
Offline Stochastic Shortest Path: Learning, Evaluation and Towards Optimality (Supplementary material)

Ming Yin^{*1,2}

Wenjing Chen^{*3}

Mengdi Wang⁴

Yu-Xiang Wang¹

¹Department of Computer Science, UC Santa Barbara

²Department of Statistics and Applied Probability, UC Santa Barbara

³Department of Electrical and Computer Engineering, Texas A&M University

⁴Department of Electrical and Computer Engineering, Princeton University

1 GENERAL BELLMAN EQUATION FOR A FIXED POLICY

In this section, we prove Proposition ?? . In particular, the first part of the proposition for V^* has been covered in Bertsekas and Tsitsiklis [1991]. Therefore, we only consider the second part, which is a Bellman equation for fixed policy. Moreover, we do not constraint to proper policy and our result holds true for all the policies.

Lemma 1.1 (General Bellman equation for fixed policy π). *Let π be a fixed policy, proper or improper and cost $c \geq 0$ for the SSP. Then the following Bellman equations hold:*

$$Q^\pi(s, a) = c(s, a) + P_{s,a}V^\pi, \quad V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)}[Q^\pi(s, a)]. \quad (1)$$

Proof of Lemma 1.1. By definition of Q^π , we have

$$Q^\pi(s, a) = \lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right].$$

We can rewrite term $\mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right]$ as

$$\begin{aligned} \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right] &= c(s, a) + \sum_{s'} \mathbb{P}(s' \mid s, a) \mathbb{E}_\pi \left[\sum_{h=1}^T c(s_h, a_h) \mid s_1 = s' \right] \\ &= c(s, a) + \sum_{s'} \mathbb{P}(s' \mid s, a) \left\{ \mathbb{E}_\pi \left[\sum_{h=0}^{T-1} c(s_h, a_h) \mid s_0 = s' \right] \right\}, \end{aligned}$$

where the first equality is by law of total expectation. The second equality follows from the fact that the transition kernel P is **homogeneous** in SSP.

Define the sequence $V_T(s) := \left\{ \mathbb{E}_\pi \left[\sum_{h=0}^{T-1} c(s_h, a_h) \mid s_0 = s \right] \right\}$. Since for any state-action pair (s, a) , $c(s, a) \geq 0$, we know that the sequence $\{V_T(s)\}_{T=1}^\infty$ is non-decreasing. It implies that $\lim_{T \rightarrow \infty} V_T(s)$ exists. ($\lim_{T \rightarrow \infty} V_T(s)$ either diverges to $+\infty$ or converges to a positive number.) It follows that (the following switching the order of limit and summation is valid since the summation is finite sum)

$$\lim_{T \rightarrow \infty} \sum_{s'} \mathbb{P}(s' \mid s, a) \left\{ \mathbb{E}_\pi \left[\sum_{h=0}^{T-1} c(s_h, a_h) \mid s_0 = s' \right] \right\} = \sum_{s'} \mathbb{P}(s' \mid s, a) \lim_{T \rightarrow \infty} \left\{ \mathbb{E}_\pi \left[\sum_{h=0}^{T-1} c(s_h, a_h) \mid s_0 = s' \right] \right\}. \quad (2)$$

^{*}Equal contribution.

Combine the above two equalities together, we can get

$$\begin{aligned} Q^\pi(s, a) &= c(s, a) + \sum_{s'} \mathbb{P}(s'|s, a) \lim_{T \rightarrow \infty} \left\{ \mathbb{E}_\pi \left[\sum_{h=0}^{T-1} c(s_h, a_h) \mid s_0 = s' \right] \right\} \\ &= c(s, a) + \sum_{s'} \mathbb{P}(s'|s, a) V^\pi(s'). \end{aligned} \tag{3}$$

From the definition of value function, we have (where the second line uses law of total expectation)

$$\begin{aligned} V^\pi(s) &= \lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s \right] \\ &= \lim_{T \rightarrow \infty} \mathbb{E}_{a_0} \left[\mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 \right] \mid s_0 = s \right] \\ &= \lim_{T \rightarrow \infty} \sum_a \pi(a|s) \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right] \end{aligned}$$

Similar to $\lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s \right]$, we can prove that $\lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right]$ exists. Then we have

$$\begin{aligned} V^\pi(s) &= \lim_{T \rightarrow \infty} \sum_a \pi(a|s) \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right] \\ &= \sum_a \pi(a|s) \lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{h=0}^T c(s_h, a_h) \mid s_0 = s, a_0 = a \right] = \sum_a \pi(a|s) Q^\pi(s, a). \end{aligned}$$

□

Remark 1.2. Essentially, the above proof only requires $c(s, a) \geq 0$. Moreover, even if the general Bellman equation holds, it does not imply $c^\pi + P^\pi(\cdot)$ is a contraction (i.e. doing value iteration for general policy π might not converge to V^π).

2 RESULTS FOR GENERAL STOCHASTIC SHORTEST PATH PROBLEM

Lemma 2.1. For any two contraction mapping T_1 and T_2 that are monotone (i.e. for any vector greater than $V \geq V'$, it holds $T_1 V \geq T_1 V'$ and $T_2 V \geq T_2 V'$) on the metric space $\mathbb{R}^{S'}$. Suppose V_1 and V_2 are the fixed points for T_1 and T_2 respectively. If we have $T_1(V)(s) \geq T_2(V)(s)$ for any $s \in S'$, then we have $V_1(s) \geq V_2(s)$ for any $s \in S'$.

Proof. First we have $T_1 V_1 \geq T_2 V_1$. Since V_1 is the fixed point of T_1 , we know $V_1 := T_1 V_1 \geq T_2 V_1$. By monotone property with recursion, we have that

$$V_1 \geq (T_2)^k V_1. \tag{4}$$

Since V_2 is the fixed point of T_2 , we have

$$\lim_{k \rightarrow \infty} (T_2)^k V_1 = V_2.$$

Combine the above inequalities together we can get $V_1 \geq V_2$. □

3 CONVERGENCES FOR ALGORITHM ??

Lemma 3.1. $\widehat{T}^\pi : \mathbb{R}^S \times \{0\} \rightarrow \mathbb{R}^S \times \{0\}$ is a contraction mapping, i.e., $\forall V_1, V_2 \in \mathbb{R}^{S'}$, $V_1(g) = V_2(g) = 0$, we have

$$\left\| \widehat{T}^\pi V_1 - \widehat{T}^\pi V_2 \right\|_\infty \leq \rho \|V_1 - V_2\|_\infty, \tag{5}$$

Here $\rho := \max_{\substack{s,a \\ s \neq g}} \left(\frac{n_{s,a}}{n_{s,a}+1} \right) < 1$ and $\widehat{T}^\pi V(s) = \langle \pi(\cdot|s), \widehat{c}(s, \cdot) + \widetilde{P}_{s, \cdot} V \rangle$.

Proof of Lemma 3.1. We first prove the result for state g . Since g is a zero-cost absorbing state, we have for any $a \in \mathcal{A}$, $\widehat{c}(g, a) = 0$ and $\widetilde{P}_{g,a}V = V(g)$. Then for any $V \in \mathbb{R}^S \times \{0\}$, $V(g) = 0$, we have

$$\widehat{\mathcal{T}}^\pi V(g) = \langle \pi(\cdot|g), \widehat{c}(g, \cdot) + \widetilde{P}_{g,\cdot}V \rangle = 0. \quad (6)$$

Therefore $\widehat{\mathcal{T}}^\pi V_1(g) - \widehat{\mathcal{T}}^\pi V_2(g) = 0 \leq \rho \|V_1 - V_2\|_\infty$. Next we only need to prove for all state $\forall s \neq g$. Indeed,

$$\begin{aligned} |\widehat{\mathcal{T}}^\pi V_1(s) - \widehat{\mathcal{T}}^\pi V_2(s)| &= |\langle \pi(\cdot|s), \widetilde{P}_{s,\cdot}(V_1 - V_2) \rangle| \\ &\leq \max_a |\widetilde{P}_{s,a}(V_1 - V_2)| \\ &= \max_a \left| \sum_{s' \neq g} \widetilde{P}(s'|s, a)(V_1(s') - V_2(s')) \right| \\ &\leq \max_a \left(\frac{n_{s,a}}{n_{s,a} + 1} \right) \left| \sum_{s' \neq g} \widehat{P}(s'|s, a)(V_1(s') - V_2(s')) \right| \\ &\leq \max_a \left(\frac{n_{s,a}}{n_{s,a} + 1} \right) \|V_1 - V_2\|_\infty. \end{aligned} \quad (7)$$

where the second inequality is due to $V_1(g) = V_2(g) = 0$ and the third inequality is by the definition of \widetilde{P} . Take the supremum over s , we get

$$\left\| \widehat{\mathcal{T}}^\pi V_1 - \widehat{\mathcal{T}}^\pi V_2 \right\|_\infty \leq \max_{\substack{s,a \\ s \neq g}} \left(\frac{n_{s,a}}{n_{s,a} + 1} \right) \|V_1 - V_2\|_\infty \quad (8)$$

□

Lemma 3.2. $\forall \pi \in \Pi_{proper}$, define $\widehat{V}^\pi := \lim_{i \rightarrow \infty} V^{(i)}$ (Note by Lemma 3.1 this limit always exists since \widehat{V}^π is the fixed point of $\widehat{\mathcal{T}}^\pi$ and $V^{(i+1)} = \widehat{\mathcal{T}}^\pi V^{(i)}$). Then (recall $\rho := \max_{\substack{s,a \\ s \neq g}} \left(\frac{n_{s,a}}{n_{s,a} + 1} \right) < 1$)

$$\|\widehat{V}^\pi\|_\infty \leq \max_{\substack{s,a \\ s \neq g}} n(s, a) + 1.$$

Proof of Lemma 3.2. Recall the definition, $\widehat{V}^\pi = \widehat{\mathcal{T}}^\pi \widehat{V}^\pi$

$$\left\| \widehat{V}^\pi \right\|_\infty = \left\| \widehat{\mathcal{T}}^\pi \widehat{V}^\pi \right\|_\infty \quad (9)$$

$$\leq \max_{s, s \neq g} |\langle \pi(\cdot|s), \widehat{c}(s, \cdot) \rangle| + \max_{s, s \neq g} |\langle \pi(\cdot|s), \widetilde{P}_{s,\cdot} \widehat{V}^\pi \rangle| \quad (10)$$

$$\leq \max_{s, s \neq g} |\langle \pi(\cdot|s), \widehat{c}(s, \cdot) \rangle| + \max_{\substack{s,a \\ s \neq g}} |\widetilde{P}_{s,a} \widehat{V}^\pi|$$

$$\leq 1 + \max_{\substack{s,a \\ s \neq g}} |\widetilde{P}_{s,a} \widehat{V}^\pi|$$

$$\leq 1 + \max_{\substack{s,a \\ s \neq g}} \left\{ \left(\frac{n_{s,a}}{n_{s,a} + 1} \right) \left| \sum_{s' \neq g} \widehat{P}(s'|s, a) \widehat{V}^\pi(s') \right| \right\}$$

$$\leq 1 + \rho \left\| \widehat{V}^\pi \right\|_\infty. \quad (11)$$

□

The first inequality follows from $\widehat{\mathcal{T}}^\pi \widehat{V}^\pi(g) = 0$ and triangle inequality. Since $\rho < 1$, we can get $\left\| \widehat{V}^\pi \right\|_\infty \leq \frac{1}{1-\rho}$. From the definition of ρ , we can conclude the proof.

Lemma 3.3. $\left\| \widehat{V}^\pi - V^{(i)} \right\|_\infty \leq \frac{\epsilon_{OPE}}{1-\rho}$, where $\rho := \max_{\substack{s,a \\ s \neq g}} \left(\frac{n_{s,a}}{n_{s,a} + 1} \right) < 1$ as in Lemma 3.1 and $V^{(i)}$ is the output of Algorithm ??.

Proof of Lemma 3.3. Using definition $\widehat{V}^\pi := \lim_{j \rightarrow \infty} V^{(j)}$ and the telescoping sum we obtain

$$\left\| \widehat{V}^\pi - V^{(i)} \right\|_\infty \leq \sum_{j=i}^{\infty} \left\| V^{(j+1)} - V^{(j)} \right\|_\infty \leq \left\| V^{(i+1)} - V^{(i)} \right\|_\infty \sum_{j=0}^{\infty} \rho^j \leq \frac{\epsilon_{\text{OPE}}}{1-\rho},$$

where the second inequality uses $\widehat{\mathcal{T}}^\pi$ is a ρ -contraction. \square

Remark 3.4. Throughout the paper, we denote the number of state-action visitation as either $n_{s,a}$ or $n(s,a)$. They represent the same quantity.

Lemma 3.5. For any $V(\cdot) \in \mathbb{R}^S$ satisfying $V(g) = 0$,

$$|(\widetilde{P}_{s,a} - \widehat{P}_{s,a})V| \leq \frac{\|V\|_\infty}{n(s,a) + 1}, \quad |\text{Var}(\widetilde{P}_{s,a}, V) - \text{Var}(\widehat{P}_{s,a}, V)| \leq \frac{2\|V\|_\infty^2 S}{n(s,a) + 1}. \quad (12)$$

Proof. See Lemma 12 of Tarbouriech et al. [2021]. \square

4 SOME KEY LEMMAS

4.1 HIGH PROBABILITY EVENT

We define the “good property” event $\mathcal{E} := \mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3 \cap \mathcal{E}_4 \cap \mathcal{E}_5$ according the following (where $\iota := \log(10S^2A/\delta)$)

$$\begin{aligned} \mathcal{E}_1 &:= \left\{ \forall (s, a, s') \in (\mathcal{S} \times \mathcal{A} \times \mathcal{S}), \forall n(s, a) \geq 1 : |P(s'|s, a) - \widehat{P}(s'|s, a)| \leq \sqrt{\frac{2P(s'|s, a)\iota}{n(s, a)}} + \frac{2\iota}{3n(s, a)} \right\} \\ \mathcal{E}_2 &:= \left\{ \forall (s, a) \in (\mathcal{S} \times \mathcal{A}), \forall n(s, a) \geq 1 : |(P_{s,a} - \widehat{P}_{s,a})V| \leq \sqrt{\frac{2\text{Var}(P_{s,a}, V)\iota}{n(s, a)}} + \frac{2\|V\|_\infty\iota}{3n(s, a)} \right\} \\ \mathcal{E}_3 &:= \left\{ \forall (s, a) \in (\mathcal{S} \times \mathcal{A}), \forall n(s, a) \geq 1 : |(P_{s,a} - \widehat{P}_{s,a})V| \leq \sqrt{\frac{2\text{Var}(\widehat{P}_{s,a}, V)\iota}{n(s, a)}} + \frac{7\|V\|_\infty\iota}{3n(s, a)} \right\} \\ \mathcal{E}_4 &:= \left\{ \forall (s, a) \in (\mathcal{S} \times \mathcal{A}), \forall n(s, a) \geq 1 : |\widehat{c}(s, a) - c(s, a)| \leq \sqrt{\frac{2\text{Var}_c(s, a)\iota}{n(s, a)}} + \frac{2\iota}{3n(s, a)} \right\} \\ \mathcal{E}_5 &:= \left\{ \forall (s, a) \in (\mathcal{S} \times \mathcal{A}), \forall n(s, a) \geq 1 : |\widehat{c}(s, a) - c(s, a)| \leq \sqrt{\frac{2\widehat{c}(s, a)\iota}{n(s, a)}} + \frac{7\iota}{3n(s, a)} \right\}. \end{aligned} \quad (13)$$

Lemma 4.1. The event \mathcal{E} holds for any V that is independent from \widehat{P} with probability $1 - \frac{\delta}{2}$.

Proof. From the empirical Bernstein’s inequality given in Lemma 12.4, we have that for each fixed (s, a) , the event $|(P_{s,a} - \widehat{P}_{s,a})V| \leq \sqrt{\frac{2\text{Var}(\widehat{P}_{s,a}, V)\iota}{n(s, a)}} + \frac{7\|V\|_\infty\iota}{3n(s, a)}$ holds with probability $1 - \frac{\delta}{10S^2A}$. By taking a union bound, we have that event \mathcal{E}_3 holds with probability $1 - \frac{\delta}{10S}$. Similarly, we have event \mathcal{E}_5 holds with probability $1 - \frac{\delta}{10S}$. By applying the standard Bernstein’s inequality in Lemma 12.3 and taking union bound over (s, a, s') , we can get that event $\mathcal{E}_1, \mathcal{E}_2$ and \mathcal{E}_4 holds with probability $1 - \frac{\delta}{10}$. Since \mathcal{E} is the intersection of the above events, we can prove the lemma by taking a union bound again over all of the five events. \square

4.2 VALUE DECOMPOSITION LEMMA

Lemma 4.2. Suppose $\widehat{V}^\pi := \lim_{j \rightarrow \infty} V^{(j)}$ where $V^{(j)} = \widehat{\mathcal{T}}^\pi V^{(j-1)}$ for all j , then we have the following suboptimality decomposition for any initial state \bar{s} :

$$\widehat{V}^\pi(\bar{s}) - V^\pi(\bar{s}) = \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h, \bar{s}}^\pi(s, a) \{(\widehat{c} - c)(s, a) + (\widetilde{P}_{s,a} - P_{s,a})\widehat{V}^\pi\} \quad (14)$$

Proof of Lemma 4.2. We prove this lemma by recursion. First, we have for any $h \geq 0$,

$$\begin{aligned}
\sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s) (\widehat{V}^\pi(s) - V^\pi(s)) &= \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s) \sum_a \pi(a|s) (\widehat{Q}^\pi(s,a) - Q^\pi(s,a)) \\
&= \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) (\widehat{Q}^\pi(s,a) - Q^\pi(s,a)) \\
&= \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \{(\widehat{c} - c)(s,a) + (\widetilde{P}_{s,a} \widehat{V}^\pi - P_{s,a} V^\pi)\} \\
&= \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \{(\widehat{c} - c)(s,a) + (\widetilde{P}_{s,a} - P_{s,a}) \widehat{V}^\pi + P_{s,a} (\widehat{V}^\pi - V^\pi)\} \\
&= \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \{(\widehat{c} - c)(s,a) + (\widetilde{P}_{s,a} - P_{s,a}) \widehat{V}^\pi\} + \sum_{\substack{s \\ s \neq g}} \xi_{h+1,\bar{s}}^\pi(s) (\widehat{V}^\pi - V^\pi)(s),
\end{aligned}$$

where the third equality uses both Bellman equations and empirical Bellman equations and the last equality follows from the fact that $\xi_{h+1}(s') = \sum_{s,a} \xi_h(s,a) P(s'|s,a)$ and $\widehat{V}^\pi(g) = V^\pi(g) = 0$. By recursion, we have that

$$\widehat{V}^\pi(\bar{s}) - V^\pi(\bar{s}) = \sum_{h=0}^H \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \{(\widehat{c} - c)(s,a) + (\widetilde{P}_{s,a} - P_{s,a}) \widehat{V}^\pi\} + \sum_{\substack{s \\ s \neq g}} \xi_{H+1,\bar{s}}^\pi(s) (\widehat{V}^\pi(s) - V^\pi(s)), \quad (15)$$

for all H . Then we have

$$\left| \sum_{\substack{s \\ s \neq g}} \xi_{H+1,\bar{s}}^\pi(s) (\widehat{V}^\pi(s) - V^\pi(s)) \right| \leq \sum_{\substack{s \\ s \neq g}} \xi_{H+1,\bar{s}}^\pi(s) \cdot \left\| \widehat{V}^\pi - V^\pi \right\|_\infty \leq P_{\bar{s}}^\pi(s_{H+1} \neq g) \cdot \left\| \widehat{V}^\pi - V^\pi \right\|_\infty.$$

Since π is proper, we have $\|V^\pi\|_\infty \leq \infty$ and by Lemma 4.3 $\lim_{H \rightarrow +\infty} P_{\bar{s}}^\pi(s_H \neq g) = 0$. From Lemma 3.2, we have $\|\widehat{V}^\pi\|_\infty \leq \infty$. It follows that

$$\lim_{H \rightarrow +\infty} \sum_{\substack{s \\ s \neq g}} \xi_{H+1,\bar{s}}^\pi(s) (\widehat{V}^\pi(s) - V^\pi(s)) = 0 \quad (16)$$

By taking H to infinity in Equation (15), we conclude the proof. \square

4.3 KEY LEMMAS: ARRIVAL TIME DECOMPOSITION AND DEPENDENCE IMPROVEMENT FOR SSP

Below we present two lemmas for SSP problem, which is the key for obtaining tight instance-dependent bounds.

Lemma 4.3 (Arrival time decomposition). *Let $T_{\bar{s}}^\pi$ be the expected time of arrival to goal state g when applying proper policy π and starting from \bar{s} , then $T_{\bar{s}}^\pi = \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a)$. Moreover, $T_{\bar{s}}^\pi < \infty$ for all \bar{s} .*

Proof of Lemma 4.3. Denote T to be the random variable of arrival time to goal state g when applying proper policy π , starting from \bar{s} . Then $\mathbb{E}[T] = T_{\bar{s}}^\pi$. Furthermore, since T is non-negative integral variable, it holds $\mathbb{E}[T] = \sum_{h=0}^{\infty} \mathbb{P}(T > h)$.

Then we have

$$\begin{aligned}
T_{\bar{s}}^\pi &= \mathbb{E}_{P,\pi} T = \sum_{h=0}^{\infty} \mathbb{P}_{P,\pi}(T > h) \\
&= \sum_{h=0}^{\infty} \mathbb{P}_{P,\pi}(s_1 \neq g, s_2 \neq g, \dots, s_h \neq g) \\
&\stackrel{(i)}{=} \sum_{h=0}^{\infty} \mathbb{P}_{P,\pi}(s_h \neq g) = \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a),
\end{aligned}$$

where equality (i) follows from the fact that g is an absorbing state, so we can only reach a state which is not a goal state if all the previous steps are not goal state and vice versa.

Lastly, since π is proper, $T_{\bar{s}}^{\pi} < \infty$ for all \bar{s} and this implies $\lim_{h \rightarrow \infty} \mathbb{P}_{P, \pi}(s_h \neq g) = 0$. \square

The next lemma is the key for achieving optimal rate.

Lemma 4.4 (Dependency Improvement). *For any probability transition matrix P , policy π , and any cost function $c \in [0, 1]$, we use $\xi_h^{\pi}(s, a)$ to denote the probability of visiting (s, a) associated with $\widehat{SSP}(P, \pi)$. Suppose $V \in \mathbb{R}^{S+1}$ is any value function satisfying order property (where $V(g) = 0$), i.e., $V(s) \geq \sum_a \pi(a|s)P_{s,a}V$ for all $s \in \mathcal{S}$, then we have*

$$\sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s, a) \text{Var}(P_{s,a}, V) \leq 2 \|V\|_{\infty} \sum_{\substack{s \\ s \neq g}} \xi_0(s) V(s) \leq 2 \|V\|_{\infty}^2. \quad (17)$$

Proof of Lemma 4.4.

$$\begin{aligned} & \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s, a) \text{Var}(P_{s,a}, V) = \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s, a) \{P_{s,a}(V)^2 - (P_{s,a}V)^2\} \\ &= \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h+1}^{\pi}(s) V^2(s) - \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s, a) (P_{s,a}V)^2 \\ &\leq \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\pi}(s) V^2(s) - \sum_{h=0}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s, a) (P_{s,a}V)^2 \\ &\stackrel{(i)}{=} \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\pi}(s) \{V^2(s) - \sum_a \pi(a|s) (P_{s,a}V)^2\} \\ &\stackrel{(ii)}{\leq} \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\pi}(s) \{V^2(s) - (\sum_a \pi(a|s) P_{s,a}V)^2\} \\ &= \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\pi}(s) \{(V(s) - \sum_a \pi(a|s) P_{s,a}V)(V(s) + \sum_a \pi(a|s) P_{s,a}V)\} \\ &\stackrel{(iii)}{\leq} 2 \|V\|_{\infty} \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\pi}(s) (V(s) - \sum_a \pi(a|s) P_{s,a}V) \\ &= 2 \|V\|_{\infty} \left[\sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\pi}(s) V(s) - \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h+1}^{\pi}(s) V(s) \right] \\ &= 2 \|V\|_{\infty} \sum_{\substack{s \\ s \neq g}} \xi_0(s) V(s) \leq 2 \|V\|_{\infty}^2, \end{aligned}$$

where (i) follows from the fact that $\xi(s, a) = \xi(s)\pi(a|s)$, (ii) uses the Jensen's inequality and the fact that $f(x) = x^2$ is a convex function. (iii) uses the ordering condition. \square

5 CRUDE EVALUATION BOUND

Theorem 5.1. *Denote $d_m := \min\{\sum_{h=1}^{\infty} \xi_h^{\mu}(s, a) : s.t. \sum_{h=1}^{\infty} \xi_h^{\mu}(s, a) > 0\}$, and T_s^{π} to be the expected time to hit g when starting from s . Define $\bar{T}^{\pi} = \max_{\bar{s} \in \mathcal{S}} T_{\bar{s}}^{\pi}$. Then when $n \geq \max\{\frac{49S_L}{9d_m}, 64(\bar{T}^{\pi})^2 \frac{S_L}{d_m}, C \cdot \log(SA/\delta) / \sum_{h=1}^{\infty} \xi_h^{\mu}(s, a)\}$, we*

have with probability $1 - \delta$, (here $\iota = O(\log(SA/\delta))$)

$$\left\| \widehat{V}^\pi - V^\pi \right\|_\infty \leq O \left(\frac{\bar{T}^\pi \sqrt{\max_{s,a} \text{Var}_c(s,a)} \iota + \sqrt{\bar{T}^\pi \|V^\pi\|_\infty^2 \iota}}{\sqrt{n \cdot d_m}} \right) + O \left(\frac{\bar{T}^\pi \|V^\pi\|_\infty \cdot \iota}{n \cdot d_m} \right).$$

Proof. We denote $\bar{s} \in \mathcal{S}$ to be any initial state. From Lemma 4.2, we have (here $\xi_{h,\bar{s}}^\pi(s,a)$ is the marginal state-action probability when starting from state \bar{s} and following π)

$$\begin{aligned} |V^\pi(\bar{s}) - \widehat{V}^\pi(\bar{s})| &= \left| \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \{(\widehat{c} - c)(s,a) + (\tilde{P}_{s,a} - P_{s,a}) \widehat{V}^\pi\} \right| \\ &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) |(\widehat{c} - c)(s,a)| + \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) |(\tilde{P}_{s,a} - P_{s,a}) (\widehat{V}^\pi - V^\pi)| \\ &\quad + \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) |(\tilde{P}_{s,a} - P_{s,a}) V^\pi| \end{aligned}$$

We bound the above three parts one by one. First of all, by Bernstein inequality, Lemma 12.6 and union bound, with probability $1 - \delta$,

$$\begin{aligned} \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) |(\widehat{c} - c)(s,a)| &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \left[2 \sqrt{\frac{\text{Var}_c(s,a) \iota}{n(s,a)}} + \frac{4\iota}{3n(s,a)} \right] \\ &\stackrel{(i)}{\leq} \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \left[2 \sqrt{\frac{2 \text{Var}_c(s,a) \iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s,a)}} + \frac{8\iota}{3n \sum_{h=1}^{\infty} \xi_h^\mu(s,a)} \right] \\ &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \left[2 \sqrt{\frac{2 \max_{s,a} \text{Var}_c(s,a) \iota}{n \cdot d_m}} + \frac{8\iota}{3n \cdot d_m} \right] \\ &\leq \bar{T}_{\bar{s}}^\pi \left[2 \sqrt{\frac{2 \max_{s,a} \text{Var}_c(s,a) \iota}{n \cdot d_m}} + \frac{8\iota}{3n \cdot d_m} \right] \end{aligned}$$

(i) uses Lemma 12.6 and the last inequality uses Lemma 4.3.

For the second part, note

$$\begin{aligned} &\sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) |(\tilde{P}_{s,a} - P_{s,a}) (\widehat{V}^\pi - V^\pi)| \\ &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \left[\sqrt{\frac{2S \cdot \text{Var}(P_{s,a}, \widehat{V}^\pi - V^\pi) \iota}{n(s,a)}} + \frac{4 \|\widehat{V}^\pi - V^\pi\|_\infty S \iota}{3n(s,a)} + \frac{\|\widehat{V}^\pi - V^\pi\|_\infty}{n(s,a) + 1} \right] \\ &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \left[\sqrt{\frac{4S \cdot \text{Var}(P_{s,a}, \widehat{V}^\pi - V^\pi) \iota}{n \sum_{h=1}^{\infty} \xi_h^\pi(s,a)}} + \frac{14 \|\widehat{V}^\pi - V^\pi\|_\infty S \iota}{3n \sum_{h=1}^{\infty} \xi_h^\pi(s,a)} \right] \\ &\leq \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{\sum_{h=0}^{\infty} \xi_{h,\bar{s}}^\pi(s,a)}{\sum_{h=0}^{\infty} \xi_h^\mu(s,a)}} \sqrt{4S \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^\pi(s,a) \|\widehat{V}^\pi - V^\pi\|_\infty^2 \frac{\iota}{n}} + \frac{14 \|\widehat{V}^\pi - V^\pi\|_\infty S \iota}{3n} \sum_{s,a, s \neq g} \frac{d_{\bar{s}}^\pi(s,a)}{d^\mu(s,a)} \end{aligned}$$

$$\begin{aligned}
&\leq \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d_{\bar{s}}^{\pi}(s,a)}{d^{\mu}(s,a)} \cdot 4ST_{\bar{s}}^{\pi} \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty}^2 \cdot \frac{\iota}{n}} + \frac{14 \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} S\iota}{3n} \sum_{s,a,s \neq g} \frac{d_{\bar{s}}^{\pi}(s,a)}{d^{\mu}(s,a)} \\
&\leq 4 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d_{\bar{s}}^{\pi}(s,a)}{d^{\mu}(s,a)} \cdot ST_{\bar{s}}^{\pi} \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty}^2 \cdot \frac{\iota}{n}} \leq 4T_{\bar{s}}^{\pi} \sqrt{\frac{S\iota}{n \cdot d_m}} \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty},
\end{aligned}$$

where the first inequality uses Lemma 10.3 and Lemma 3.5, the second inequality uses Lemma 12.6, the third inequality uses $\text{Var}(\cdot) \leq \|\cdot\|_{\infty}^2$, the fourth and fifth inequality use Lemma 4.3 and the last inequality follows from the condition $n \geq \frac{49S\iota}{9d_m} \geq \frac{49S\iota}{9T_{\bar{s}}^{\pi}} \sum_{s,a,s \neq g} \frac{d_{\bar{s}}^{\pi}(s,a)}{d^{\mu}(s,a)}$.

For the third part, we have

$$\begin{aligned}
&\sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^{\pi}(s,a) |(\tilde{P}_{s,a} - P_{s,a})V^{\pi}| \\
&\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^{\pi}(s,a) \left[2\sqrt{\frac{\text{Var}(P_{s,a}, V^{\pi})\iota}{n(s,a)}} + \frac{4\|V^{\pi}\|_{\infty}\iota}{3n(s,a)} + \frac{\|V^{\pi}\|_{\infty}}{n(s,a)} \right] \\
&\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_{h,\bar{s}}^{\pi}(s,a) \left[2\sqrt{\frac{2\text{Var}(P_{s,a}, V^{\pi})\iota}{n \sum_{h=1}^{\infty} \xi_h^{\mu}(s,a)}} + \frac{14\|V^{\pi}\|_{\infty}\iota}{3n \sum_{h=1}^{\infty} \xi_h^{\mu}(s,a)} \right] \\
&= \sqrt{8 \sum_{\substack{s,a \\ s \neq g}} \frac{\sum_{h=1}^{\infty} \xi_{h,\bar{s}}^{\pi}(s,a)}{\sum_{h=1}^{\infty} \xi_h^{\mu}(s,a)} \cdot \sum_{s,a,s \neq g} \sum_{h=1}^{\infty} \xi_{h,\bar{s}}^{\pi}(s,a) \text{Var}(P_{s,a}, V^{\pi}) \frac{\iota}{n}} + \frac{14\|V^{\pi}\|_{\infty}\iota}{3n} \cdot \sum_{s,a,s \neq g} \frac{d_{\bar{s}}^{\pi}(s,a)}{d^{\mu}(s,a)}} \\
&\leq \sqrt{8 \frac{T_{\bar{s}}^{\pi}}{d_m} \cdot \frac{\|V^{\pi}\|_{\infty}^2 \iota}{n}} + \frac{14\|V^{\pi}\|_{\infty}\iota}{3n} \cdot \frac{T_{\bar{s}}^{\pi}}{d_m}.
\end{aligned}$$

where the first inequality uses Lemma 3.5 and Bernstein inequality, the second inequality uses Lemma 12.6 and the last one uses Lemma 4.3. Recall $\bar{T}^{\pi} = \max_{\bar{s} \in \mathcal{S}} T_{\bar{s}}^{\pi}$, then combine all the three parts together and take the max over \bar{s} , we can derive

$$\begin{aligned}
\left(1 - 4\bar{T}^{\pi} \sqrt{\frac{S\iota}{nd_m}}\right) \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} &\leq T_{\bar{s}}^{\pi} \left[2\sqrt{\frac{2 \max_{s,a} \text{Var}_c(s,a)\iota}{n \cdot d_m}} + \frac{8\iota}{3n \cdot d_m} \right] + \sqrt{8 \frac{T_{\bar{s}}^{\pi}}{d_m} \cdot \frac{\|V^{\pi}\|_{\infty}^2 \iota}{n}} + \frac{14\|V^{\pi}\|_{\infty}\iota}{3n} \frac{T_{\bar{s}}^{\pi}}{d_m} \\
&\leq O\left(\frac{\bar{T}^{\pi} \sqrt{\max_{s,a} \text{Var}_c(s,a)\iota} + \sqrt{\bar{T}^{\pi} \|V^{\pi}\|_{\infty}^2 \iota}}{\sqrt{n \cdot d_m}}\right) + O\left(\frac{\bar{T}^{\pi} \|V^{\pi}\|_{\infty} \cdot \iota}{n \cdot d_m}\right),
\end{aligned}$$

therefore it implies (by applying the condition $n \geq 64(\bar{T}^{\pi})^2 \frac{S\iota}{d_m}$)

$$\|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} \leq O\left(\frac{\bar{T}^{\pi} \sqrt{\max_{s,a} \text{Var}_c(s,a)\iota} + \sqrt{\bar{T}^{\pi} \|V^{\pi}\|_{\infty}^2 \iota}}{\sqrt{n \cdot d_m}}\right) + O\left(\frac{\bar{T}^{\pi} \|V^{\pi}\|_{\infty} \cdot \iota}{n \cdot d_m}\right)$$

□

6 PROOF OF THEOREM ??

Theorem 6.1 (Restatement of Theorem ??). *Denote $d_m := \min\{\sum_{h=1}^{\infty} \xi_h^{\mu}(s,a) : s.t. \sum_{h=1}^{\infty} \xi_h^{\mu}(s,a) > 0\}$, and T_s^{π} to be the expected time to hit g when starting from s . Define $\bar{T}^{\pi} = \max_{\bar{s} \in \mathcal{S}} T_{\bar{s}}^{\pi}$. Then when $n \geq \max\{\frac{49S\iota}{9d_m}, 64(\bar{T}^{\pi})^2 \frac{S\iota}{d_m}, C \cdot \iota/d_m\}$,*

we have with probability $1 - \delta$,

$$|V^{(i)}(s_{\text{init}}) - V^\pi(s_{\text{init}})| \leq 4 \sum_{s,a,s \neq g} d^\pi(s,a) \sqrt{\frac{2\text{Var}_{P_{s,a}}[V^\pi + c]}{n \cdot d^\mu(s,a)}} + \tilde{O}\left(\frac{1}{n}\right) + \frac{\epsilon_{\text{OPE}}}{1-\rho}.$$

where the \tilde{O} absorbs Polylog term and even higher order term.

Proof. Recall that we start from the initial state s_{init} . Then by Lemma 3.3,

$$\left| V^\pi(s_{\text{init}}) - V^{(i)}(s_{\text{init}}) \right| \leq \left| V^\pi(s_{\text{init}}) - \widehat{V}^\pi(s_{\text{init}}) \right| + \left| \widehat{V}^\pi(s_{\text{init}}) - V^{(i)}(s_{\text{init}}) \right| \leq \left| V^\pi(s_{\text{init}}) - \widehat{V}^\pi(s_{\text{init}}) \right| + \frac{\epsilon_{\text{OPE}}}{1-\rho}, \quad (18)$$

it remains to bound $|V^\pi(s_{\text{init}}) - \widehat{V}^\pi(s_{\text{init}})|$.

From Lemma 4.2, we have

$$\begin{aligned} \left| V^\pi(s_{\text{init}}) - \widehat{V}^\pi(s_{\text{init}}) \right| &= \left| \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) \{(\widehat{c} - c)(s,a) + (\tilde{P}_{s,a} - P_{s,a})\widehat{V}^\pi\} \right| \\ &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) |(\widehat{c} - c)(s,a)| + \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) |(\tilde{P}_{s,a} - P_{s,a})(\widehat{V}^\pi - V^\pi)| \\ &\quad + \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) |(\tilde{P}_{s,a} - P_{s,a})V^\pi| \end{aligned}$$

We bound the above three parts one by one. First of all, by Bernstein inequality and Lemma 12.6 together with union bound, with probability $1 - \delta$,

$$\begin{aligned} \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) |(\widehat{c} - c)(s,a)| &\leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) \left[2\sqrt{\frac{\text{Var}_c(s,a)\iota}{n(s,a)}} + \frac{4\iota}{3n(s,a)} \right] \\ &\stackrel{\text{(i)}}{\leq} \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s,a) \left[2\sqrt{\frac{2\text{Var}_c(s,a)\iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s,a)}} + \frac{8\iota}{3n \sum_{h=1}^{\infty} \xi_h^\mu(s,a)} \right] \\ &\stackrel{\text{(ii)}}{\leq} \sqrt{\frac{\sum_{\substack{s,a \\ s \neq g}} \sum_{h=1}^{\infty} \xi_h^\pi(s,a)}{\sum_{\substack{s,a \\ s \neq g}} \sum_{h=1}^{\infty} \xi_h^\mu(s,a)}} \sqrt{\frac{\sum_{\substack{s,a \\ s \neq g}} \sum_{h=1}^{\infty} \xi_h^\pi(s,a) \text{Var}_c(s,a)\iota}{n}} + \left(\sum_{s,a,s \neq g} \frac{d^\pi(s,a)}{d^\mu(s,a)} \right) \frac{8\iota}{3n} \\ &\stackrel{\text{(iii)}}{\leq} \sqrt{\frac{\sum_{\substack{s,a \\ s \neq g}} \frac{d^\pi(s,a)}{d^\mu(s,a)} \cdot \frac{V^\pi(s_{\text{init}}) \cdot \iota}{n}}{\sum_{\substack{s,a \\ s \neq g}} \frac{d^\pi(s,a)}{d^\mu(s,a)}}} + \left(\sum_{s,a,s \neq g} \frac{d^\pi(s,a)}{d^\mu(s,a)} \right) \frac{8\iota}{3n}, \end{aligned} \quad (19)$$

(i) uses Lemma 12.6 and (ii) uses Cauchy-Schwartz inequality. (iii) uses the fact that $\text{Var}_c(s,a) \leq \mathbb{E}C(s,a)^2 \leq c(s,a)$ since the realization $C(s,a) \in [0, 1]$ and the definition of $V^\pi(s_{\text{init}})$.

For the second part, note

$$\begin{aligned}
& \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h(s,a) |(\tilde{P}_{s,a} - P_{s,a})(\widehat{V}^{\pi} - V^{\pi})| \\
& \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h(s,a) \left[\sqrt{\frac{2S \cdot \text{Var}(P_{s,a}, \widehat{V}^{\pi} - V^{\pi})_{\ell_{s,a}}}{n(s,a)}} + \frac{4 \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} S_{\ell_{s,a}}}{3n(s,a)} + \frac{\|\widehat{V}^{\pi} - V^{\pi}\|_{\infty}}{n(s,a) + 1} \right] \\
& \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h(s,a) \left[\sqrt{\frac{4S \cdot \text{Var}(P_{s,a}, \widehat{V}^{\pi} - V^{\pi})_{\ell_{s,a}}}{n \sum_{h=1}^{\infty} \xi_h^{\pi}(s,a)}} + \frac{14 \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} S_{\ell_{s,a}}}{3n \sum_{h=1}^{\infty} \xi_h^{\pi}(s,a)} \right] \\
& \leq \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{\sum_{h=0}^{\infty} \xi_h^{\pi}(s,a)}{\sum_{h=0}^{\infty} \xi_h^{\mu}(s,a)}} \sqrt{4S \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h(s,a) \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty}^2 \frac{\ell}{n} + \frac{14 \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} S_{\ell}}{3n} \sum_{s,a,s \neq g} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)}} \\
& \leq \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)} \cdot 4ST^{\pi} \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty}^2 \cdot \frac{\ell}{n} + \frac{14 \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty} S_{\ell}}{3n} \sum_{s,a,s \neq g} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)}} \\
& \leq 4 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)} \cdot ST^{\pi} \|\widehat{V}^{\pi} - V^{\pi}\|_{\infty}^2 \cdot \frac{\ell}{n}}
\end{aligned} \tag{20}$$

where the first inequality uses Lemma 10.3 and Lemma 3.5, the second inequality uses Lemma 12.6, the third inequality uses $\text{Var}(\cdot) \leq \|\cdot\|_{\infty}^2$ and CS inequality, the fourth inequality use Lemma 4.3 and the last inequality follows from the condition $n \geq \frac{49S_{\ell}}{9d_m} \geq \frac{49S_{\ell}}{9T^{\pi}} \sum_{s,a,s \neq g} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)}$.

For the third part, we have

$$\begin{aligned}
& \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s,a) |(\tilde{P}_{s,a} - P_{s,a})V^{\pi}| \\
& \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s,a) \left[2\sqrt{\frac{\text{Var}(P_{s,a}, V^{\pi})_{\ell}}{n(s,a)}} + \frac{4 \|V^{\pi}\|_{\infty} \ell}{3n(s,a)} + \frac{\|V^{\pi}\|_{\infty}}{n(s,a)} \right] \\
& \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^{\pi}(s,a) \left[2\sqrt{\frac{2\text{Var}(P_{s,a}, V^{\pi})_{\ell}}{n \sum_{h=1}^{\infty} \xi_h^{\mu}(s,a)}} + \frac{14 \|V^{\pi}\|_{\infty} \ell}{3n \sum_{h=1}^{\infty} \xi_h^{\mu}(s,a)} \right] \\
& (\leq \sqrt{8 \sum_{\substack{s,a \\ s \neq g}} \frac{\sum_{h=1}^{\infty} \xi_h^{\pi}(s,a)}{\sum_{h=1}^{\infty} \xi_h^{\mu}(s,a)} \cdot \sum_{s,a,s \neq g} \sum_{h=1}^{\infty} \xi_h^{\pi}(s,a) \text{Var}(P_{s,a}, V^{\pi}) \frac{\ell}{n} + \frac{14 \|V^{\pi}\|_{\infty} \ell}{3n} \cdot \sum_{s,a,s \neq g} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)}}} \\
& \leq \sqrt{8 \sum_{\substack{s,a \\ s \neq g}} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)} \cdot \frac{\|V^{\pi}\|_{\infty}^2 \ell}{n} + \frac{14 \|V^{\pi}\|_{\infty} \ell}{3n} \cdot \sum_{s,a,s \neq g} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)}}} \\
& \leq 4 \sqrt{2 \sum_{\substack{s,a \\ s \neq g}} \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)} \cdot \frac{\|V^{\pi}\|_{\infty}^2 \ell}{n}},
\end{aligned} \tag{21}$$

where the first inequality uses Lemma 10.3 and Lemma 3.5, the second inequality uses Lemma 12.6. The third inequality uses the Cauchy-Schwartz inequality.

Combine Equation (19), (20) and (21) together, we obtain

$$\begin{aligned}
& \left| V^\pi(s_{\text{init}}) - \widehat{V}^\pi(s_{\text{init}}) \right| \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s, a) \left[2\sqrt{\frac{2\text{Var}(P_{s,a}, V^\pi)\iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)}} + \frac{14 \|V^\pi\|_\infty \iota}{3n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)} \right] \\
& + \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s, a) \left[2\sqrt{\frac{2\text{Var}_c(s, a)\iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)}} + \frac{8\iota}{3n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)} \right] \\
& + 4 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d^\pi(s, a)}{d^\mu(s, a)} \cdot ST^\pi \|\widehat{V}^\pi - V^\pi\|_\infty^2 \cdot \frac{\iota}{n}} \\
& \leq 4 \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^\pi(s, a) \sqrt{\frac{2[\text{Var}(P_{s,a}, V^\pi) + \text{Var}_c(s, a)]\iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)}} + \sum_{\substack{s,a \\ s \neq g}} \frac{22 \sum_{h=1}^{\infty} \xi_h^\pi(s, a) \|V^\pi\|_\infty \iota}{3n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)} \\
& + 4 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d^\pi(s, a)}{d^\mu(s, a)} \cdot ST^\pi \cdot \frac{(\bar{T}^\pi)^2 \cdot \max_{s,a} \text{Var}_c(s, a)\iota + \bar{T}^\pi \|V^\pi\|_\infty^2 \iota}{n \cdot d_m} \cdot \frac{\iota}{n}} + \tilde{O}\left(\frac{1}{n^{3/2}}\right) \\
& = 4 \sum_{s,a, s \neq g} d^\pi(s, a) \sqrt{\frac{2\text{Var}_{P_{s,a}}[V^\pi + c]}{nd^\mu(s, a)}} + \tilde{O}\left(\frac{1}{n}\right)
\end{aligned}$$

where the only inequality uses Theorem 5.1. Combining this with (18) we finish the proof of Theorem ??.

□

7 PREPARATIONS FOR PROVING OFFLINE LEARNING SSP

Throughout the whole section, we denote $\iota = O(\log(SA/\delta))$. All the results apply to the construction of Algorithm ?. In particular, we use \bar{V} to denote the limit of $V^{(i)}$ (by letting $\epsilon_{\text{OPL}} = 0$). This limit exists, as guaranteed by Lemma 7.6.

7.1 AUXILIARY LEMMAS

Lemma 7.1. Denote the limit of sequence $V^{(i)}$ in Algorithm ?? as \bar{V} , we have that

$$\|\bar{V}\|_\infty \leq \tilde{B}$$

Proof. First of all, by Lemma 7.6, we know \bar{V} exists. Next, from the Algorithm ??, we can get that

$$Q^{(i+1)}(s, a) = \min\{\widehat{c}(s, a) + \tilde{P}'_{s,a} V^{(i)} + b_{s,a}(V^{(i)}), \tilde{B}\} \leq \tilde{B} \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall i \in \mathbb{N}$$

and thus

$$V^{(i+1)}(s) = \min_a Q^{(i+1)}(s, a) \leq \tilde{B} \quad \forall s \in \mathcal{S} \times \mathcal{A}, \forall i \in \mathbb{N}.$$

It implies that $\bar{V}(s) = \lim_{i \rightarrow \infty} V^{(i)}(s) \leq \tilde{B}$.

□

Lemma 7.2. For any $V(\cdot) \in \mathbb{R}^{\mathcal{S}}$ satisfying $V(g) = 0$,

$$|(\tilde{P}'_{s,a} - \widehat{P}_{s,a})V| \leq \frac{\|V\|_\infty}{n_{\max} + 1}, \quad |\text{Var}(\tilde{P}'_{s,a}, V) - \text{Var}(\widehat{P}_{s,a}, V)| \leq \frac{2\|V\|_\infty^2}{n_{\max} + 1}. \quad (22)$$

Proof. The proof is similar to Lemma 12 in Tarbouriech et al. [2021]. We include the proof for completeness. Since $V(g) = 0$, for all state $s \neq g$, we have $\tilde{P}'_{s,a} V = \sum_{s', s' \neq g} \left(\frac{n_{\max}}{n_{\max} + 1}\right) \widehat{P}(s'|s, a) V(s') = \left(\frac{n_{\max}}{n_{\max} + 1}\right) \widehat{P}_{s,a} V$

$$|(\tilde{P}'_{s,a} - \widehat{P}_{s,a})V| = \left|\left(\frac{n_{\max}}{n_{\max} + 1}\right) \widehat{P}_{s,a} V - \widehat{P}_{s,a} V\right|$$

$$= \frac{|\sum_{s', s' \neq g} \widehat{P}(s'|s, a)V(s')|}{n_{max} + 1} \leq \frac{\|V\|_\infty}{n_{max} + 1}.$$

Then we prove the second inequality. Similarly,

$$\begin{aligned} |\text{Var}(\widetilde{P}'_{s,a}, V) - \text{Var}(\widehat{P}_{s,a}, V)| &= |\widetilde{P}'_{s,a}V^2 - (\widetilde{P}'_{s,a}V)^2 - \widehat{P}_{s,a}V^2 + (\widehat{P}_{s,a}V)^2| \\ &= |(\frac{n_{max}}{n_{max} + 1})\widehat{P}_{s,a}V^2 - (\frac{n_{max}}{n_{max} + 1})\widehat{P}_{s,a}V^2 - \widehat{P}_{s,a}V^2 + (\widehat{P}_{s,a}V)^2| \\ &= |-(\frac{1}{n_{max} + 1})\{\widehat{P}_{s,a}V^2 - (\widehat{P}_{s,a}V)^2\} + \frac{n_{max}}{(n_{max} + 1)^2}(\widehat{P}_{s,a}V)^2| \\ &= (\frac{1}{n_{max} + 1})|\text{Var}(\widehat{P}_{s,a}, V)| + |\frac{n_{max}}{(n_{max} + 1)^2}(\widehat{P}_{s,a}V)^2| \leq \frac{2\|V\|_\infty^2}{n_{max} + 1}. \end{aligned}$$

□

Lemma 7.3. Let $T_{\max} = \max_i T_i$ and $n > O(T_{\max}^2 \log(SA/\delta)/d_m^2)$. If in addition $n \geq \frac{S\iota}{2d_m}$, with probability at least $1 - \delta$, we have that for any state action pair (s, a) ,

$$|(P_{s,a} - \widetilde{P}'_{s,a})\bar{V}| \leq \frac{\widetilde{B}}{n(s, a)} + \frac{16\widetilde{B}\iota}{3n(s, a)} + 2\sqrt{\frac{\text{Var}(\widetilde{P}', \bar{V})\iota}{n(s, a)}} + 6\sqrt{\frac{S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty$$

Proof. First, we can bound term $(P_{s,a} - \widetilde{P}'_{s,a})\bar{V}$.

$$|(P_{s,a} - \widetilde{P}'_{s,a})\bar{V}| \leq |(\widehat{P}_{s,a} - \widetilde{P}'_{s,a})\bar{V}| + |(P_{s,a} - \widehat{P}_{s,a})(\bar{V} - V^*)| + |(P_{s,a} - \widehat{P}_{s,a})V^*| \quad (23)$$

Then we bound the above three terms one by one. From Lemma 7.1 and Lemma 7.2, we have

$$|(\widehat{P}_{s,a} - \widetilde{P}'_{s,a})\bar{V}| \leq \frac{\|\bar{V}\|_\infty}{n_{max} + 1} \leq \frac{\widetilde{B}}{n_{max} + 1}. \quad (24)$$

For the second term, we have

$$\begin{aligned} |(P_{s,a} - \widehat{P}_{s,a})(\bar{V} - V^*)| &\leq \sqrt{\frac{2S\text{Var}(P_{s,a}, \bar{V} - V^*)\iota}{n(s, a)}} + \frac{2\|\bar{V} - V^*\|_\infty S\iota}{3n(s, a)} \\ &\leq \sqrt{\frac{2S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty + \frac{\|\bar{V} - V^*\|_\infty S\iota}{n(s, a)}, \end{aligned} \quad (25)$$

where the first inequality holds because of lemma 10.3. For the last term, we have that

$$\begin{aligned} |(P_{s,a} - \widehat{P}_{s,a})V^*| &\stackrel{(i)}{\leq} \sqrt{\frac{2\text{Var}(\widehat{P}_{s,a}, V^*)\iota}{n(s, a)}} + \frac{7\|V^*\|_\infty \iota}{3n(s, a)} \\ &\stackrel{(ii)}{\leq} 2\sqrt{\frac{\text{Var}(\widehat{P}_{s,a}, V^* - \bar{V})\iota}{n(s, a)}} + 2\sqrt{\frac{\text{Var}(\widehat{P}_{s,a}, \bar{V})\iota}{n(s, a)}} + \frac{7\widetilde{B}\iota}{3n(s, a)} \\ &\leq 2\sqrt{\frac{\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty + 2\sqrt{\frac{\text{Var}(\widehat{P}_{s,a}, \bar{V})\iota}{n(s, a)}} + \frac{7\widetilde{B}\iota}{3n(s, a)} \\ &\stackrel{(iii)}{\leq} 2\sqrt{\frac{\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty + 2\sqrt{\frac{\text{Var}(\widetilde{P}'_{s,a}, \bar{V})\iota}{n(s, a)}} + \frac{2\sqrt{2\iota}\|\bar{V}\|_\infty}{n(s, a)} + \frac{7\widetilde{B}\iota}{3n(s, a)} \\ &\leq 2\sqrt{\frac{\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty + 2\sqrt{\frac{\text{Var}(\widetilde{P}'_{s,a}, \bar{V})\iota}{n(s, a)}} + \frac{16\widetilde{B}\iota}{3n(s, a)}, \end{aligned} \quad (26)$$

where (i) holds under event \mathcal{E}_3 . (ii) holds because of $\text{Var}(X + Y) \leq 2\text{Var}(X) + 2\text{Var}(Y)$. (iii) comes from Lemma 7.2. Both (ii) and (iii) uses the result that $\sqrt{a + b} \leq \sqrt{a} + \sqrt{b}$ when $a \geq 0$ and $b \geq 0$. Combine the above inequalities together, we can get

$$(P_{s,a} - \tilde{P}'_{s,a})\bar{V} \leq \frac{\tilde{B}}{n(s,a)} + \frac{16\tilde{B}\iota}{3n(s,a)} + 2\sqrt{\frac{\text{Var}(\tilde{P}', \bar{V})\iota}{n(s,a)}} + \left(\sqrt{\frac{2S\iota}{n(s,a)}} + \frac{S\iota}{n(s,a)} + 2\sqrt{\frac{\iota}{n(s,a)}}\right) \|\bar{V} - V^*\|_\infty.$$

Since with probability $1 - \delta$, by Lemma 12.2 $n(s,a) \geq \frac{1}{2}nd_m$. When $n \geq \frac{S\iota}{2d_m}$, we have

$$\begin{aligned} \left(\sqrt{\frac{2S\iota}{n(s,a)}} + \frac{S\iota}{n(s,a)} + 2\sqrt{\frac{\iota}{n(s,a)}}\right) \|\bar{V} - V^*\|_\infty &\leq (\sqrt{2} + 2 + 2)\sqrt{\frac{S\iota}{n(s,a)}} \|\bar{V} - V^*\|_\infty \\ &\leq 6\sqrt{\frac{S\iota}{n(s,a)}} \|\bar{V} - V^*\|_\infty \end{aligned}$$

□

Lemma 7.4. Define function $f : \mathbb{R}^{S'} \times \mathbb{R}^{S'} \times \mathbb{R} \rightarrow \mathbb{R}$ as $f(p, v, n) = pv + \max\{2\sqrt{\frac{\text{Var}(p,v)\iota}{n}}, 4\frac{\tilde{B}\iota}{n}\}$, if $\|v\|_\infty \leq \tilde{B}$ and $v(g) = 0$, then we have $(\frac{\partial f}{\partial v})(s) \geq 0$ and $\sum_{s,s \neq g} (\frac{\partial f}{\partial v})(s) \leq 1 - p(g)^2$.

Proof.

$$\begin{aligned} \left(\frac{\partial f}{\partial v}\right)(s) &= p(s) + \mathbb{I}\left\{2\sqrt{\frac{\text{Var}(p,v)\iota}{n}} \geq 4\frac{\tilde{B}\iota}{n}\right\} 2\sqrt{\frac{\iota}{n}} \frac{\partial(\sqrt{\text{Var}(p,v)})}{\partial v(s)} \\ &= p(s) + \mathbb{I}\left\{2\sqrt{\frac{\text{Var}(p,v)\iota}{n}} \geq 4\frac{\tilde{B}\iota}{n}\right\} 2\sqrt{\frac{\iota}{n}} \frac{p(s)(v(s) - pv)}{\sqrt{\text{Var}(p,v)}} \end{aligned}$$

Simplifying the above equation, we can get

$$\left(\frac{\partial f}{\partial v}\right)(s) \geq \min\{p(s), p(s) - \frac{p(s)(pv - v(s))}{\tilde{B}}\}$$

(27)

Since $|pv - v(s)| \leq \tilde{B}$, we have $p(s) - \frac{p(s)(pv - v(s))}{\tilde{B}} \geq 0$. Then we would have $(\frac{\partial f}{\partial v})(s) \geq 0$. For the second part, we have

1. Case I: $2\sqrt{\frac{\text{Var}(p,v)\iota}{n}} \geq 4\frac{\tilde{B}\iota}{n}$, we have

$$\begin{aligned} \sum_{s,s \neq g} \left(\frac{\partial f}{\partial v}\right)(s) &= \sum_{s,s \neq g} \left\{p(s) + 2\sqrt{\frac{\iota}{n\text{Var}(p,v)}} p(s)(v(s) - pv)\right\} \\ &\leq \sum_{s,s \neq g} p(s) + 2\sqrt{\frac{\iota}{n\text{Var}(p,v)}} \left\{ \sum_{s,s \neq g} p(s)v(s) - \sum_{s,s \neq g} p(s) \left(\sum_{s',s' \neq g} p(s')v(s') \right) \right\} \\ &\leq \sum_{s,s \neq g} p(s) + 2\sqrt{\frac{\iota}{n\text{Var}(p,v)}} \left[\sum_{s,s \neq g} p(s)v(s) \right] \left(1 - \sum_{s,s \neq g} p(s)\right) \\ &\leq \sum_{s,s \neq g} p(s) + \frac{[\sum_{s,s \neq g} p(s)v(s)](1 - \sum_{s,s \neq g} p(s))}{\tilde{B}} \\ &\leq \sum_{s,s \neq g} p(s) + \left(\sum_{s,s \neq g} p(s) \right) \left(1 - \sum_{s,s \neq g} p(s)\right) = 1 - p(g)^2 \end{aligned}$$

2. Case II: $2\sqrt{\frac{\text{Var}(p,v)\iota}{n}} \leq 4\frac{\tilde{B}\iota}{n}$, we have

$$\sum_{s,s \neq g} \left(\frac{\partial f}{\partial v}\right)(s) = \sum_{s,s \neq g} p(s) = 1 - p(g) \leq 1 - p(g)^2$$

Combine this inequality with (27), we complete the proof. \square

Lemma 7.5. *When $n \geq \max\{\frac{4B_* - 2c_{\min}}{c_{\min}d_{\max}}, \frac{26^2 \times 2S\iota(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{B_*}+1)^4 S\iota\bar{T}^*\bar{T}}{B_*(\sqrt{B_*}+1)^2 d_m}, O(T_{\max}^2 \log(SA/\delta)/d_m^2)\}$, $\bar{\pi}$ is a proper policy (Recall $\bar{\pi}$ is the output of Algorithm ??).*

Proof. By definition we need to show that $T^{\bar{\pi}}(s) < \infty$ for any $s \in \mathcal{S}$. We prove this by contradiction. Suppose $T^{\bar{\pi}}(s) = \infty$, then we have that there exists at least one state e such that the expected visiting times of state e is infinite, i.e., $\exists e \in \mathcal{S}$, such that $T_{e,e}^{\bar{\pi}} = \infty$. In this case, e is a (positive) recurrent state in the finite Markov Chain induced by policy $\bar{\pi}$. Denote the communication class which e belongs as \mathcal{S}_0 . Since the state space is finite, we have that every state in the communication class \mathcal{S}_0 is recurrent. From the finite Markov Chain theory, we know that the communication class \mathcal{S}_0 is closed. In other words, $\forall x \in \mathcal{S}_0$, and $\forall y \in \mathcal{S} \setminus \mathcal{S}_0$, we have $P(y|x, \bar{\pi}(x)) = 0$. Thus with probability 1, we have $\sum_{i=1}^n \sum_{j=1}^{T_i} \mathbb{I}(s_j^{(i)} = x, a_j^{(i)} = \bar{\pi}(x), s_{j+1}^{(i)} = y) = 0$. This implies that $\hat{P}(y|x, \bar{\pi}(x)) = \frac{\sum_{i=1}^n \sum_{j=1}^{T_i} \mathbb{I}(s_j^{(i)} = x, a_j^{(i)} = \bar{\pi}(x), s_{j+1}^{(i)} = y)}{n(x, \bar{\pi}(x))} = 0$. By definition of the estimated transition matrix \tilde{P}' , we have $\tilde{P}'(g|x, \bar{\pi}(x)) = \frac{1}{n_{\max} + 1}$. It follows that

$$\sum_{h=0}^{\infty} \sum_{s \in \mathcal{S}_0} \tilde{\xi}_{h,e}(s) = \sum_{h=0}^{\infty} \left(\frac{n_{\max}}{n_{\max} + 1}\right)^h = n_{\max} + 1. \quad (28)$$

Then we have

$$\sum_{h=0}^{\infty} \sum_{s \in \mathcal{S}_0} \tilde{\xi}_{h,e}(s) = n_{\max} + 1 \geq \frac{1}{2}nd_{\max} + 1, \quad (29)$$

where the last inequality holds with probability $1 - \delta$ by Lemma 12.6. Define $\tilde{V}^{\bar{\pi}}(e) = \sum_{h=0}^{\infty} \mathbb{E}_{\tilde{P}', \bar{\pi}}[\hat{c}(s_h, a_h) | s_0 = e]$, then

$$\begin{aligned} \tilde{V}^{\bar{\pi}}(e) &\geq \sum_{h=0}^{\infty} \mathbb{E}_{\tilde{P}', \bar{\pi}}[c_{\min}] \\ &\geq \sum_{h=0}^{\infty} \mathbb{E}_{\tilde{P}', \bar{\pi}}[c_{\min} \mathbb{I}(s_h \in \mathcal{S}_0)] = c_{\min} \sum_{h=0}^{\infty} \sum_{s \in \mathcal{S}_0} \tilde{\xi}_{h,e}(s) \geq c_{\min} \left(\frac{1}{2}nd_{\max} + 1\right). \end{aligned}$$

Apply Lemma 1.1 to the SSP problem with $M := \langle \mathcal{S}, \mathcal{A}, \tilde{P}', \hat{c}, e, g \rangle$, we can get

$$\tilde{V}^{\bar{\pi}}(s) = \tilde{P}'_{s, \bar{\pi}(s)} \tilde{V}^{\bar{\pi}} + \hat{c}(s, \bar{\pi}(s)) := \tilde{\mathcal{T}}' \tilde{V}^{\bar{\pi}}(s)$$

Since $\bar{V}(s) = \tilde{P}'_{s, \bar{\pi}(s)} \bar{V} + \hat{c}(s, \bar{\pi}(s)) + b_{s, \bar{\pi}(s)}(\bar{V}) = \tilde{\mathcal{T}} \bar{V}(s)$, from Lemma 2.1 we have $\bar{V}(s) \geq \tilde{V}^{\bar{\pi}}(s)$ (since $b(\bar{V})$ is non-negative and both $\tilde{\mathcal{T}}, \tilde{\mathcal{T}}'$ are monotone operators). Then we get

$$\bar{V}(e) \geq \tilde{V}^{\bar{\pi}}(e) \geq c_{\min} \left(\frac{1}{2}nd_{\max} + 1\right). \quad (30)$$

From Lemma 8.1 (note Lemma 8.1 only bounds \bar{V} and V^* and has nothing to do with $\bar{\pi}$), we have that with probability $1 - \delta$, when $n \geq \max\{\frac{26^2 \times 2S\iota(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{B_*}+1)^4 S\iota\bar{T}^*\bar{T}}{B_*(\sqrt{B_*}+1)^2 d_m}\}$, which implies $n \geq \frac{900\bar{T}^* \iota(\sqrt{B_*}+1)^2}{B_* d_m}$, we have $\bar{V}(e) \leq V^*(e) + B_* \leq 2B_*$. Combine this inequality with (30), we can get $n \leq \frac{4B_* - 2c_{\min}}{c_{\min}d_{\max}}$, which contradicts with the assumptions in the lemma. \square

7.2 CONVERGENCE OF PESSIMISTIC VALUE ITERATION IN ALGORITHM ??

Define the operator $\tilde{\mathcal{T}}$ as $\tilde{\mathcal{T}}(V)(s) = \min_a \left\{ \min\{\hat{c}(s, a) + \tilde{P}_{s,a}V + b_{s,a}(V), \tilde{B}\} \right\}$. First, we prove that $\tilde{\mathcal{T}}$ is a contraction mapping.

Lemma 7.6. $\tilde{\mathcal{T}} : \mathbb{R}^{\mathcal{S}} \times \{0\} \rightarrow \mathbb{R}^{\mathcal{S}} \times \{0\}$ is a contraction mapping, i.e., $\forall V_1, V_2 \in \mathbb{R}^{\mathcal{S}}, V_1(g) = V_2(g) = 0$, we have

$$\left\| \tilde{\mathcal{T}}V_1 - \tilde{\mathcal{T}}V_2 \right\|_{\infty} \leq \gamma \|V_1 - V_2\|_{\infty}, \quad (31)$$

where $\gamma := 1 - \frac{1}{(1 + \max_{s,a} n(s,a))^2}$.

Proof of Lemma 7.6. First, we prove the result for $s = g$. Since $b_{g,a}(V) = 0$, then we have $\tilde{\mathcal{T}}(V)(g) = 0$ and thus $\tilde{\mathcal{T}}(V_1)(g) - \tilde{\mathcal{T}}(V_2)(g) = 0$. When $s \neq g$, we have

$$\begin{aligned}
|\tilde{\mathcal{T}}V_1(s) - \tilde{\mathcal{T}}V_2(s)| &\leq \max_a |\min\{\hat{c}(s, a) + \tilde{P}'_{s,a} V_1 + b_{s,a}(V_1), \tilde{B}\} - \min\{\hat{c}(s, a) + \tilde{P}'_{s,a} V_2 + b_{s,a}(V_2), \tilde{B}\}| \\
&\stackrel{(i)}{\leq} \max_a |\{\hat{c}(s, a) + \tilde{P}'_{s,a} V_1 + b_{s,a}(V_1)\} - \{\hat{c}(s, a) + \tilde{P}'_{s,a} V_2 + b_{s,a}(V_2)\}| \\
&= \max_a |f(\tilde{P}'_{s,a}, V_1) - f(\tilde{P}'_{s,a}, V_2)| \\
&\stackrel{(ii)}{=} \max_a \sum_s |(\frac{\partial f}{\partial v})(\tilde{P}'_{s,a}, \theta(V_1 - V_2) + V_2), V_1 - V_2| \\
&\leq \max_a (\sum_{s', s' \neq g} |(\frac{\partial f}{\partial v}(s'))(\tilde{P}'_{s,a}, \theta(V_1 - V_2) + V_2)|) \|V_1 - V_2\|_\infty \\
&\stackrel{(iii)}{\leq} \max_a \{1 - \tilde{P}'_{s,a}(g)^2\} \|V_1 - V_2\|_\infty,
\end{aligned}$$

where (i) comes from Lemma 12.7. (ii) is due to the mean value theorem. (iii) uses the result in Lemma 7.4. Then we have

$$\left\| \tilde{\mathcal{T}}V_1(s) - \tilde{\mathcal{T}}V_2(s) \right\|_\infty \leq \max_{s,a} \{1 - \tilde{P}'_{s,a}(g)^2\} \|V_1 - V_2\|_\infty \leq \left\{ 1 - \frac{1}{(1 + \max_{s,a} n(s, a))^2} \right\} \|V_1 - V_2\|_\infty$$

□

We then introduce the following two regret decomposition lemma.

7.3 REGRET DECOMPOSITION LEMMA FOR POLICY OPTIMIZATION

Lemma 7.7. *Suppose \bar{V} is the limit of the sequence $V^{(i)}$ in Algorithm ??, we have the following decomposition lemma.*

$$\bar{V} - V^* \leq \sum_{h=0}^{\infty} \sum_{s \neq g} \xi_h^*(s) \{(\tilde{P}'_{s, \pi^*(s)} - P_{s, \pi^*(s)})\bar{V} + \hat{c}(s, \pi^*(s)) - c(s, \pi^*(s)) + b_{s, \pi^*(s)}(\bar{V})\} \quad (32)$$

Proof.

$$\bar{V} - V^* = \sum_{s, s \neq g} \xi_0^*(s) (\bar{V}(s) - V^*(s))$$

Since for any $h \in \mathbf{N}$, we have

$$\begin{aligned}
\sum_{s, s \neq g} \xi_h^*(s) (\bar{V}(s) - V^*(s)) &\stackrel{(i)}{\leq} \sum_{s, s \neq g} \xi_h^*(s) (\bar{Q}(s, \pi^*(s)) - Q^*(s, \pi^*(s))) \\
&\stackrel{(ii)}{=} \sum_{s, s \neq g} \xi_h^*(s) \{ \tilde{P}'_{s, \pi^*(s)} \bar{V} - P_{s, \pi^*(s)} V^* + \hat{c}(s, \pi^*(s)) - c(s, \pi^*(s)) + b_{s, \pi^*(s)}(\bar{V}) \} \\
&= \sum_{s, s \neq g} \xi_h^*(s) \{ (\tilde{P}'_{s, \pi^*(s)} - P_{s, \pi^*(s)}) \bar{V} + \hat{c}(s, \pi^*(s)) - c(s, \pi^*(s)) + b_{s, \pi^*(s)}(\bar{V}) \} \\
&\quad + \sum_{s, s \neq g} \xi_h^*(s) P_{s, \pi^*(s)} (\bar{V} - V^*) \\
&= \sum_{s, s \neq g} \xi_h^*(s) \{ (\tilde{P}'_{s, \pi^*(s)} - P_{s, \pi^*(s)}) \bar{V} + \hat{c}(s, \pi^*(s)) - c(s, \pi^*(s)) + b_{s, \pi^*(s)}(\bar{V}) \} \\
&\quad + \sum_{s, s \neq g} \xi_{h+1}^*(s) (\bar{V} - V^*)(s)
\end{aligned}$$

(i) follows from the fact that $\bar{V}(s) = \min_a \bar{Q}(s, a) \leq \bar{Q}(s, \pi^*(s))$. (ii) uses the fact that \bar{V} is the limit of $V^{(i)}$ and we have $\bar{Q}(s, a) = \hat{c}(s, a) + \tilde{P}'_{s,a} \bar{V} + b_{s,a}(\bar{V})$. By recursion over time step h , we can get

$$\begin{aligned} \bar{V} - V^* &\leq \sum_{h=0}^H \sum_{s, s \neq g} \xi_h^*(s) \{ (\tilde{P}'_{s, \pi^*(s)} - P_{s, \pi^*(s)}) \bar{V} + \hat{c}(s, \pi^*(s)) - c(s, \pi^*(s)) + b_{s, \pi^*(s)}(\bar{V}) \} \\ &\quad + \sum_{s, s \neq g} \xi_{H+1}^*(s) (\bar{V} - V^*), \end{aligned} \quad (33)$$

Since π^* is a proper policy, we have that for any $s \neq g$, $\lim_{H \rightarrow \infty} \xi_{H+1}^*(s) = 0$. Also, \bar{V} is bounded by \tilde{B} and V^* is bounded by B_* . Thus let H goes to infinity, we can complete the proof. \square

Lemma 7.8. When $n \geq \max\left\{\frac{4B_* - 2c_{\min}}{c_{\min} d_{max}}, \frac{26^2 \times 2S_t(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{\tilde{B}+1})^4 S_t \bar{T}^* \tilde{T}}{B_*(\sqrt{B_*+1})^2 d_m}, O(T_{\max}^2 \log(SA/\delta)/d_m^2)\right\}$, we have

$$V^{\bar{\pi}} - \bar{V} = \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\bar{\pi}}(s) \{ (P_{s, \bar{\pi}(s)} - \tilde{P}'_{s, \bar{\pi}(s)}) \bar{V} + c(s, \bar{\pi}(s)) - \hat{c}(s, \bar{\pi}(s)) - b_{s, \bar{\pi}(s)}(\bar{V}) \}. \quad (34)$$

Proof. First of all, by the condition we have $\bar{\pi}$ is a proper policy by Lemma 7.5. We prove by recursion formula.

$$\begin{aligned} \sum_{s, s \neq g} \xi_h^{\bar{\pi}}(s) (V^{\bar{\pi}} - \bar{V}) &= \sum_{s, s \neq g} \xi_h^{\bar{\pi}}(s) \{ P_{s, \bar{\pi}(s)} V^{\bar{\pi}} + c(s, \bar{\pi}(s)) - \tilde{P}'_{s, \bar{\pi}(s)} \bar{V} - \hat{c}(s, \bar{\pi}(s)) - b_{s, \bar{\pi}(s)}(\bar{V}) \} \\ &= \sum_{s, s \neq g} \xi_h^{\bar{\pi}}(s) \{ (P_{s, \bar{\pi}(s)} - \tilde{P}'_{s, \bar{\pi}(s)}) \bar{V} + P_{s, \bar{\pi}(s)} (V^{\bar{\pi}} - \bar{V}) + c(s, \bar{\pi}(s)) - \hat{c}(s, \bar{\pi}(s)) - b_{s, \bar{\pi}(s)}(\bar{V}) \} \\ &= \sum_{s, s \neq g} \xi_h^{\bar{\pi}}(s) \{ (P_{s, \bar{\pi}(s)} - \tilde{P}'_{s, \bar{\pi}(s)}) \bar{V} + c(s, \bar{\pi}(s)) - \hat{c}(s, \bar{\pi}(s)) - b_{s, \bar{\pi}(s)}(\bar{V}) \} \\ &\quad + \sum_{s, s \neq g} \xi_{h+1}^{\bar{\pi}}(s) (V^{\bar{\pi}} - \bar{V}), \end{aligned}$$

where the first inequality uses the Bellman equation for policy $\bar{\pi}$, which follows from Lemma 1.1. By recursion, we have

$$\begin{aligned} V^{\bar{\pi}} - \bar{V} &= \sum_{h=0}^H \sum_{s, s \neq g} \xi_h^{\bar{\pi}}(s) \{ (P_{s, \bar{\pi}(s)} - \tilde{P}'_{s, \bar{\pi}(s)}) \bar{V} + c(s, \bar{\pi}(s)) - \hat{c}(s, \bar{\pi}(s)) - b_{s, \bar{\pi}(s)}(\bar{V}) \} \\ &\quad + \sum_{s, s \neq g} \xi_H^{\bar{\pi}}(s) (V^{\bar{\pi}} - \bar{V}). \end{aligned}$$

From Lemma 7.5, we have that when $n \geq \max\left\{\frac{4B_* - 2c_{\min}}{c_{\min} d_{max}}, \frac{26^2 \times 2S_t(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{\tilde{B}+1})^4 S_t \bar{T}^* \tilde{T}}{B_*(\sqrt{B_*+1})^2 d_m}, O(T_{\max}^2 \log(SA/\delta)/d_m^2)\right\}$, $\bar{\pi}$ is a proper policy. Thus $\|V^{\bar{\pi}}\|_{\infty} < +\infty$, and for any state $s, s \neq g$, $\lim_{H \rightarrow \infty} \xi_H^{\bar{\pi}}(s) = 0$. Let H goes to infinity, we can prove the lemma. \square

8 CRUDE OPTIMIZATION BOUND

In this section, we give a rough bound for $\bar{V} - V^*$.

Theorem 8.1. Denote $d_m := \min\{\sum_{h=1}^{\infty} \xi_h^{\mu}(s, a) : s.t. \sum_{h=1}^{\infty} \xi_h^{\mu}(s, a) > 0\}$ and $T_{\max} = \max_i T_i$. Let T_s^* be the expected time to hit g when starting from s with the optimal policy and denote $\bar{T}^{\pi} = \max_s T_s^{\pi}$. Then when $n \geq \max\left\{\frac{26^2 \times 2S_t(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{\tilde{B}+1})^4 S_t \bar{T}^* \tilde{T}}{B_*(\sqrt{B_*+1})^2 d_m}, O(T_{\max}^2 \log(SA/\delta)/d_m^2)\right\}$, we have with probability $1 - \delta$,

$$\|\bar{V} - V^*\|_{\infty} \leq 30 \sqrt{\frac{\bar{T}^* B_* t}{n d_m}} (\sqrt{B_*} + 1) \quad (35)$$

Proof. From Lemma 7.7, we have (by choosing $\xi_0 = \mathbf{1}[s_0 = \bar{s}]$)

$$|\bar{V}(\bar{s}) - V^*(\bar{s})| \leq \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \{ (\tilde{P}_{s,\pi^*(s)} - P_{s,\pi^*(s)}) \bar{V} + \tilde{c}(s, \pi^*(s)) - c(s, \pi^*(s)) + b_{s,\pi^*(s)}(\bar{V}) \}$$

For the first term, we can bound it by Lemma 7.3

$$|(\tilde{P}_{s,a} - P_{s,a}) \bar{V}| \leq \frac{\tilde{B}}{n(s,a)} + \frac{16\tilde{B}\iota}{3n(s,a)} + 2\sqrt{\frac{\text{Var}(\tilde{P}', \bar{V})\iota}{n(s,a)}} + 6\sqrt{\frac{S\iota}{n(s,a)}} \|V - V^*\|_{\infty}.$$

Conditioned on the event \mathcal{E}_5 , we have

$$|\tilde{c}(s,a) - c(s,a)| \leq \sqrt{\frac{2\tilde{c}(s,a)\iota}{n(s,a)}} + \frac{7\iota}{3n(s,a)}$$

Combine the above inequalities together, we can get

$$\begin{aligned} |\bar{V}(\bar{s}) - V^*(\bar{s})| &\leq \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ \frac{2\tilde{B}}{n(s,\pi^*(s))} + \frac{32\tilde{B}\iota}{3n(s,\pi^*(s))} + 2\sqrt{\frac{\text{Var}(\tilde{P}', \bar{V})\iota}{n(s,\pi^*(s))}} + 6\sqrt{\frac{S\iota}{n(s,\pi^*(s))}} \|V - V^*\|_{\infty} \right. \\ &\quad \left. + 2\sqrt{\frac{2\tilde{c}(s,\pi^*(s))\iota}{n(s,\pi^*(s))}} + \frac{14\iota}{3n(s,\pi^*(s))} + \max\left\{ 2\sqrt{\frac{\text{Var}(\tilde{P}, V)\iota}{n(s,\pi^*(s))}}, 4\frac{\tilde{B}\iota}{n(s,\pi^*(s))} \right\} \right. \\ &\quad \left. + 180\sqrt{\frac{3\tilde{T}\tilde{B}S}{2n(s,\pi^*(s))n_{\min}}} (\sqrt{\tilde{B}} + 1)\iota \right\} \\ &\stackrel{(i)}{\leq} \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ \frac{2\tilde{B}}{n(s,\pi^*(s))} + \frac{44\tilde{B}\iota}{3n(s,\pi^*(s))} + 4\sqrt{\frac{\text{Var}(\tilde{P}', \bar{V})\iota}{n(s,\pi^*(s))}} + 6\sqrt{\frac{S\iota}{n(s,\pi^*(s))}} \|\bar{V} - V^*\|_{\infty} \right. \\ &\quad \left. + 2\sqrt{\frac{2\tilde{c}(s,\pi^*(s))\iota}{n(s,\pi^*(s))}} + \frac{14\iota}{3n(s,\pi^*(s))} + 180\sqrt{\frac{3\tilde{T}\tilde{B}S}{2n(s,\pi^*(s))n_{\min}}} (\sqrt{\tilde{B}} + 1)\iota \right\}, \end{aligned}$$

where (i) uses the inequality that $\max\{a, b\} \leq a + b$. For notation simplicity, we define

$$b_0(s,a) := 180\sqrt{\frac{3\tilde{T}\tilde{B}S}{2n(s,a)n_{\min}}} (\sqrt{\tilde{B}} + 1)\iota.$$

First, we bound the variance term

$$\begin{aligned} \text{Var}(\tilde{P}'_{s,a}, \bar{V}) &\stackrel{(i)}{\leq} \text{Var}(\tilde{P}_{s,a}, \bar{V}) + \frac{2\|\bar{V}\|_{\infty}^2}{n_{\max} + 1} \\ &\stackrel{(ii)}{\leq} \frac{3}{2}\text{Var}(P_{s,a}, \bar{V}) + \frac{2\|\bar{V}\|_{\infty}^2 S\iota}{n(s,a)} + \frac{2\|\bar{V}\|_{\infty}^2}{n_{\max} + 1} \\ &\stackrel{(iii)}{\leq} 3\text{Var}(P_{s,a}, \bar{V} - V^*) + 3\text{Var}(P_{s,a}, V^*) + \frac{2\|\bar{V}\|_{\infty}^2 S(\iota + 1)}{n(s,a)} \\ &\stackrel{(iv)}{\leq} 3\|\bar{V} - V^*\|_{\infty}^2 + 3\text{Var}(P_{s,a}, V^*) + \frac{2\tilde{B}^2 S(\iota + 1)}{n(s,a)}, \end{aligned}$$

where (i) follows from Lemma 7.2. (ii) uses Lemma 10.1. (iii) uses the fact that $\text{Var}(X + Y) \leq 2\text{Var}(X) + 2\text{Var}(Y)$. (iv) uses the fact that $\|\bar{V}\|_{\infty} \leq \tilde{B}$. Then we can have

$$|\bar{V}(\bar{s}) - V^*(\bar{s})| \leq \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ \frac{2\tilde{B}}{n(s,\pi^*(s))} + \frac{44\tilde{B}\iota}{3n(s,\pi^*(s))} + 4\sqrt{\frac{3\iota}{n(s,\pi^*(s))}} \|\bar{V} - V^*\|_{\infty} \right\}$$

$$\begin{aligned}
& + 4\sqrt{\frac{3\text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{n(s, \pi^*(s))}} + 4\sqrt{\frac{2S\iota(\iota+1)}{n^2(s, \pi^*(s))}} \tilde{B} \\
& + 6\sqrt{\frac{S\iota}{n(s, \pi^*(s))}} \|\bar{V} - V^*\|_\infty + 2\sqrt{\frac{2\tilde{c}(s, \pi^*(s))\iota}{n(s, \pi^*(s))}} + \frac{14\iota}{3n(s, \pi^*(s))} + b_0(s, a) \} \\
\leq & \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ \frac{27 \max\{\tilde{B}, 1\} \sqrt{S}\iota}{n(s, \pi^*(s))} + 13\sqrt{\frac{S\iota}{n(s, \pi^*(s))}} \|\bar{V} - V^*\|_\infty \right. \\
& \left. + 4\sqrt{\frac{3\text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{n(s, \pi^*(s))}} + 2\sqrt{\frac{2\tilde{c}(s, \pi^*(s))\iota}{n(s, \pi^*(s))}} + b_0(s, a) \right\} \\
\leq & \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ \frac{31 \max\{\tilde{B}, 1\} \sqrt{S}\iota}{n(s, \pi^*(s))} + 13\sqrt{\frac{S\iota}{n(s, \pi^*(s))}} \|\bar{V} - V^*\|_\infty \right. \\
& \left. + 4\sqrt{\frac{3\text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{n(s, \pi^*(s))}} + 2\sqrt{\frac{3c(s, \pi^*(s))\iota}{n(s, \pi^*(s))}} + b_0(s, a) \right\},
\end{aligned}$$

where (i) uses the assumption $\iota \geq 1$ and that $S \geq 1$. (ii) holds because of Lemma 10.2.

Then we have

$$\begin{aligned}
|\bar{V}(\bar{s}) - V^*(\bar{s})| & \stackrel{(i)}{\leq} \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ \frac{62 \max\{\tilde{B}, 1\} \sqrt{S}\iota}{nd_m} + 13\sqrt{2} \sqrt{\frac{S\iota}{nd_m}} \|\bar{V} - V^*\|_\infty \right. \\
& \left. + 4\sqrt{\frac{6\text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{nd_m}} + 2\sqrt{\frac{6c(s, \pi^*(s))\iota}{nd_m}} + \bar{b}_0 \right\} \\
& \stackrel{(ii)}{\leq} \frac{62T_{\bar{s}}^* \max\{\tilde{B}, 1\} \sqrt{S}\iota}{nd_m} + 13\sqrt{\frac{2S\iota}{nd_m}} \|\bar{V} - V^*\|_\infty T_{\bar{s}}^* + T_{\bar{s}}^* \bar{b}_0 \\
& \quad + \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ 4\sqrt{\frac{6\text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{nd_m}} + 2\sqrt{\frac{6c(s, \pi^*(s))\iota}{nd_m}} \right\},
\end{aligned}$$

where $\bar{b}_0 := 180\sqrt{\frac{6\tilde{T}\tilde{B}S}{n^2d_m^2}}(\sqrt{\tilde{B}} + 1)\iota$. (i) holds with probability $1 - \delta$ because of Lemma 12.6. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have $n(s, a) \geq \frac{1}{2}nd(s, a) \geq \frac{1}{2}nd_m$. In particular, $n_{\min} \geq \frac{1}{2}nd_m$. Since

$$\begin{aligned}
\sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ 4\sqrt{\frac{6\text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{nd_m}} \right\} & \stackrel{(i)}{\leq} 4\sqrt{\frac{\sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \text{Var}(P_{s,\pi^*(s)}, V^*)\iota}{nd_m}} \\
& \stackrel{(ii)}{\leq} 4\sqrt{T_{\bar{s}}^*} \sqrt{\frac{12\iota}{nd_m}} \|V^*\|_\infty,
\end{aligned}$$

where (i) uses the Cauchy-Schwartz inequality. (ii) uses the result in Lemma 4.3 and Lemma 10.1. Similarly, we have

$$\begin{aligned}
\sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) \left\{ 2\sqrt{\frac{6c(s, \pi^*(s))\iota}{nd_m}} \right\} & \leq 2\sqrt{\frac{\sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_{h,\bar{s}}^*(s) c(s, \pi^*(s))\iota}{nd_m}} \\
& \leq 2\sqrt{T_{\bar{s}}^*} \sqrt{\frac{6\|V^*\|_\infty \iota}{nd_m}}.
\end{aligned}$$

Combine the above together, we get

$$|\bar{V}(\bar{s}) - V^*(\bar{s})| \leq \frac{62T_{\bar{s}}^* \max\{\tilde{B}, 1\} \sqrt{S}\iota}{nd_m} + 13\sqrt{\frac{2S\iota}{nd_m}} \|\bar{V} - V^*\|_\infty T_{\bar{s}}^* + T_{\bar{s}}^* \bar{b}_0$$

$$\begin{aligned}
& + 4\sqrt{\bar{T}_s^*} \sqrt{\frac{12\iota}{nd_m}} \|V^*\|_\infty + 2\sqrt{\bar{T}_s^*} \sqrt{\frac{6\|V^*\|_\infty \iota}{nd_m}} \\
& \leq \frac{62\bar{T}^* \max\{\tilde{B}, 1\} \sqrt{S\iota}}{nd_m} + 13\sqrt{\frac{2S\iota}{nd_m}} \|\bar{V} - V^*\|_\infty \bar{T}^* + \bar{T}^* \bar{b}_0 \\
& + 8\sqrt{\bar{T}^*} \sqrt{\frac{3\iota}{nd_m}} \|V^*\|_\infty + 2\sqrt{\bar{T}^*} \sqrt{\frac{6\|V^*\|_\infty \iota}{nd_m}},
\end{aligned}$$

It implies that

$$\begin{aligned}
(1 - 13\sqrt{\frac{2S\iota}{nd_m}} \bar{T}^*) \|\bar{V} - V^*\|_\infty & \leq \frac{62 \max\{\tilde{B}, 1\} \sqrt{S\iota} \bar{T}^*}{nd_m} + \bar{T}^* \bar{b}_0 \\
& + 14\sqrt{\frac{\bar{T}^* B_* \iota}{nd_m}} (\sqrt{B_*} + 1). \tag{36}
\end{aligned}$$

Since $n \geq \frac{26^2 \times 2S\iota(\bar{T}^*)^2}{d_m}$,

$$\begin{aligned}
\|\bar{V}(s) - V^*(s)\|_\infty & \leq \frac{124 \max\{\tilde{B}, 1\} \sqrt{S\iota} \bar{T}^*}{nd_m} + 2\bar{T}^* \bar{b}_0 \\
& + 21\sqrt{\frac{\bar{T}^* B_* \iota}{nd_m}} (\sqrt{B_*} + 1) \\
& \leq \frac{124 \max\{\tilde{B}, 1\} \sqrt{S\iota} \bar{T}^*}{nd_m} + 360\bar{T}^* \sqrt{\frac{6\tilde{T}\tilde{B}S}{n^2 d_m^2}} (\sqrt{\tilde{B}} + 1)\iota \\
& + 28\sqrt{\frac{\bar{T}^* B_* \iota}{nd_m}} (\sqrt{B_*} + 1) \\
& \leq 720\bar{T}^* \sqrt{\frac{6\tilde{T}S}{n^2 d_m^2}} (\sqrt{\tilde{B}} + 1)^2 \iota + 28\sqrt{\frac{\bar{T}^* B_* \iota}{nd_m}} (\sqrt{B_*} + 1)
\end{aligned}$$

When $n \geq \frac{10^6(\sqrt{\tilde{B}}+1)^4 S\iota \bar{T}^* \tilde{T}}{B_*(\sqrt{B_*}+1)^2 d_m}$, we have

$$\|\bar{V}(s) - V^*(s)\|_\infty \leq 30\sqrt{\frac{\bar{T}^* B_* \iota}{nd_m}} (\sqrt{B_*} + 1)$$

□

9 PROOF OF THEOREM ??

In this section, we provide the proof of Theorem ?. However, before that, we first present a lemma that guarantees pessimism.

Lemma 9.1. *When $n \geq \max\{\frac{26^2 \times 2S\iota(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{\tilde{B}}+1)^4 S\iota \bar{T}^* \tilde{T}}{B_*(\sqrt{B_*}+1)^2 d_m}, O(T_{\max}^2 \log(SA/\delta)/d_m^2)\}$ (where $T_{\max} = \max_i T_i$), with probability at least $1 - \delta$, we have that for any state action pair (s, a) ,*

$$c(s, a) - \hat{c}(s, a) + (P_{s,a} - \tilde{P}'_{s,a})\bar{V} - b_{s,a}(\bar{V}) \leq 0$$

Proof. Applying the result in Theorem 8.1, we can get

$$6\sqrt{\frac{S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty \leq 180\sqrt{\frac{\bar{T}^* B_* S}{n(s, a)nd_m}} (\sqrt{B_*} + 1)\iota.$$

Combine the above inequality with Lemma 7.3 implies that

$$(P_{s,a} - \tilde{P}'_{s,a})\bar{V} \leq \frac{\tilde{B}}{n(s,a)} + \frac{16\tilde{B}\iota}{3n(s,a)} + 2\sqrt{\frac{\text{Var}(\tilde{P}', \bar{V})\iota}{n(s,a)}} + 180\sqrt{\frac{\bar{T}^*B_*S}{n(s,a)nd_m}}(\sqrt{B_*} + 1)\iota.$$

Conditioned on the event \mathcal{E}_5 , then we have

$$c(s,a) - \hat{c}(s,a) + (P_{s,a} - \tilde{P}'_{s,a})\bar{V} - b_{s,a}(\bar{V}) \leq 180\sqrt{\frac{\bar{T}^*B_*S}{n(s,a)nd_m}}(\sqrt{B_*} + 1)\iota - 180\sqrt{\frac{3\tilde{T}\tilde{B}S}{2n(s,a)n_{\min}}}(\sqrt{\tilde{B}} + 1)\iota.$$

Applying the Chernoff bound given in Lemma 12.6, we have that with probability $1 - \delta$, $n(s,a) < \frac{3}{2}nd^\mu(s,a)$ for any state action pair (s,a) . Thus $n_{\min} := \min_{s,a,n(s,a)>0} n(s,a) < \frac{3}{2}n(\min_{n(s,a)>0} d^\mu(s,a))$. For any $(s,a) \in \mathcal{S} \times \mathcal{A}$, if we have $d^\mu(s,a) > 0$, by the Lemma 12.6 we have that with probability $1 - \delta$, $n(s,a) > \frac{nd^\mu(s,a)}{2} > 0$, which implies

$$\{(s,a) \in \mathcal{S} \times \mathcal{A} : d^\mu(s,a) > 0\} \subseteq \{(s,a) \in \mathcal{S} \times \mathcal{A} : n(s,a) > 0\}.$$

Then we can get $\min_{n(s,a)>0} d^\mu(s,a) \leq \min_{d^\mu(s,a)>0} d^\mu(s,a) = d_m$ and thus $n_{\min} \leq \frac{3}{2}nd_m$. Because $\bar{T}^* \leq \tilde{T}$ and $B_* \leq \tilde{B}$, we can prove the result in the Lemma. \square

Now we are ready to introduce the final proof.

Theorem 9.2. *Given Assumption ?? and Assumption ?. When $n \geq n_0$, the suboptimality bound of the output policy $\bar{\pi}$ can be upper bounded as follows with probability $1 - \delta$ (where $\iota = O(\log(SA/\delta))$),*

$$V^{\bar{\pi}}(s_{\text{init}}) - V^*(s_{\text{init}}) \leq 8 \sum_{s,a,s \neq g} d^*(s,a) \sqrt{\frac{3\text{Var}_{P_{s,a}}[V^* + c]\iota}{n \cdot d^\mu(s,a)}} + \tilde{O}\left(\frac{1}{n}\right), \quad (37)$$

where we define $n_0 := n \geq \max\left\{\frac{4B_* - 2c_{\min}}{c_{\min}d_{max}}, \frac{26^2 \times 2S\iota(\bar{T}^*)^2(\sqrt{B_*} + 1)^2}{d_m}, \frac{10^6(\sqrt{\tilde{B}} + 1)^4 S\iota\bar{T}^*\tilde{T}}{B_*(\sqrt{B_*} + 1)^2 d_m}, O(T_{\max}^2 \log(SA/\delta)/d_m^2)\right\}$.

Proof.

$$V^{\bar{\pi}}(s) - V^*(s) = (V^{\bar{\pi}}(s) - \bar{V}(s)) + (\bar{V}(s) - V^*(s)). \quad (38)$$

From Lemma 7.8, we have with probability $1 - \delta$ and when $n \geq \max\left\{\frac{4B_* - 2c_{\min}}{c_{\min}d_{max}}, \frac{26^2 \times 2S\iota(\bar{T}^*)^2}{d_m}, \frac{10^6(\sqrt{\tilde{B}} + 1)^4 S\iota\bar{T}^*\tilde{T}}{B_*(\sqrt{B_*} + 1)^2 d_m}\right\}$,

$$V^{\bar{\pi}} - \bar{V} = \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^{\bar{\pi}}(s) \{(P_{s,\bar{\pi}(s)} - \tilde{P}'_{s,\bar{\pi}(s)})\bar{V} + c(s,\bar{\pi}(s)) - \hat{c}(s,\bar{\pi}(s)) - b_{s,\bar{\pi}(s)}(\bar{V})\}. \quad (39)$$

Next by Lemma 9.1 with probability $1 - \delta$, we have

$$(P_{s,\bar{\pi}(s)} - \tilde{P}'_{s,\bar{\pi}(s)})\bar{V} + c(s,\bar{\pi}(s)) - \hat{c}(s,\bar{\pi}(s)) - b_{s,\bar{\pi}(s)}(\bar{V}) \leq 0. \quad (40)$$

Thus combine (39) and (40) we have $V^{\bar{\pi}} - \bar{V} \leq 0$ by pessimism. For the term $\bar{V} - V^*$, we apply Lemma 7.7 to obtain:

$$\bar{V} - V^* \leq \sum_{h=0}^{\infty} \sum_{\substack{s \\ s \neq g}} \xi_h^*(s) \{(\tilde{P}'_{s,\pi^*(s)} - P_{s,\pi^*(s)})\bar{V} + \hat{c}(s,\pi^*(s)) - c(s,\pi^*(s)) + b_{s,\pi^*(s)}(\bar{V})\}. \quad (41)$$

From Lemma 7.3, we have

$$|(P_{s,a} - \tilde{P}'_{s,a})\bar{V}| \leq \frac{\tilde{B}}{n(s,a)} + \frac{16\tilde{B}\iota}{3n(s,a)} + 2\sqrt{\frac{\text{Var}(\tilde{P}'_{s,a}, \bar{V})\iota}{n(s,a)}} + 6\sqrt{\frac{S\iota}{n(s,a)}} \|\bar{V} - V^*\|_{\infty} \quad (42)$$

Conditioned on the event \mathcal{E}_5 , we have

$$|\widehat{c}(s, a) - c(s, a)| \leq \sqrt{\frac{2\widehat{c}(s, a)\iota}{n(s, a)}} + \frac{7\iota}{3n(s, a)}$$

Combine the above inequalities together, we can get

$$\begin{aligned} & (\widetilde{P}'_{s,a} - P_{s,a})\bar{V} + \widehat{c}(s, a) - c(s, a) + b_{s,a}(\bar{V}) \\ & \leq 2\sqrt{\frac{2\widehat{c}(s, a)\iota}{n(s, a)}} + \frac{14\iota}{3n(s, a)} + \frac{2\widetilde{B}}{n(s, a)} + \frac{32\widetilde{B}\iota}{3n(s, a)} + 4\sqrt{\frac{\text{Var}(\widetilde{P}'_{s,a}, \bar{V})\iota}{n(s, a)}} \\ & \quad + \frac{4\widetilde{B}\iota}{n(s, a)} + 6\sqrt{\frac{S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty + 180\sqrt{\frac{3\widetilde{T}\widetilde{B}S}{2n(s, a)n_{\min}}}(\sqrt{\widetilde{B}} + 1)\iota \\ & \leq 2\sqrt{\frac{2\widehat{c}(s, a)\iota}{n(s, a)}} + 4\sqrt{\frac{\text{Var}(\widetilde{P}'_{s,a}, \bar{V})\iota}{n(s, a)}} + \widetilde{O}\left(\frac{(\widetilde{B} + 1)\iota}{n(s, a)}\right) + 6\sqrt{\frac{S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty + 180\sqrt{\frac{3\widetilde{T}\widetilde{B}S}{2n(s, a)n_{\min}}}(\sqrt{\widetilde{B}} + 1)\iota \\ & \leq 2\sqrt{\frac{3c(s, a)\iota}{n(s, a)}} + 4\sqrt{\frac{3\text{Var}(P_{s,a}, V^*)\iota}{n(s, a)}} + \widetilde{O}\left(\frac{(\widetilde{B} + 1)\sqrt{S}\iota}{n(s, a)}\right) \\ & \quad + \widetilde{O}\left(\sqrt{\frac{S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty\right) + 180\sqrt{\frac{3\widetilde{T}\widetilde{B}S}{2n(s, a)n_{\min}}}(\sqrt{\widetilde{B}} + 1)\iota \end{aligned}$$

Plug the above into (41), then we bound all the terms one by one. First,

$$\sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s, a) \left(2\sqrt{\frac{3c(s, a)\iota}{n(s, a)}} \right) \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s, a) \left[2\sqrt{\frac{6c(s, a)\iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)}} \right] = 2 \sum_{\substack{s,a \\ s \neq g}} d^*(s, a) \sqrt{\frac{6c(s, a)\iota}{nd^\mu(s, a)}}. \quad (43)$$

For the second term, first we have

$$\begin{aligned} \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s, a) \left(4\sqrt{\frac{3\text{Var}(P_{s,a}, V^*)\iota}{n(s, a)}} \right) & \leq \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s, a) \left[4\sqrt{\frac{6\text{Var}(P_{s,a}, V^*)\iota}{n \sum_{h=1}^{\infty} \xi_h^\mu(s, a)}} \right] \\ & = 4 \sum_{\substack{s,a \\ s \neq g}} d^*(s, a) \sqrt{\frac{6\text{Var}(P_{s,a}, V^*)\iota}{nd^\mu(s, a)}} \end{aligned} \quad (44)$$

From Chernoff bound given in Lemma 12.6, we have with probability $1 - \delta$, we have

$$\widetilde{O}\left(\sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s, a) \frac{(\widetilde{B} + 1)\sqrt{S}\iota}{n(s, a)}\right) \leq \widetilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} \frac{d^\pi(s, a)}{d^\mu(s, a)} \cdot \frac{(\widetilde{B} + 1)\sqrt{S}\iota}{n}\right). \quad (45)$$

Similarly, we have

$$\begin{aligned} \widetilde{O}\left(\sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s, a) \sqrt{\frac{S\iota}{n(s, a)}} \|\bar{V} - V^*\|_\infty\right) & \leq \widetilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} d^*(s, a) \sqrt{\frac{S\iota}{nd^\mu(s, a)}} \|\bar{V} - V^*\|_\infty\right) \\ & \stackrel{(i)}{\leq} \widetilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} d^*(s, a) \sqrt{\frac{\bar{T}^* B_* S}{n^2 d^\mu(s, a) d_m}} (\sqrt{B_*} + 1)\iota\right), \end{aligned} \quad (46)$$

where (i) uses the Crude optimization bound given in Theorem 8.1. For the last term, we have

$$\tilde{O}\left(\sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s,a) \sqrt{\frac{\tilde{T}\tilde{B}S}{n(s,a)n_{\min}}} (\sqrt{\tilde{B}}+1)\iota\right) \leq \tilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{\tilde{T}\tilde{B}S}{n^2 d^\mu(s,a)d_m}} (\sqrt{\tilde{B}}+1)\iota\right), \quad (47)$$

where the inequality comes from Lemma 12.6 again. Combine the inequalities (43), (44), (45), (46) and (45) together, we have

$$\begin{aligned} \bar{V}(s_{\text{init}}) - V^*(s_{\text{init}}) &\stackrel{\text{(i)}}{\leq} 2 \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{6c(s,a)\iota}{nd^\mu(s,a)}} + 4 \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{6\text{Var}(P_{s,a}, V^*)\iota}{nd^\mu(s,a)}} \\ &\quad + \tilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{\tilde{T}\tilde{B}S}{n^2 d^\mu(s,a)d_m}} (\sqrt{\tilde{B}}+1)\iota\right) + \tilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} \frac{d^\pi(s,a)}{d^\mu(s,a)} \cdot \frac{(\tilde{B}+1)\sqrt{S}\iota}{n}\right) \\ &\stackrel{\text{(ii)}}{\leq} 2 \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{6c(s,a)\iota}{nd^\mu(s,a)}} + 4 \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{6\text{Var}(P_{s,a}, V^*)\iota}{nd^\mu(s,a)}} \\ &\quad + \tilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{\tilde{T}S}{n^2 d^\mu(s,a)d_m}} (\tilde{B}+1)\iota\right) \\ &\leq 8 \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{3\text{Var}(P_{s,a}, V^* + c)\iota}{nd^\mu(s,a)}} + \tilde{O}\left(\sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \sqrt{\frac{\tilde{T}S}{n^2 d^\mu(s,a)d_m}} (\tilde{B}+1)\iota\right), \quad (48) \end{aligned}$$

where the inequality (i) uses the assumption that $\tilde{T}^* \leq \tilde{T}$ and $B_* \leq \tilde{B}$. (ii) uses the fact that $\frac{d^\pi(s,a)}{d^\mu(s,a)} \leq \frac{d^\pi(s,a)}{\sqrt{d^\mu(s,a)d_m}}$ and that $\tilde{B} + \sqrt{\tilde{B}} \leq 2(\tilde{B}+1)$. The last inequality comes from $\sqrt{a} + \sqrt{b} \leq \sqrt{2a+2b}$. \square

Based on Theorem 8.1, we can also get the Proposition below.

Proposition 9.3. *When $n \geq n_0$ (where n_0 is defined the same as Theorem 9.2), then the suboptimality incurred by the limit of the output policy $\bar{\pi}$ can be upper bounded as (with probability $1 - \delta$)*

$$V^{\bar{\pi}}(s_{\text{init}}) - V^*(s_{\text{init}}) \leq 8 \left(\sum_{\substack{s,a \\ s \neq g}} \frac{d^*(s,a)}{d^\mu(s,a)} \cdot \frac{6\iota}{n} \right) \cdot (B_* + 1) + \tilde{O}\left(\frac{1}{n}\right). \quad (49)$$

Proof. By Theorem ??,

$$\begin{aligned} V^{\bar{\pi}}(s_{\text{init}}) - V^*(s_{\text{init}}) &\leq 8 \sum_{s,a,s \neq g} d^*(s,a) \sqrt{\frac{3\text{Var}_{P_{s,a}}[V^* + c]\iota}{n \cdot d^\mu(s,a)}} + \tilde{O}\left(\frac{1}{n}\right) \\ &\stackrel{\text{(i)}}{\leq} 8 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d^*(s,a)}{n \cdot d^\mu(s,a)}} \cdot \sqrt{3 \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \text{Var}_{P_{s,a}}[V^* + c]\iota}, \quad (50) \end{aligned}$$

where (i) uses the Cauchy-Schwartz inequality. Since

$$\begin{aligned} \sum_{\substack{s,a \\ s \neq g}} d^*(s,a) \text{Var}_{P_{s,a}}[V^* + c] &= \sum_{h=1}^{\infty} \sum_{\substack{s,a \\ s \neq g}} \xi_h^*(s,a) (\text{Var}(P_{s,a}, V^*) + c(s,a)) \\ &\stackrel{\text{(i)}}{\leq} 2 \|V^*\|_\infty^2 + V^*(s_0) \end{aligned}$$

$$\leq 2B_\star^2 + B_\star, \quad (51)$$

where (i) comes from Lemma 10.1 and the definition of value function. Plug (51) into (50), we obtain

$$\begin{aligned} V^{\bar{\pi}}(s_{\text{init}}) - V^\star(s_{\text{init}}) &\leq 8 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{d^\star(s,a)}{n \cdot d^\mu(s,a)}} \cdot \sqrt{6(B_\star^2 + B_\star)\iota} + \tilde{O}\left(\frac{1}{n}\right) \\ &\leq 8 \sqrt{\sum_{\substack{s,a \\ s \neq g}} \frac{6d^\star(s,a) \cdot \iota}{d^\mu(s,a) \cdot n}} \cdot (B_\star + 1) + \tilde{O}\left(\frac{1}{n}\right), \end{aligned} \quad (52)$$

which completes the proof. \square

10 PROPERTIES OF TRANSITION MATRIX ESTIMATE \hat{P}

Lemma 10.1. *For any $V(\cdot) \in \mathbb{R}^S$ satisfying $V(g) = 0$, i.e. (13), and suppose event \mathcal{E}_1 holds. we have*

$$\begin{aligned} \text{Var}(\hat{P}_{s,a}, V) &\leq \frac{3}{2} \text{Var}(P_{s,a}, V) + \frac{2\|V\|_\infty^2 S \iota}{n(s,a)} \\ \text{Var}(P_{s,a}, V) &\leq 2\text{Var}(\hat{P}_{s,a}, V) + \frac{4\|V\|_\infty^2 S \iota}{n(s,a)} \end{aligned} \quad (53)$$

Proof. From the event \mathcal{E}_1 , we have

$$|\hat{P}(s'|s,a) - P(s'|s,a)| \leq \sqrt{\frac{2P(s'|s,a)\iota}{n(s,a)}} + \frac{2\iota}{3n(s,a)} \leq \frac{P(s'|s,a)}{2} + \frac{5\iota}{3n(s,a)},$$

where the second inequality uses $\sqrt{ab} \leq \frac{a+b}{2}$ with $a = \frac{2\iota}{n(s,a)}$, $b = P(s'|s,a)$. Thus we have

$$\begin{aligned} \hat{P}(s'|s,a) &\leq \frac{3P(s'|s,a)}{2} + \frac{5\iota}{3n(s,a)} \leq \frac{3P(s'|s,a)}{2} + \frac{2\iota}{n(s,a)} \\ P(s'|s,a) &\leq 2\hat{P}(s'|s,a) + \frac{4\iota}{3n(s,a)} \leq 2\hat{P}(s'|s,a) + \frac{4\iota}{n(s,a)}. \end{aligned} \quad (54)$$

For the first inequality, we have

$$\begin{aligned} \text{Var}(\hat{P}_{s,a}, V) &= \hat{P}_{s,a}(V - \hat{P}_{s,a}V)^2 \leq \hat{P}_{s,a}(V - P_{s,a}V)^2 \\ &\leq \sum_{s'} \left(\frac{3P(s'|s,a)}{2} + \frac{2\iota}{n(s,a)} \right) (V(s') - P_{s,a}V)^2 \\ &\leq \frac{3}{2} \text{Var}(P_{s,a}, V) + \frac{2\|V\|_\infty^2 S \iota}{n(s,a)}, \end{aligned}$$

here the first inequality is due to $\hat{P}_{s,a}V := \text{argmin}_z \sum_{s'} \hat{P}_{s,a}(s')(V(s') - z)^2$, and the last term has $S + 1$ due to the extra state g . For the second part, we have

$$\begin{aligned} \text{Var}(P_{s,a}, V) &= P_{s,a}(V - P_{s,a}V)^2 \leq P_{s,a}(V - \hat{P}_{s,a}V)^2 \\ &\leq \sum_{s'} \left(2\hat{P}(s'|s,a) + \frac{4\iota}{n(s,a)} \right) (V(s') - \hat{P}_{s,a}V)^2 \\ &\leq 2\text{Var}(\hat{P}_{s,a}, V) + \frac{4\|V\|_\infty^2 (S+1)\iota}{n(s,a)}, \end{aligned}$$

\square

Lemma 10.2. *With probability at least $1 - \delta$, we have*

$$\begin{aligned} c(s, a) &\leq 2\widehat{c}(s, a) + \frac{10\iota}{3n(s, a)} \\ \widehat{c}(s, a) &\leq \frac{3}{2}c(s, a) + \frac{5\iota}{3n(s, a)} \end{aligned}$$

Proof. Conditioned on event \mathcal{E}_4 , we have

$$\begin{aligned} |c(s, a) - \widehat{c}(s, a)| &\leq \sqrt{\frac{2c(s, a)\iota}{n(s, a)}} + \frac{2\iota}{3n(s, a)} \\ &\leq \frac{\iota}{n(s, a)} + \frac{1}{2}c(s, a) + \frac{2\iota}{3n(s, a)} \\ &\leq \frac{5\iota}{3n(s, a)} + \frac{1}{2}c(s, a), \end{aligned}$$

where the first inequality uses the assumption that $c(s, a) \in [0, 1]$. The second inequality follows from the result that $\sqrt{ab} \leq \frac{a+b}{2}$. Simplify the above inequality, we can conclude the proof. \square

Lemma 10.3. *With probability $1 - \delta$, for all $V(\cdot) \in \mathbb{R}^{S'}$ such that $\|V\|_\infty < \infty$, we have for all state-action pair (s, a)*

$$(\widehat{P}_{s,a} - P_{s,a})V \leq \sqrt{\frac{2S\text{Var}(P_{s,a}, V)\iota}{n(s, a)}} + \frac{2\|V\|_\infty S\iota}{3n(s, a)}, \quad (55)$$

where $\iota = O(\log(SA/\delta))$.

Proof. Suppose the event \mathcal{E}_1 holds. Then we have (deterministically)

$$\begin{aligned} |(\widehat{P}_{s,a} - P_{s,a})V| &\stackrel{(i)}{=} |(\widehat{P}_{s,a} - P_{s,a})(V - P_{s,a}V\mathbf{1}_S)| \\ &\leq \left(\sqrt{\frac{2P(\cdot|s, a)\iota}{n(s, a)}} + \frac{2\iota}{3n(s, a)}\right) |V - P_{s,a}V\mathbf{1}_S| \\ &\leq \sqrt{\frac{2P_{s,a}\iota}{n(s, a)}} |V - P_{s,a}V\mathbf{1}_S| + \frac{2S\|V\|_\infty\iota}{3n(s, a)} \\ &\stackrel{(ii)}{\leq} \left(\sqrt{S} \sqrt{\frac{2P_{s,a}(V - P_{s,a}V\mathbf{1}_S)^2\iota}{n(s, a)}}\right) + \frac{2S\|V\|_\infty\iota}{3n(s, a)} \\ &\leq \sqrt{\frac{2S\text{Var}(P_{s,a}, V)\iota}{n(s, a)}} + \frac{2\|V\|_\infty S\iota}{3n(s, a)}, \end{aligned}$$

where (i) follows from the fact that $P_{s,a}V$ is a scalar, which implies that $(\widehat{P}_{s,a} - P_{s,a})(P_{s,a}V)\mathbf{1}_S = (P_{s,a}V) \sum_{s'} (\widehat{P}(s'|s, a) - P(s'|s, a)) = 0$. (ii) uses the Cauchy-Schwartz inequality. Lastly, \mathcal{E}_1 fails with probability only δ (by Lemma 4.1). \square

11 MINIMAX LOWER BOUND FOR OFFLINE SSP

In this section, we provide the minimax lower bound for offline stochastic shortest path problem. Concretely, we consider the family of problems satisfying bounded partial coverage, *i.e.* $\max_{s, a, s \neq g} \frac{d^{\pi^*}(s, a)}{d^\mu(s, a)} \leq C^*$, where $d^\pi(s, a) = \sum_{h=0}^{\infty} \xi_h^\pi(s, a) < \infty$ for all s, a (excluding g) for any proper policy π . Formally, we have the following result:

Theorem 11.1 (Restatement of Theorem ??). *We define the following family of SSPs:*

$$\text{SSP}(C^*) = \{(s_{\text{init}}, \mu, P, c) \mid \max_{s,a,s \neq g} \frac{d^{\pi^*}(s,a)}{d^\mu(s,a)} \leq C^*\},$$

where $d^\pi(s,a) = \sum_{h=0}^{\infty} \xi_h^\pi(s,a)$. Then for any $C^* \geq 1$, $\|V^*\|_\infty = B_\star > 1$, it holds (for some universal constant c)

$$\inf_{\hat{\pi} \text{ proper}} \sup_{(s_{\text{init}}, \mu, P, c) \in \text{SSP}(C^*)} \mathbb{E}_{\mathcal{D}}[V^{\hat{\pi}}(s_{\text{init}}) - V^*(s_{\text{init}})] \geq c \cdot B_\star \sqrt{\frac{SC^*}{n}}.$$

The proof of Theorem 11.1 relies on the hard instances construction that is similar to Rashidinejad et al. [2021]. However, we need to incorporate the absorbing state g and assign the transition of initial state s_{init} carefully to make sure the optimal proper policy exists.

Proof of Theorem 11.1. We create hard instances of SSPs as follows: we split $S - 1$ states (except s_{init}) into $S' = (S - 1)/2$ groups, and denote it as $\mathcal{S} = \{s_{\text{init}}\} \cup \{s_1^j, s_+^j\}_{j=1}^{S'}$. For s_1^j , $j = 1, \dots, S'$, there are two actions a_1, a_2 and for states s_{init}, s_+^j and goal state g there is only one default action a_d (therefore the only choice is always optimal for those states). Concretely,

- For state s_{init} , it transitions to s_1^j ($j = 1, \dots, S'$) uniformly with probability $1/S'$, i.e. $P(s_1^j | s_{\text{init}}, a_d) = 1/S'$;
- For each state s_1^j , it satisfies

$$P(s_+^j | s_1^j, a_1) = P(g | s_1^j, a_1) = 1/2; \quad P(s_+^j | s_1^j, a_2) = \frac{1}{2} + v_j \delta; \quad P(g | s_1^j, a_2) = \frac{1}{2} - v_j \delta.$$

where $v_j \in \{+1, -1\}$ and δ to be specified later.

- For s_+^j , it satisfies

$$P(s_+^j | s_+^j, a_d) = q, \quad P(g | s_+^j, a_d) = 1 - q,$$

where $q = 1 - \frac{1}{B_\star}$ and g is absorbing.

- the cost function satisfies (regardless of actions):

$$c(s_{\text{init}}) = c(s_1^j) = c(s_+^j) = 1, \quad c(g) = 0.$$

It is easy to check this is a SSP. Moreover, it is clear when $v_j = 1$, the optimal action at s_1^j is a_1 and if $v_j = -1$ the optimal action is a_2 . Note by straightforward calculation we have that $\|V^*\|_\infty \leq 2B_\star$.

We consider the family of SSP instances \mathcal{P} to satisfy Lemma 12.1, i.e. it satisfies $|\mathcal{P}| \geq e^{S'/8}$ and for any two instances in \mathcal{P} , $\|v_i - v_j\|_1 \geq S'/2$. Also, it suffices to consider all the deterministic learning algorithms, as stochastic output policies are randomized versions over deterministic ones (c.f. Krishnamurthy et al. [2016]). Then we have the following lemma:

Lemma 11.2. *For any (deterministic) policy π and any two different transition probabilities $P_1, P_2 \in \mathcal{P}$, it holds:*

$$V_{P_1}^\pi(s_{\text{init}}) - V_{P_1}^{\pi^*}(s_{\text{init}}) + V_{P_2}^\pi(s_{\text{init}}) - V_{P_2}^{\pi^*}(s_{\text{init}}) \geq \delta B_\star / 2.$$

Proof of Lemma 11.2. Since s_{init} uniformly transitions to S' states (w.r.t the default action a_d), therefore for any policy π ,

$$V_{P_1}^\pi(s_{\text{init}}) - V_{P_1}^{\pi^*}(s_{\text{init}}) = 1 + \frac{1}{S'} \sum_{i=1}^{S'} V_{P_1}^\pi(s_1^i) - \left(1 + \frac{1}{S'} \sum_{i=1}^{S'} V_{P_1}^{\pi^*}(s_1^i)\right) = \frac{1}{S'} \sum_{i=1}^{S'} (V_{P_1}^\pi(s_1^i) - V_{P_1}^{\pi^*}(s_1^i)).$$

Case 1. If $v_j = 1$, then

$$P(s_+^i | s_1^i, a_2) = \frac{1}{2} + \delta, \quad P(g | s_1^i, a_2) = \frac{1}{2} - \delta$$

and in this case $\pi^*(s_1^i) = a_1$.

If $\pi(s_1^i) = a_2$, then

$$V^\pi(s_1^i) = 1 + \left(\frac{1}{2} + \delta\right)V^\pi(s_+^i) + \left(\frac{1}{2} - \delta\right)V^\pi(g) = 1 + \left(\frac{1}{2} + \delta\right)V^\pi(s_+^i),$$

and this implies

$$\begin{aligned} V^\pi(s_1^i) - V^{\pi^*}(s_1^i) &= \left(\frac{1}{2} + \delta\right)V^\pi(s_+^i) - \frac{1}{2}V^{\pi^*}(s_+^i) \\ &\geq \delta V^\pi(s_+^i) = \delta(1 + q + q^2 + \dots) = \delta \cdot \frac{1}{1 - q} = \delta B_\star. \end{aligned}$$

If $\pi(s_1^i) = a_1$, then $V^\pi(s_1^i) - V^{\pi^*}(s_1^i) \geq 0$. Therefore, in this case, one has

$$V^\pi(s_1^i) - V^{\pi^*}(s_1^i) \geq \delta B_\star \cdot \mathbf{1}[\pi(s_1^i) \neq \pi^*(s_1^i)].$$

Case 2. If $v_j = -1$, then

$$P(s_+^i | s_1^i, a_2) = \frac{1}{2} - \delta, \quad P(g | s_1^i, a_2) = \frac{1}{2} + \delta$$

and in this case $\pi^*(s_1^i) = a_2$.

If $\pi(s_1^i) = a_1$, then

$$V^\pi(s_1^i) = 1 + \frac{1}{2} \cdot V^\pi(s_+^i) + \frac{1}{2}V^\pi(g) = 1 + \frac{1}{2}V^\pi(s_+^i),$$

and this implies

$$\begin{aligned} V^\pi(s_1^i) - V^{\pi^*}(s_1^i) &= \frac{1}{2}V^\pi(s_+^i) - \left(\frac{1}{2} - \delta\right)V^{\pi^*}(s_+^i) \\ &\geq \delta V^{\pi^*}(s_+^i) = \delta(1 + q + q^2 + \dots) = \delta \cdot \frac{1}{1 - q} = \delta B_\star. \end{aligned}$$

If $\pi(s_1^i) = a_2$, then $V^\pi(s_1^i) - V^{\pi^*}(s_1^i) \geq 0$. Therefore, in this case, we still have

$$V^\pi(s_1^i) - V^{\pi^*}(s_1^i) \geq \delta B_\star \cdot \mathbf{1}[\pi(s_1^i) \neq \pi^*(s_1^i)].$$

Combine the above two cases, we have

$$\begin{aligned} &V_{P_1}^\pi(s_{\text{init}}) - V_{P_1}^{\pi^*}(s_{\text{init}}) + V_{P_2}^\pi(s_{\text{init}}) - V_{P_2}^{\pi^*}(s_{\text{init}}) \\ &\geq \frac{1}{S'} \sum_{i=1}^{S'} (V_{P_1}^\pi(s_1^i) - V_{P_1}^{\pi^*}(s_1^i)) + \frac{1}{S'} \sum_{i=1}^{S'} (V_{P_2}^\pi(s_1^i) - V_{P_2}^{\pi^*}(s_1^i)) \\ &\geq \frac{1}{S'} \delta B_\star \sum_{i=1}^{S'} (\mathbf{1}[\pi(s_1^i) \neq \pi_{P_1}^*(s_1^i)] + \mathbf{1}[\pi(s_1^i) \neq \pi_{P_2}^*(s_1^i)]) \\ &\geq \frac{\delta B_\star}{S'} \sum_{i=1}^{S'} \mathbf{1}[\pi_{P_1}^*(s_1^i) \neq \pi_{P_2}^*(s_1^i)] \end{aligned} \tag{56}$$

Lastly, by Lemma 12.1, $\sum_{i=1}^{S'} \mathbf{1}[\pi_{P_1}^*(s_1^i) \neq \pi_{P_2}^*(s_1^i)] = \|v_{P_1} - v_{P_2}\|_1 \geq S'/2$, and plug this back to (56) we obtain the result. \square

Now we construct the behavior policy μ such that the data trajectories generated from the induced distribution $\mu \circ P$ suffice for the lower bound. Since only s_1^i has two actions, we specify below:

$$\mu(a_2 | s_1^i) = 1/C^\star, \quad \mu(a_1 | s_1^i) = 1 - 1/C^\star, \quad \forall i \in \{1, \dots, S'\}.$$

First, we examine this choice belongs to $\text{SSP}(C^*)$. Indeed, the only case where a_2 is the suboptimal action for all s_1^i ($i \in \{1, \dots, S'\}$) is when $v_1, \dots, v_{S'}$ all equal 1. We can eliminate this SSP from \mathcal{P} and the property of Lemma 12.1 still holds. Then, for some i_0 such that $v_{i_0} = -1$ (a_2 is the optimal action for this state), we have

$$d^*(s_1^{i_0}, a_2) = d^*(s_1^{i_0}) \cdot 1 = \frac{1}{S'}, \quad d^\mu(s_1^{i_0}, a_2) = d^\mu(s_1^{i_0})\mu(a_2|s_1^{i_0}) = \frac{1}{S'C^*},$$

therefore $d^*(s_1^{i_0}, a_2)/d^\mu(s_1^{i_0}, a_2) = C^*$ and this $(s_{\text{init}}, \mu, P, c) \in \text{SSP}(C^*)$.

Recall n is the number of episodes. Now apply Fano's inequality (Lemma 12.5) (where each whole trajectory is considered one single data point over the distribution $\mu \circ P$ therefore $\mathcal{D} := \{(s_1^{(i)}, a_1^{(i)}, c_1^{(i)}, s_2^{(i)}, \dots, s_{T_i}^{(i)})\}_{i=1, \dots, n}$ consists of n i.i.d. samples) and Lemma 11.2, we have¹

$$\inf_{\hat{\pi}} \sup_{(s_{\text{init}}, \mu, P, c) \in \mathcal{P}} \mathbb{E}_{\mathcal{D}}[V^{\hat{\pi}}(s_{\text{init}}) - V^*(s_{\text{init}})] \geq \frac{\delta B^*}{4} \left(1 - \frac{n \cdot \max_{i \neq j} \text{KL}(\mu \circ P_i || \mu \circ P_j) + \log 2}{\log |\mathcal{P}|} \right).$$

Note by the choice of \mathcal{P} , $\log |\mathcal{P}| \geq S'/8$, therefore it remains to bound $\max_{i \neq j} \text{KL}(\mu \circ P_i || \mu \circ P_j)$. By definition, we have

$$\text{KL}(\mu \circ P_1 || \mu \circ P_2) = \frac{1}{S'} \sum_{i=1}^{S'} \sum_{\tau_{s_1^i}} \mathbb{P}_1(\tau_{s_1^i}) \log \frac{\mathbb{P}_1(\tau_{s_1^i})}{\mathbb{P}_2(\tau_{s_1^i})},$$

where $\tau_{s_1^i}$ corresponds to all the possible trajectories starting from s_1^i . Then there are the following several cases:²

- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_1 \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_1 \rightarrow g) = (1 - \frac{1}{C^*})\frac{1}{2}$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_2 \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_2 \rightarrow g) = \frac{1}{C^*} \cdot (\frac{1}{2} - v_i\delta)$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_1 \rightarrow s_+^i \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_1 \rightarrow s_+^i \rightarrow g) = (1 - \frac{1}{C^*})\frac{1}{2}(1 - q)$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_2 \rightarrow s_+^i \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_2 \rightarrow s_+^i \rightarrow g) = \frac{1}{C^*} \cdot (\frac{1}{2} + v_i\delta)(1 - q)$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_1 \rightarrow s_+^i \rightarrow s_+^i \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_1 \rightarrow s_+^i \rightarrow s_+^i \rightarrow g) = (1 - \frac{1}{C^*})\frac{1}{2}q(1 - q)$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_2 \rightarrow s_+^i \rightarrow s_+^i \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_2 \rightarrow s_+^i \rightarrow s_+^i \rightarrow g) = \frac{1}{C^*}(\frac{1}{2} + v_i\delta)q(1 - q)$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_1 \rightarrow (s_+^i)_{\times k} \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_1 \rightarrow (s_+^i)_{\times k} \rightarrow g) = (1 - \frac{1}{C^*})\frac{1}{2}q^{k-1}(1 - q)$;
- If $\tau_{s_1^i} = \{s_1^i \rightarrow a_2 \rightarrow (s_+^i)_{\times k} \rightarrow g\}$, then $\mathbb{P}(s_1^i \rightarrow a_2 \rightarrow (s_+^i)_{\times k} \rightarrow g) = \frac{1}{C^*}(\frac{1}{2} + v_i\delta)q^{k-1}(1 - q)$;

Note for path $\tau_{s_1^i}$ that chooses action a_1 , $\mathbb{P}_1(\tau_{s_1^i}) = \mathbb{P}_2(\tau_{s_1^i})$ which implies $\mathbb{P}_1(\tau_{s_1^i}) \log \frac{\mathbb{P}_1(\tau_{s_1^i})}{\mathbb{P}_2(\tau_{s_1^i})} = 0$, so we only need to sum over the paths that choose a_2 . In particular, we have

$$\begin{aligned} & \sum_{\tau_{s_1^i}} \mathbb{P}_1(\tau_{s_1^i}) \log \frac{\mathbb{P}_1(\tau_{s_1^i})}{\mathbb{P}_2(\tau_{s_1^i})} = \frac{1}{C^*} \cdot (\frac{1}{2} - v_i^{P_1}\delta) \log \frac{\frac{1}{C^*} \cdot (\frac{1}{2} - v_i^{P_1}\delta)}{\frac{1}{C^*} \cdot (\frac{1}{2} - v_i^{P_2}\delta)} \\ & + \sum_{k=0}^{\infty} \frac{1}{C^*} (\frac{1}{2} + v_i^{P_1}\delta) q^{k-1} (1 - q) \log \frac{\frac{1}{C^*} (\frac{1}{2} + v_i^{P_1}\delta) q^{k-1} (1 - q)}{\frac{1}{C^*} (\frac{1}{2} + v_i^{P_2}\delta) q^{k-1} (1 - q)} \\ & = \frac{1}{C^*} \cdot (\frac{1}{2} - v_i^{P_1}\delta) \log \frac{\frac{1}{2} - v_i^{P_1}\delta}{\frac{1}{2} - v_i^{P_2}\delta} + \frac{1}{C^*} \cdot (\frac{1}{2} + v_i^{P_1}\delta) \log \frac{\frac{1}{2} + v_i^{P_1}\delta}{\frac{1}{2} + v_i^{P_2}\delta} \\ & \leq \frac{1}{C^*} \cdot (\frac{1}{2} - v_i^{P_1}\delta) \log \frac{\frac{1}{2} - v_i^{P_1}\delta}{\frac{1}{2} + v_i^{P_1}\delta} + \frac{1}{C^*} \cdot (\frac{1}{2} + v_i^{P_1}\delta) \log \frac{\frac{1}{2} + v_i^{P_1}\delta}{\frac{1}{2} - v_i^{P_1}\delta} \\ & = \frac{1}{C^*} 2\delta \log \frac{\frac{1}{2} + \delta}{\frac{1}{2} - \delta} = \frac{1}{C^*} 2\delta \log(1 + \frac{2\delta}{\frac{1}{2} - \delta}) \leq \frac{4\delta^2}{C^*}, \end{aligned}$$

¹Note here we drop $\hat{\pi}$ is proper as the theorem statement did. We can do this since, for all the instances in \mathcal{P} , any policy is proper.

²We omit the subscript j in P_j here and only uses \mathbb{P} to denote $\mu \circ P$ for the moment.

where the first inequality comes from when $v_i^{P_1} = v_i^{P_2}$ then the term is simply 0 and the second to the last equality holds true regardless of whether $v_i = 1$ or $v_i = -1$. The last inequality comes from $\log(1+x) \leq x$ for all $x > -1$ (here $0 < \delta < \frac{1}{2}$).

Plug above back into the definition we obtain

$$\max_{i \neq j} \text{KL}(\mu \circ P_i \| \mu \circ P_j) \leq 4\delta^2 / C^*,$$

and as long as

$$\frac{4n\delta^2}{C^*S'/8} \leq \frac{1}{2}$$

e.g. if we choose $\delta = \frac{1}{16} \sqrt{\frac{C^*S}{n}}$ (recall $S' = (S-1)/2$), then we have

$$\inf_{\hat{\pi}} \sup_{(s_{\text{init}}, \mu, P, c) \in \mathcal{P}} \mathbb{E}_{\mathcal{D}}[V^{\hat{\pi}}(s_{\text{init}}) - V^*(s_{\text{init}})] \geq \frac{\delta B_*}{4} \frac{1}{2} = \frac{1}{128} B_* \sqrt{\frac{C^*S}{n}}.$$

This completes the proof. □

12 TECHNICAL LEMMAS

Lemma 12.1 (Gilbert-Varshamov). *There exists a subset \mathcal{V} of $\{-1, 1\}^S$ such that*

- $|\mathcal{V}| \geq 2^{S/8}$;
- for any two different $v_i, v_j \in \mathcal{V}$, it holds $\|v_i - v_j\|_1 \geq S/2$.

Lemma 12.2 (Generalized Chernoff bound). *Suppose X_1, \dots, X_n are independent random variables taking values in $[a, b]$. Let $X = \sum_{i=1}^n X_i$ denote their sum and let $\mu = E[X_i]$. Then for any $\delta > 0$,*

$$\mathbb{P}[X < (1-\theta)n\mu] \leq e^{-2\theta^2 n\mu^2 / (b-a)^2} \quad \text{and} \quad \mathbb{P}[X \geq (1+\theta)pn] \leq e^{-2\theta^2 n\mu^2 / (b-a)^2}.$$

This result can be found in Sums of independent bounded random variables Section of https://en.wikipedia.org/wiki/Chernoff_bound.

Lemma 12.3 (Bernstein's Inequality). *Let x_1, \dots, x_n be independent bounded random variables such that $\mathbb{E}[x_i] = 0$ and $|x_i| \leq \xi$ with probability 1. Let $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}[x_i]$, then with probability $1 - \delta$ we have*

$$\frac{1}{n} \sum_{i=1}^n x_i \leq \sqrt{\frac{2\sigma^2 \cdot \log(1/\delta)}{n}} + \frac{2\xi}{3n} \log(1/\delta)$$

Lemma 12.4 (Empirical Bernstein's Inequality [Maurer and Pontil, 2009]). *Let x_1, \dots, x_n be i.i.d random variables such that $|x_i| \leq \xi$ with probability 1. Let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\widehat{V}_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$, then with probability $1 - \delta$ we have*

$$\left| \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E}[x] \right| \leq \sqrt{\frac{2\widehat{V}_n \cdot \log(2/\delta)}{n}} + \frac{7\xi}{3n} \log(2/\delta).$$

Lemma 12.5 (Generalized Fano's inequality). *Let $L : \Theta \times \mathcal{A} \rightarrow \mathbb{R}_+$ be any loss function, and there exist $\theta_1, \dots, \theta_m \in \Theta$ such that*

$$L(\theta_i, a) + L(\theta_j, a) \geq \Delta, \quad \forall i \neq j \in [m], a \in \mathcal{A}.$$

Then it holds

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} L(\theta, \hat{\theta}) \geq \frac{\Delta}{2} \left(1 - \frac{n \cdot \max_{i \neq j} \text{KL}(\mathbb{P}_{\theta_i} \| \mathbb{P}_{\theta_j}) + \log 2}{\log m} \right),$$

where n is the number of i.i.d. samples sampled from the distribution \mathbb{P}_{θ} .

Proof of Lemma 12.5. The proof come from the combination of Lemma 1 and Lemma 3 of [Han and Fischer-Hwang \[2019\]](#). \square

Lemma 12.6 (Chernoff Bound for Stochastic Shortest Path). *Recall by definition*

$$n(s, a) = \sum_{i=1}^n \sum_{h=1}^{T_i} \mathbf{1}[s_h^{(i)} = s, a_h^{(i)} = a] = \sum_{i=1}^n \sum_{h=1}^{\infty} \mathbf{1}[s_h^{(i)} = s, a_h^{(i)} = a].$$

Let $T_{\max} = \max_i T_i$ and recall $d_m := \min\{\sum_{h=0}^{\infty} \xi_h^\mu(s, a) : s.t. \sum_{h=0}^{\infty} \xi_h^\mu(s, a) > 0\}$. When $n > C \cdot T_{\max}^2 \log(SA/\delta)/d_m^2$, with probability $1 - \delta$, for all $s, a \in \mathcal{S} \times \mathcal{A}$,

$$\frac{1}{2}n \cdot \sum_{h=1}^{\infty} \xi_h^\mu(s, a) \leq n(s, a) \leq \frac{3}{2}n \cdot \sum_{h=1}^{\infty} \xi_h^\mu(s, a).$$

Proof of Lemma 12.6. Indeed, denote $n_t(s, a) = \sum_{i=1}^n \sum_{h=1}^t \mathbf{1}[s_h^{(i)} = s, a_h^{(i)} = a]$, then

$$\mathbb{E}[n_t(s, a)] = \sum_{i=1}^n \sum_{h=1}^t \mathbb{E}[\mathbf{1}[s_h^{(i)} = s, a_h^{(i)} = a]] = \sum_{i=1}^n \sum_{h=1}^t \xi_h^\mu(s, a) = n \sum_{h=1}^t \xi_h^\mu(s, a).$$

Now define $X_{i,t} = \sum_{h=1}^t \mathbf{1}[s_h^{(i)} = s, a_h^{(i)} = a]$, then by $T_{\max} = \max_i T_i$ we have $0 \leq X_{i,t} \leq T_{\max}$ for all i, t since T_{\max} denotes the maximum length of trajectory. Then apply Lemma 12.2 (where we pick $\theta = \frac{1}{2}$) to $n_t(s, a)$ and $\sum_{h=1}^t \xi_h^\mu(s, a)$ and union bound over s, a , we have with probability $1 - \delta$, for any fixed t ,

$$\mathbb{P} \left[\frac{1}{2}n \cdot \sum_{h=1}^t \xi_h^\mu(s, a) \leq n_t(s, a) \leq \frac{3}{2}n \cdot \sum_{h=1}^t \xi_h^\mu(s, a), \forall s, a \right] \geq 1 - \delta$$

Next note $n_t(s, a) \rightarrow n(s, a)$ almost surely, and $\sum_{h=1}^t \xi_h^\mu(s, a) \rightarrow \sum_{h=1}^{\infty} \xi_h^\mu(s, a)$ almost surely, and that a.s. convergence implies convergence in distribution, we have

$$\begin{aligned} & \mathbb{P} \left[\frac{1}{2}n \cdot \sum_{h=1}^{\infty} \xi_h^\mu(s, a) \leq n(s, a) \leq \frac{3}{2}n \cdot \sum_{h=1}^{\infty} \xi_h^\mu(s, a), \forall s, a \right] \\ &= \lim_{t \rightarrow \infty} \mathbb{P} \left[\frac{1}{2}n \cdot \sum_{h=1}^t \xi_h^\mu(s, a) \leq n_t(s, a) \leq \frac{3}{2}n \cdot \sum_{h=1}^t \xi_h^\mu(s, a), \forall s, a \right] \geq 1 - \delta \end{aligned}$$

\square

Lemma 12.7. For any $a, b, c \in \mathbb{R}$, we have

$$|\min\{a, b\} - \min\{a, c\}| \leq |b - c|. \quad (57)$$

Proof. 1. Case I: $a \leq b$ and $a \leq c$, $|\min\{a, b\} - \min\{a, c\}| = 0$.

2. Case II: $a \geq b$ and $a \geq c$, $|\min\{a, b\} - \min\{a, c\}| = |b - c|$.

3. Case III: $b < a < c$ or $c < a < b$, $|\min\{a, b\} - \min\{a, c\}| \leq \max\{|a - b|, |a - c|\} \leq |b - c|$.

\square

References

Dimitri P Bertsekas and John N Tsitsiklis. An analysis of stochastic shortest path problems. *Mathematics of Operations Research*, 16(3):580–595, 1991.

Yanjun Han and Irena Fischer-Hwang. Multiple hypothesis testing: Tree, fano and assoaud. 2019. URL <https://theinformaticists.com/2019/09/16/lecture-8-multiple-hypothesis-testing-tree-fano-and-assoaud/>.

Akshay Krishnamurthy, Alekh Agarwal, and John Langford. PAC reinforcement learning with rich observations. In *Advances in Neural Information Processing Systems*, pages 1840–1848, 2016.

Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*, 2009.

Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *arXiv preprint arXiv:2103.12021*, 2021.

Jean Tarbouriech, Runlong Zhou, Simon S Du, Matteo Pirota, Michal Valko, and Alessandro Lazaric. Stochastic shortest path: Minimax, parameter-free and towards horizon-free regret. *Advances in neural information processing systems*, 2021.