

## SUPPLEMENTARY MATERIAL

### A UNFAIRNESS OF IMPORTANCE SAMPLING

Suppose we want to use importance sampling to select the better of two policies,  $\pi_e$  and  $\pi_b$ , where we have prior data collected from  $\pi_b$ , in a MAB with two actions  $a_1$  and  $a_2$ , with rewards and probabilities as described in Table 3. Notice that  $V^{\pi_b} = p + (1-p)r$  and  $V^{\pi_e} = 1$ , so a fair policy selection algorithm should choose  $\pi_e$  at least half the time since  $p + (1-p)r < 1$ . If we draw only a single sample from  $\pi_b$ , we get that with probability  $1-p$ , IS would select  $\pi_b$  over  $\pi_e$ . Thus as long as  $p < 0.5$ , IS will be unfair. Furthermore, notice that as we decrease  $p$ , the gap between the performance of the policies increases, yet the probability that IS chooses the right policy only decreases!

Now suppose we draw  $n$  samples from  $\pi_b$ . Notice that as long as  $\tau_2$  is never sampled, IS will choose  $\pi_b$ , since in that case  $\hat{V}_{\text{IS}}^{\pi_3} = 0$ .  $\pi_b$  will never sample  $\tau_2$  with probability  $(1-p)^n$ . Thus IS is unfair as long as  $(1-p)^n \geq 0.5$ , or as long as  $p \leq 1 - 0.5^{1/n} \approx \ln(2)/n$  for large  $n$ .

It may appear that this unfairness is not a big problem when we have a reasonable number of samples, but the practical significance of this problem becomes more pronounced in more realistic domains where we have a large number of possible trajectories, or equivalently, a long horizon. For example consider a domain where there are only two actions and the agent must take 50 sequential actions and receives a reward only at the end of a trajectory. Furthermore, consider that the only valuable trajectory is to take a particular action for the entire trajectory (analogous to  $a_2$  above). In this case, we would need over  $10^{14}$  samples just to get IS to be fair!

Table 3: Domain in Supplementary Material A. Rewards are deterministic. The bottom two rows give the probability distributions for  $\pi_e$  and  $\pi_b$  over the two actions.

	$a_1$ ( $R = r < 1$ )	$a_2$ ( $R = 1$ )
$\pi_e$	0	1
$\pi_b$	$1-p$	$p$

### B NON-TRANSITIVITY

**Theorem B.1** (Non-Transitivity of  $\succ_{\text{MC},n}$ ). *The relation induced by  $\succ_{\text{MC},n}$  is non-transitive. Specifically, there exists policies  $\pi_1$ ,  $\pi_2$ , and  $\pi_3$  where*  

$$\Pr(\hat{V}_{\text{MC},n}^{\pi_1} > \hat{V}_{\text{MC},n}^{\pi_2}) = \Pr(\hat{V}_{\text{MC},n}^{\pi_2} > \hat{V}_{\text{MC},n}^{\pi_3}) =$$

$\Pr(\hat{V}_{\text{MC},n}^{\pi_3} > \hat{V}_{\text{MC},n}^{\pi_1}) = 1 - \phi \approx 0.618$  where  $\phi = \frac{\sqrt{5}+1}{2}$  is the golden ratio. Moreover, it is possible that for any policy  $\pi$ , there is another policy  $\pi'$  where  $\succ_{\text{MC},n}(\pi', \pi) = \text{True}$ .

*Proof of Theorem B.1.* Consider a multi-armed bandit where there are three actions:  $a_1$ , which gives a reward of  $n+1$  with probability  $p$  and a reward of 0 with probability  $1-p$ ,  $a_2$  which always gives a reward of 1, and  $a_3$  which gives a reward of  $n^2 + n + 1$  with probability  $1-q$  and a reward of 0.5 with probability  $q$ . Suppose policies  $\pi_1$ ,  $\pi_2$ , and  $\pi_3$  always choose action  $a_1$ ,  $a_2$ , and  $a_3$  respectively. Now suppose we want to estimate the three policies with  $n$  on-policy samples from each. We have that  $\pi_1$  gives a higher reward than  $\pi_2$  whenever we get the large reward at least once, which happens with probability  $1 - (1-p)^n$ . Thus

$$\Pr(\hat{V}_{\text{MC},n}^{\pi_1} > \hat{V}_{\text{MC},n}^{\pi_2}) = 1 - (1-p)^n$$

Furthermore, clearly  $\pi_2$  gives a larger reward than  $\pi_3$  whenever all samples of  $\pi_2$  give a reward of 0.5, which happens with probability  $q^n$ . Now, finally we see that  $\pi_3$  gives a larger reward than  $\pi_1$  whenever it gives at least one sample with a large reward or when both of them give only samples of their small rewards, which happens with probability  $(1-q^n) + q^n(1-p)^n$ , so

$$\Pr(\hat{V}_{\text{MC},n}^{\pi_3} > \hat{V}_{\text{MC},n}^{\pi_1}) = (1-q^n) + q^n(1-p)^n$$

Now let  $p = 1 - (2-\phi)^{1/n}$  and  $q = (\phi-1)^{1/n}$ , where  $\phi = \frac{\sqrt{5}+1}{2} \approx 1.618$  is the golden ratio. Thus we have that:

$$\Pr(\hat{V}_{\text{MC},n}^{\pi_1} > \hat{V}_{\text{MC},n}^{\pi_2}) = \phi - 1$$

$$\Pr(\hat{V}_{\text{MC},n}^{\pi_2} > \hat{V}_{\text{MC},n}^{\pi_3}) = \phi - 1$$

$$\Pr(\hat{V}_{\text{MC},n}^{\pi_3} > \hat{V}_{\text{MC},n}^{\pi_1}) = (2-\phi) + (\phi-1)(2-\phi) = \phi - 1$$

We now show that for this multi-armed bandit, there is no optimal policy with respect to  $\succ_{\text{MC},n}$ . A policy in this setting is simply a distribution over  $a_1$ ,  $a_2$ , and  $a_3$ . Equivalently, we can view any policy as a mix of the policies  $\pi_1$ ,  $\pi_2$ , and  $\pi_3$ . Suppose a policy  $\pi$  executes  $\pi_1$  with probability  $p$ ,  $\pi_2$  with probability  $q$ , and  $\pi_3$  with probability  $r$ . If  $p$  is the largest of the probabilities, then notice that  $\Pr(\hat{V}_{\text{MC},n}^{\pi_3} > \hat{V}_{\text{MC},n}^{\pi}) = p(\phi-1) \geq q(\phi-1) = \Pr(\hat{V}_{\text{MC},n}^{\pi} > \hat{V}_{\text{MC},n}^{\pi_3})$ . We can make a similar argument if  $q$  or  $r$  are the largest probabilities. Thus, there is no optimal policy.  $\square$

### C FAIRNESS PROOFS

**Theorem 6.1.** *Using the on-policy Monte Carlo estimator for policy selection when we have  $n$  samples*

from each of policies  $\pi_1$  and  $\pi_2$  is fair provided that  $V_{\text{Max}}^{\pi_1} + V_{\text{Max}}^{\pi_2} \leq |V^{\pi_1} - V^{\pi_2}| \sqrt{\frac{2n}{\ln 2}}$ . We can guarantee strict fairness if the inequality above is strict.

*Proof of Theorem 6.1.* Suppose without loss of generality that  $V^{\pi_1} > V^{\pi_2}$ . Let  $\hat{V}_i^\pi = \sum_{t=1}^{T_i} R_{i,t}$  (i.e., the estimate of the value of policy  $\pi$  using only  $\tau_i^\pi$ ). Now let

$$X_i = \hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2}$$

Note that the range of  $X_i$  is  $[-V_{\text{Max}}^{\pi_2}, V_{\text{Max}}^{\pi_1}]$ . Let  $\omega$  be the difference between the upper and lower bounds of  $X_i$ , that is,  $\omega = V_{\text{Max}}^{\pi_1} + V_{\text{Max}}^{\pi_2}$ . Because all  $\tau_i^{\pi_1}$  and  $\tau_i^{\pi_2}$  are independent of  $\tau_j^{\pi_1}$  and  $\tau_j^{\pi_2}$  for all  $i \neq j$ , we know that  $X_i$  is independent of  $X_j$  for all  $i \neq j$ . Thus we can use Hoeffding's inequality to find that:

$$\begin{aligned} \Pr(\bar{X} \leq 0) &= \Pr(\bar{X} - \mathbf{E}[\bar{X}] \leq -\mathbf{E}[\bar{X}]) \\ &\leq \exp\left(\frac{-2n\mathbf{E}[\bar{X}]^2}{\omega^2}\right) \end{aligned}$$

Note that

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n \hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2} = \hat{V}_{\text{MC},n}^{\pi_1} - \hat{V}_{\text{MC},n}^{\pi_2}$$

and  $\mathbf{E}[\bar{X}] = V^{\pi_1} - V^{\pi_2}$ . Thus, if we want to guarantee

$$\Pr\left(\hat{V}_{\text{MC},n}^{\pi_1} - \hat{V}_{\text{MC},n}^{\pi_2} \leq 0\right) \leq \delta$$

we can simply guarantee

$$\exp\left(\frac{-2n(V^{\pi_1} - V^{\pi_2})^2}{\omega^2}\right) \leq \delta$$

Solving for  $\omega$ , we must have that:

$$\omega \leq (V^{\pi_1} - V^{\pi_2}) \sqrt{\frac{2n}{\ln(1/\delta)}}$$

Substituting  $\delta = 0.5$ , we can thus guarantee that  $\Pr\left(\hat{V}_{\text{MC},n}^{\pi_1} - \hat{V}_{\text{MC},n}^{\pi_2} > 0\right) \geq 0.5$ , which guarantees fairness when  $V^{\pi_1} > V^{\pi_2}$ . Since we do not actually know which policy has a greater value, we can guarantee fairness with the following condition:

$$\omega \leq |V^{\pi_1} - V^{\pi_2}| \sqrt{\frac{2n}{\ln 2}}$$

□

**Theorem 6.2.** *Using importance sampling for policy selection when we have  $n$  samples from the behavior policy is fair with respect to  $\succ_V$ , provided that*

$$w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2} \leq |V^{\pi_1} - V^{\pi_2}| \sqrt{\frac{2n}{\ln 2}}$$

*We can guarantee strict fairness if the inequality above is strict.*

*Proof of Theorem 6.2.* Let  $\hat{V}_i^\pi = \sum_{t=1}^{T_i} w_i^{\pi_e} R_{i,t}$  (i.e., the estimate of the value of policy  $\pi$  using only  $\tau_i$ ). Now let

$$X_i = \hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2}$$

Note that the range of  $X_i$  is  $[-w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}, w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1}]$ . Let  $\omega$  be the difference between the upper and lower bounds of  $X_i$ , that is,  $\omega = w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}$ . Because  $\tau_i$  and  $\tau_j$  are independent for all  $i \neq j$ , we know that  $X_i$  is independent of  $X_j$  for all  $i \neq j$ . The rest of the proof follows exactly as in the proof of Theorem 6.1. □

**Theorem 6.3.** *For any two policies  $\pi_1$  and  $\pi_2$ , behavior policy  $\pi_b$ , and for all  $k \in \{1, 2, 3, \dots\}$ , using importance sampling for policy selection when we have  $n$  samples from the behavior policy is fair with respect to  $\succ_{\text{MC},kn}$  provided that there exists  $\epsilon > 0$  and  $\delta < 0.5$  such that  $\Pr\left(|\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2}| \geq \epsilon\right) \geq 1 - \delta$  and  $(w_{\text{Max}}^{\pi_1} + 1)V_{\text{Max}}^{\pi_1} + (w_{\text{Max}}^{\pi_2} + 1)V_{\text{Max}}^{\pi_2} \leq \epsilon \sqrt{\frac{2n}{\ln \frac{1-\delta}{0.5-\delta}}}$ . Importance sampling in this setting is transitively fair, provided that  $\delta \leq 0.25$ .*

*Proof of Theorem 6.3.* Suppose without loss of generality that  $\Pr\left(\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} \geq \epsilon\right) \geq 1 - \delta$ . Recall that IS uses trajectories  $\tau_1, \dots, \tau_n \sim \pi_b$ . Consider additional random samples  $\tau_1^{\pi_1}, \dots, \tau_{kn}^{\pi_1} \sim \pi_1$  and  $\tau_1^{\pi_2}, \dots, \tau_{kn}^{\pi_2} \sim \pi_2$ . Note that these samples are all independent from each other. For  $i \in \{1, 2, \dots, n\}$ , let

$$\hat{V}_i^\pi = \frac{1}{k} \sum_{j=1}^k \sum_{t=1}^{T_{j,i}} R_{j,i,t}$$

(i.e., the estimate of the value of policy  $\pi$  using only samples  $\tau_i^\pi, \tau_{2i}^\pi, \dots, \tau_{ki}^\pi$ ). Furthermore let  $\hat{V}_{\text{IS},i}^\pi = \sum_{t=1}^{T_i} w_{i,t} R_{i,t}$ . Now let

$$X_i = (\hat{V}_{\text{IS},i}^{\pi_1} - \hat{V}_{\text{IS},i}^{\pi_2}) - (\hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2})$$

Notice that the range of  $X_i$  is  $[-V_{\text{Max}}^{\pi_2} - w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1}, V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}]$ . Let  $\omega$  be the difference between the upper and lower bounds of  $X_i$ , that is,  $\omega = (w_{\text{Max}}^{\pi_1} + 1)V_{\text{Max}}^{\pi_1} + (w_{\text{Max}}^{\pi_2} + 1)V_{\text{Max}}^{\pi_2}$ .

Notice that

$$\begin{aligned} \bar{X} &= \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n (\hat{V}_{\text{IS},i}^{\pi_1} - \hat{V}_{\text{IS},i}^{\pi_2}) - (\hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2}) \\ &= (\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2}) - (\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2}) \end{aligned}$$

and

$$\mathbf{E}[\bar{X}] = (V^{\pi_1} - V^{\pi_2}) - (V^{\pi_1} - V^{\pi_2}) = 0$$

Thus we have that

$$\Pr\left(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} \leq 0\right) = \Pr\left(\bar{X} \leq -(\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2})\right)$$

$$= \Pr\left(\bar{X} - \mathbf{E}[\bar{X}] \leq -(\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2})\right)$$

Thus, we can use Hoeffding's inequality to find that:

$$\Pr\left(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} \leq 0 \mid \hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} \geq \epsilon\right) \leq \exp\left(\frac{-2n\epsilon^2}{\omega^2}\right)$$

So if we want to guaranteee

$$\Pr\left(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} \leq 0\right) \leq \gamma$$

we can simply guaranteee

$$\exp\left(\frac{-2n(\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2})^2}{\omega^2}\right) \leq \gamma$$

Solving for  $\omega$ , we must have that:

$$\omega \leq (\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2}) \sqrt{\frac{2n}{\ln(1/\gamma)}}$$

Notice that

$$\begin{aligned} & \Pr\left(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} \geq 0\right) \\ & \geq \Pr\left(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} \leq 0 \mid \hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} \geq \epsilon\right) \\ & \quad \times \Pr\left(\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} \geq \epsilon\right) \\ & \geq (1-\gamma)(1-\delta) \end{aligned}$$

If we set  $\gamma = \frac{0.5-\delta}{1-\delta}$ , we have that

$$\Pr\left(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} \geq 0\right) \geq 0.5$$

which is what we want.

As long as  $\delta \leq 0.25$ , we have that the IS is transitively fairness since any fair algorithm with respect to  $\succ_{\text{MC},kn}$  is transitively fair whenever  $\Pr(|\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2}| > 0) \geq 0.75$   $\square$

## D MULTIPLE COMPARISONS

Here we present an algorithm that uses pairwise comparisons to select amongst  $k \geq 2$  policies (Algorithm 3). This algorithm can take as input either of the off-policy fair policy selection algorithms above (or some variant thereof).

**Theorem D.1.** *For any finite set of  $k$  policies  $\Pi$ , behavior policy  $\pi_b$ ,  $p = 0.5$ , and fair off-policy policy selection algorithm FPS, Algorithm 3 is a strictly fair policy selection algorithm with when we have  $n$  samples drawn from  $\pi_b$ .*

---

### Algorithm 3 Off-Policy FPS for $k$ policies

---

```

input  $\Pi, V_{\text{Max}}^{\Pi}, \epsilon, p$ 
       $\mathcal{T} = \{\tau_1, \tau_2, \dots, \tau_n \sim \pi_b\}, \text{FPS}$ 
       $\delta \leftarrow (1-p)/(2k+3)$ 
       $\pi^* \leftarrow \Pi.\text{next}$ 
      Eliminated  $\leftarrow \emptyset$ 
      CurrBeat  $\leftarrow \emptyset$ 
      repeat
         $\pi' \leftarrow (\Pi \setminus \text{CurrBeat}).\text{next}$ 
        winner  $\leftarrow \text{FPS}(\pi^*, \pi', V_{\text{Max}}^{\pi^*}, V_{\text{Max}}^{\pi'}, \epsilon, \delta, \mathcal{T})$ 
        if winner ==  $\pi^*$  then
          Eliminated  $\leftarrow \text{Eliminated} \cup \{\pi'\}$ 
          CurrBeat  $\leftarrow \text{CurrBeat} \cup \{\pi'\}$ 
        else if winner ==  $\pi'$  then
           $\pi^* \leftarrow \pi'$ 
          Eliminated  $\leftarrow \text{Eliminated} \cup \{\pi^*\}$ 
          CurrBeat  $\leftarrow \text{CurrBeat} \cup \{\pi^*\}$ 
        else
           $\pi^* \leftarrow (\Pi \setminus \text{Eliminated}).\text{next}$ 
          Eliminated  $\leftarrow \text{Eliminated} \cup \{\pi^*, \pi'\}$ 
          CurrBeat  $\leftarrow \emptyset$ 
        end if
      until len(Eliminated) ==  $k-1$  or len(CurrBeat) ==  $k$ 
      if len(Eliminated) ==  $k-1$  then
        return  $\pi^*$ 
      else
        return No Fair Comparison
      end if

```

---

*Proof of Theorem D.1.* The algorithm essentially applies an algorithm for finding the maximum element of a set, with the exception that whenever it cannot make a fair comparison between two policies, it will eliminate both of those policies from consideration of being better than all other policies with respect to the better-than function. The algorithm must return `No Fair Comparison` if and only if every policy is eliminated. Notice that we only eliminate a policy when it is not returned by `FPS` or when `No Fair Comparison` is returned, which is correct. Notice that until there are  $k - 1$  policies that are eliminated, at every comparison at least one policy is eliminated and the last of those comparisons must include the only remaining non-eliminated policy. Afterwards it takes at most  $k - 2$  comparisons with the final policy (comparing it to every other policy other than the one it was already compared to) to determine that no fair comparison is possible, making a total of  $2k - 3$  comparisons.

On the other hand, the algorithm must return a policy when that policy is outputted by `FPS` from comparisons with every other policy, which is exactly what it does (i.e. when the `CurrBeated` set includes  $k - 1$  policies). The maximum number of comparisons it takes to output a policy is the number of comparisons it takes to eliminate  $k - 1$  other policies plus the number of comparisons it takes to beat  $k - 2$  policies using the same argument as above, making a total of  $2k - 3$  comparisons. Thus, if we let  $\delta = (1 - p)/(2k - 3)$  all of the comparisons made by `FPS` will simultaneously hold with probability  $1 - (2k - 3)\delta = p$ , and so setting  $p = 0.5$  ensures fairness.  $\square$

## E SAFETY THEOREMS AND PROOFS

In this section, we will prove the theorems for ensuring safety when using importance sampling for policy selection.

**Theorem 6.4.** *For any two policies  $\pi_1$  and  $\pi_2$ , behavior policy  $\pi_b$ ,  $\omega = w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}$ , and  $\delta \leq 0.5$ , Algorithm 2 is a safe policy selection algorithm with respect to  $\succ_V$  with probability  $1 - \delta$ .*

*Proof of Theorem 6.4.* Let  $\hat{V}_i^\pi = \sum_{t=1}^{T_i} w_i^{\pi_e} R_{i,t}$  (i.e., the estimate of the value of policy  $\pi$  using only  $\tau_i$ ). Now let

$$X_i = \hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2}$$

Note that the range of  $X_i$  is  $[-w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}, w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1}]$ . Because  $\tau_i$  and  $\tau_j$  are independent for all  $i \neq j$ , we know that  $X_i$  is independent of  $X_j$  for all  $i \neq j$ . Thus we can

use Hoeffding's inequality to find that:

$$\Pr \left( \bar{X} - \mathbf{E} [\bar{X}] \geq -\omega \sqrt{\frac{\ln(1/\gamma)}{2n}} \right) \geq 1 - \gamma$$

and

$$\Pr \left( \bar{X} - \mathbf{E} [\bar{X}] \leq \omega \sqrt{\frac{\ln(1/\gamma)}{2n}} \right) \geq 1 - \gamma$$

where  $\omega = w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}$ . Note that

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n \hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2} = \hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2}$$

and  $\mathbf{E} [\bar{X}] = V^{\pi_1} - V^{\pi_2}$ . Thus, substituting  $(1 - p)/2$  for  $\gamma$ , we have that the following two statements hold with probability at least  $1 - 2\gamma = p$ ,

$$V^{\pi_1} - V^{\pi_2} \geq \hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} - \sqrt{\frac{\ln(2/(1-p))}{2n}}$$

and

$$V^{\pi_1} - V^{\pi_2} \leq \hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} + \sqrt{\frac{\ln(2/(1-p))}{2n}}$$

Thus the probability that  $V^{\pi_1} - V^{\pi_2} < 0$  but  $\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} - \sqrt{\frac{\ln(2/(1-p))}{2n}} > 0$  or  $V^{\pi_1} - V^{\pi_2} > 0$  but  $\hat{V}_{\text{IS}}^{\pi_1} + \hat{V}_{\text{IS}}^{\pi_2} - \sqrt{\frac{\ln(2/(1-p))}{2n}} < 0$  is less than  $p$ , which means for  $p = 1 - \delta$ , we will output the worse policy according to  $\succ_V$  with probability at most  $\delta$ , which is exactly what we need.  $\square$

**Theorem 6.5.** *For any two policies  $\pi_1$  and  $\pi_2$ , where  $\Pr \left( |\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2}| \geq 0 \right) \geq 1 - \delta_{\text{MC}}$  any behavior policy  $\pi_b$ ,  $\omega = (w_{\text{Max}}^{\pi_1} + 1)V_{\text{Max}}^{\pi_1} + (w_{\text{Max}}^{\pi_2} + 1)V_{\text{Max}}^{\pi_2}$ ,  $p = 1 - \delta_{\text{MC}}\delta$  for some  $\delta \leq 0.5$ , and for all  $k \in \{1, 2, 3, \dots\}$ , Algorithm 2 is a safe policy selection algorithm with respect to  $\succ_{\text{MC},kn}$  with probability  $1 - \delta$  when we have  $n$  samples drawn from  $\pi_b$ . It is a transitively fair policy selection algorithm whenever  $\delta_{\text{MC}} \leq 0.25$ .*

*Proof of Theorem 6.5.* Recall that Algorithm 2 receives as input  $\tau_1, \dots, \tau_n \sim \pi_b$ . Consider additional random samples  $\tau_1^{\pi_1}, \dots, \tau_{kn}^{\pi_1} \sim \pi_1$  and  $\tau_1^{\pi_2}, \dots, \tau_{kn}^{\pi_2} \sim \pi_2$ . Note that these samples are all independent from each other. For  $i \in \{1, 2, \dots, n\}$ , let

$$\hat{V}_i^\pi = \frac{1}{k} \sum_{j=1}^k \sum_{t=1}^{T_{ji}} R_{j,i,t}$$

(i.e., the estimate of the value of policy  $\pi$  using only samples  $\tau_i^\pi, \tau_{2i}^\pi, \dots, \tau_{ki}^\pi$ ). Furthermore let  $\hat{V}_{IS,i}^\pi = \sum_{t=1}^{T_i} w_{i,t} R_{i,t}$ . Now let

$$X_i = (\hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2}) - (\hat{V}_{IS,i}^{\pi_1} - \hat{V}_{IS,i}^{\pi_2})$$

Notice that the range of  $X_i$  is  $[-V_{\text{Max}}^{\pi_2} - w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1}, V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}]$ . Thus we can use Hoeffding's inequality to find that:

$$\Pr\left(\bar{X} - \mathbf{E}[\bar{X}] \geq -\omega \sqrt{\frac{\ln(1/\delta)}{2n}}\right) \geq 1 - \gamma$$

and

$$\Pr\left(\bar{X} - \mathbf{E}[\bar{X}] \leq \omega \sqrt{\frac{\ln(1/\delta)}{2n}}\right) \geq 1 - \gamma$$

where  $\omega = (w_{\text{Max}}^{\pi_1} + 1)V_{\text{Max}}^{\pi_1} + (w_{\text{Max}}^{\pi_2} + 1)V_{\text{Max}}^{\pi_2}$ .

Notice that

$$\begin{aligned} \bar{X} &= \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n (\hat{V}_i^{\pi_1} - \hat{V}_i^{\pi_2}) - (\hat{V}_{IS,i}^{\pi_1} - \hat{V}_{IS,i}^{\pi_2}) \\ &= (\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2}) - (\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2}) \end{aligned}$$

and

$$\mathbf{E}[\bar{X}] = (V^{\pi_1} - V^{\pi_2}) - (V^{\pi_1} - V^{\pi_2}) = 0$$

Thus, substituting  $(1-p)/2$  for  $\gamma$ , we have that the following two statements hold with probability at least  $1 - 2\gamma = p$ ,

$$\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} \geq \hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} - \sqrt{\frac{\ln(2/(1-p))}{2n}}$$

and

$$\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} \leq \hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} + \sqrt{\frac{\ln(2/(1-p))}{2n}}$$

Thus the probability that  $\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} < 0$  but  $(\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2}) - \omega \sqrt{\frac{\ln(2/(1-p))}{2n}} > 0$  or  $\hat{V}_{\text{MC},kn}^{\pi_1} - \hat{V}_{\text{MC},kn}^{\pi_2} > 0$  but  $\hat{V}_{\text{IS}}^{\pi_1} - \hat{V}_{\text{IS}}^{\pi_2} + \omega \sqrt{\frac{\ln(2/(1-p))}{2n}} < 0$  is less than  $1 - p$ . Now suppose without loss of generality that  $P(\hat{V}_{\text{MC},kn}^{\pi_1} > \hat{V}_{\text{MC},kn}^{\pi_2}) = \delta_{\text{MC}} > 0.5$ . The probability that we output  $\pi_2$  is at most  $p/\delta_{\text{MC}}$ . So if  $p = 1 - \delta_{\text{MC}}\delta$ , we output the worse policy with respect to  $\succ_{\text{MC},kn}$  with probability at most  $\delta$ , which is exactly what we need.

As long as  $\delta_{\text{MC}} \leq 0.25$ , we have that the algorithm is transitively safe since any fair algorithm with respect to  $\succ_{\text{MC},kn}$  is transitively fair whenever  $\Pr(|\hat{V}_{\text{MC}}^{\pi_1} - \hat{V}_{\text{MC}}^{\pi_2}| > 0) \geq 0.75$   $\square$

**Theorem 6.6.** *There exists policies  $\pi_1, \pi_2$ , and behavior policy  $\pi_b$  for which Algorithm 2 with inputs as described in Theorem 6.4 is not a safe policy selection algorithm with respect to  $\succ_{\text{MC},1}$  with  $p = 0.5$  when we have a single sample drawn from  $\pi_b$ .*

*Proof of Theorem 6.6.* Consider a world where there are three trajectories:  $\tau_1$  with reward 0.0001,  $\tau_2$  with reward 0.0002, and  $\tau_3$  with reward 1. We want to select between two policies:  $\pi_1$ , which places probability 1 on  $\tau_2$  and  $\pi_2$  which places probability 0.51 on  $\tau_1$  and probability 0.49 on  $\tau_3$ . When we only have one sample from each policy,  $\Pr(\hat{V}_{\text{MC},1}^{\pi_1} > \hat{V}_{\text{MC},1}^{\pi_2}) = 0.51 > 0.5$ , but clearly  $V^{\pi_1} \ll V^{\pi_2}$ . Now consider using IS with behavior policy  $\pi_b$  which places probability 0.48 on  $\tau_1$  and probability 0.01 on  $\tau_2$  and probability 0.51 on  $\tau_3$ . If we apply Algorithm 2 with the inputs to guarantee that the algorithm is safe with respect  $\succ_V$  (as given in Theorem 6.4), we find that whenever  $\pi_b$  samples from  $\tau_3$ ,

$$\begin{aligned} &V_{\text{IS}}^{\pi_1} - V_{\text{IS}}^{\pi_2} + (w_{\text{Max}}^{\pi_1} V_{\text{Max}}^{\pi_1} + w_{\text{Max}}^{\pi_2} V_{\text{Max}}^{\pi_2}) \sqrt{\frac{\ln 4}{2n}} \\ &= 0(1) - \frac{0.49}{0.51}(1) + \left(\frac{1}{0.01}(0.0002) + \frac{0.51}{0.48}(1)\right) \sqrt{\frac{\ln 4}{2}} \\ &\approx -0.060 < 0 \end{aligned}$$

Since this event occurs with probability 0.51, we find that Algorithm 2 returns  $\pi_2$  more than half the time, indicating that Algorithm 2 is not a safe policy with respect to  $\succ_{\text{MC},1}$   $\square$