

CS292F StatRL Lecture 16

Offline RL with function approximation

Instructor: Yu-Xiang Wang

Spring 2021

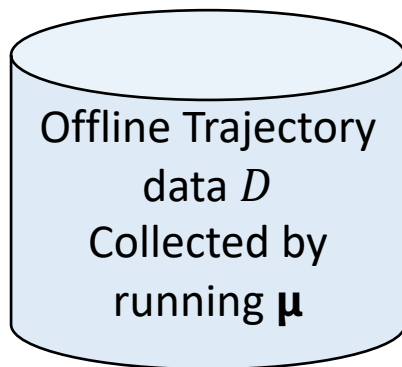
UC Santa Barbara

Plan for the remaining sessions

- Two lectures next week
 - Q&A session for me to help you guys with your HWs and projects
 - Will not be recorded
- Jun 2 lecture will be made into a **4-hour-long mini-symposium** of project presentations.

Recap: Offline Reinforcement Learning, aka. Batch RL

- Task 1: Offline Policy Evaluation. (OPE)

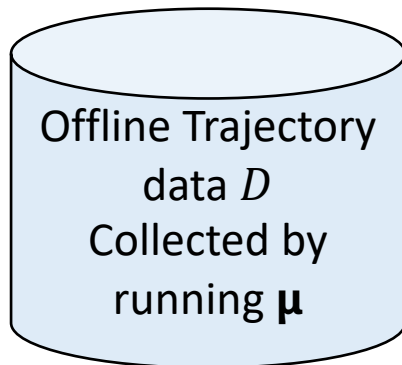


Task: design OPE
methods

Evaluate fixed Target
Policy π

**Via
Uniform
OPE**

- Task 2: Offline Policy Learning. (OPL)



Task: design OPO
methods

Find near optimal
Policy $\hat{\pi}^*$

Recap: Offline Policy Learning

- Model-based approach

- ERM to maximize the model-based OPE
- Possible extensions: Optimism and Pessimism

$$\text{ERM: } \underset{\pi}{\operatorname{argmax}} \hat{V}^{\pi}$$

$$\underset{\pi}{\operatorname{argmax}} \hat{V}^{\pi} + b^{\pi}$$

$$\underset{\pi}{\operatorname{argmax}} \hat{V}^{\pi} - b^{\pi}$$

- Model-free approaches

- Variance-Reduced Value Iteration
- Fitted Q-Learning. (We will talk about this today!)

$$(\text{fit}, \text{bias}, w.)^{21}$$

In the tabular setting, the problem is (almost) completely solved.

$$(1-\gamma)^{-1}$$

$$N = n \cdot (1-\gamma)^{-1}$$

Offline Policy Learning

	H-horizon, Stationary	H-Horizon, Nonstationary	Infinite horizon γ -discounted
Upper bound	$\sqrt{\frac{H^2}{n d_m}}$	$\sqrt{\frac{H^3}{n d_m}}$	$\sqrt{\frac{(1-\gamma)^{-3}}{N d_m}}$
Lower bound	$\sqrt{\frac{H^2}{n d_m}}$	$\sqrt{\frac{H^3}{n d_m}}$	$\sqrt{\frac{(1-\gamma)^{-3}}{N d_m}}$

Remaining open research threads:

- More adaptive bounds: More explicit dependence on importance weights.
- Reward free / Task-Agnostic settings
- Function approximation settings

OPE

$$\frac{1}{n} \sum_{i=1}^n \frac{(r(s_i) + \gamma V_{\pi}(s_i) - V_{\pi}(s_i))^2}{d(s_i)} \text{Var}[V_{\pi}(s_i) | s_i]$$

$$(s, a, s')^n$$

This lecture

- Function approximation in RL in general
- Theoretical analysis of Fitted Q Iteration for offline RL in function approximation setting

Borrowed some ideas / materials from Nan Jiang.
FQI analysis from AJKS Ch 15.

Recap: Large MDPs

- State space can be exponentially large

- Planning horizon H is large

- Typical solutions:

- use features to denote state (or state-action pairs)
- use function approximation of various quantities

$$s \rightarrow \phi(s) \in \mathbb{R}^d$$

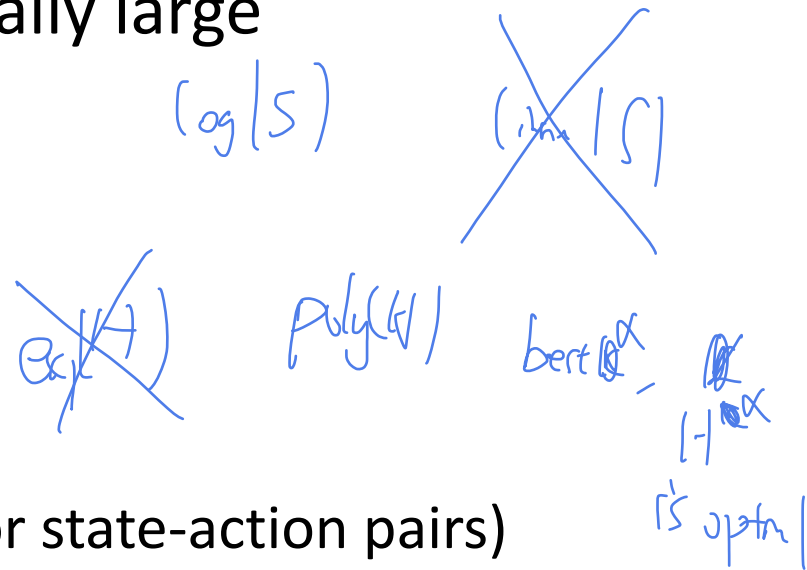
$$(s,a) \rightarrow \phi(s,a) \in \mathbb{R}^d$$

$$f_0 \in \mathcal{F}_H$$

$$f_\theta(\phi(s)) \approx V^*(s)$$

$$f_\theta(\phi(s,a)) \approx Q^*(s,a)$$

$$\theta \in \Theta$$



Ideally, if we have access to the MDP exactly then we could solve

- We could run value iterations

$$f_t \leftarrow \mathcal{T} f_{t-1}$$

$\mathcal{F} = \{f: (S, A) \rightarrow \mathbb{R}\}$
 $\text{if } Q^* \in \mathcal{F}$

$$\mathcal{T} f(s, a) := \underbrace{r(s, a)} + \underbrace{\gamma \mathbb{E}_{s' \sim P(\cdot | s, a)}} \max_{a' \in \mathcal{A}} f(s', a').$$

γ - contraction
 $\|f_c - f^*\| \leq \gamma \|f - f^*\|$

- The issue is that is the updated function still within the function class?

- Add a projection

$$f_t \leftarrow \Pi_F \mathcal{T} f_{t-1}$$

Fitted Q-Iteration

$$f_t = \arg \min_{f \in \mathcal{F}} \sum_{(s,a,r,s') \in D} \left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s',a') \right) \right)^2$$

(constant)

$\mathbb{T} f_{t-1} = \mathbb{T} f_{t-1}$
if deterministic transition

- Stochastic semigradient updates

Treat as constant; don't pass gradient

$$\begin{aligned} \theta &\leftarrow \theta - \frac{\alpha}{2} \cdot \nabla_{\theta} \left(f_{\theta}(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{\theta}(s',a') \right) \right)^2 \\ &= \theta - \alpha \left(f_{\theta}(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{\theta}(s',a') \right) \right) \nabla_{\theta} f_{\theta}(s,a) \end{aligned}$$

Same as Q-learning in tabular / linear function approximation case.

Very similar to DQN if we use a neural network function approximation

Questions about convergence?

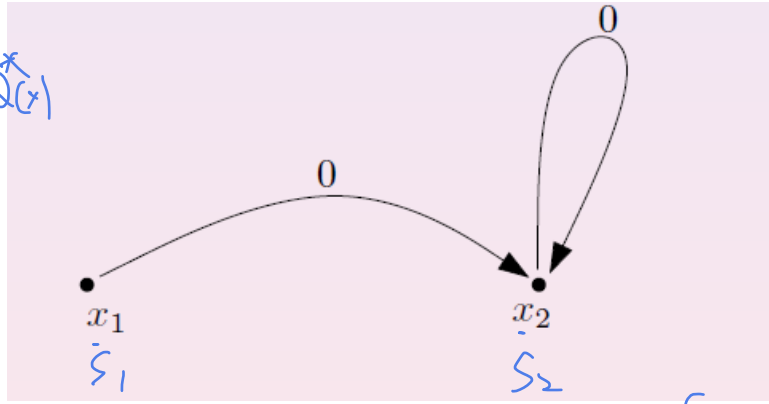
- If realizable, then we know that Q^* is a fixed point.

$$TQ^* = Q^*$$

- But convergence guarantee is not guaranteed in general.
 - Even if it is a linear function approximation.
 - Even if it is realizable.
 - Even if we know the MDP

2.1 Counter-example for least-square regression [Tsitsiklis and van Roy, 1996]

An MDP with two states x_1, x_2 , 1-d features for the two states: $f_{x_1} = 1, f_{x_2} = 2$. Linear Function approximation with $\tilde{V}_\theta(x) = \theta f_x$.



Handwritten notes:

$$\phi(s_1) = 1 \text{ CR}^1$$

$$\phi(s_2) = 2 \text{ CR}^1$$

$$V^* = 0$$

Handwritten note:

$$-B \leq \theta \leq B$$

initialize at θ

Handwritten note:

$$(x_1, a, x_2, 0)$$

credit: course notes from Shipra Agrawal

$$\begin{aligned} \theta_k &:= \arg \min_{\theta} \frac{1}{2} (\theta - \text{target}_1)^2 + (2\theta - \text{target}_2)^2 \\ &= \arg \min_{\theta} \frac{1}{2} (\theta - \gamma \theta^{k-1} f_{x_2})^2 + (2\theta - \gamma \theta^{k-1} f_{x_2})^2 \\ &= \arg \min_{\theta} \frac{1}{2} (\theta - \gamma 2\theta^{k-1})^2 + (2\theta - \gamma 2\theta^{k-1})^2 \end{aligned}$$

Handwritten note:

$$(x_2, a, x_2, 0)$$

$$(\theta - \gamma 2\theta^{k-1}) + 2(2\theta - \gamma 2\theta^{k-1}) = 0 \Rightarrow \underline{5\theta = 6\gamma\theta^{k-1}}$$

$$\underline{\theta_k = \frac{6}{5}\gamma\theta_{k-1}}$$

This diverges if $\gamma \geq 5/6$.

Direct Bellman Residual Minimization

minimize $\|f - \mathcal{T}f\|$ over $f \in F$

$$\sum_{s,a} c_{s,a} (f(s,a) - (r(s,a) + \gamma \max_{a'} f(s',a'))) \quad \leftarrow$$

- We do not have access to the transition kernel
- So what we can minimize is the following

$$\min_{f \in F} \sum_{(s,a) \sim \mu} \mathbb{E} \left[\left(f(s,a) - (r + \gamma \max_{a'} f(s',a')) \right)^2 \right]$$

Handwritten notes: $(s,a) \sim \mu$ is annotated with $(s,a) \sim R(s,a)$ and $s' \sim P(s,a)$. The term $\max_{a'} f(s',a')$ is annotated with $\max_{a'} f(s',a')$.

Are they equivalent, as the dataset goes to infinity?

Double sampling issue

$$\mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - \left(r + \gamma \max_{a'} f(s', a') \right) \right)^2 \right]$$

$\left(\mathbb{E} f(s,a) - \left(r + \gamma \max_{a'} \mathbb{E} f(s', a') \right) \right)$

$$= \mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - (\mathcal{T}f)(s,a) \right)^2 \right] + \mathbb{E}_{(s,a) \sim \mu} \left[\left((\mathcal{T}f)(s,a) - \left(r + \gamma \max_{a'} f(s', a') \right) \right)^2 \right]$$

$\mathbb{E} \left[\frac{f(s,a) - \mathbb{E} f(s', a')}{r + \gamma \max_{a'} f(s', a')} \right]$

- Workaround #1: If you can get **two samples** from the same s, a then you could do stochastic approximation for both.

$$\mathbb{E} \left(\left(f(s,a) - \left(r_A + \gamma \max_{a' \in \mathcal{A}} f(s'_A, a') \right) \right) \left(f(s,a) - \left(r_B + \gamma \max_{a' \in \mathcal{A}} f(s'_B, a') \right) \right) \right)$$

Workaround #2: Solve a saddle point problem instead (Antos et al. 08)

- Idea: let us estimate the second term and subtract it away.

$$\arg \min_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \left(\mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f(s', a') \right) \right)^2 - \left(g(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f(s', a') \right) \right)^2 \right] \right)$$

Handwritten annotations: $\mathbb{E} \left(r + \gamma \max_{a' \in \mathcal{A}} f(s', a') \right) \forall f$ with arrows pointing to the corresponding terms in the equation.

- If function class G is sufficiently expressive, then it can make the second term 0 for all f.

Quick checkpoint

- Idea: Approximate Q^* function
 - Minimize the best approximation error by
 - Value iterations?
 - Direct Bellman Residual Minimization?
- Still an active area of research in both theory and practice
- Standard techniques seem to work (especially when you have a simulator and can restart..)

Function approximations of other quantities

- Approximating the occupancy measure of logging policy

$$E_{(S,a) \sim d^{\pi}(S)} \frac{\pi(a|s)}{u(a|s)} \cdot V(S,a) = V^{\pi}$$

- Approximation the occupancy measure of target policy

- Approximation of the importance weights

$$\max_{W \in \mathcal{F}} E_{S \sim \tau} W(S,a) V(S,a)$$

← Invar. Constraints, W

Practical Course at HKUST 2021
Dual/Dual

Remaining part of the lecture

- A theoretical analysis of Fitted Q-Iterations
- Under a number of additional conditions to make it tractable

Fitted Q Iterations: Problem Setup

- Infinite Horizon Discounted MDP

$$\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \gamma, P, r, \rho\}$$

$$\sup_{s,a} r(s,a) \in [0, 1]$$

$$\mathcal{D} := \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$$

- Goal: find nearly optimal policy

$$V^* - V^{\pi} < \epsilon$$

do $\mathcal{M}, \mathcal{S}_0$

$s_t \sim P(\cdot | s_t, a_t)$
 $r_t \sim r(s_t, a_t)$

$a_t \sim \mu(a | s_t)$

iid: by \mathcal{S}_i in $\mathcal{D}(s)$

$d_{\gamma}^{\pi} = \sum_{t=0}^{\infty} \gamma^t |V(s_t) - V^*|$

Assumptions

$$|\mathcal{F}| < +\infty$$

1. Realizability

We assume \mathcal{F} is rich enough such that $Q^* \in \mathcal{F}$.

2. Uniform concentrability

$$d^{\mu}(s,a) > d_m \text{ for all } s,a$$

$$\forall \pi, h, x, a : \frac{d^{\pi}(s,a)}{\mu(s,a)} \leq C.$$

3. Bellman completeness

We assume that for any $f \in \mathcal{F}$, $\mathcal{T}f \in \mathcal{F}$.

Recap: the FQI algorithm

- Initialize at any function f_0 in the function class \mathcal{F}

$$\text{FQI: } f_t \in \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \left(f(s'_i, a_i) - r_i - \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s_i, a_i) \right)^2.$$

- Natural policy of the approximate Q^* function

$$\pi_k(s) := \operatorname{argmax}_a f_k(s, a), \forall s.$$

$$V^k - V^{\pi_k} \xrightarrow{(\leftarrow) \leftarrow \text{or}} \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$

FQI “works” under these conditions

$$\boxed{|V|=n \cdot H} \quad \sqrt{\frac{H^3}{N \ln n}}$$

$$\log |F| = \sqrt{S H \log d}$$

$$A \approx A^S H$$

Theorem 15.4 (AJKS): Assume Assumptions 1--3, with probability at least $1 - \delta$

$$V^* - V^{\pi_k} \leq \mathcal{O} \left(\frac{1}{(1-\gamma)^3} \sqrt{\frac{\mathcal{O}(\log(|F|/\delta))}{n}} \right) + \frac{2\gamma^k}{(1-\gamma)^2}$$

Uniform Convergence

+ Approx. error

$$(1-\gamma)^{-3} \simeq H \quad \text{Statistical}$$

$$\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$

$$\log |F| \simeq d$$

Optimization error



$\int |S| \cdot H^3$ suboptimal when applying it to the Tabular Setting

Sketch of the proof

- Use contraction of properties of Bellman updates
- Use uniform convergence
 - the approximate Bellman update using the offline dataset is similar to the actual exact Bellman update
 - simultaneously for all functions within the function class

- By Performance Difference Lemma

$$(1 - \gamma)(V^* - V^{\pi_k}) = \mathbb{E}_{s \sim d^{\pi_k}} [-A^*(s, \pi_k(s))]$$

$$\pi_k = \operatorname{argmax}_a (r_c(\cdot, a))$$

$$= \mathbb{E}_{s \sim d^{\pi_k}} [Q^*(s, \pi_k^*(s)) - Q^*(s, \pi_k(s))]$$

$$\leq \mathbb{E}_{s \sim d^{\pi_k}} [Q^*(s, \pi_k^*(s)) - r_c(s, \pi_k^*(s)) + r_c(s, \pi_k^*(s)) - r_c(s, \pi_k(s)) + r_c(s, \pi_k(s)) - Q^*(s, \pi_k(s))]$$

≤ 0

$$\| \cdot \|_{2, d^{\pi_k}} \implies \mathbb{E}_{s \sim d^{\pi_k}} [Q^* - r_c | (s, \pi_k^*(s))] + \mathbb{E}_{s \sim d^{\pi_k}} [r_c - Q^* | (s, \pi_k(s))]$$

$$\leq \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [(Q^* - r_c)^2 | (s, \pi_k^*(s))]} + \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [(Q^* - r_c)^2 | (s, \pi_k(s))]}$$

$$\| Q^* - r_c \|_{2, d^{\pi_k^*}}$$

$$\| Q^* - r_c \|_{2, d^{\pi_k}}$$

$$\underline{(\mathbb{E} X)^2 \leq \mathbb{E} X^2}$$

$$\underline{\left(\max_a Q^* - \max_a f_{k-1} \right)^2 \leq \max_a (Q^* - f_{k-1})^2}$$

$$\frac{\nu}{m} \leq C$$

$$\underline{\|Q^* - f_k\|_{2,\nu}} \leq \|Q^* - \mathcal{T}f_{k-1}\|_{2,\nu} + \|f_k - \mathcal{T}f_{k-1}\|_{2,\nu}$$

$$\leq \gamma \sqrt{\mathbb{E}_{s,a \sim \nu} \left[\left(\mathbb{E}_{s' \sim P(\cdot|s,a)} \max_a Q^*(s',a) - \max_a f_{k-1}(s',a) \right)^2 \right]} + \|f_k - \mathcal{T}f_{k-1}\|_{2,\nu}$$

$$\Rightarrow \leq \gamma \sqrt{\mathbb{E}_{s,a \sim \nu, s' \sim P(\cdot|s,a)} \max_a (Q^*(s',a) - f_{k-1}(s',a))^2} + \sqrt{C} \|f_k - \mathcal{T}f_{k-1}\|_{2,\mu}$$

$$\uparrow \|Q^* - f_{k-1}\|_{2,\nu'}$$

$$\underline{\nu'(s',a')} = \sum_{s,a} \nu(s,a) P(s'|s,a) \mathbf{1}\{a' = \operatorname{argmax}_a (Q^*(s',a) - f_{k-1}(s',a))^2\}$$

Roll in with ν , then take action a'

Recursive application

$$\|Q^* - f_k\|_{2,\nu} \leq \gamma \|Q^* - f_{k-1}\|_{2,\nu'} + \sqrt{C} \|f_k - \mathcal{T}f_{k-1}\|_{2,\mu}.$$

for over $t = k-1, \dots, 0$

$$\|Q^* - f_k\|_{2,\nu} \leq \sqrt{C} \sum_{t=0}^{k-1} \gamma^t \|f_{k-t} - \mathcal{T}f_{k-t-1}\|_{2,\mu} + \gamma^k \|Q^* - f_0\|_{2,\tilde{\nu}},$$

$$\gamma \|Q^* - f_0\|_{2,\tilde{\nu}} \leq \gamma^k V_{\max}.$$

$$\begin{aligned} &\leq \gamma^k V_{\max} \\ &V_{\max} \leq \frac{1}{1-\gamma} \end{aligned}$$

Reducing to a uniform convergence problem

- Apply Uniform convergence

$$\sqrt{C} \sum_{t=0}^{k-1} \gamma^t \|f_{k-t} - \mathcal{T}f_{k-t-1}\|_{2,\mu} \leq \mathcal{O} \left(\sqrt{C} \sum_{t=0}^{k-1} \gamma^k \sqrt{\frac{V_{\max}^2 \ln(|\mathcal{F}|/\delta)}{n}} \right) = \mathcal{O} \left(\frac{V_{\max} \sqrt{C \ln(|\mathcal{F}|/\delta)}}{(1-\gamma)\sqrt{n}} \right)$$

Handwritten annotations:
 - Blue circles around f_{k-t-1} and $\mathcal{T}f_{k-t-1}$ in the summand.
 - Blue arrows pointing from the circles to the text "Bellman Operator by solving FQI" and "True Bellman Operator \mathcal{T} ".
 - Blue underlines under the summation and the final asymptotic expression.

- Finally

$$\|Q^* - f_k\|_{2,\nu} = \mathcal{O} \left(\frac{V_{\max} \sqrt{C \ln(|\mathcal{F}|/\delta)}}{(1-\gamma)\sqrt{n}} \right) + \gamma^k V_{\max},$$

Handwritten annotation: A blue underline under the entire expression.

Uniform convergence


Lemma 15.5 (Least Square Generalization Error). Given $f \in \mathcal{F}$, denote $\hat{f}_f := \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n (f(s_i, a_i) - r_i - \gamma \max_{a'} f(s'_i, a'))^2$. With probability at least $1 - \delta$, for all $f \in \mathcal{F}$, we have:

$$\mathbb{E}_{s, a \sim \mu} \left(\hat{f}_f(s, a) - \mathcal{T}f(s, a) \right)^2 = \mathcal{O} \left(\frac{V_{\max}^2 \ln \left(\frac{|\mathcal{F}|}{\delta} \right)}{n} \right).$$

- Implies that

$$\mathbb{E}_{s, a \sim \mu} \left(f_t(s, a) - \mathcal{T}f_{t-1}(s, a) \right)^2 = \mathcal{O} \left(\frac{V_{\max}^2 \ln \left(\frac{|\mathcal{F}|}{\delta} \right)}{n} \right).$$

Proof sketch

- Bernstein + Union bound 

$$z_i^f := \left(f(s_i, a_i) - r_i - \gamma \max_{a' \in \mathcal{A}} f'(s'_i, a') \right)^2 - \underbrace{\left(\mathcal{T} f'(s_i, a_i) - r_i - \gamma \max_{a' \in \mathcal{A}} f'(s'_i, a') \right)^2}.$$

- Boundedness

$$\underbrace{|z_i^f|} \leq V_{\max}^2,$$

- Small variance

$$\underbrace{V_{\max}^2} \mathbb{E}_{s, a \sim \mu} \left[\underbrace{\left(f(s, a) - \mathcal{T} f'(s, a) \right)^2} \right]$$

A self-bounding trick

$$\mathbb{E}_{s,a \sim \mu} \left(\hat{f}(s,a) - \mathcal{T}f'(s,a) \right)^2 \leq \sqrt{\frac{8V_{\max}^2 \mathbb{E}_{s,a \sim \mu} \left[(\hat{f}(s,a) - \mathcal{T}f'(s,a))^2 \right] \ln(|\mathcal{F}|/\delta)}{n}} + \frac{4V_{\max}^2 \ln(|\mathcal{F}|/\delta)}{3n}.$$

- Solve a quadratic equation

$$\mathbb{E}_{s,a \sim \mu} \left(\hat{f}(s,a) - \mathcal{T}f'(s,a) \right)^2 \leq \left(\sqrt{2} + \sqrt{10/3} \right)^2 \frac{V_{\max}^2 \ln(|\mathcal{F}|/\delta)}{n}.$$

Thank you very much!