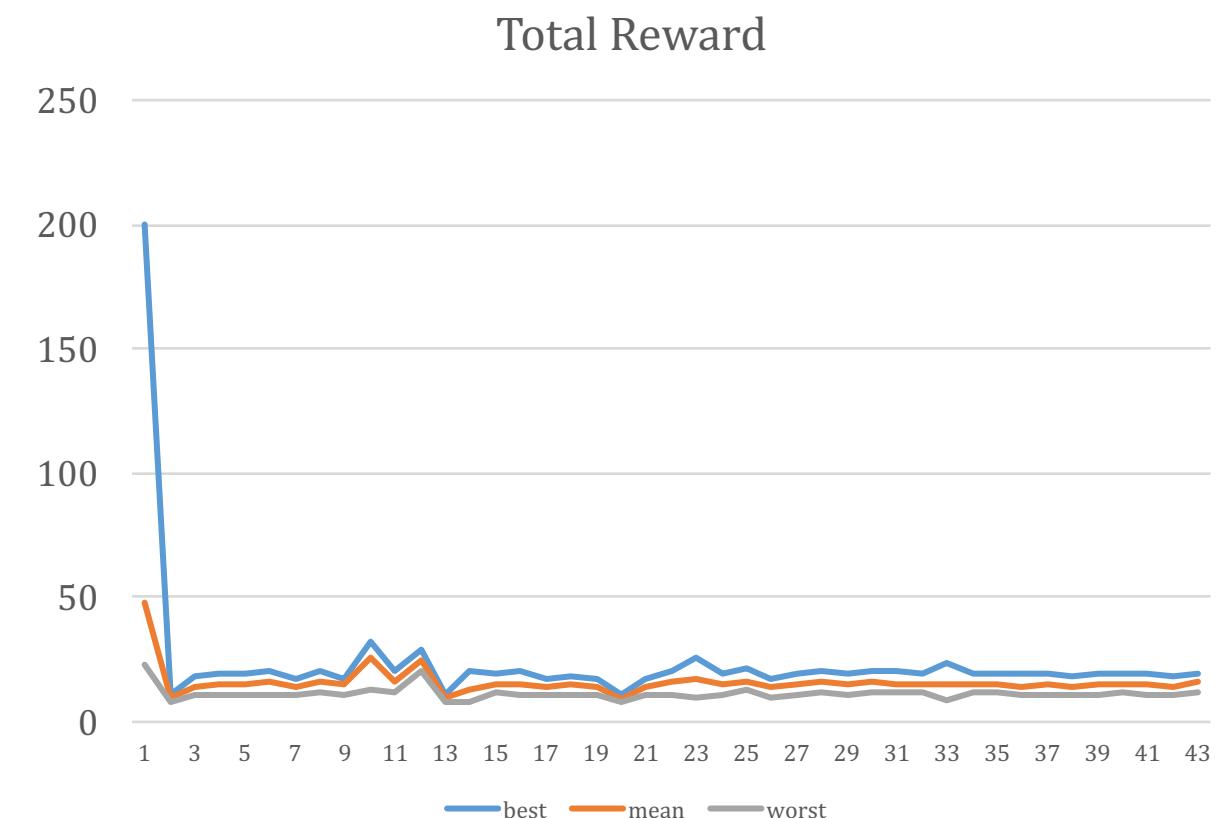
Reward basis function : RBF (ANN cartpole + RBF)

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : RBF (ANN cartpole + RBF)

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



# IRL + LSTD-mu

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100 (100 x ? samples) (**LSTD-mu**)

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

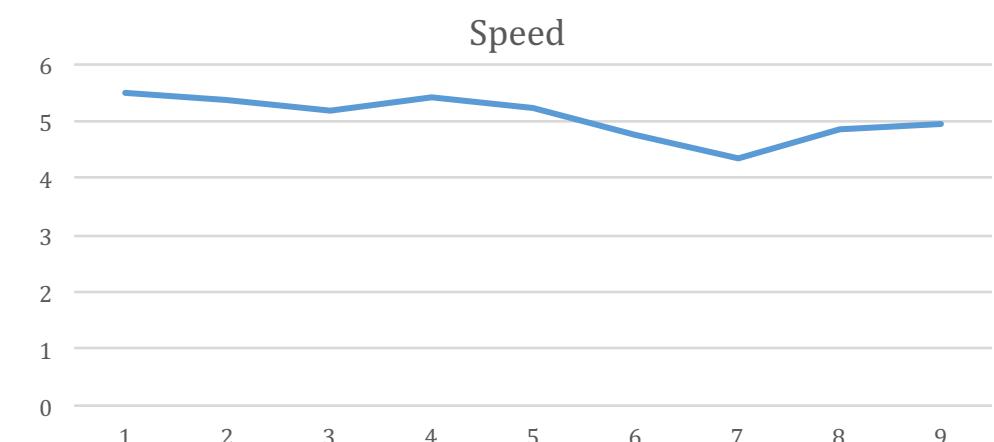
Reward basis function : ANN cartpole + RBF

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



Average speed : x5.06 (compared to MC approach with # Epi 100)

# IRL + LSTD-mu

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 1 (100) (1 x ? samples) (**LSTD-mu**)

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

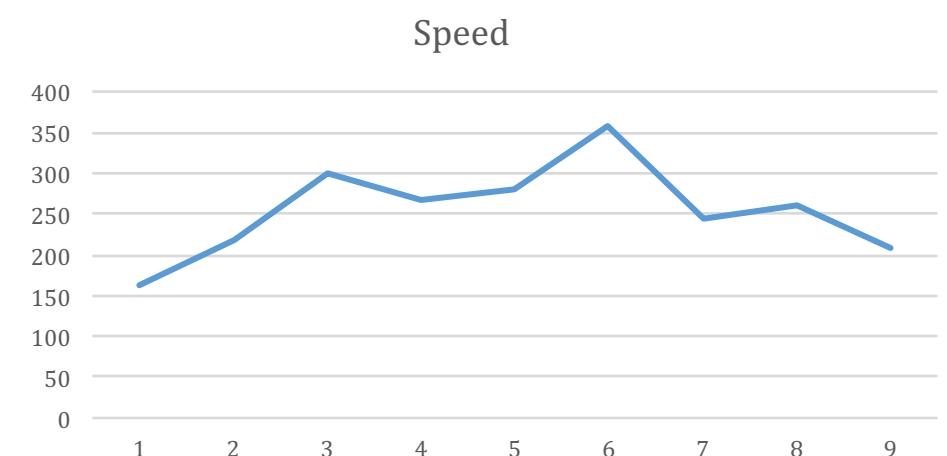
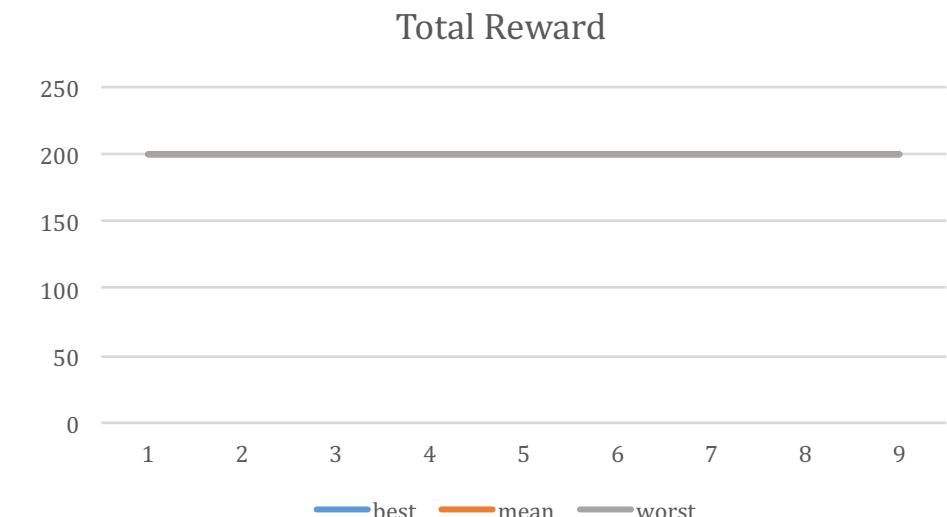
Reward basis function : ANN cartpole + RBF

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : ANN cartpole + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



Average speed : x255 (compared to MC approach with # Epi 100)

# Given information, Expert Trajectories.

---

Expert Trajectory,

$$T_i = \{(s_t, a_t, s_{t+1}) \mid 1 \leq t \leq \text{Done}\}$$

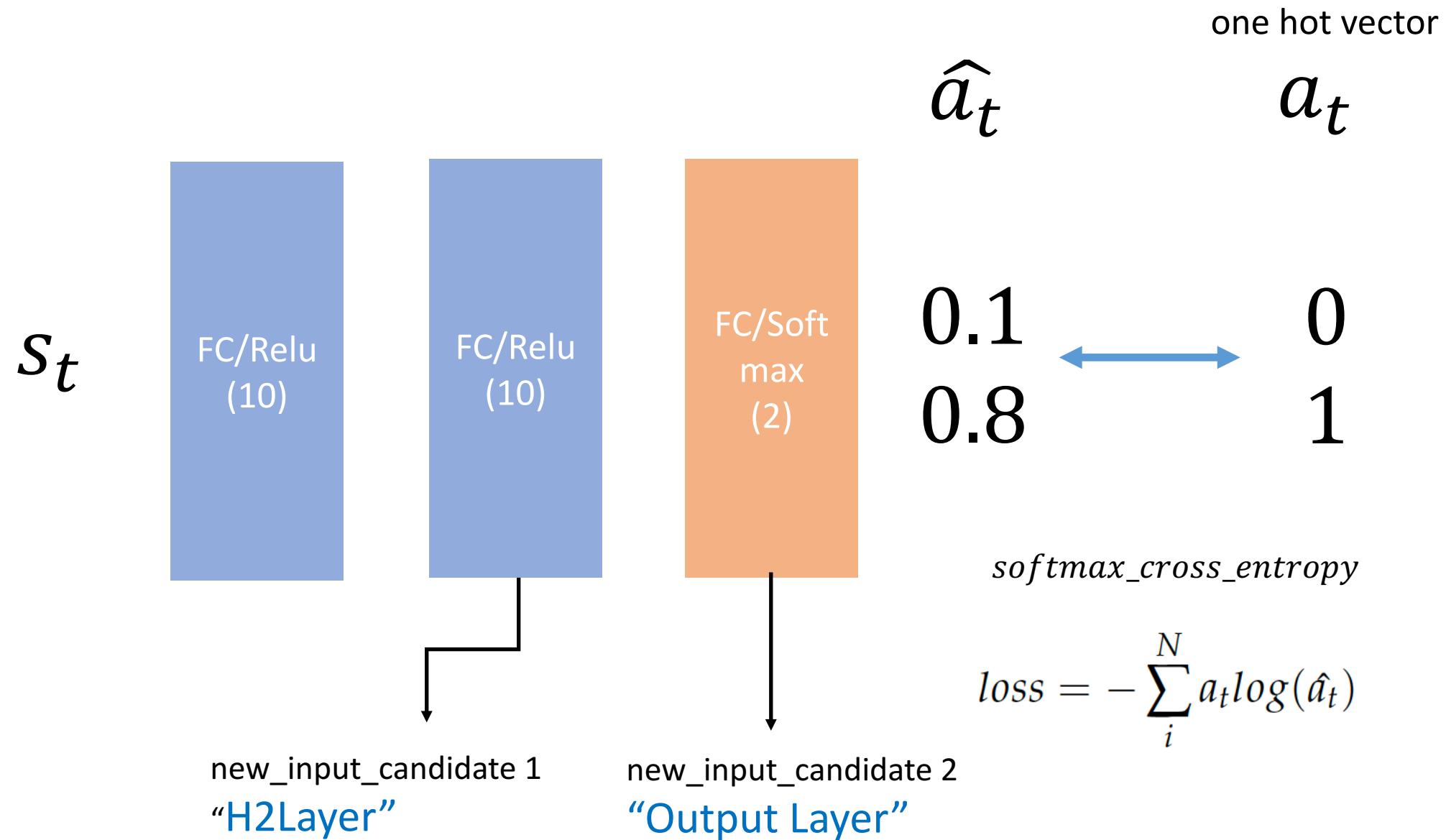
Expert Trajectories,

$$T_{collection} = \{T_i, 1 \leq i \leq n\}$$

Supervised Action Estimator? (not DQN but quite similar)

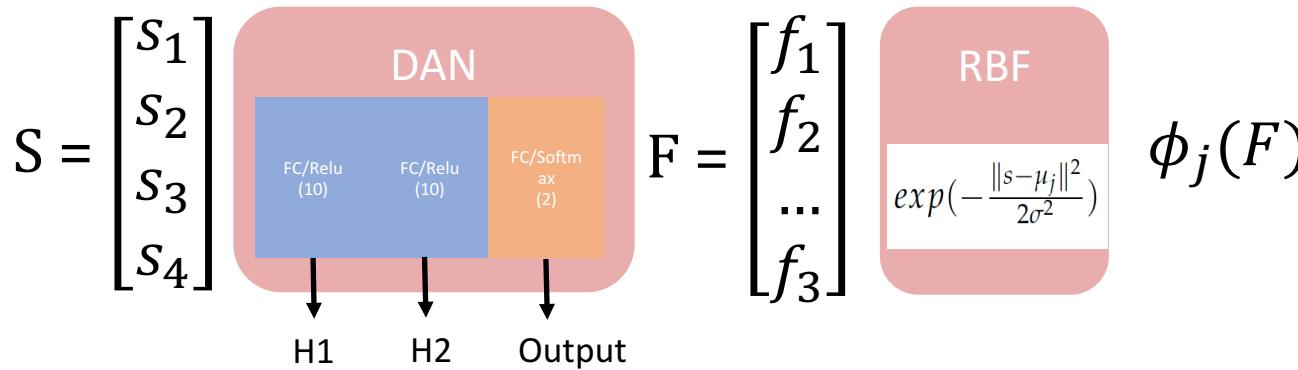
$$s_t(\textit{input}) \rightarrow a_t(\textit{output})?$$

# Deep Action Network(DAN)



# DAN + RBF as basis function

DAN + RBF

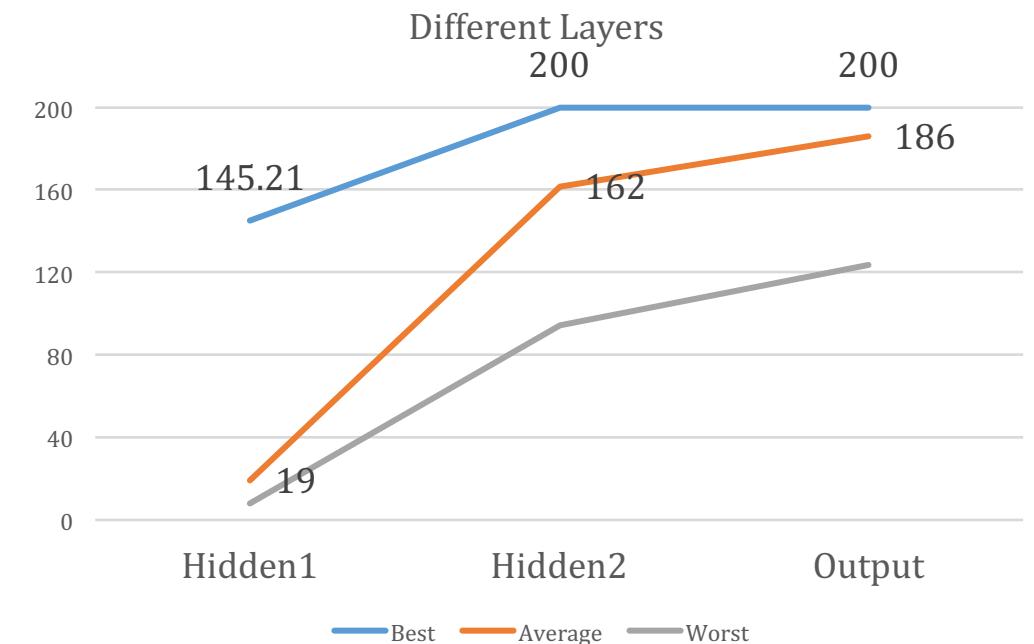


# Episode for training  $w^\pi$ (LSTD-Q) : 100 (2000 samples)

# Evaluation for each policy( $w^\pi$ ) : 1000

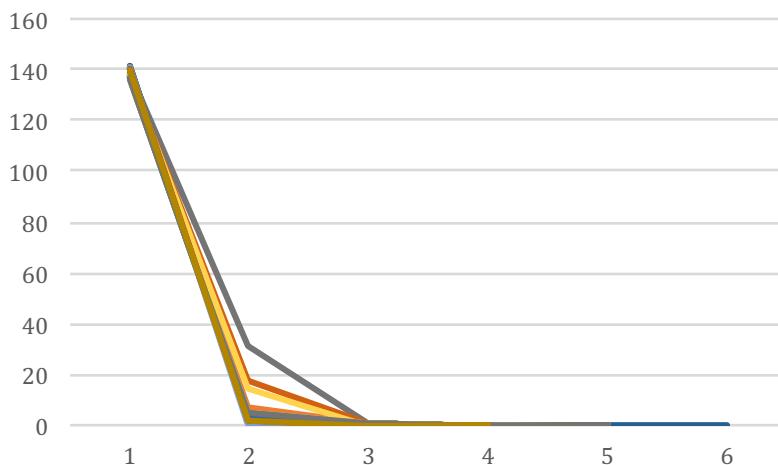
# Trials(Experiments) : 100

# basis feature(P) = 10

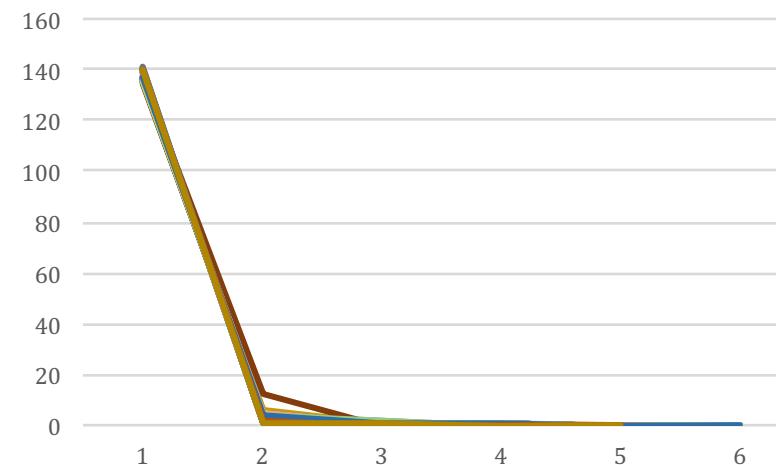


# DAN + RBF as basis function

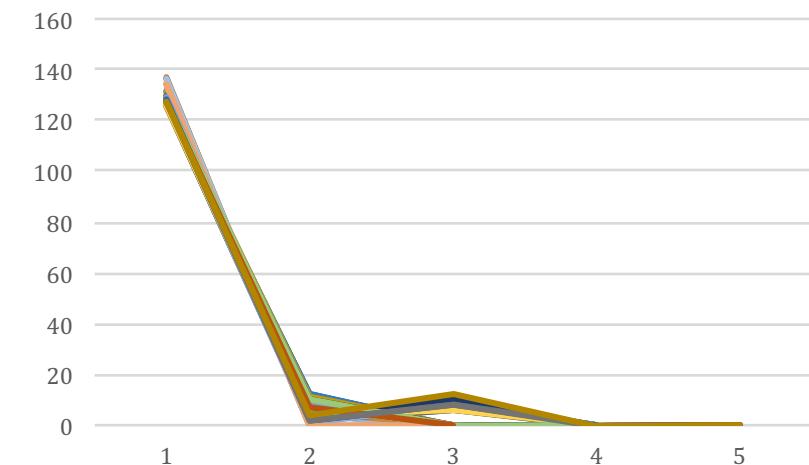
DAN(hidden 1)



DAN(hidden 2)



DAN(output)



$$t = \|w^{\pi_{j-1}} - w^{\pi_j}\|_2$$

# DAN + IRL + MC approach

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

Reward basis function : RBF (Tensorflow conflict)

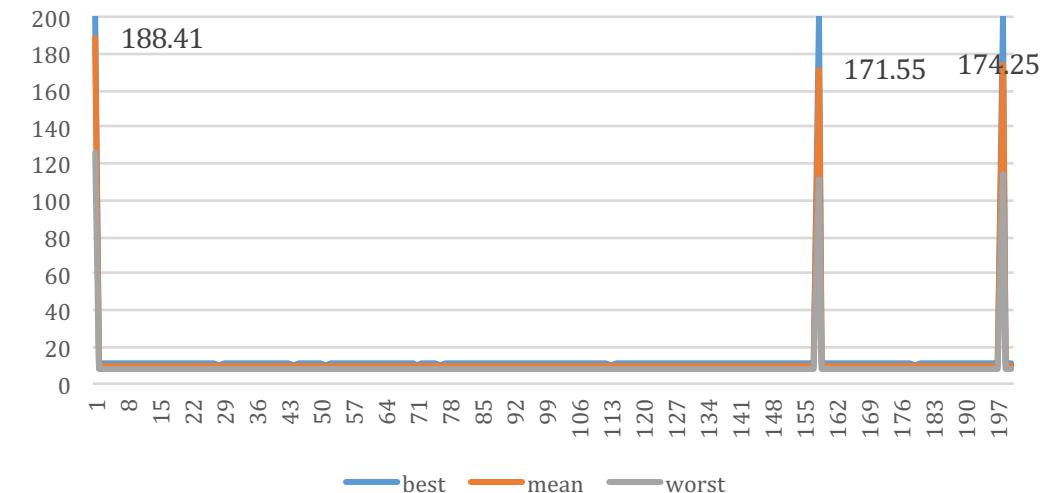
Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

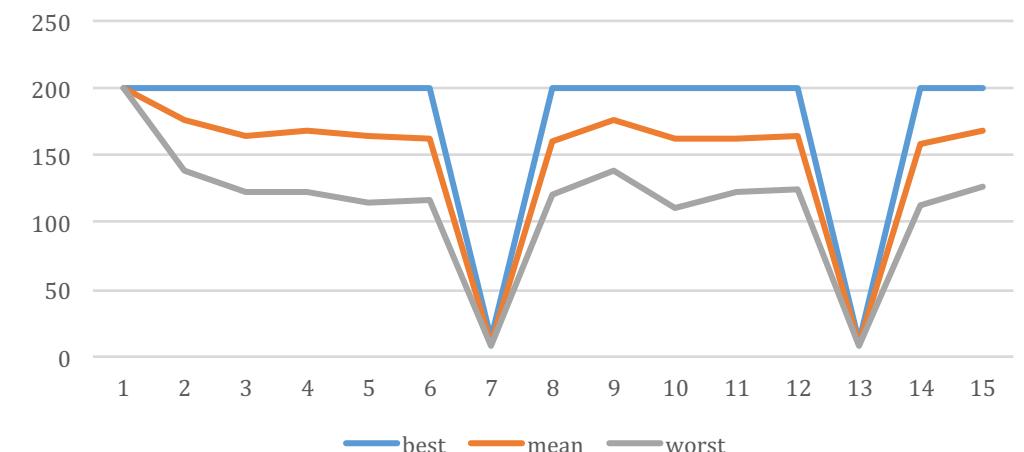
LSTD-Q basis function : DAN(output layer) + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : {5, 10}

LSTDQ basis dimension( $\phi_{LSTD}$ ) = 5



LSPI basis dimension( $\phi_{LSTD}$ ) = 10



# DAN + IRL + MC approach

# of Expert Episode : 100

# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

$$R(s) = \theta^\top \phi_{reward}(s)$$

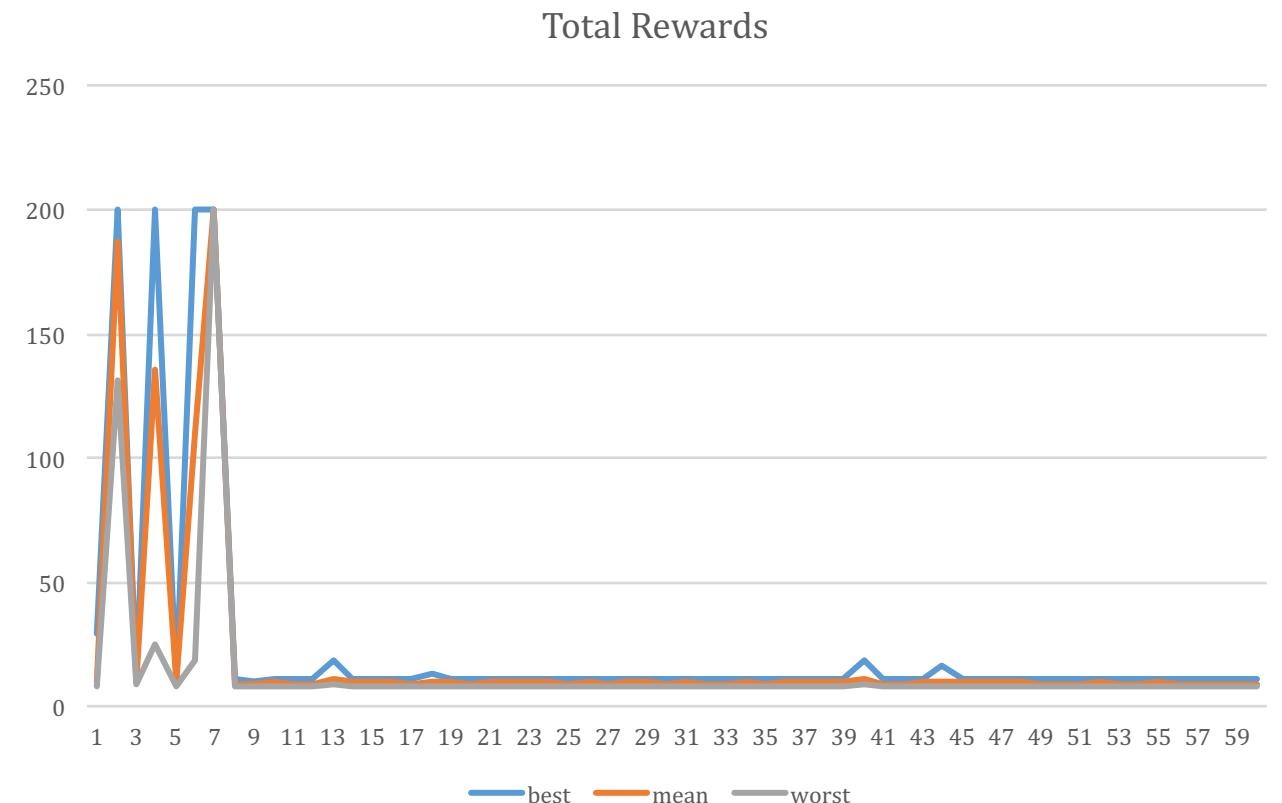
Reward basis function : RBF (Tensorflow conflict)

Reward basis feature dim( $\phi_{reward}$ ) : 5

$$\hat{Q}^\pi = \Phi_{LSTD} w^\pi$$

LSTD-Q basis function : DAN(hidden2) + RBF

LSTD-Q basis dimension( $\phi_{LSTD}$ ) : 5



# ANN for cartpole vs. Deep Action Network (for LSTD-Q basis func.)

- ANN for cartpole

label : should be provided how to act, according to danger situation(state)

only Cartpole problem

- Deep Action Network(DAN)

label : expert trajectory ( $s, a$ )

more general

# of Expert Episode : 100

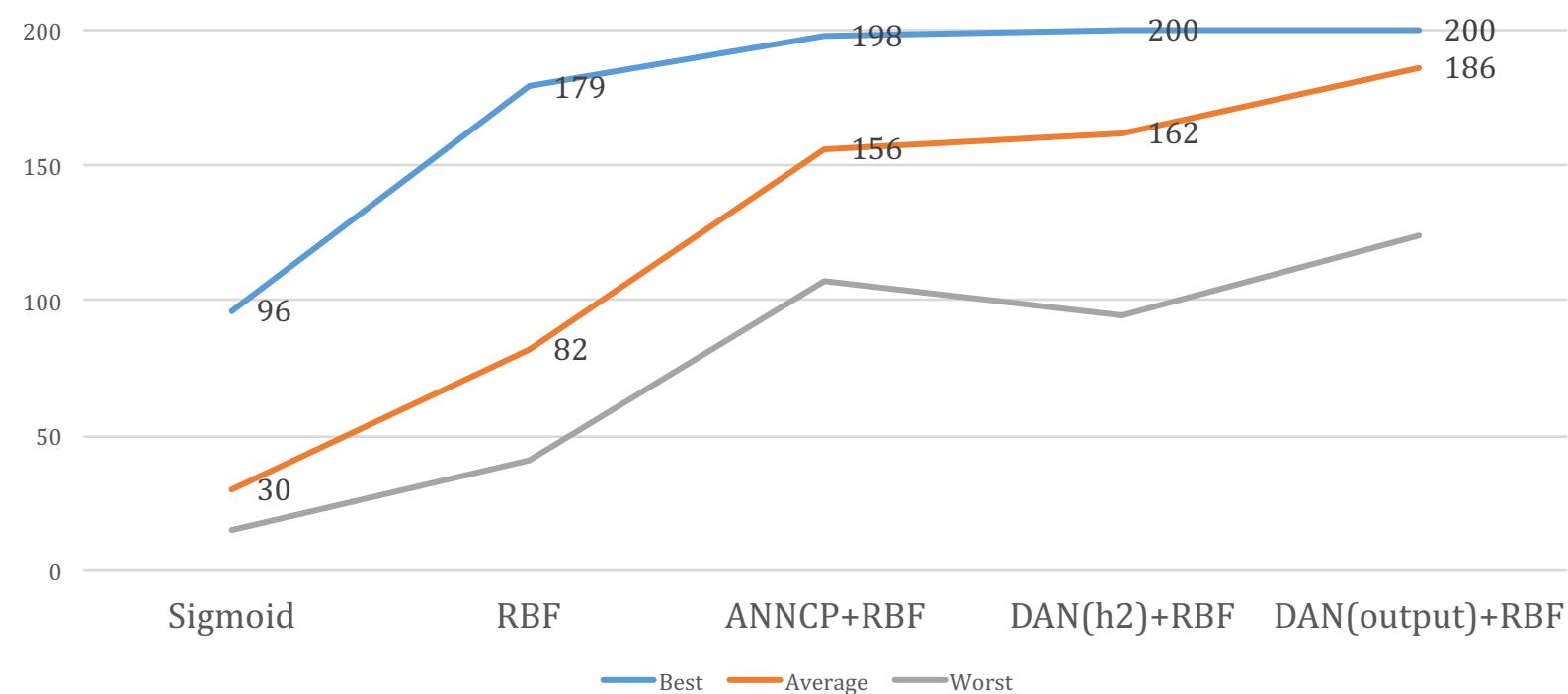
# of Episode for training  $w^\pi$ (LSTD-Q) : 100

# of Episode for calculating  $\mu^\pi$  : 100

# of Evaluation of each policy( $w^\pi$ ) : 100

# basis feature(P) = 10

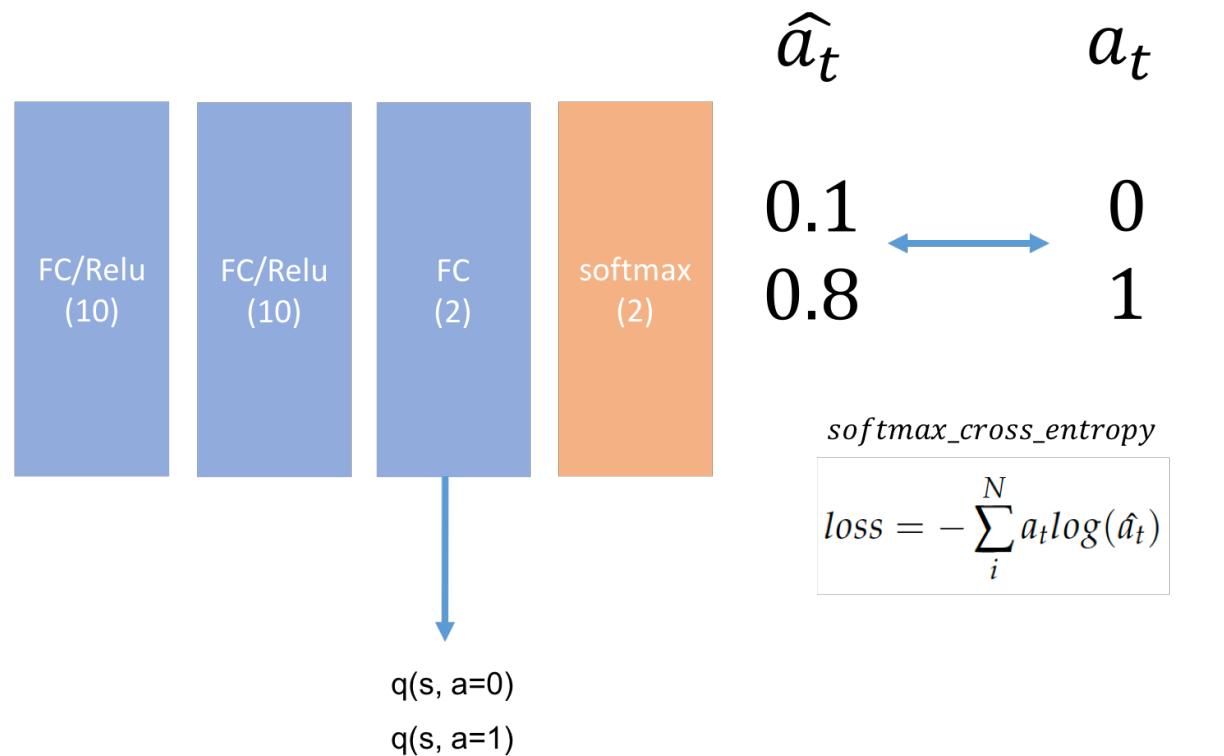
Different Basis function



# DAN + DQN

## Deep Action Network

Dataset =  $\{(s_{i,t}, a_{i,t})\} i \in \text{Trajectories}, 1 \leq t \leq \text{Done}$



Now, we know  $a_t$  (only for  $s_t$ )

$s_t \rightarrow a_t \rightarrow s_{t+1} \rightarrow a_{t+1} \dots$

Dataset =  $\{(s_{i,t}, a_{i,t}, s_{i,t+1})\} i \in \text{Traj}, 1 \leq t \leq \text{Done}\}$

$$q(s, a) = r(s, a) + \gamma q(s', \pi(s')) \quad (\text{by Bellman equation})$$

$$r(s, a) = q(s, a) - \gamma q(s', \pi(s'))$$

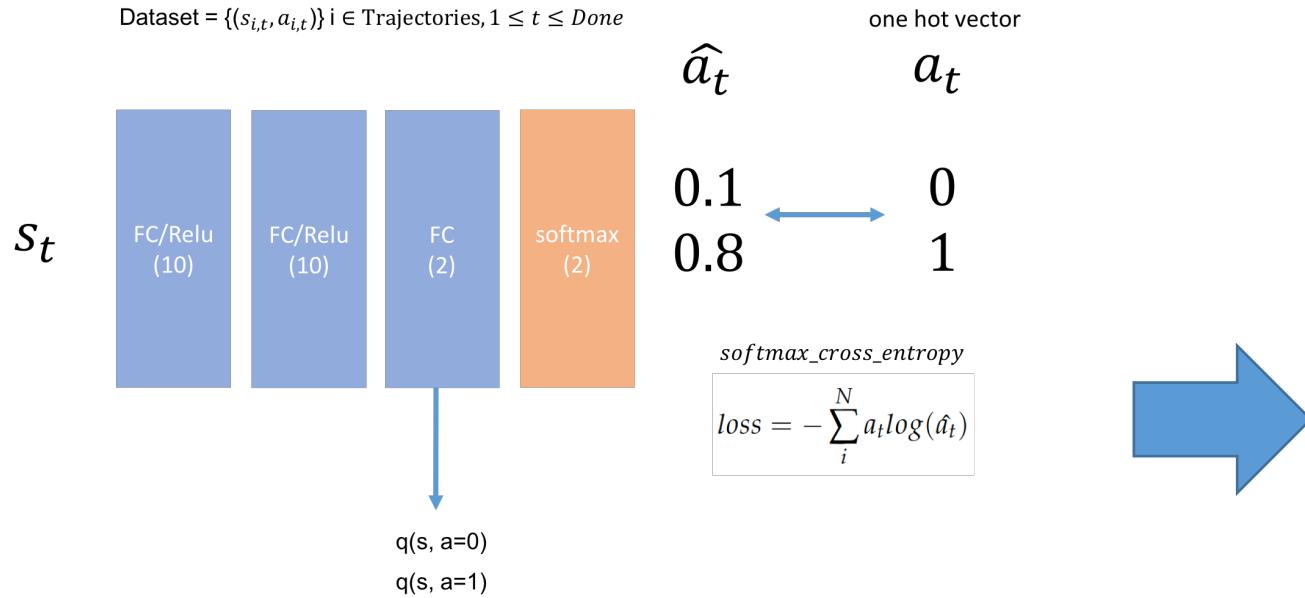
$$r(s, a) = q(s, a) - \gamma \max_{a'} q(s', a')$$

known

known

# DAN + DQN

## 1. DAN : Estimate Action and get $q(s,a)$



## 3. MDP solver(control solver)

### a. DQN

$$q(s, a) = r(s, a) + \gamma q(s', \pi(s')) \quad (\text{by Bellman equation})$$

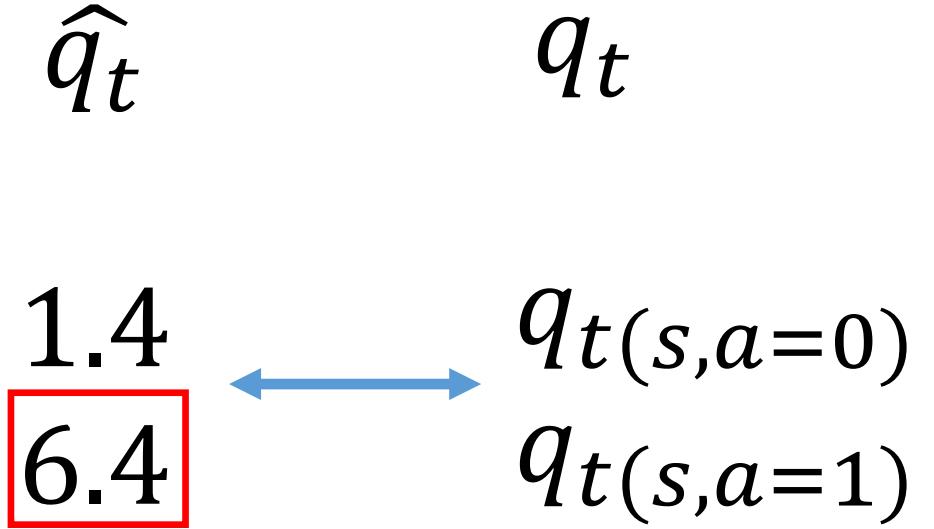
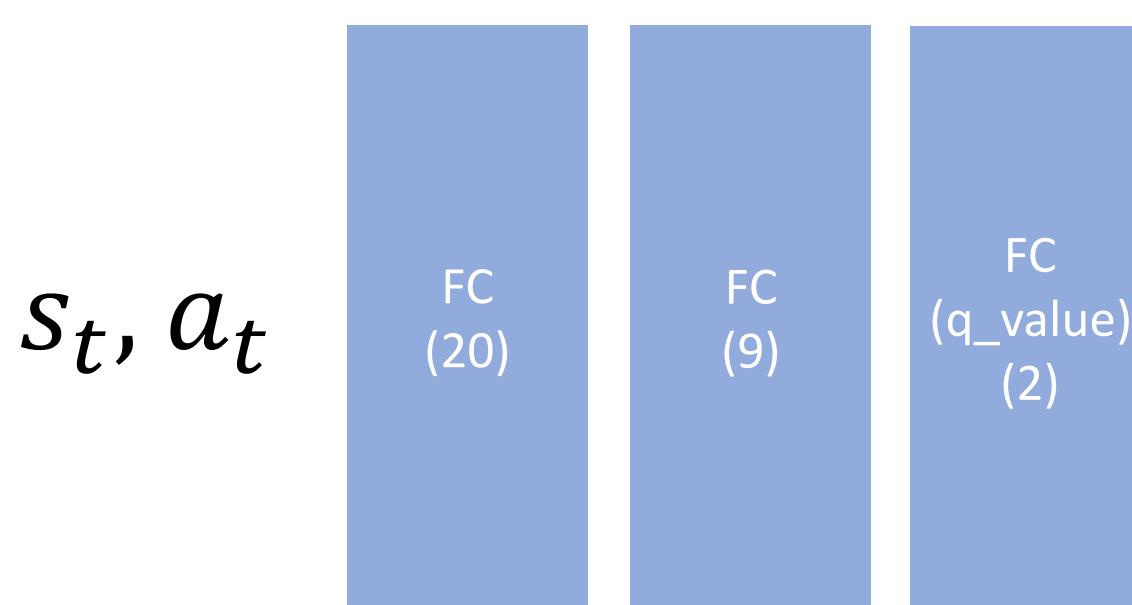
$$r(s, a) = q(s, a) - \gamma q(s', \pi(s'))$$

$$r(s, a) = q(s, a) - \gamma \max_{a'} q(s', a')$$

known                      known

# Deep Q Network(DQN)

Dataset =  $\{(s_{i,t}, a_{i,t}, r_{i,t})\} i \in \text{Trajectories}, 1 \leq t \leq Done$



$$Loss(w) = \sum (\hat{q}_t - q_t)^2$$

$$q_t(s, a) = r_t(s, a) - \gamma \max_{a'} q_t(s', a')$$

$$r_t(s, a) = q_t(s, a) - \gamma \max_{a'} q_t(s', a')$$

from DAN

# DAN->DRN->DQN (Previous Results)

```

Deep Action Network Training iteration 1810, 0.03702983260154724
Deep Action Network Training iteration 1820, 0.015642434358596802
Deep Action Network Training iteration 1830, 0.010453480295836926
Deep Action Network Training iteration 1840, 0.050940681248903275
Deep Action Network Training iteration 1850, 0.011584173887968063
Deep Action Network Training iteration 1860, 0.006984895560890436
Deep Action Network Training iteration 1870, 0.020961787551641464
Deep Action Network Training iteration 1880, 0.024781856685876846
Deep Action Network Training iteration 1890, 0.039983492344617844
Deep Action Network Training iteration 1900, 0.04822589084506035
Deep Action Network Training iteration 1910, 0.01002182811498642
Deep Action Network Training iteration 1920, 0.012212877161800861
Deep Action Network Training iteration 1930, 0.05930352210998535
Deep Action Network Training iteration 1940, 0.05499160289764404
Deep Action Network Training iteration 1950, 0.040910910815000534
Deep Action Network Training iteration 1960, 0.013394010253250599
Deep Action Network Training iteration 1970, 0.03657654672861099
Deep Action Network Training iteration 1980, 0.021590299904346466
Deep Action Network Training iteration 1990, 0.030237708240747452
Testing Deep Action Network... 5 times
2019-01-09 10:17:15.416 Python[4242:642329] ApplePersistenceIgnoreStat
will not be touched. New state will be written to /var/folders/fw/st6g
...0000g.../T/_org.python.python savedState
Test DAN 0 : 30 timesteps
Test DAN 1 : 27 timesteps
Test DAN 2 : 31 timesteps
Test DAN 3 : 48 timesteps
Test DAN 4 : 24 timesteps

```

Expert trajectories were wrong.

```

Deep Reward Network iteration : 900, loss : 0.49980152130126954
Deep Reward Network iteration : 910, loss : 0.4997991180419922
Deep Reward Network iteration : 920, loss : 0.4997987747192383
Deep Reward Network iteration : 930, loss : 0.49979644775390625
Deep Reward Network iteration : 940, loss : 0.4997944641113281
Deep Reward Network iteration : 950, loss : 0.49979454040527344
Deep Reward Network iteration : 960, loss : 0.49979267120361326
Deep Reward Network iteration : 970, loss : 0.49979217529296877
Deep Reward Network iteration : 980, loss : 0.49979118347167967
Deep Reward Network iteration : 990, loss : 0.49979095458984374
DQN Training episode 0 timestep 0 reward -0.37572297159576416
DQN Training episode 0 timestep 1 reward -0.36829471588134766
DQN Training episode 0 timestep 2 reward -0.3732357621192932
DQN Training episode 0 timestep 3 reward -0.3663175404071808
DQN Training episode 0 timestep 4 reward -0.36510223150253296
DQN Training episode 0 timestep 5 reward -0.3578903079032898
DQN Training episode 0 timestep 6 reward -0.3613246977329254
DQN Training episode 0 timestep 7 reward -0.3559782803058624
DQN Training episode 0 timestep 8 reward -0.35866567492485046
DQN Training episode 0 timestep 9 reward -0.35467588901519775
DQN Training episode 0 timestep 10 reward -0.3565685749053955
DQN Training episode 0 timestep 11 reward -0.3541203737258911
DQN Training episode 0 timestep 12 reward -0.3483473062515259
DQN Training episode 0 timestep 13 reward -0.3440174162387848
DQN Training episode 0 timestep 14 reward -0.35596293210983276
DQN Training Episode 0 finished after 15 timesteps
average reward : -0.3854454053299768

```

fortunately, reward are same!

# DAN->DQN

## DAN

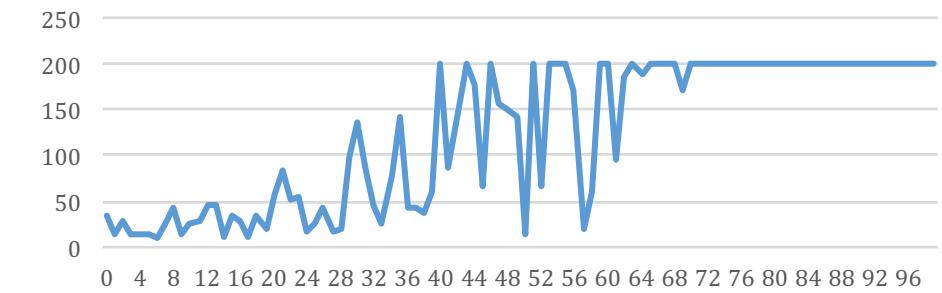
```

Deep Action Network Training iteration 9960, 0.5171039593219757
Deep Action Network Training iteration 9970, 0.5489578205347061
Deep Action Network Training iteration 9980, 0.5222935244441033
Deep Action Network Training iteration 9990, 0.5158185619115829
saveing our trained weights!!
Testing Deep Action Network... 100 times
Test DAN 0 : 200 timesteps
Test DAN 1 : 200 timesteps
Test DAN 2 : 200 timesteps
Test DAN 3 : 200 timesteps
Test DAN 4 : 200 timesteps
Test DAN 5 : 200 timesteps
Test DAN 6 : 200 timesteps
Test DAN 7 : 200 timesteps
Test DAN 8 : 200 timesteps
Test DAN 9 : 200 timesteps
Test DAN 10 : 200 timesteps
Test DAN 11 : 200 timesteps
Test DAN 12 : 200 timesteps
Test DAN 13 : 200 timesteps
Test DAN 14 : 200 timesteps
Test DAN 15 : 200 timesteps
Test DAN 16 : 200 timesteps
Test DAN 17 : 200 timesteps
Test DAN 18 : 200 timesteps
Test DAN 19 : 150 timesteps

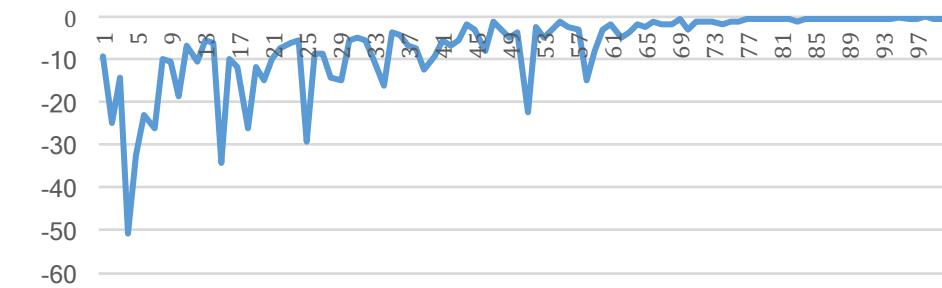
```

average reward : 194.85

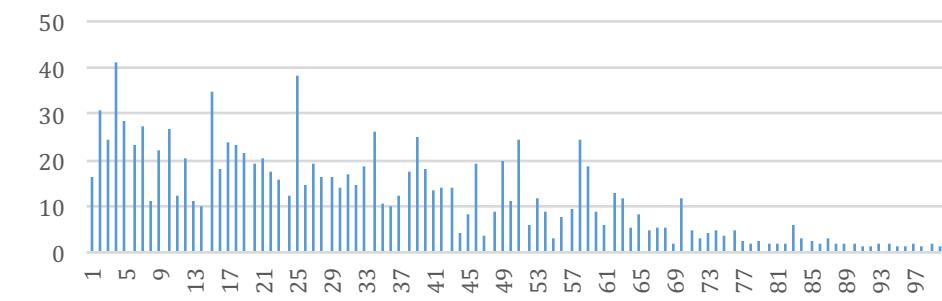
Total Rewards



average of reward



standard deviation of reward



# What I implemented,

---

1. LSPI
2. basis function : sigmoid, RBF, ANNCP, DAN(h2, output)
3. Batch, off-policy and model-free Apprenticeship learning(LSTD-mu)
4. Deep Action Network(DAN) + DQN

# Next steps, (in order of importance)

---

1. Fix LSPI problem(oscillation)
2. Analyze Deep Action, Reward Network
3. Other simulator (CartPole might be easy)
4. IRL + DQN (as MDP solver) or apply other Deep RL
5. Projection method → Quadratic programming (reward bounding part)
6. Apply deep features for LSTD-mu

# Backup slide

---

# Backup slide : Radial Basis Functions

RBFs are the natural generalization of coarse coding(binary features {0, 1}) to continuous-valued features([0, 1]).

A typical RBF feature has a Gaussian response dependent on the distance between the state,  $s$ , and the feature's center state,  $c_i$ , and the feature's width,  $\sigma_i$ .

$$x_i(s) \doteq \exp\left(-\frac{\|s - c_i\|^2}{2\sigma_i^2}\right)$$

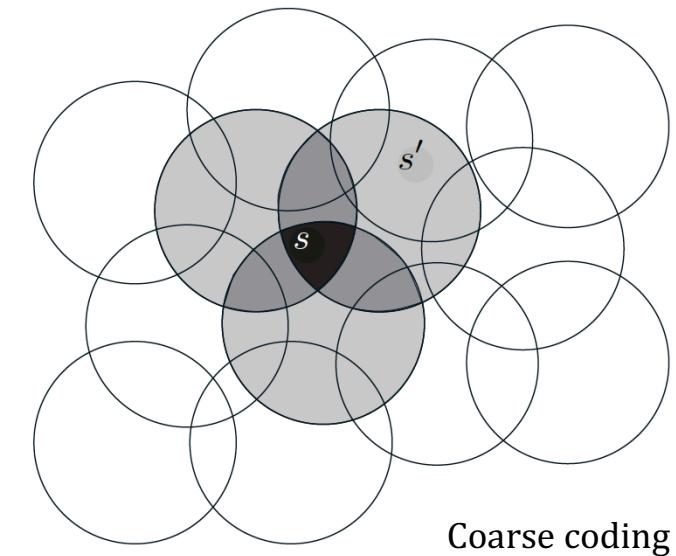
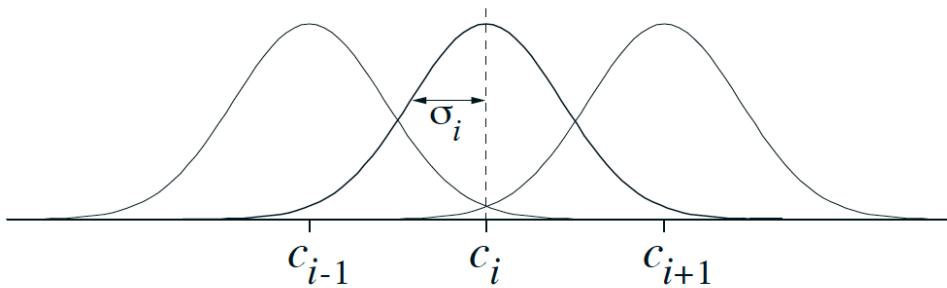


Figure 9.13: One-dimensional radial basis functions.

Advantage : vary smoothly, and differentiable.