

A. Details of Coutherexamples

In this section we provide details of computing the variance in Figure 1. For each MDP, there are totally four possible trajectories (product of two actions and two steps), and the probabilities of them under behavior policy are all 1/4. We list the return of different estimators for those four trajectories, then compute the variance of the estimators.

	Probabilities of path	Example 1a			Example 1b			Example 1c		
		IS	PDIS	MIS	IS	PDIS	MIS	IS	PDIS	MIS
a_1, a_1	0.25	1.44	1.2	1.2	0	0	0	0	0	0
a_1, a_2	0.25	1.92	2.16	2.0	1.44	1.44	1.2	0.96	0.96	0.8
a_2, a_1	0.25	0.96	0.8	0.8	0.64	0.8	0.8	0.96	0.8	0.8
a_2, a_2	0.25	1.28	1.44	1.6	1.92	1.76	2.0	1.28	1.44	1.6
Expectation		1.4	1.4	1.4	1	1	1	0.8	0.8	0.8
Variance		0.12	0.2448	0.2	0.5424	0.4528	0.52	0.2304	0.2688	0.32

Table 1: Importance sampling returns and the variance. See figure 1 for the problem structure.

B. Proof of Lemma 1

Proof. In this proof, we use τ to denote the trajectory without reward: $\tau_{1:t} = \{s_k, a_k\}_{k=1}^t$. Since $\mathbb{E}(\rho_{1:t}|s_t, a_t) = \mathbb{E}(\rho_{1:t-1}|s_t, a_t)\rho_t$, we only need to prove that $\mathbb{E}(\rho_{1:t-1}|s_t, a_t) = \frac{d^\pi(s_t)}{d^\mu(s_t)}$.

$$\mathbb{E}(\rho_{1:t-1}|s_t, a_t) = \int \prod_{k=1}^{t-1} \frac{\pi(s_k, a_k)}{\mu(s_k, a_k)} p_\mu(\tau_{1:t-1}|s_t, a_t) d\tau_{1:t-1} \quad (1)$$

$$= \int \frac{p_\pi(\tau_{1:t-1})}{p_\mu(\tau_{1:t-1})} p_\mu(\tau_{1:t-1}|s_t, a_t) d\tau_{1:t-1} \quad (2)$$

$$= \int \frac{p_\pi(\tau_{1:t-1})}{p_\mu(\tau_{1:t-1})} \frac{p_\mu(\tau_{1:t-1}) p(s_t|\tau_{1:t-1}) \mu(a_t|s_t)}{p_\mu(s_t, a_t)} d\tau_{1:t-1} \quad (3)$$

$$= \int \frac{p_\pi(\tau_{1:t-1})}{p_\mu(\tau_{1:t-1})} \frac{p_\mu(\tau_{1:t-1}) p(s_t|s_{t-1}, a_{t-1}) \mu(a_t|s_t)}{d_t^\mu(s_t) \mu(a_t|s_t)} d\tau_{1:t-1} \quad (4)$$

$$= \frac{1}{d_t^\mu(s_t)} \int p(s_t|s_{t-1}, a_{t-1}) p_\pi(\tau_{1:t-1}) d\tau_{1:t-1} \quad (5)$$

$$= \frac{1}{d_t^\mu(s_t)} \int p(s_t|\tau_{1:t-1}) p_\pi(\tau_{1:t-1}) d\tau_{1:t-1} \quad (6)$$

$$= \frac{d_t^\pi(s_t)}{d_t^\mu(s_t)} \quad (7)$$

□

C. Proofs for Finite Horizon Case

C.1. Proof of Lemma 2

Proof. Since $\mathbb{E}(\sum_t \mathbb{E}(Y_t|X_t)) = \mathbb{E}(\sum_t Y_t)$, we just need to compute the difference between the second moment of $\sum_t Y_t$ and $\sum_t \mathbb{E}(Y_t|X_t)$:

$$\mathbb{E}\left(\sum_t \mathbb{E}(Y_t|X_t)\right)^2 = \mathbb{E}\left(\sum_t (\mathbb{E}(Y_t|X_t))^2 + 2 \sum_{t < k} \mathbb{E}(Y_t|X_t)\mathbb{E}(Y_k|X_k)\right) \quad (8)$$

$$= \sum_t \mathbb{E}(\mathbb{E}(Y_t|X_t))^2 + 2 \sum_{t < k} \mathbb{E}(\mathbb{E}(Y_t|X_t)\mathbb{E}(Y_k|X_k)) \quad (9)$$

$$\leq \sum_t \mathbb{E}(\mathbb{E}(Y_t^2|X_t)) + 2 \sum_{t < k} \mathbb{E}(\mathbb{E}(Y_t|X_t)\mathbb{E}(Y_k|X_k)) \quad (10)$$

$$= \sum_t \mathbb{E}(Y_t^2) + 2 \sum_{t < k} \mathbb{E}(\mathbb{E}(Y_t|X_t)\mathbb{E}(Y_k|X_k)) \quad (11)$$

$$\mathbb{E}\left(\sum_t Y_t\right)^2 = \mathbb{E}\left(\sum_t Y_t^2 + 2 \sum_{t < k} Y_t Y_k\right) \quad (12)$$

$$= \sum_t \mathbb{E}(Y_t^2) + 2 \sum_{t < k} \mathbb{E}(Y_t Y_k) \quad (13)$$

Thus we finished the proof by taking the difference between $\mathbb{E}(\sum_t Y_t)^2$ and $\mathbb{E}(\sum_t \mathbb{E}(Y_t|X_t))^2$. \square

C.2. Proof of Theorem 1

Proof. Let $\tau_{1:t}$ be the first t steps in a trajectory: $(s_1, a_1, r_1, \dots, s_t, a_t, r_t)$, then $\rho_{1:t} r_t = \mathbb{E}(\rho_{1:T} r_t | \tau_{1:t})$. To prove the inequality between the variance of importance sampling and per decision importance sampling, we apply Lemma 2 to the variance, letting $Y_t = r_t \rho_{1:T}$ and $X_t = \tau_{1:t}$. Then it is sufficient to show that for any $1 \leq t < k \leq T$,

$$\mathbb{E}(r_t r_k \rho_{1:T} \rho_{1:T}) = \mathbb{E}(Y_t Y_k) \geq \mathbb{E}(\mathbb{E}(Y_t|X_t)\mathbb{E}(Y_k|X_k)) = \mathbb{E}(r_t r_k \rho_{1:t} \rho_{1:k}) \quad (14)$$

To prove that, it is sufficient to show $\mathbb{E}(r_t r_k \rho_{1:T} \rho_{1:T} | \tau_{1:t}) \geq \mathbb{E}(r_t r_k \rho_{1:t} \rho_{1:k} | \tau_{1:t})$. Since

$$\mathbb{E}(r_t r_k \rho_{1:t} \rho_{1:k} | \tau_{1:t}) = r_t \rho_{1:t}^2 \mathbb{E}(r_k \rho_{t+1:k} | \tau_{1:t}) \quad (15)$$

$$= r_t \rho_{1:t}^2 \mathbb{E}(r_k \rho_{t+1:T} | \tau_{1:t}) \quad (16)$$

$$= r_t \rho_{1:t}^2 \mathbb{E}(r_k \rho_{t+1:T} | \tau_{1:t}) \mathbb{E}(\rho_{t+1:T} | \tau_{1:t}) \quad (17)$$

$$(18)$$

Given $\tau_{1:t}$, r_k and $\rho_{t+1:T}$ can be viewed as r_{k-t+1} and $\rho_{1:T-t+1}$ on a new trajectory. Then according to the statement of theorem, $r_{k-t+1} \rho_{1:T-t+1}$ and $\rho_{1:T-t+1}$ are positively correlated. Now we can upper bound $\mathbb{E}(r_t r_k \rho_{1:t} \rho_{1:k} | \tau_{1:t})$ by:

$$r_t \rho_{1:t}^2 \mathbb{E}(r_k \rho_{t+1:T} | \tau_{1:t}) \mathbb{E}(\rho_{t+1:T} | \tau_{1:t}) \leq r_t \rho_{1:t}^2 \mathbb{E}(r_k \rho_{t+1:T} \rho_{t+1:T} | \tau_{1:t}) \quad (19)$$

$$= \mathbb{E}(r_t r_k \rho_{1:T} \rho_{1:T} | \tau_{1:t}) \quad (20)$$

This implies $\mathbb{E}(r_t r_k \rho_{1:T} \rho_{1:T}) \geq \mathbb{E}(r_t r_k \rho_{1:t} \rho_{1:k})$ by taking expectation over $\tau_{1:t}$, and finish the proof. \square

C.3. Proof of Theorem 2

Proof. Using lemma 2 by $Y_t = \rho_{1:t} r_t$ and $X_t = s_t, a_t, r_t$, we have that the variance of \hat{v}_{SIS} is smaller than the variance of \hat{v}_{PDIS} if for any $t < k$:

$$\mathbb{E}[\rho_{1:t} \rho_{0:k} r_t r_k] \geq \mathbb{E}[\mathbb{E}(\rho_{1:t} | s_t, a_t) \mathbb{E}(\rho_{0:k} | s_k, a_k) r_t r_k] \quad (21)$$

$$= \mathbb{E}\left[\frac{d_t^\pi(s, a)}{d_t^\mu(s, a)} \frac{d_k^\pi(s, a)}{d_k^\mu(s, a)} r_t r_k\right] \quad (22)$$

The second line follows from Lemma 1 to simplify $\mathbb{E}(\rho_{0:t}|s_t, a_t)$. To show that, we will transform the above equation into an expression about two covariances. To proceed we subtracting $\mathbb{E}(\rho_{1:t}r_t)\mathbb{E}(\rho_{1:k}r_k)$ from both sides, and note that the resulting left hand side is simply the covariance:

$$\begin{aligned} \text{Cov}[\rho_{1:t}r_t, \rho_{0:k}r_k] &= \mathbb{E}[\rho_{1:t}\rho_{1:k}r_t r_k] - \mathbb{E}(\rho_{1:t}r_t)\mathbb{E}(\rho_{1:k}r_k) \\ &\geq \mathbb{E}\left[\frac{d_t^\pi(s, a)}{d_t^\mu(s, a)}\frac{d_k^\pi(s, a)}{d_k^\mu(s, a)}r_t r_k\right] - \mathbb{E}(\rho_{1:t}r_t)\mathbb{E}(\rho_{1:k}r_k) \end{aligned} \quad (23)$$

We now expand the second term in the right hand side

$$\mathbb{E}(\rho_{1:t}r_t)\mathbb{E}(\rho_{1:k}r_k) = \mathbb{E}(r_t\mathbb{E}(\rho_{1:t}|s_t, a_t))\mathbb{E}(r_k\mathbb{E}(\rho_{1:k}|s_k, a_k)) \quad (24)$$

$$= \mathbb{E}\left[\frac{d_t^\pi(s, a)}{d_t^\mu(s, a)}r_t\right]\mathbb{E}\left[\frac{d_k^\pi(s, a)}{d_k^\mu(s, a)}r_k\right] \quad (25)$$

This shows that both sides of 23 are covariances. The result then follows under the assumption of the proof. \square

D. Proofs for infinite horizon case

D.1. Proof of Theorem 3

Proof. We can write the log of likelihood ratio as sum of random variables on a Markov chain,

$$\log \rho_{1:T} = \sum_{t=1}^T \log \rho_t = \sum_{t=1}^T \log \left(\frac{\pi(a_t|s_t)}{\mu(a_t|s_t)} \right) \quad (26)$$

By the strong law of large number on Markov chain (Breiman, 1960):

$$\frac{1}{T} \log \rho_{1:T} = \frac{1}{T} \sum_{i=1}^T \log \left(\frac{\pi(a_i|s_i)}{\mu(a_i|s_i)} \right) \xrightarrow{a.s.} \mathbb{E}_{d^\mu} \log \left(\frac{\pi(a_i|s_i)}{\mu(a_i|s_i)} \right) = -c \quad (27)$$

If $\pi \neq \mu$, the strict concavity of log function implies that:

$$c = \mathbb{E}_{d^\mu} \log \left(\frac{\pi(a|s)}{\mu(a|s)} \right) < \log \mathbb{E}_{d^\mu} \left(\frac{\pi(a|s)}{\mu(a|s)} \right) = 0 \quad (28)$$

Thus $\frac{1}{T} \log \rho_{1:T} \xrightarrow{a.s.} c$ and $\rho_{1:T}^{1/T} \xrightarrow{a.s.} e^{-c}$. Since $r_t \leq 1$, $|\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t|^{1/T} \leq \rho_{1:T}^{1/T} T^{1/T}$. Since $T^{1/T} \rightarrow 1$, $\overline{\lim}_{T \rightarrow \infty} |\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t|^{1/T} < e^{-c}$. \square

D.2. Proof of Corollary 1

Proof. $\rho_{1:T} \xrightarrow{a.s.} 0$ directly follows from $\rho_{1:T}^{1/T} \xrightarrow{a.s.} e^{-c}$ in Theorem 3. For $\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t$, if there exist $\epsilon > 0$ such that $\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t > \epsilon$ for any T, then:

$$\overline{\lim}_{T \rightarrow \infty} \left| \rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t \right|^{1/T} \geq \overline{\lim}_{T \rightarrow \infty} \epsilon^{1/T} = 1$$

This contradicts $e^{-c} > \overline{\lim}_{T \rightarrow \infty} |\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t|^{1/T}$ So $\overline{\lim}_{T \rightarrow \infty} \rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t \leq 0$, which implies that $\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t \xrightarrow{a.s.} 0$. \square

D.3. Proof of Lemma 3

Proof. Let $f(s, a) = \log \frac{\pi(s, a)}{\mu(s, a)}$. According to Assumption 3, $|f(s, a)| < \infty$. Since $B(s, a) \geq 1$, $\frac{|f(s, a)|}{\sqrt{B(s, a)}} < \infty$. Since f^2 and B are both finite, $\mathbb{E}_{d^\mu} f^2 < \infty$ and $\mathbb{E}_{d^\mu} B < \infty$. Now we satisfy the condition of Lemma 3 in (Glynn and Olvera-Cravioto, 2019): in the proof of Lemma 3 in (Glynn and Olvera-Cravioto, 2019) they used their Assumption

165 i) Harris Chain, which is our Assumption 1, their Assumption vii) $\|f\|_{sqrTV}$ bounded (whic is satisfied by our bound
166 on B in Assumption 2), which is explained by f is bounded and $\sqrt{B} \geq 1$, and finally their assumption iv), which is
167 our assumption 2. The only difference is we assume a ‘‘petite’’ K which is a slight generalization of the ‘‘small’’ set K
168 (See discussion in (Meyn and Tweedie, 2012, Section 5)). The proof in Meyn and Glynn 1996 also used petite (which is
169 the part where Glynn and Olvera-Cravioto need assumption iv)). This assumption (drift condition) is often necessary for
170 quantitative analysis of general state Markov Chains. The geometric ergodicity for general state MC is also defined with a
171 petite/small set. By Thm 15.0.1 in Meyn and Tweedie the drift property is equivalent to geometric ergodicity. According
172 to Lemma 3 in (Glynn and Olvera-Cravioto, 2019), whose proof is similar with Theorem 2.3 in (Glynn et al., 1996), we
173 have that there exist a solution \hat{f} to the following Poisson’s equation:

$$174 \hat{f}(s, a) - \mathbb{E}_{\cdot|s,a} \hat{f}(s', a') = f(s, a) - \mathbb{E}_{d^\mu} f(s, a) \quad (29)$$

176 satisfying $|\hat{f}(s, a)| < c_1 \sqrt{B(s, a)}$ for some constant c_1 . Following from the Poisson’s equation we have:

$$177 \log \rho_{1:T} + Tc = \sum_{t=1}^T (f(s_t, a_t) - \mathbb{E}_{d^\mu} f(s, a)) \quad (30)$$

$$181 = \sum_{t=1}^T \left(\hat{f}(s_t, a_t) - \mathbb{E}_{s', a' | s_t, a_t} \hat{f}(s', a') \right) \quad (31)$$

$$184 = \hat{f}(s_1, a_1) - \hat{f}(s_{T+1}, a_{T+1}) + \sum_{t=2}^{T+1} \left(\hat{f}(s_t, a_t) - \mathbb{E}_{s', a' | s_{t-1}, a_{t-1}} \hat{f}(s', a') \right) \quad (32)$$

187 $\left(\hat{f}(s_t, a_t) - \mathbb{E}_{s', a' | s_{t-1}, a_{t-1}} \hat{f}(s', a') \right)$ are martingale differences. The absolute value of difference is upper bounded by
188 $2\|\hat{f}\|_\infty \leq 2c_1 \sqrt{\|B\|_\infty}$. \square

190 D.4. Proof of Theorem 4

192 **Lemma 1.** *If $\mathbb{E}_\mu[\rho^2|s] \leq M_\rho^2$ for any s , $\mathbb{E}[\rho_{0:k}^2] \leq M_\rho^{2k}$*

194 *Proof.*

$$196 \mathbb{E}[\rho_{0:t}^2] = \mathbb{E} \left[\prod_{i=1}^k \rho_i^2 \right] \quad (33)$$

$$199 = \mathbb{E} \left[\left(\prod_{i=1}^{k-1} \rho_i^2 \right) \mathbb{E}_{s_k, a_k} [\rho_k^2 | s_1, a_1, s_2, \dots, s_{k-1}, a_{k-1}] \right] \quad (34)$$

$$202 \leq \mathbb{E} \left[\left(\prod_{i=1}^{k-1} \rho_i^2 \right) M_\rho^2 \right] \quad (35)$$

$$205 = M_\rho^2 \mathbb{E} \left[\prod_{i=1}^{k-1} \rho_i^2 \right] \quad (36)$$

$$207 \dots \quad (37)$$

$$209 = M_\rho^{2k} \quad (38)$$

212 *Proof.* Define $Y = \rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t$ and $Z = \mathbf{1}(Y > v^\pi/2)$, then $v^\pi = \mathbb{E}(Y)$. By the law of total variance,

$$214 \text{Var}(Y) = \text{Var}(\mathbb{E}(Y|Z)) + \mathbb{E}(\text{Var}(Y|Z)) \quad (39)$$

$$215 \geq \text{Var}(\mathbb{E}(Y|Z)) \quad (40)$$

$$216 = \mathbb{E}(\mathbb{E}(Y|Z))^2 - (v^\pi)^2 \quad (41)$$

$$218 \geq \Pr(Y > v^\pi/2) (\mathbb{E}(Y|Y > v^\pi/2))^2 - (v^\pi)^2 \quad (42)$$

219

Now we are going to lower bound $\mathbb{E}(Y|Y > v^\pi/2)$. We can rewrite $\mathbb{E}(Y) = v^\pi$ as:

$$v^\pi = \mathbb{E}(Y) = \mathbb{E}(\mathbb{E}(Y|Z)) \quad (43)$$

$$= \Pr(Y > v^\pi/2)\mathbb{E}(Y|Y > v^\pi/2) + \Pr(Y \leq v^\pi/2)\mathbb{E}(Y|Y \leq v^\pi/2) \quad (44)$$

$$\leq \Pr(Y > v^\pi/2)\mathbb{E}(Y|Y > v^\pi/2) + 1 \times v^\pi/2 \quad (45)$$

So $\mathbb{E}(Y|Y > v^\pi/2) \geq \frac{v^\pi}{2\Pr(Y > v^\pi/2)}$. Substitute this into the RHS of Equation 42:

$$\text{Var}(Y) \geq \frac{(v^\pi)^2}{4\Pr(Y > v^\pi/2)} - (v^\pi)^2 \quad (46)$$

Now we are going to upper bound $\Pr(Y > v^\pi/2)$. Recall that we define $c = \mathbb{E}_{d^\mu} D_{\text{KL}}(\mu||\pi) = -\mathbb{E}_{d^\mu} \log\left(\frac{\pi(a|s)}{\mu(a|s)}\right)$. Now we define $c(T) = -\mathbb{E}_{d_{1:T}^\mu} \log\left(\frac{\pi(a|s)}{\mu(a|s)}\right) = -\frac{1}{T}\mathbb{E}_\mu[\log \rho_{1:T}]$.

$$\Pr(Y > v^\pi/2) \quad (47)$$

$$= \Pr(\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t > v^\pi/2) \leq \Pr(\rho_{1:T} T > v^\pi/2) \quad (48)$$

$$= \Pr\left(\rho_{1:T} > \frac{v^\pi}{2T}\right) \quad (49)$$

$$= \Pr(\log \rho_{1:T} > \log v^\pi - \log(2T)) \quad (50)$$

$$= \Pr\left(\frac{\log \rho_{1:T}}{T} > \frac{\log v^\pi - \log(2T)}{T}\right) \quad (51)$$

$$= \Pr\left(\frac{\log \rho_{1:T}}{T} + c + \frac{\hat{f}(s_{T+1}, a_{T+1}) - \hat{f}(s_1, a_1)}{T} > c + \frac{\log v^\pi - \log(2T) + \hat{f}(s_{T+1}, a_{T+1}) - \hat{f}(s_1, a_1)}{T}\right) \quad (52)$$

Since $\log v^\pi$ is a constant, $\hat{f}(s_{T+1}, a_{T+1}) - \hat{f}(s_1, a_1)$ could be upper bounded by constant $2c_1\sqrt{\|B\|_\infty}$, and $\lim_{T \rightarrow \infty} \frac{\log(2T)}{T} = 0$, we know that $\lim_{T \rightarrow \infty} \frac{\log v^\pi - \log(2T) + \hat{f}(s_{T+1}, a_{T+1}) - \hat{f}(s_1, a_1)}{T} = 0$. So there exists a constant $T_0 > 0$ such that for all $T > T_0$,

$$\frac{\log v^\pi - \log(2T) + \hat{f}(s_{T+1}, a_{T+1}) - \hat{f}(s_1, a_1)}{T} > -\frac{c}{2}$$

Therefore for all $T > T_0$:

$$\Pr(Y > v^\pi/2) \leq \Pr\left(\frac{\log \rho_{1:T}}{T} + c + \frac{\hat{f}(s_{T+1}, a_{T+1}) - \hat{f}(s_1, a_1)}{T} > c/2\right)$$

According to Lemma 3, and Azuma's inequality (Azuma, 1967), we have:

$$\Pr(Y > v^\pi/2) \leq \exp\left(\frac{-Tc^2}{8c_1^2\|B\|_\infty}\right)$$

Thus we can lower bound the variance of importance sampling estimator Y :

$$\text{Var}(Y) \geq \frac{(v^\pi)^2}{4} \exp\left(\frac{Tc^2}{8c_1^2\|B\|_\infty}\right) - (v^\pi)^2 \quad (53)$$

If the one step likelihood ratio is upper bounded by U_ρ , then the variance of importance sampling estimator can be upper bounded by:

$$\text{Var}(\hat{v}_{\text{IS}}) = \mathbb{E}[Y^2] - (v^\pi)^2 = \mathbb{E}\left[\rho_{0:T}^2 \left(\sum_{t=1}^T \gamma^{t-1} r_t\right)^2\right] - (v^\pi)^2 \quad (54)$$

$$\leq T^2 \mathbb{E}[\rho_{0:T}^2] - (v^\pi)^2 \quad (55)$$

$$\leq T^2 U_\rho^{2T} - (v^\pi)^2 \quad (56)$$

Following from lemma 1, the variance term can also be upper bounded by:

$$\text{Var}(\hat{v}_{\text{IS}}) = \mathbb{E}[Y^2] - (v^\pi)^2 = \mathbb{E} \left[\rho_{0:T}^2 \left(\sum_{t=1}^T \gamma^{t-1} r_t \right)^2 \right] - (v^\pi)^2 \quad (57)$$

$$\leq T^2 \mathbb{E} [\rho_{0:T}^2] - (v^\pi)^2 \quad (58)$$

$$\leq T^2 M_\rho^{2T} - (v^\pi)^2 \quad (59)$$

□

D.5. Proof of Theorem 5

Proof. Let $Y_t = \rho_{1:t} \gamma^{t-1} r_t$. For the upper bound:

$$\text{Var}(\hat{v}_{\text{PDIS}}) = \mathbb{E} \left(\left(\sum_{t=1}^T Y_t \right)^2 \right) - (v^\pi)^2 \quad (60)$$

$$\leq \mathbb{E} \left(T \sum_{t=1}^T Y_t^2 \right) - (v^\pi)^2 \quad (61)$$

$$= T \sum_{t=1}^T \mathbb{E}(Y_t^2) - (v^\pi)^2 \quad (62)$$

$$= T \sum_{t=1}^T \mathbb{E}(\rho_{0:t}^2 \gamma^{2t-2} r_t^2) - (v^\pi)^2 \quad (63)$$

$$\leq T \sum_{t=1}^T U_\rho^{2t} \gamma^{2t-2} \mathbb{E}_\mu[(r_t)^2] - (v^\pi)^2 \quad (64)$$

Or it can also be bounded as:

$$\text{Var}(\hat{v}_{\text{PDIS}}) \leq T \sum_{t=1}^T \mathbb{E}(\rho_{0:t}^2 \gamma^{2t-2} r_t^2) - (v^\pi)^2 \quad (65)$$

$$= T \sum_{t=1}^T \gamma^{2t-2} \mathbb{E}(\rho_{0:t}^2) - (v^\pi)^2 \quad (66)$$

$$\leq T \sum_{t=1}^T \gamma^{2t-2} M_\rho^{2t} - (v^\pi)^2 \quad (67)$$

The last step follows from lemma 1. For the lower bound, we notice that $Y_t \geq 0$ for any t , then:

$$\mathbb{E} \left(\left(\sum_{t=1}^T Y_t \right)^2 \right) \geq \mathbb{E} \left(\sum_{t=0}^T Y_t^2 \right) = \sum_{t=1}^T \mathbb{E}(Y_t^2) \quad (68)$$

For each t , we will follow a similar proof as how to lower bound part in Theorem 4:

$$\mathbb{E}(Y_t^2) = \mathbb{E} \left(\mathbb{E}(Y_t^2 | \mathbb{1}(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2)) \right) \quad (69)$$

$$\geq \mathbb{E} \left(\mathbb{E}(Y_t | \mathbb{1}(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2)) \right)^2 \quad (70)$$

$$\geq \Pr(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \left(\mathbb{E}(Y_t | Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \right)^2 \quad (71)$$

Notice that $\mathbb{E}(Y_t) = \gamma^{t-1} \mathbb{E}_\pi(r_t)$,

$$\gamma^{t-1} \mathbb{E}_\pi(r_t) = \mathbb{E}(Y_t) \quad (72)$$

$$= \Pr(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \mathbb{E}(Y_t | Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) + \Pr(Y_t \leq \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \mathbb{E}(Y_t | Y_t \leq \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \quad (73)$$

$$\leq \Pr(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \mathbb{E}(Y_t | Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) + \gamma^{t-1} \mathbb{E}_\pi(r_t)/2 \quad (74)$$

So we can lower bound the $\mathbb{E}(Y_t^2)$:

$$\mathbb{E}(Y_t | Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2) \geq \frac{\gamma^{t-1} \mathbb{E}_\pi(r_t)}{2 \Pr(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2)} \quad (75)$$

$$\mathbb{E}(Y_t^2) \geq \frac{\gamma^{2t-2} (\mathbb{E}_\pi(r_t))^2}{4 \Pr(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2)} \quad (76)$$

Now we are going to upper bound the tail probability $\Pr(Y_t > \gamma^{t-1} \mathbb{E}_\pi(r_t)/2)$:

$$\Pr\left(Y_t | Y_t > \frac{\gamma^{t-1} \mathbb{E}_\pi(r_t)}{2}\right) \quad (77)$$

$$= \Pr\left(\rho_{1:t} \gamma^{t-1} r_t > \frac{\gamma^{t-1} \mathbb{E}_\pi(r_t)}{2}\right) \quad (78)$$

$$\leq \Pr\left(\rho_{1:t} > \frac{\mathbb{E}_\pi(r_t)}{2}\right) \quad (79)$$

$$= \Pr(\log \rho_{1:t} > \log \mathbb{E}_\pi(r_t) - \log 2) \quad (80)$$

$$= \Pr\left(\frac{1}{t} \log \rho_{1:t} > \frac{\mathbb{E}_\pi(r_t) - \log 2}{t}\right) \quad (81)$$

$$= \Pr\left(\frac{1}{t} \log \rho_{1:t} + c + \frac{\hat{f}(s_{t+1}, a_{t+1}) - \hat{f}(s_1, a_1)}{T} > c + \frac{\mathbb{E}_\pi(r_t) - \log 2 + \hat{f}(s_{t+1}, a_{t+1}) - \hat{f}(s_1, a_1)}{t}\right) \quad (82)$$

Since $|\mathbb{E}_\pi(r_t) - \log 2 + \hat{f}(s_{t+1}, a_{t+1}) - \hat{f}(s_1, a_1)|$ is bounded, there exist some $T_0 > 0$ such that if $t > T_0$, we can lower bound the right hand side in the probability by $c/2$. Then for $t > T_0$, by Azuma's inequality (Azuma, 1967),

$$\Pr\left(Y_t | Y_t > \frac{\gamma^{t-1} \mathbb{E}_\pi(r_t)}{2}\right) \leq \Pr\left(\frac{\log \rho_{1:t}}{t} + c + \frac{\hat{f}(s_{t+1}, a_{t+1}) - \hat{f}(s_1, a_1)}{t} > \frac{c}{2}\right) \quad (83)$$

$$\leq \exp\left(\frac{-tc^2}{8c_1^2 \|B\|_\infty}\right) \quad (84)$$

So we have that for $t > T_0$:

$$\mathbb{E}(Y_t^2) \geq \frac{\gamma^{2t-2} \mathbb{E}_\pi(r_t)}{4} \exp\left(\frac{tc^2}{8c_1^2 \|B\|_\infty}\right)$$

For $0 < t \leq T_0$, $\mathbb{E}(Y_t^2) \geq 0$ completes the proof. \square

D.6. Proof of Corollary 2

Proof. First, $\gamma \geq \exp\left(\frac{-c^2}{16c_1^2 \|B\|_\infty}\right)$ indicate $\left(\frac{c^2}{8c_1^2 \|B\|_\infty} + 2 \log \gamma\right) > 0$. This is necessary for the second condition to hold since $r_t < 1$. The second condition $\mathbb{E}_\pi(r_t) = \Omega\left(\exp\left(\frac{-tc^2}{8c_1^2 \|B\|_\infty} - 2t \log \gamma + \epsilon t/2\right)\right)$ implies that there exist a $T_1 > 0$ and a constant $C > 0$ such that $(\mathbb{E}_\pi(r_t))^2 \geq C \left(\exp\left(\frac{-tc^2}{8c_1^2 \|B\|_\infty} - 2t \log \gamma + \epsilon t\right)\right)$, for any $t > T_1$. Then let $T > \max\{T_1, T_0\}$, where T_0 is the constant in Theorem 5:

$$\text{Var}\left(\sum_{t=T_0}^T \rho_{1:t} \gamma^{t-1} r_t\right) \geq \sum_{t=1}^T \frac{\gamma^{2t-2} (\mathbb{E}_\pi(r_t))^2}{4} \exp\left(\frac{tc^2}{8c_1^2 \|B\|_\infty}\right) - (v^\pi)^2 \quad (85)$$

$$\geq \frac{\gamma^{2T-2} (\mathbb{E}_\pi(r_T))^2}{4} \exp\left(\frac{Tc^2}{8c_1^2 \|B\|_\infty}\right) - (v^\pi)^2 \quad (86)$$

$$\geq \frac{\gamma^{-2} C}{4} \exp(\epsilon T) - (v^\pi)^2 = \Omega(\exp \epsilon T) \quad (87)$$

\square

D.7. Proof of Corollary 3

Proof. If $U_\rho\gamma \leq 1$, $U_\rho^t\gamma^{t-1}\mathbb{E}_\pi(r_t) \leq 1/\gamma$ for any t since $r_t \in [0, 1]$. If $U_\rho\gamma \lim (\mathbb{E}_\mu(r_T))^{1/T} < 1$, let $\delta = 1 - U_\rho\gamma \lim (\mathbb{E}_\mu(r_T))^{1/T} > 0$. There exist a $T_0 > 0$ such that for all $t > T_0$, $U_\rho\gamma(\mathbb{E}_\pi(r_t))^{1/t} \leq U_\rho\gamma(\lim (\mathbb{E}_\mu(r_T))^{1/T} + \delta/2(U_\rho\gamma)) = 1 - \delta/2 < 1$. Therefore in both case, for all $T > T_0$, $U_\rho^t\gamma^{t-1}\mathbb{E}_\mu(r_T) \leq 1/\gamma$:

$$\text{Var}\left(\sum_{t=1}^T \rho_{1:t}\gamma^{t-1}r_t\right) \leq T \sum_{t=1}^T U_\rho^t\gamma^{t-1}\mathbb{E}_\mu(r_T) \leq T \sum_{t=1}^{T_0} U_\rho^t\gamma^{t-1}\mathbb{E}_\mu(r_T) + T \sum_{t=T_0+1}^T U_\rho^t\gamma^{t-1}\mathbb{E}_\mu(r_T) \quad (88)$$

$$\leq TT_0 \frac{U_\rho^{T_0} - 1}{U_\rho - 1} + 2T^2 \frac{1}{\gamma} \quad (89)$$

Since T_0 is a constant, the variance is $O(T^2)$. \square

D.8. Proof of Theorem 6

Proof.

$$\text{Var}\left(\sum_{t=1}^T \frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t\right) \quad (90)$$

$$= \sum_{t=1}^T \text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t\right) + 2 \sum_{t < k} \text{Cov}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t, \frac{d_k^\pi(s_k, a_k)}{d_k^\mu(s_k, a_k)} \gamma^{k-1} r_k\right) \quad (91)$$

$$\leq \sum_{t=1}^T \text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t\right) + \sum_{t < k} 2 \sqrt{\text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t\right) \text{Var}\left(\frac{d_k^\pi(s_k, a_k)}{d_k^\mu(s_k, a_k)} \gamma^{k-1} r_k\right)} \quad (92)$$

$$\leq \sum_{t=1}^T \text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t\right) + \sum_{t < k} \left(\text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t\right) + \text{Var}\left(\frac{d_k^\pi(s_k, a_k)}{d_k^\mu(s_k, a_k)} \gamma^{k-1} r_k\right)\right) \quad (93)$$

$$= T \sum_{t=1}^T \gamma^{2t-2} \text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} r_t\right) \quad (94)$$

$$\leq T \sum_{t=1}^T \gamma^{2t-2} \text{Var}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)}\right) \quad (95)$$

$$= T \sum_{t=1}^T \gamma^{2t-2} \left(\mathbb{E}\left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)}\right)^2 - 1\right) \quad (96)$$

\square

D.9. Proof of Corollary 4

Lemma 4. If $d_t^\mu(s_t)$ and $d_t^\pi(s_t)$ are asymptotically equi-continuous, $\frac{d^\pi(s)}{d^\mu(s)} \leq U_s$, and $\frac{\pi(a|s)}{\mu(a|s)} \leq U_\rho$, then,

$$\lim_t \mathbb{E}_{s_t, a_t \sim d_t^\mu} \left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)}\right)^2 = \mathbb{E}_{s, a \sim d^\mu} \left(\frac{d^\pi(s, a)}{d^\mu(s, a)}\right)^2$$

Proof. According to the law of large number on Markov chain (Breiman, 1960), the distribution of d_t^μ converge to the stationary distribution d^μ in distribution. By the Lemma 1 in (Boos et al., 1985), $d_t^\mu(s, a)$ converge to $d^\mu(s, a)$ pointwisely, $d_t^\pi(s, a)$ converge to $d^\pi(s, a)$ pointwisely. So $\frac{d_t^\pi(s)}{d_t^\mu(s)}$ converge to $\frac{d^\pi(s)}{d^\mu(s)}$ pointwisely.

$$\mathbb{E}_{s_t, a_t \sim d_t^\mu} \left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right)^2 - \mathbb{E}_{s, a \sim d^\mu} \left(\frac{d^\pi(s, a)}{d^\mu(s, a)} \right)^2 \quad (97)$$

$$= \int_{s, a} \frac{(d_t^\pi(s, a))^2}{d_t^\mu(s, a)} \mathbf{d}s \mathbf{d}a - \int_{s, a} \frac{(d^\pi(s, a))^2}{d^\mu(s, a)} \mathbf{d}s \mathbf{d}a \quad (98)$$

$$= \int_{s, a} \frac{(d_t^\pi(s))^2 (\pi(a|s))^2}{d_t^\mu(s) \mu(a|s)} - \frac{(d^\pi(s))^2 (\pi(a|s))^2}{d^\mu(s) \mu(a|s)} \mathbf{d}s \mathbf{d}a \quad (99)$$

$$\leq U_\rho \int_{s, a} \left| \frac{(d_t^\pi(s))^2}{d_t^\mu(s)} - \frac{(d^\pi(s))^2}{d^\mu(s)} \right| \mathbf{d}s \mathbf{d}a \quad (100)$$

$$\leq U_\rho \int_{s, a} \left| \frac{d^\pi(s) (d_t^\pi(s) - d^\pi(s))}{d^\mu(s)} + d_t^\pi(s) \frac{d_t^\pi(s)}{d_t^\mu(s)} - d_t^\pi(s) \frac{d^\pi(s)}{d^\mu(s)} \right| \mathbf{d}s \mathbf{d}a \quad (101)$$

$$\leq U_\rho \int_{s, a} \left| \frac{d^\pi(s) (d_t^\pi(s) - d^\pi(s))}{d^\mu(s)} \right| + d_t^\pi(s) \left| \frac{d_t^\pi(s)}{d_t^\mu(s)} - \frac{d^\pi(s)}{d^\mu(s)} \right| \mathbf{d}s \mathbf{d}a \quad (102)$$

$$\leq U_\rho U_s d_{\text{TV}}(d_t^\pi, d^\pi) + U_\rho \int_{s, a} \left| \frac{d_t^\pi(s)}{d_t^\mu(s)} - \frac{d^\pi(s)}{d^\mu(s)} \right| \mathbf{d}s \mathbf{d}a \quad (103)$$

By the law of large number on Markov chain (Breiman, 1960), $d_{\text{TV}}(d_t^\pi, d^\pi) \rightarrow 0$. Since U_ρ and U_s are constant, and $\frac{d_t^\pi(s)}{d_t^\mu(s)} \rightarrow \frac{d^\pi(s)}{d^\mu(s)}$, the right hand side of equation above converge to zero, which completes the proof. \square

Proof of Corollary 4:

Proof. Since $\frac{d^\pi(s, a)}{d^\mu(s, a)}$ is bounded by $U_\rho U_s$ and then $\mathbb{E}_{s, a \sim d^\mu} \left(\frac{d^\pi(s, a)}{d^\mu(s, a)} \right)^2$ is bounded by $U_\rho^2 U_s^2$. Following from Lemma 4, there exist $T_0 > 0$ such that for all $t > T_0$, $\mathbb{E}_{s_t, a_t \sim d_t^\mu} \left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right)^2 \leq 2\mathbb{E}_{s, a \sim d^\mu} \left(\frac{d^\pi(s, a)}{d^\mu(s, a)} \right)^2 \leq 2U_\rho^2 U_s^2$. Then by Theorem 6, for $T > T_0$

$$\begin{aligned} \text{Var} \left(\sum_{t=1}^T \frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \gamma^{t-1} r_t \right) &\leq T \sum_{t=1}^T \gamma^{t-1} \mathbb{E} \left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right)^2 \\ &\leq T \sum_{t=1}^{T_0} \gamma^{t-1} \mathbb{E} \left(\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right)^2 + 2T(T - T_0)U_\rho^2 U_s^2 \\ &= O(T^2) \end{aligned}$$

\square

D.10. Proof of Corollary 5

Now we consider an type of approximate MIS estimators, which plug an approximate density ratio into the MIS estimator. More specifically, we consider it use a function $w_t(s_t, a_t)$ to approximate density ratio $\frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)}$, and construct the estimator as:

$$\hat{v}_{\text{ASIS}} = \sum_{t=1}^T w_t(s, a) \gamma^{t-1} r_t \quad (104)$$

This approximate MIS estimator is often biased based on the choice of $w_t(s, a)$, so we consider the upper bound of their mean square error with respect to T and the error of the ratio estimator.

Theorem 7. \hat{v}_{ASIS} with w_t such that where $\mathbb{E}_\mu \left(w_t(s_t, a_t) - \frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right)^2 \leq \epsilon_w$

$$\text{MSE}(\hat{v}_{\text{ASIS}}) \leq 2\text{Var}(\hat{v}_{\text{SIS}}) + 2T^2 \epsilon_w \quad (105)$$

495 *Proof.*

$$496 \text{MSE} \left(\sum_{t=1}^T w_t(s_t, a_t) \gamma^{t-1} r_t \right) = \mathbb{E} \left(\sum_{t=1}^T w_t(s_t, a_t) \gamma^{t-1} r_t - v^\pi \right)^2 \quad (106)$$

$$497 = \mathbb{E} \left(\sum_{t=1}^T w_t(s_t, a_t) \gamma^{t-1} r_t - \hat{v}_{\text{SIS}} + \hat{v}_{\text{SIS}} - v^\pi \right)^2 \quad (107)$$

$$500 \leq 2\mathbb{E} \left(\sum_{t=1}^T w_t(s_t, a_t) \gamma^{t-1} r_t - \hat{v}_{\text{SIS}} \right)^2 + 2\mathbb{E} (\hat{v}_{\text{SIS}} - v^\pi)^2 \quad (108)$$

$$503 \leq 2\mathbb{E} \left(\sum_{t=1}^T \gamma^{t-1} r_t \left(w_t(s_t, a_t) - \frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right) \right)^2 + 2\text{Var}(\hat{v}_{\text{SIS}}) \quad (109)$$

$$506 \leq 2T \sum_{t=1}^T \gamma^{2t-2} \mathbb{E} \left(w_t(s_t, a_t) - \frac{d_t^\pi(s_t, a_t)}{d_t^\mu(s_t, a_t)} \right)^2 + 2\text{Var}(\hat{v}_{\text{SIS}}) \quad (110)$$

$$511 \leq 2\text{Var}(\hat{v}_{\text{SIS}}) + 2T \sum_{t=1}^T \gamma^{2t-2} \epsilon_w \quad (111)$$

$$512 \leq 2\text{Var}(\hat{v}_{\text{SIS}}) + 2T^2 \epsilon_w \quad (112)$$

513 □

514 **Proof of Corollary 5:**

515 *Proof.* By Theorem 7 we have that the MSE is bounded by

$$516 2\text{Var}(\hat{v}_{\text{SIS}}) + 2T^2 \epsilon_w \quad (113)$$

517 According to Corollary 4:

$$518 2\text{Var}(\hat{v}_{\text{SIS}}) + 2T^2 \epsilon_w = O(T^2) + 2T^2 \epsilon_w = O(T^2(1 + \epsilon_w)) \quad (114)$$

519 □

520 **E. Return-Conditional IS estimators**

521 A natural extension of the conditional importance sampling estimators is to condition on the observed returns G_t . Precisely we examine the general conditional importance sampling estimator:

$$522 G_t \mathbb{E} [\rho_{1:t} | \phi_t] \quad , \quad (115)$$

523 and consider when $\phi_t = G_t$. An analytic expression for $\mathbb{E} [\rho_{1:t} | G_t]$ is not available, but we can model this as a regression problem to predict $\mathbb{E} [\rho_{1:t} | G_t]$ given an input G_t . A natural approach is to use ordinary least squares (OLS) estimator to estimate $\mathbb{E} [\rho_{1:t} | \phi_t]$ viewing ϕ_t (or any other statistics G_t) as an *input* and $\rho_{1:t}$ as an *output*. While tempting at first glance, we show that this approach produces exactly the same estimates of the expected return as that of the crude importance sampling estimator.

524 We start by considering the OLS problem associated with the conditional weights in which we want to find a $\hat{\theta}$ such that $\phi_t^\top \hat{\theta} \approx \mathbb{E} [\rho_{1:t} | \phi_t]$. Let $\Phi \in \mathbb{R}^{n \times 2}$ be the *design* matrix containing the observed returns $G_t^{(i)}$ after t steps and $Y \in \mathbb{R}^n$ be the vector of importance ratios $\rho_{1:t}^{(i)}$ for each rollout i :

$$525 Y = \begin{bmatrix} \rho_t^{(0)} \\ \vdots \\ \rho_t^{(N)} \end{bmatrix}, \quad \Phi = \begin{bmatrix} G_t^{(0)} & 1 \\ \vdots & \vdots \\ G_t^{(N)} & 1 \end{bmatrix} .$$

526

The OLS estimator for the return-conditional weights is then $\hat{Y} = \Phi \hat{\theta}$ and where $\hat{\theta} \in \mathbb{R}^2$ is defined as:

$$\hat{\theta} = (\Phi^\top \Phi)^{-1} \Phi^\top Y .$$

We can now use the approximate return-conditional weights to form a Monte Carlo estimate of the expected return under the target policy:

$$\hat{v}_{\text{RCIS}} \equiv \frac{1}{N} \sum_{i=0}^N G_t^{(i)} \hat{Y}^{(i)} = \frac{1}{N} [1, 0] \Phi^\top \hat{Y} , \quad (116)$$

where $\hat{Y}^{(i)} = [G_t^{(i)}, 0]^\top \hat{\theta}$ and the equality follows from the fact that $\Phi^\top Y = [\sum_{i=1}^n \rho_t^{(i)} G_t^{(i)}, \sum_{i=1}^n \rho_t^{(i)}]^\top$. Using this observation, we can also express the crude importance sampling estimator with the linear combination $\Phi^\top Y$, where Y now consists of the *true* weights:

$$\hat{v}_{\text{IS}} \equiv \frac{1}{N} [1, 0] \Phi^\top Y . \quad (117)$$

Note that equation (116) differs from (117) only in the term $\hat{Y} = \Phi \hat{\theta} = \Phi (\Phi^\top \Phi)^{-1} \Phi^\top Y$ and upon closer inspection, we find that:

$$\Phi^\top \hat{e} = \Phi^\top Y - \Phi^\top \hat{Y} = \Phi^\top (Y - \Phi (\Phi^\top \Phi)^{-1} \Phi^\top Y) = \Phi^\top (I - H) Y = \mathbf{0} ,$$

where \hat{e} is residual vector $Y - \hat{Y}$ and $H = \Phi (\Phi^\top \Phi)^{-1} \Phi^\top$ is the *hat* matrix. Hence, it follows that the estimate of the expected return made under the crude importance sampling estimator must be identical to the extended estimator which uses approximate return-conditional weights:

$$\hat{v}_{\text{IS}} - \hat{v}_{\text{RCIS}} = \frac{1}{n} [1, 0] \Phi^\top Y - \frac{1}{n} [1, 0] \Phi^\top \hat{Y} = \frac{1}{n} [1, 0] (\Phi^\top Y - \Phi^\top \hat{Y}) = [1, 0] \mathbf{0} = 0 .$$

This analysis can be generalized to any conditional importance sampling estimator for which G_t can be expressed as a linear combination of ϕ_t . For example, rather than conditioning on the final return, we could condition on the return so far (the sum of returns to the present) and use $\phi_t = [r_1, r_2, \dots, r_t]$ with the coefficient vector $[1, 1, \dots, 1, 0]$. Similarly, this negative result carries to reward-conditional weights if the immediate reward r_t can be expressed as linear combination of ϕ_t , including if ϕ_t is simply the immediate reward.

References

- Kazuoki Azuma. Weighted sums of certain dependent random variables. *Tohoku Mathematical Journal, Second Series*, 19(3):357–367, 1967.
- Dennis D Boos et al. A converse to scheffe’s theorem. *The Annals of Statistics*, 13(1):423–427, 1985.
- Leo Breiman. The strong law of large numbers for a class of markov chains. *The Annals of Mathematical Statistics*, 31(3):801–803, 1960.
- Peter W. Glynn and Mariana Olvera-Cravioto. Likelihood ratio gradient estimation for steady-state parameters. *Stochastic Systems*, 9(2):83–100, June 2019.
- Peter W Glynn, Sean P Meyn, et al. A liapounov bound for solutions of the poisson equation. *The Annals of Probability*, 24(2):916–931, 1996.
- Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.