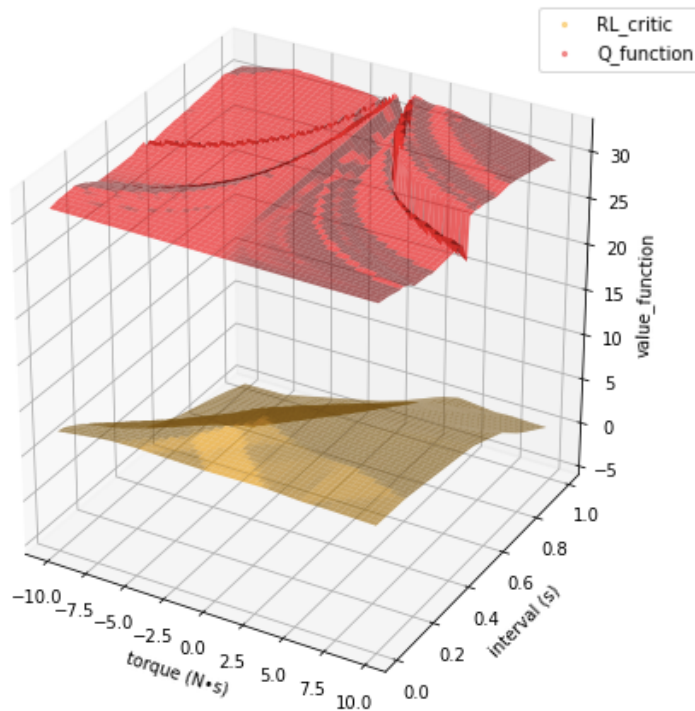# Weekly Report

M2 Ibuki Takeuchi

- Low accuracy of $Q$-function approximation

Graph of $Q(s, a)$ for fixed $s$
(function of $a$)



- $Q$-function are learned with supervised learning (Least Square)

$$\min_{\omega} \frac{1}{N} \sum_{(s,a)\in E} \{Q(s, a|\omega) - \underline{(r(s, a) + \gamma Q(s, \pi(s)|\omega))}\}^2$$

teacher data

- Low accuracy may come from..
  - Optimization algorithm
  - Data bias in $E$

    (This may be the reason)

# Weekly Report

M2 Ibuki Takeuchi

- Check as if optimization is well conducted
    - Define $\omega_{RL}$ be the $Q$-function parameter learned with RL

    - Loss function for data set $E$ should be minimized by $\omega_{RL}$

$$Loss(\omega) = \frac{1}{N} \sum_{(s,a) \in E} \{Q(s,a|\omega) - (r(s,a) + \gamma Q(s, \pi(s)|\omega))\}^2$$

    - By comparing the loss function with some $\omega$
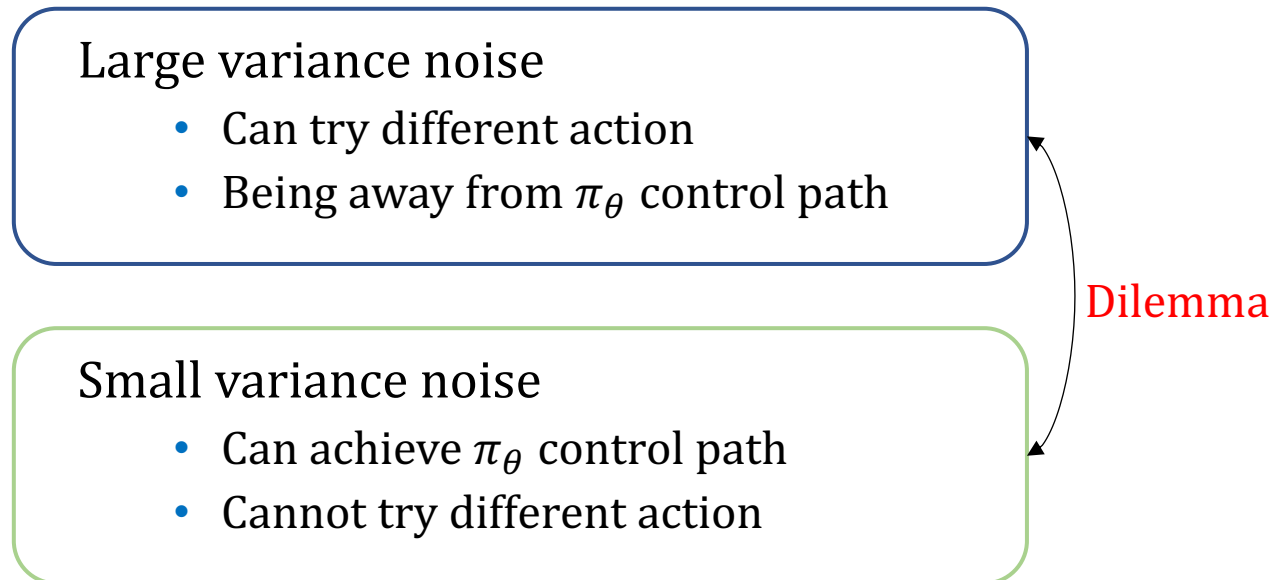
$$Loss(\omega) = 17, Loss(\omega_{RL}) = 0.02$$

    - I want to check again on what should be compared, with Kashima Sensei

# Weekly Report

M2 Ibuki Takeuchi

- Overcome data bias
  - Various experiences of $(s, a)$ are needed
  - Teacher's data are collected by agent's experience

  - How does the agent collect data?
    - Data exploration:

$$a = \pi_\theta(s) + e$$

    - Store data of $(s, a)$ to set $E$

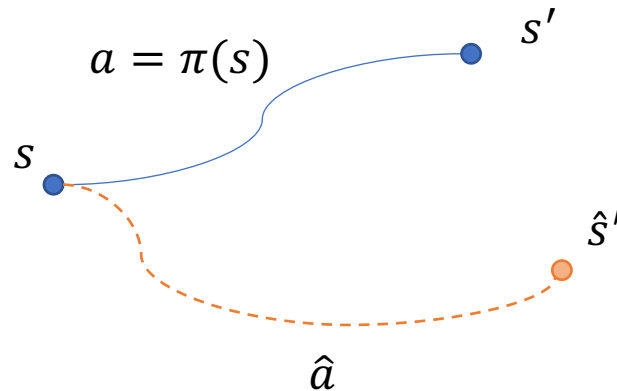    - If the variance of $e$ is large, that of $(s, a)$ become large

# Weekly Report

M2 Ibuki Takeuchi

- DDPG requires agent to experience <span style="color:red">states</span> on control path w.r.t. current policy $\pi_\theta$

- Exploration-Exploitation Dilemma

Large variance noise
- Can try different action
- Being away from $\pi_\theta$ control path

Dilemma

Small variance noise
- Can achieve $\pi_\theta$ control path
- Cannot try different action

# Weekly Report

M2 Ibuki Takeuchi

- Adaptive noise scaling
  - By changing action...

    if the change of next state is large $\Longleftrightarrow \frac{\partial s'}{\partial a}$ is Large

    $\rightarrow$ small noise

    if the change of next state is small $\Longleftrightarrow \frac{\partial s'}{\partial a}$ is Small

    $\rightarrow$ large noise

$a = \pi(s)$

$s'$

$s$

$\hat{s}'$

$\hat{a}$

# Weekly Report

M2 Ibuki Takeuchi

- Try following noise scaling

$$\frac{c}{\|g\| + c} \times \mathcal{N}(0,1) \ : c \text{ is hyper parameter}$$

- Because I have not summarized my consideration, I want to report the result on colloquium next week

Control System Theory Group

M2 Ibuki Takeuchi

- Try noise

$$\frac{c}{\|g\| + c} \times \mathcal{N}(0,1) \ : c \text{ is hyper parameter}$$



Control System Theory Group