

Electronic Companion of “Thompson Sampling with Information Relaxation Penalties”

Seungki Min

KAIST

`skmin@kaist.ac.kr`

Costis Maglaras

Columbia University

`c.maglaras@gsb.columbia.edu`

Ciamac C. Moallemi

Columbia University

`ciamac@gsb.columbia.edu`

Current Revision: February 2022

A. An Illustrative Example

Let us consider a Bernoulli MAB with eight periods ($T = 8$) and three arms ($K = 3$) with the following priors:

$$\mu_1 \sim \text{Beta}(3, 1), \quad \mu_2 \sim \text{Beta}(1, 1), \quad \mu_3 \sim \text{Beta}(1, 3), \quad (69)$$

where $R_{a,n} \sim \text{Bernoulli}(\mu_a)$ for each $a \in \{1, 2, 3\}$ and $n \in \{1, 2, \dots, 8\}$. Given this prior belief, the predictive mean reward of each arm is $\bar{\mu}_1 = \mathbb{E}_{\mu_1 \sim \text{Beta}(3,1)}[\mu_1] = \frac{3}{4}$, $\bar{\mu}_2 = \frac{1}{2}$, and $\bar{\mu}_3 = \frac{1}{4}$, respectively. As an illustrative example, we examine a particular instance where the true outcome ω is given as follows:

	True means $\mu_a(\theta_a)$	Rewards $R_{a,n}$							
		$n = 1$	2	3	4	5	6	7	8
Arm 1 ($a = 1$)	0.235	0	1	1	1	0	0	0	0
Arm 2 ($a = 2$)	0.443	1	0	0	1	1	1	1	0
Arm 3 ($a = 3$)	0.787	1	1	1	1	0	0	1	1

Table 8: An example of the outcome in a Bernoulli MAB with $K = 3$ and $T = 8$.

If we consider only the priors, arm 1 is best since $\bar{\mu}_1$ is largest among $(\bar{\mu}_1, \bar{\mu}_2, \bar{\mu}_3)$. If, however,

*The authors wish to thank Daniel Russo, Martin Haugh, David Brown, Jim Smith, and anonymous reviewers for helpful comments. A preliminary version of this paper appeared in the conference proceedings Advances in Neural Information Processing Systems 32 (NeurIPS 2019) (Min et al., 2019).

we have full information about the parameter values, arm 3 is best since μ_3 is largest among (μ_1, μ_2, μ_3) .

A.1. Inner Problems Induced by Different Penalty Functions

No penalty. To clarify the role of penalties, we first consider the case of zero penalty, i.e., $z_t \equiv 0$, which was not discussed in §3. With zero penalty, the DM at any time earns the current realized reward without adjustment. The clairvoyant DM, who is informed of the outcome ω , can find the best action sequence for this particular outcome ω . Recall that $R_{a,n}$ is defined to be the reward from the n^{th} pull of arm a , not the reward from arm a at time n , and so the DM is not allowed to skip any of the reward realizations and the total reward does not depend on the order of pulls. As depicted in the table below, the optimal solution is to pull arm 1 four times, arm 2 once, and arm 3 three times, which yields a total reward of 7.

	Payoffs under zero penalty								Maximal payoff
	$n = 1$	2	3	4	5	6	7	8	
Arm 1	0	1	1	1	0	0	0	0	7
Arm 2	1	0	0	1	1	1	1	0	
Arm 3	1	1	1	1	0	0	1	1	

TS penalty. Next, let us examine the penalty $z_t^{\text{TS}}(\mathbf{a}_{1:t}, \omega) \triangleq r_t(\mathbf{a}_{1:t}, \omega) - \mu_{a_t}(\theta_{a_t})$ under which the DM earns μ_a whenever playing arm a . The hindsight optimal action sequence is to pull arm 3 (the arm with the largest mean reward μ_a) eight times in a row and the DM can earn a total reward of $T \times \mu_3 = 6.296$ at most.

	Payoffs under z_t^{TS}								Maximal payoff
	$n = 1$	2	3	4	5	6	7	8	
Arm 1	.235	.235	.235	.235	.235	.235	.235	.235	6.296
Arm 2	.443	.443	.443	.443	.443	.443	.443	.443	
Arm 3	.787	.787	.787	.787	.787	.787	.787	.787	

IRS.FH penalty. When the penalties are given by $z_t^{\text{IRS.FH}}(\mathbf{a}_{1:t}, \omega) \triangleq r_t(\mathbf{a}_{1:t}, \omega) - \hat{\mu}_{a_t, T-1}(\omega)$, the DM earns $\hat{\mu}_{a, T-1}(\omega)$ whenever playing arm a . Recall that $\hat{\mu}_{a, T-1}(\omega)$ is the Bayesian estimate on mean reward of arm a after observing reward realizations $R_{a,1}, \dots, R_{a, T-1}$. In this particular example, we have $(\hat{\mu}_{1, T-1}, \hat{\mu}_{2, T-1}, \hat{\mu}_{3, T-1}) = \left(\frac{6}{11}, \frac{6}{9}, \frac{6}{11}\right)$ and the maximal payoff is $T \times \hat{\mu}_{2, T-1} = 5.333$, which can be obtained by playing arm 2 throughout the entire time horizon.

IRS.V-Zero penalty. Finally, let us focus on $z_t^{\text{IRS.V-ZERO}}(\mathbf{a}_{1:t}, \omega) \triangleq r_t(\mathbf{a}_{1:t}, \omega) - \hat{\mu}_{a_t, n_{t-1}(\mathbf{a}_{1:t-1}, a_t)}$ under which the DM earns $\hat{\mu}_{a, n-1}(\omega)$ from the n^{th} pull of arm a . Since the payoff from an arm changes over time as the Bayesian estimate evolves, playing only one arm is no longer optimal, unlike in the previous two cases. It can be easily verified that the optimal allocation is to play arm 1 six times and arm 2 two times, as visualized in the table below.

	Payoffs under $z_t^{\text{IRS.FH}}$								Maximal payoff	
	$n = 1$	2	3	4	5	6	7	8		
Arm 1	6/11	6/11	6/11	6/11	6/11	6/11	6/11	6/11	6/11	5.333
Arm 2	6/9	6/9	6/9	6/9	6/9	6/9	6/9	6/9	6/9	
Arm 3	6/11	6/11	6/11	6/11	6/11	6/11	6/11	6/11	6/11	

	Payoffs under $z_t^{\text{IRS.V-ZERO}}$								Maximal payoff
	$n = 1$	2	3	4	5	6	7	8	
Arm 1	3/4	3/5	4/6	5/7	6/8	6/9	6/10	6/11	5.314
Arm 2	1/2	2/3	2/4	2/5	3/6	4/7	5/8	6/9	
Arm 3	1/4	2/5	3/6	4/7	5/8	5/9	5/10	6/11	

IRS.V-EMax and the ideal penalty. Regarding the penalty functions $z_t^{\text{IRS.V-EMAX}}$ and z_t^{ideal} , we cannot visualize the optimal solution with a table since the total payoff depends on the detailed sequence of pulls and not only the number of pulls. While omitting the visual proof of optimality, we have that the action sequence $\mathbf{a}_{1:8}^* = (1, 2, 2, 1, 1, 1, 1, 1)$ achieves the maximal payoff of 5.806 under $z_t^{\text{IRS.V-EMAX}}$, and $\mathbf{a}_{1:8}^* = (1, 1, 1, 1, 1, 1, 1, 1)$ achieves the maximal payoff of 6.063 under z_t^{ideal} . In particular for z_t^{ideal} , the maximal payoff depends only on the prior belief \mathbf{y} and the time horizon T , irrespective of the outcome¹ ω .

We have so far illustrated how the different penalty functions induce the different inner problems and the different best actions given the same outcome ω . The readers may notice from the above examples that, as the penalty function becomes more complicated, the hindsight best action sequence becomes less dependent on a particular realization of ω . Instead, it becomes more dependent on the prior belief.

A.2. IRS Performance Bounds

The maximal payoffs above are calculated for a particular outcome given by Table 8. Recall that the IRS performance bound W^z is defined as the expected value of the maximal payoff where the expectation is taken with respect to the randomness of outcome ω over its prior distribution $\mathcal{I}(T, \mathbf{y})$. We can obtain this value by simulation, i.e., by solving a bunch of inner problems with respect to the randomly generated outcomes $\omega^{(1)}, \omega^{(2)}, \dots, \omega^{(S)}$ and taking the average of the maximal values. For this particular Bernoulli MAB setting ($T = 8$ with given priors), we obtain the following performance bounds:

W^0	W^{TS}	$W^{\text{IRS.FH}}$	$W^{\text{IRS.V-ZERO}}$	$W^{\text{IRS.V-EMAX}}$	$W^{\text{ideal}} = V^*$
6.805	6.429	6.279	6.111	6.075	6.063

¹For details, see the proof of the strong duality theorem in §C.1. While the maximal value does not depend on ω , the optimal action sequence still depends on ω . More specifically, it is the sequence of actions that the (non-anticipating) Bayesian optimal policy will take when ω is sequentially revealed.

We observe that the performance bounds are monotone, i.e., $W^0 > W^{\text{TS}} > W^{\text{IRS.FH}} > W^{\text{IRS.V-ZERO}} > W^{\text{IRS.V-EMAX}} > W^{\text{ideal}} = V^*$, which is consistent with Theorem 2.

A.3. Illustration of the IRS Policy (IRS.V-Zero)

We illustrate how the policy $\pi^{\text{IRS.V-ZERO}}$ makes decisions sequentially when the true outcome ω is the one specified in Table 8. At $t = 1$, it first synthesizes a future scenario based on the prior belief (i.e., sampling $\tilde{\omega}_1 \sim \mathcal{I}(\mathbf{y}_0)$) and finds the best action sequence in the presence of penalties $z_t^{\text{IRS.V-ZERO}}$ in the belief that the sampled outcome $\tilde{\omega}_1$ is the ground truth. The following table shows an example in which $\pi^{\text{IRS.V-ZERO}}$ plays arm 1.

$t = 1$	Priors \mathbf{y}_0	Payoffs with respect to $\tilde{\omega}_1 \sim \mathcal{I}(\mathbf{y}_0)$								Action
		$n = 1$	2	3	4	5	6	7	8	
Arm 1	Beta(3, 1)	3/4	4/5	5/6	6/7	7/8	7/9	8/10	9/11	$a_1 = 1$
Arm 2	Beta(1, 1)	1/2	1/3	1/4	1/5	1/6	1/7	2/8	3/9	
Arm 3	Beta(1, 3)	1/4	1/5	1/6	1/7	1/8	1/9	1/10	2/11	

As a result of the first action ($a_1 = 1$), we observe that $R_{1,1} = 0$ (encoded in the true outcome ω) and the associated belief is updated from Beta(3, 1) to Beta(3, 2) according to Bayes' rule. In order to make the next decision a_2 at time $t = 2$, $\pi^{\text{IRS.V-ZERO}}$ simulates an outcome for the remaining time horizon, i.e., $\tilde{\omega}_2 \sim \mathcal{I}(\mathbf{y}_1)$, independently of the outcome $\tilde{\omega}_1$ used at $t = 1$. Again, $\pi^{\text{IRS.V-ZERO}}$ finds the best action sequence for this new scenario and takes its first action.² The table below shows an instance of $\tilde{\omega}_2$ in which the policy will pull arm 2.

$t = 2$	Priors \mathbf{y}_1	Payoffs with respect to $\tilde{\omega}_2 \sim \mathcal{I}(\mathbf{y}_1)$							Action
		$n = 1$	2	3	4	5	6	7	
Arm 1	Beta(3, 2)	3/5	4/6	4/7	4/8	4/9	5/10	5/11	$a_2 = 2$
Arm 2	Beta(1, 1)	1/2	2/3	3/4	3/5	4/6	4/7	5/8	
Arm 3	Beta(1, 3)	1/4	1/5	1/6	1/7	1/8	1/9	1/10	

We can update the prior of arm 2 as a new reward realization $R_{2,1} = 1$ is revealed. In the following decision epochs $t = 3, 4, \dots$, the policy repeats the same decision-making procedure – (i) samples $\tilde{\omega}_t \sim \mathcal{I}(\mathbf{y}_{t-1})$, (ii) solves the inner problem, and (iii) plays the best arm that the optimal solution suggests – while updating the priors as the true reward realizations are revealed sequentially.

The following table illustrates the last decision epoch. As there remains one time period only, the policy $\pi^{\text{IRS.V-ZERO}}$ tries to maximize $\hat{\mu}_{a,0}(\tilde{\omega}_7) = \bar{\mu}_a(\mathbf{y}_7)$, which is the expected mean reward given the prior at that moment. Such a decision is totally myopic, but it is Bayesian optimal.

²In case of IRS.V-ZERO, we select the arm with the largest pull allocation as a first action.

$t = 8$	Priors \mathbf{y}_7	Payoffs with respect to $\tilde{\omega}_7 \sim \mathcal{I}(\mathbf{y}_7)$	Action
		$n = 1$	
Arm 1	Beta(6, 3)	6/9	$a_8 = 1$
Arm 2	Beta(2, 2)	2/4	
Arm 3	Beta(1, 3)	1/4	

B. Algorithms in Detail

B.1. Implementation of IRS.V-Zero

We provide a pseudo-code of the policy $\pi^{\text{IRS.V-ZERO}}$ introduced in §3.3. The same logic can be directly used to compute the performance bound $W^{\text{IRS.V-ZERO}}$ if the sampled outcome $\tilde{\omega}$ is replaced with the true outcome ω .

Algorithm 5: Arm selection rule of $\pi^{\text{IRS.V-ZERO}}$ when remaining time is T and current belief is \mathbf{y}

Function IRS.V-Zero(T, \mathbf{y})

```

1   $\tilde{\theta}_a \sim \mathcal{P}_a(y_a), \tilde{R}_{a,n} \sim \mathcal{R}_a(\tilde{\theta}), \quad \forall n \in \{1, \dots, T\}, \forall a \in \{1, \dots, K\}$ 
2  for  $a = 1, \dots, K$  do
3       $\tilde{y}_{a,0} \leftarrow y_a, \tilde{S}_{a,0} \leftarrow 0$ 
4      for  $n = 1, \dots, T$  do
5           $\tilde{S}_{a,n} \leftarrow \tilde{S}_{a,n-1} + \bar{\mu}_a(\tilde{y}_{a,n-1})$ 
6           $\tilde{y}_{a,n} \leftarrow \mathcal{U}_a(\tilde{y}_{a,n-1}, \tilde{R}_{a,n})$ 
7      end
8       $\tilde{M}_{0,0} \leftarrow 0, \tilde{M}_{0,n} \leftarrow -\infty, \forall n \in \{1, \dots, T\}$ 
9      for  $a = 1, \dots, K$  do
10         for  $n = 0, \dots, T$  do
11              $\tilde{M}_{a,n} \leftarrow \max_{0 \leq m \leq n} \{\tilde{M}_{a-1,n-m} + \tilde{S}_{a,m}\}$ 
12              $\tilde{L}_{a,n} \leftarrow \operatorname{argmax}_{0 \leq m \leq n} \{\tilde{M}_{a-1,n-m} + \tilde{S}_{a,m}\}$ 
13         end
14     end
15      $\tau \leftarrow T$ 
16     for  $a = K, \dots, 1$  do
17          $\tilde{n}_a^* \leftarrow \tilde{L}_{a,\tau}$ 
18          $\tau \leftarrow \tau - \tilde{n}_a^*$ 
19     end
20 return  $\operatorname{argmax}_a \tilde{n}_a^*$ 

```

B.2. Implementation of IRS.V-EMax

We use the notation $\mathbf{y}_t(\mathbf{n}_{1:K}, \omega)$ to denote the belief as a function of pull counts $\mathbf{n}_{1:K} \triangleq (n_1, \dots, n_K) \in \mathbb{N}_0^K$, based on the observation that the belief is completely determined by how many times each of the arms has been pulled, $\mathbf{n}_{1:K}$, irrespective of the specific sequence in which the arms have been pulled. Given the pull counts $\mathbf{n}_{1:K}$, we define the payoff of pulling arm a one more time after pulling the individual arms n_1, \dots, n_K times respectively: with $t = \sum_{a=1}^K n_a$, the effective payoff associated with arm a at time t is

$$r^z(\mathbf{n}_{1:K}, a, \omega) \triangleq \hat{\mu}_{a, n_a}(\omega) - W^{\text{TS}}(T - t - 1, \mathbf{y}_{t+1}(\mathbf{n}_{1:K} + \mathbf{e}_a, \omega)) + W^{\text{TS}}(T - t - 1, \mathbf{y}_t(\mathbf{n}_{1:K}, \omega)), \quad (70)$$

where $\mathbf{e}_a \in \mathbb{N}_0^K$ is a basis vector such that the a^{th} component is one and the others are zero. Note that we used the fact that $\mathbb{E} \left[W^{\text{TS}}(T - t, \mathbf{y}_t) \middle| H_{t-1} \right] = W^{\text{TS}}(T - t, \mathbf{y}_{t-1})$.

Consider a subproblem of [\(6\)](#) that maximizes the total payoff given the number of pulls $\mathbf{n}_{1:K}$ across all the arms: with $t = \sum_{a=1}^K n_a$, we get

$$M(\mathbf{n}_{1:K}, \omega) \triangleq \max_{\mathbf{a}_{1:t} \in \mathcal{A}^t} \left\{ \sum_{s=1}^t r_s(\mathbf{a}_{1:s}, \omega) - z_s^{\text{IRS.V-EMAX}}(\mathbf{a}_{1:s}, \omega); \sum_{s=1}^t \mathbf{1}\{a_s = a\} = n_a, \forall a \right\}. \quad (71)$$

Consequently, the maximal value $M(\mathbf{n}_{1:K}, \omega)$ should satisfy the following Bellman equation:

$$M(\mathbf{n}_{1:K}, \omega) = \max_{a \in \mathcal{A}: n_a \geq 1} \{M(\mathbf{n}_{1:K} - \mathbf{e}_a, \omega) + r^z(\mathbf{n}_{1:K} - \mathbf{e}_a, a, \omega)\}, \quad (72)$$

i.e., when letting a^* be the maximizer of [\(72\)](#), it is optimal to play arm a^* after making the best effort within the allocation $\mathbf{n}_{1:K} - \mathbf{e}_{a^*}$. For all feasible counts $\mathbf{n}_{1:K}$'s such that $\sum_{a=1}^K n_a \leq T$, we can compute $M(\mathbf{n}_{1:K}, \omega)$'s by sequentially solving [\(72\)](#) in an appropriate order. By doing so, we can obtain the maximal value of the original inner problem [\(6\)](#) by evaluating

$$\max_{\mathbf{n}_{1:K} \in N_T} \{M(\mathbf{n}_{1:K}, \omega)\}, \quad (73)$$

where $N_T \triangleq \{(n_1, \dots, n_K) \in \mathbb{N}_0^K : \sum_{a=1}^K n_a = T\}$, and the performance bound $W^{\text{IRS.V-EMAX}}$ is the expected value of [\(73\)](#) with respect to the random realization of ω . The optimal action sequence $\mathbf{a}_{1:T}^*$ can be obtained by tracking $M(\mathbf{n}_{1:K}, \omega)$'s backward.

Algorithm 6: Arm selection rule of $\pi^{\text{IRS.V-ZERO}}$ when remaining time is T and current belief is \mathbf{y}

Function IRS.V-EMax(T, \mathbf{y})

```

1   $\tilde{\theta}_a \sim \mathcal{P}_a(y_a), \tilde{R}_{a,n} \sim \mathcal{R}_a(\tilde{\theta}), \quad \forall n \in \{1, \dots, T\}, \forall a \in \{1, \dots, K\}$ 
2   $\tilde{y}_{a,0} \leftarrow y_a, \tilde{y}_{a,n} \leftarrow \mathcal{U}_a(\tilde{y}_{a,n-1}, \tilde{R}_{a,n}), \quad \forall n \in \{1, \dots, T\}, \forall a \in \{1, \dots, K\}$ 
3  for each  $\mathbf{n}_{1:K} \in N_{\leq T}$  do
4  |    $\tilde{\Gamma}[\mathbf{n}_{1:K}] \leftarrow \mathbb{E}_{\tilde{\mathbf{y}}(\mathbf{n}_{1:K})} [\max_a \mu_a(\theta_a)]$ 
   end
5  for each  $\mathbf{n}_{1:K} \in N_{< T}$  do
6  |    $\tilde{r}^z[\mathbf{n}_{1:K}, a] \leftarrow \bar{\mu}_a(\tilde{y}_{a,n_a-1}) + (T - \sum_{a=1}^K n_a - 1) \times (\tilde{\Gamma}[\mathbf{n}_{1:K}] - \tilde{\Gamma}[\mathbf{n}_{1:K} + \mathbf{e}_a]), \quad \forall a \in$ 
   |    $\{1, \dots, K\}$ 
   end
7   $\tilde{M}[\mathbf{0}] \leftarrow 0$ 
8  for each  $\mathbf{n}_{1:K} \in N_{\leq T} \setminus \{\mathbf{0}\}$  in order with increasing  $\sum_{a=1}^K n_a$  do
9  |    $\tilde{M}[\mathbf{n}_{1:K}] \leftarrow \max_{a:n_a>0} \{ \tilde{M}[\mathbf{n}_{1:K} - \mathbf{e}_a] + \tilde{r}^z[\mathbf{n}_{1:K} - \mathbf{e}_a, a] \}$ 
10 |    $\tilde{A}[\mathbf{n}_{1:K}] \leftarrow \operatorname{argmax}_{a:n_a>0} \{ \tilde{M}[\mathbf{n}_{1:K} - \mathbf{e}_a] + \tilde{r}^z[\mathbf{n}_{1:K} - \mathbf{e}_a, a] \}$ 
   end
11  $\mathbf{m}_{1:K} \leftarrow \operatorname{argmax}_{\mathbf{n}_{1:K} \in N_T} \{ \tilde{M}[\mathbf{n}_{1:K}] \}$ 
12 for  $t = T, \dots, 1$  do
13 |    $\tilde{a}_t^* \leftarrow \tilde{A}[\mathbf{m}_{1:K}]$ 
14 |    $m_{\tilde{a}_t^*} \leftarrow m_{\tilde{a}_t^*} - 1$ 
   end
15 return  $\tilde{a}_1^*$ 

```

Here, $\tilde{\mathbf{y}}(\mathbf{n}_{1:K}) \triangleq (\tilde{y}_{1,n_1}, \dots, \tilde{y}_{K,n_K})$, $N_{\leq T} \triangleq \{\mathbf{n}_{1:K}; \sum_a n_a \leq T\}$, $N_{< T} \triangleq \{\mathbf{n}_{1:K}; \sum_a n_a < T\}$, and in line 8, $\mathbf{n}_{1:K}$ iterates over $N_{\leq T} \setminus \{\mathbf{0}\}$ in an order in which $\sum_{a=1}^K n_a$ is non-decreasing.

Since $|N_{\leq T}| = O(T^K)$, it requires $O(KT^K)$ operations to compute all $M(\mathbf{n}_{1:K}, \omega)$'s. However, another practical issue is the cost of computing $W^{\text{TS}}(T, \mathbf{y}) = T \times \mathbb{E}_{\mathbf{y}} [\max_a \mu_a(\theta_a)]$ which has to be evaluated $O(T^K)$ times in total. There is no simple closed-form expression in general, and it should be evaluated with numerical integration or Monte Carlo sampling.

B.3. Implementation of IRS.Index

We first prove the identity that was utilized in §3.5, and then provide the pseudo code for IRS.INDEX policy.

Proposition 2. *The optimization problem (46) can be reformulated as*

$$\max_{0 \leq n \leq T} \left\{ T \times \Gamma_0^\lambda + (T - n) \times \left(\lambda - \min_{0 \leq i \leq n} \Gamma_i^\lambda \right) + \sum_{i=1}^n \left(\hat{\mu}_{a,i-1} - \Gamma_{i-1}^\lambda \right) \right\}. \quad (74)$$

Here, the decision variable n is the total number of pulls of a stochastic arm.

Proof. Fix $m \triangleq n_T$, i.e., the total number of pulls on the stochastic arm. Note that if $a_t = 0$, then $(T - t) \times (\Gamma_{n_t}^\lambda - \Gamma_{n_{t-1}}^\lambda) = 0$ since $n_t = n_{t-1}$. The objective function can be represented as

$$\sum_{n=1}^m \hat{\mu}_{a,n-1} + (T - m) \times \lambda - \sum_{n=1}^m (T - t_n) \times (\Gamma_n^\lambda - \Gamma_{n-1}^\lambda), \quad (75)$$

where $t_n \triangleq \inf\{t; n_t \geq n\}$ represents the time at which the n^{th} pull on the stochastic arm is made. It suffices to find the optimal pulling times (t_1, \dots, t_m) with $1 \leq t_1 < t_2 < \dots < t_m \leq T$ by which $\sum_{n=1}^m (T - t_n) \times (\Gamma_n^\lambda - \Gamma_{n-1}^\lambda)$ is minimized. With $t_0 \triangleq 0$ and $t_{m+1} \triangleq T + 1$, we have

$$\sum_{n=1}^m (T - t_n) \times (\Gamma_n^\lambda - \Gamma_{n-1}^\lambda) \quad (76)$$

$$= \sum_{n=1}^m (T - t_n) \times \Gamma_n^\lambda - \sum_{n=1}^m (T - t_n) \times \Gamma_{n-1}^\lambda \quad (77)$$

$$= \sum_{n=1}^m (T - t_n) \times \Gamma_n^\lambda - \sum_{n=0}^{m-1} (T - t_{n+1}) \times \Gamma_n^\lambda \quad (78)$$

$$= \sum_{n=0}^m (T - t_n) \times \Gamma_n^\lambda - (T - t_0) \times \Gamma_0^\lambda - \sum_{n=0}^m (T - t_{n+1}) \times \Gamma_n^\lambda + (T - t_{m+1}) \times \Gamma_m^\lambda \quad (79)$$

$$= -\Gamma_m^\lambda - T \times \Gamma_0^\lambda + \sum_{n=0}^m (t_{n+1} - t_n) \times \Gamma_n^\lambda. \quad (80)$$

Consider the minimum value among $\Gamma_0^\lambda, \dots, \Gamma_m^\lambda$ and let $n^* \triangleq \operatorname{argmin}_{0 \leq n \leq m} \Gamma_n^\lambda$. In order to minimize (80), it should satisfy that $t_{n+1} - t_n = T - m + 1$ for $n = n^*$ and $t_{n+1} - t_n = 1$ for $n \neq n^*$. For such t_n 's, (75) reduces to

$$\sum_{n=1}^m \hat{\mu}_{a,n-1} + (T - m) \times \lambda - \left(-\Gamma_m^\lambda - T \times \Gamma_0^\lambda + \sum_{n=0}^m \Gamma_n^\lambda + (T - m) \times \min_{0 \leq n \leq m} \Gamma_m^\lambda \right) \quad (81)$$

$$= \sum_{n=1}^m \hat{\mu}_{a,n-1} + (T - m) \times \left(\lambda - \min_{0 \leq n \leq m} \Gamma_m^\lambda \right) + T \times \Gamma_0^\lambda - \sum_{n=0}^{m-1} \Gamma_n^\lambda. \quad (82)$$

By taking its maximum value over $m = 0, \dots, T$, we obtain (49). ■

The following pseudo code implements the arm selection rule of the IRS.INDEX policy when remaining time is T and current belief is \mathbf{y} . In line 14, the infimum can be found via the bisection method, and $\tilde{\mathbf{y}}_{a,0:T} \triangleq (\tilde{y}_{a,0}, \dots, \tilde{y}_{a,T})$ represents the sequence of beliefs under the sampled outcome.

Algorithm 7: Arm selection rule of IRS.INDEX policy when remaining time is T and current belief is \mathbf{y}

Function IRS.Single.Worth-Trying($a, T, \lambda, \tilde{\mathbf{y}}_{a,0:T}$)

```

1  |  $\tilde{\Gamma}_n^\lambda \leftarrow \mathbb{E}_{\tilde{y}_{a,n}} [\max(\mu_a(\theta_a), \lambda)], \forall n \in \{0, \dots, T\}$ 
2  |  $\tilde{S}_{a,0}^\mu \leftarrow 0, \tilde{S}_0^\Gamma \leftarrow 0, \tilde{m}_0^\Gamma \leftarrow \tilde{\Gamma}_0^\lambda$ 
3  | for  $n = 1, \dots, T$  do
4  |   |  $\tilde{S}_{a,n}^\mu \leftarrow \tilde{S}_{a,n-1}^\mu + \bar{\mu}_a(\tilde{y}_{a,n-1})$ 
5  |   |  $\tilde{S}_n^\Gamma \leftarrow \tilde{S}_{a,n-1}^\Gamma + \tilde{\Gamma}_n^\lambda$ 
6  |   |  $\tilde{m}_n^\Gamma \leftarrow \min(\tilde{m}_{n-1}^\Gamma, \tilde{\Gamma}_{n-1}^\lambda)$ 
   | end
7  |  $\tilde{\varphi}_a \leftarrow \max_{1 \leq n \leq T} \{ \tilde{S}_{a,n}^\mu + T \times \tilde{\Gamma}_0^\lambda + (T - n) \times (\lambda - \tilde{m}_n^\Gamma) - \tilde{S}_n^\Gamma \} - T \times \lambda$ 
8  | if  $\tilde{\varphi}_a \geq 0$  then
9  |   | return true
   | else
10 |   | return false
   | end

```

Function IRS.Index(T, \mathbf{y})

```

11 |  $\tilde{\theta}_a \sim \mathcal{P}_a(y_a), \tilde{R}_{a,n} \sim \mathcal{R}_a(\tilde{\theta}), \quad \forall n \in \{1, \dots, T\}, \forall a \in \{1, \dots, K\}$ 
12 |  $\tilde{y}_{a,0} \leftarrow y_a, \quad \tilde{y}_{a,n} \leftarrow \mathcal{U}_a(\tilde{y}_{a,n-1}, \tilde{R}_{a,n}), \quad \forall n \in \{1, \dots, T\}, \quad \forall a \in \{1, \dots, K\}$ 
13 | for  $a = 1, \dots, K$  do
14 |   |  $\tilde{\lambda}_a^* \leftarrow \inf \{ \lambda; \text{IRS.Single.Worth-Trying}(a, T, \lambda, \tilde{\mathbf{y}}_{a,0:T}) = \text{true} \}$ 
   | end
15 | return  $\text{argmax}_a \tilde{\lambda}_a^*$ 

```

C. Proofs for §3

Proposition 3 (Mean equivalence). *If the penalty function z_t is dual feasible, the presence of penalties does not affect the performance of a non-anticipating policy π : i.e.,*

$$\mathbb{E}_{\mathbf{y}}^\pi \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi, \omega) - z_t(\mathbf{A}_{1:t}^\pi, \omega) \right] = \mathbb{E}_{\mathbf{y}}^\pi \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi, \omega) \right] =: V(\pi, T, \mathbf{y}). \quad (83)$$

Proof. The claim immediately follows from the definition of dual feasibility and the linearity of the expectation operator. ■

C.1. Proof of Theorem 1

Despite that the results of Theorem 1 were already well established in Brown et al. (2010), we provide the detailed proof as our context is slightly different from that of Brown et al. (2010) regarding the measurability of r_t . We define an appending operator \oplus that concatenates an element into a vector so that $\mathbf{a}_{1:t} = \mathbf{a}_{1:t-1} \oplus a_t$.

Weak duality. Define the filtration for the perfect information relaxation $\mathcal{G}_t \triangleq \mathcal{F}_t \cup \sigma(\omega)$ and consider a relaxed policy space $\Pi_{\mathcal{G}} \triangleq \{\pi : A_t^\pi \text{ is } \mathcal{G}_{t-1}\text{-measurable, } \forall t\}$. Then, we have

$$V^*(T, \mathbf{y}) \triangleq \sup_{\pi \in \Pi_{\mathbb{F}}} \mathbb{E} \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi) \right] \stackrel{\text{Prop 3}}{=} \sup_{\pi \in \Pi_{\mathbb{F}}} \mathbb{E} \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi) - z_t(\mathbf{A}_{1:t}^\pi) \right] \quad (84)$$

$$\leq \sup_{\pi \in \Pi_{\mathcal{G}}} \mathbb{E} \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi) - z_t(\mathbf{A}_{1:t}^\pi) \right] = \mathbb{E} \left[\max_{\mathbf{a}_{1:T} \in \mathcal{A}^T} \sum_{t=1}^T r_t(\mathbf{a}_{1:t}) - z_t(\mathbf{a}_{1:t}) \right] \quad (85)$$

$$= W^z(T, \mathbf{y}), \quad (86)$$

where the inequality holds since $\Pi_{\mathbb{F}} \subseteq \Pi_{\mathcal{G}}$. ■

Strong duality. Fix T and \mathbf{y} . Let $V_t^{\text{in}}(\mathbf{a}_{1:t-1}, \omega)$ and $Q_t^{\text{in}}(\mathbf{a}_{1:t-1}, a, \omega)$ be, respectively, the value function and the state-action value (Q-value) function that are associated with the inner problem (*) given a particular outcome ω under the ideal penalty (22). With $V_{T+1}^{\text{in}} \equiv 0$, we have the following Bellman equation for the inner problem:

$$Q_t^{\text{in}}(\mathbf{a}_{1:t-1}, a, \omega) \triangleq r_t(\mathbf{a}_{1:t-1} \oplus a, \omega) - z_t^{\text{ideal}}(\mathbf{a}_{1:t-1} \oplus a, \omega) + V_{t+1}^{\text{in}}(\mathbf{a}_{1:t-1} \oplus a, \omega), \quad (87)$$

$$V_t^{\text{in}}(\mathbf{a}_{1:t-1}, \omega) = \max_{a \in \mathcal{A}} \left\{ Q_t^{\text{in}}(\mathbf{a}_{1:t-1}, a, \omega) \right\}. \quad (88)$$

We argue by induction to show that

$$V_t^{\text{in}}(\mathbf{a}_{1:t-1}, \omega) = V^*(T - t + 1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega)), \quad (89)$$

$$Q_t^{\text{in}}(\mathbf{a}_{1:t-1}, a, \omega) = Q^*(T - t + 1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega), a), \quad (90)$$

for all $\mathbf{a}_{1:t-1} \in \mathcal{A}^{t-1}$, $a \in \mathcal{A}$ and $t \in \{1, \dots, T+1\}$.

As a terminal case, when $t = T+1$, the claim holds trivially, since $V_{T+1}^{\text{in}}(\mathbf{a}_{1:T}, \omega) = 0 = V^*(0, \mathbf{y}_T(\mathbf{a}_{1:T}, \omega))$. Now assume that the claim holds for $t+1$: i.e., $V_{t+1}^{\text{in}}(\mathbf{a}_{1:t}, \omega) = V^*(T -$

$t, \mathbf{y}_t(\mathbf{a}_{1:t}, \omega)$ for all $\mathbf{a}_{1:t} \in \mathcal{A}^t$. For any $\mathbf{a}_{1:t-1} \in \mathcal{A}^{t-1}$ and $a \in \mathcal{A}$, then,

$$Q_t^{\text{in}}(\mathbf{a}_{1:t-1}, a, \omega) = r_t(\mathbf{a}_{1:t-1} \oplus a, \omega) - z_t^{\text{ideal}}(\mathbf{a}_{1:t-1} \oplus a, \omega) + V_{t+1}^{\text{in}}(\mathbf{a}_{1:t-1} \oplus a, \omega) \quad (91)$$

$$= \mathbb{E} [r_t(\mathbf{a}_{1:t-1} \oplus a, \omega) + V^*(T - t, \mathbf{y}_t(\mathbf{a}_{1:t-1} \oplus a, \omega)) | H_{t-1}(\mathbf{a}_{1:t-1}, \omega)] \quad (92)$$

$$\underbrace{-V^*(T - t, \mathbf{y}_t(\mathbf{a}_{1:t-1} \oplus a, \omega)) + V_{t+1}^{\text{in}}(\mathbf{a}_{1:t-1} \oplus a, \omega)}_{=0} \quad (93)$$

$$= \mathbb{E} [r_t(\mathbf{a}_{1:t-1} \oplus a, \omega) + V^*(T - t, \mathbf{y}_t(\mathbf{a}_{1:t-1} \oplus a, \omega)) | H_{t-1}(\mathbf{a}_{1:t-1}, \omega)] \quad (94)$$

$$= \mathbb{E}_{\mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega)} [R_a + V^*(T - t, \mathcal{U}(\mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega), a, R_a))] \quad (95)$$

$$= Q^*(T - t, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega), a), \quad (96)$$

where the last equality follows from the original Bellman equation (15). Consequently, we obtain

$$V_t^{\text{in}}(\mathbf{a}_{1:t-1}, \omega) = \max_{a \in \mathcal{A}} \{Q_t^{\text{in}}(\mathbf{a}_{1:t-1}, a, \omega)\} \quad (97)$$

$$= \max_{a \in \mathcal{A}} \{Q^*(T - t, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega), a)\} \quad (98)$$

$$= V^*(T - t, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega)). \quad (99)$$

Therefore the claim holds for all $t = 1, \dots, T$. In particular for $t = 1$, we have

$$V_1^{\text{in}}(\emptyset, \omega) = V^*(T, \mathbf{y}), \quad Q_1^{\text{in}}(\emptyset, a, \omega) = Q^*(T, \mathbf{y}, a), \quad \forall \omega. \quad (100)$$

Note that the maximal value of the inner problem does not depend on the outcome ω , i.e., it is deterministic with respect to the randomness of ω . As its expected value, $W^{\text{ideal}}(T, \mathbf{y}) = V^*(T, \mathbf{y})$. ■

C.2. Proof of Remark 2

We proceed on the proof of strong duality. The policy π^{ideal} solves the same inner problem with respect to a randomly sampled outcome $\tilde{\omega}$. When the remaining time is T and the current belief is \mathbf{y} , it takes an action with the largest Q-value: together with (100), it yields

$$a^{\pi^{\text{ideal}}} = \operatorname{argmax}_a Q_1^{\text{in}}(\emptyset, a, \tilde{\omega}) = \operatorname{argmax}_a Q^*(T, \mathbf{y}, a). \quad (101)$$

Therefore, at each moment, irrespective of the sampled outcome $\tilde{\omega}$, the policy π^{ideal} always takes the same action that the Bayesian optimal policy would take. Although there might be some ambiguity regarding tie breaking in argmax , it does not affect the expected performance. Therefore, $V(\pi^{\text{ideal}}, T, \mathbf{y}) = V^*(T, \mathbf{y})$. ■

C.3. Proof of Remark 3

First observe that for any non-anticipating policy $\pi \in \Pi_{\mathbb{F}}$, since A_t^π is \mathcal{F}_{t-1} -measurable, we have

$$\mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi, \omega) \right] = \mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T \mathbb{E}(r_t(\mathbf{A}_{1:t}^\pi, \omega) | \mathcal{F}_{t-1}, \boldsymbol{\theta}) \right] = \mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T \mu_{A_t^\pi}(\theta_{A_t^\pi}) \right]. \quad (102)$$

Since $\mathbb{E}[r_t(\mathbf{a}_{1:t}, \omega) | \boldsymbol{\theta}] = \mu_{a_t}(\theta_{a_t})$ for any $\mathbf{a}_{1:t} \in \mathcal{A}^t$, we further deduce that

$$\mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T z_t^{\text{TS}}(\mathbf{A}_{1:t}^\pi, \omega) \right] = \mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T r_t(\mathbf{A}_{1:t}^\pi, \omega) \right] - \mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T \mu_{A_t^\pi}(\theta_{A_t^\pi}) \right] = 0, \quad (103)$$

and thus z_t^{TS} is dual feasible.

Also observe that $\mathbb{E}[r_t(\mathbf{a}_{1:t}) | \hat{\boldsymbol{\mu}}_{T-1}] = \mathbb{E}[\mu_{a_t} | \hat{\boldsymbol{\mu}}_{T-1}] = \mathbb{E}[\mu_{a_t} | \hat{\boldsymbol{\mu}}_{T-1}, H_{t-1}]$ and $\mathbb{E}[r_t(\mathbf{a}_{1:t}) | H_{t-1}] = \mathbb{E}[\mu_{a_t} | H_{t-1}]$ for any $\mathbf{a}_{1:t} \in \mathcal{A}^t$. We can easily verify that each of penalty functions (22)–(26) has a form of

$$z_t(\mathbf{a}_{1:t}, \omega) = z_t^{\text{TS}}(\mathbf{a}_{1:t}, \omega) + w_t(\mathbf{a}_{1:t}, \omega) - \mathbb{E}[w_t(\mathbf{a}_{1:t}, \omega) | G_{t-1}(\mathbf{a}_{1:t-1}, \omega)], \quad (104)$$

for some deterministic function w_t and some relaxed information set $G_{t-1} \supseteq H_{t-1}$. By invoking Proposition 2.3 (iii) of Brown et al. (2010), we have that $z_t^{\text{IRS.FH}} - z_t^{\text{TS}}$, $z_t^{\text{IRS.V-ZERO}} - z_t^{\text{TS}}$, $z_t^{\text{IRS.V-EMAX}} - z_t^{\text{TS}}$, and $z_t^{\text{ideal}} - z_t^{\text{TS}}$ are dual feasible, and therefore so are $z_t^{\text{IRS.FH}}$, $z_t^{\text{IRS.V-ZERO}}$, $z_t^{\text{IRS.V-EMAX}}$, and z_t^{ideal} . ■

D. Proofs for §4

D.1. Notes on Regularity

Proposition 4. *If $\mathbb{E}_{\mathbf{y}} |R_{a,n}| < \infty$ for all a ,*

$$\mathbb{E}_{\mathbf{y}} |\mu_a(\theta_a)| < \infty, \quad \text{and} \quad W^{\text{TS}}(T, \mathbf{y}) < \infty, \quad \forall T \in \mathbb{N}. \quad (105)$$

Proof. By Jensen's inequality,

$$\mathbb{E}_{\mathbf{y}} |\mu_a(\theta_a)| = \mathbb{E}_{\mathbf{y}} [|\mathbb{E}(R_{a,n} | \theta_a)|] \leq \mathbb{E}_{\mathbf{y}} [\mathbb{E}(|R_{a,n}| | \theta_a)] = \mathbb{E}_{\mathbf{y}} |R_{a,n}| < \infty. \quad (106)$$

Consequently,

$$\mathbb{E}_{\mathbf{y}} \left[\max_a \mu_a(\theta_a) \right] \leq \mathbb{E}_{\mathbf{y}} \left[\sum_{a=1}^K |\mu_a(\theta_a)| \right] = \sum_{a=1}^K \mathbb{E}_{\mathbf{y}} |\mu_a(\theta_a)| < \infty. \quad (107)$$

The claim holds since $W^{\text{TS}}(T, \mathbf{y}) = T \times \mathbb{E}_{\mathbf{y}} [\max_a \mu_a(\theta_a)]$. ■

Proposition 5. If $\mathbb{E}_{\mathbf{y}} |R_{a,n}| < \infty$,

$$\lim_{n \rightarrow \infty} \hat{\mu}_{a,n}(\omega; y_a) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n R_{a,i} = \mu_a(\theta_a) \quad \text{almost surely,} \quad (108)$$

where $\hat{\mu}_{a,n}(\omega; y_a) \triangleq \mathbb{E}_{y_a} [\mu_a(\theta_a) | R_{a,1}, \dots, R_{a,n}]$.

Proof. Fix a and let $\mathcal{H}_n \triangleq \sigma(R_{a,1}, \dots, R_{a,n})$. First note that, by the strong law of large numbers, $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n R_{a,i} = \mu_a(\theta_a)$ almost surely. Therefore, $\mu_a(\theta_a)$ is measurable with respect to $\mathcal{H}_\infty \triangleq \bigcup_n \mathcal{H}_n$. Also note that $\hat{\mu}_{a,n} = \mathbb{E}(\mu_a(\theta_a) | \mathcal{H}_n)$ is a Doob martingale adapted to \mathcal{H}_n . By Levy's upward theorem, since $\mu_a(\theta_a) \in \mathcal{L}^1$ by Proposition [4](#), $\hat{\mu}_{a,n}$ converges to $\mathbb{E}(\mu_a(\theta_a) | \mathcal{H}_\infty) = \mu_a(\theta_a)$ almost surely as $n \rightarrow \infty$. \blacksquare

D.2. Proof of Proposition [1](#)

Asymptotic behavior of $\pi^{\text{IRS.FH}}$. Let $\tilde{\omega}$ be the sampled outcome used by $\pi^{\text{IRS.FH}}$. By Proposition [5](#), we have $\lim_{n \rightarrow \infty} \hat{\mu}_{a,n}(\tilde{\omega}) = \mu_a(\tilde{\theta}_a)$ for almost all $\tilde{\omega}$. This, together with the assumption that $\mu_i(\theta_i) \neq \mu_j(\theta_j)$ for $i \neq j$, since $\operatorname{argmax}_a \mu_a(\tilde{\theta}_a)$ is uniquely defined for almost all $\tilde{\omega}$, yields

$$\operatorname{argmax}_a \mu_a(\tilde{\theta}_a) = \operatorname{argmax}_a \lim_{n \rightarrow \infty} \hat{\mu}_{a,n}(\tilde{\omega}) = \lim_{n \rightarrow \infty} \operatorname{argmax}_a \hat{\mu}_{a,n}(\tilde{\omega}) \quad \text{a.s.} \quad (109)$$

Since almost-sure convergence guarantees convergence in distribution, for any $a \in \mathcal{A}$,

$$\lim_{T \rightarrow \infty} \mathbb{P} \left[A^{\text{IRS.FH}}(T, \mathbf{y}) = a \right] = \lim_{T \rightarrow \infty} \mathbb{P} \left[\operatorname{argmax}_{a'} \hat{\mu}_{a', T-1}(\tilde{\omega}) = a \right] \quad (110)$$

$$= \mathbb{P} \left[\operatorname{argmax}_{a'} \mu_{a'}(\tilde{\theta}_{a'}) = a \right] \quad (111)$$

$$= \mathbb{P} \left[A^{\text{TS}}(\mathbf{y}) = a \right]. \quad (112)$$

Note that we are not assuming that $\pi^{\text{IRS.FH}}$ and π^{TS} share the randomness. The sampled parameters used in π^{TS} are not necessarily the ones used in $\pi^{\text{IRS.FH}}$, but their distributions are identical since they are drawn from the same prior. \blacksquare

Asymptotic behavior of $\pi^{\text{IRS.V-ZERO}}$. To simplify notation, let $A_T^\circ \triangleq A^{\text{IRS.V-ZERO}}(T, \mathbf{y})$. As above, it suffices to show that $\lim_{T \rightarrow \infty} A_T^\circ = \operatorname{argmax}_{a \in \mathcal{A}} \mu_a(\tilde{\theta}_a) := A^{\text{TS}}$ for almost all sampled outcome $\tilde{\omega}$. We hide $\tilde{\omega}$ and $\tilde{\theta}_a$ from the notation for the further simplification.

Define

$$\Delta \triangleq \min_{a \neq A^{\text{TS}}} |\mu_{A^{\text{TS}}} - \mu_a| \quad \text{and} \quad M \triangleq \sup_{a \in \mathcal{A}, n \geq 0} |\hat{\mu}_{a,n}|. \quad (113)$$

We have $0 < \Delta < 2M < \infty$ almost surely since $\mu_i(\tilde{\theta}_i) \neq \mu_j(\tilde{\theta}_j)$ for $i \neq j$ and $\lim_{n \rightarrow \infty} \hat{\mu}_{a,n} = \mu_a < \infty$

almost surely for all a . In addition, there exists $N \in \mathbb{N}$ such that

$$|\hat{\mu}_{a,n} - \mu_a| < \frac{\Delta}{4}, \quad \forall n \geq N, \quad \forall a \in \mathcal{A}. \quad (114)$$

For such N , we have

$$\inf_{n \geq N} \hat{\mu}_{a^{\text{TS}},n} \geq \sup_{n \geq N} \hat{\mu}_{a,n} + \frac{\Delta}{2}, \quad \forall a \neq A^{\text{TS}}. \quad (115)$$

Note that A^{TS} , Δ , M , and N do not have the dependency on T .

To argue by contradiction, suppose that $A_T^\circ \neq A^{\text{TS}}$ for some large T such that $T \geq 2N + \frac{8MN}{\Delta} + 2$. Define the optimal allocation to the inner problem of IRS.V-ZERO for such T :

$$\mathbf{n}_{1:K}^\circ \triangleq \operatorname{argmax}_{\mathbf{n}_{1:K} \in \mathcal{N}_T} \left\{ \sum_{a=1}^K \sum_{s=1}^{n_a} \hat{\mu}_{a,s-1} \right\}, \quad (116)$$

where the ties are broken arbitrarily in $\operatorname{argmax}\{\}$. We let $n^\circ(a)$ be the a^{th} component of $\mathbf{n}_{1:K}^\circ$. According to the specified arm selection rule, we have $A_T^\circ = \operatorname{argmax}_a n^\circ(a)$ and hence $n^\circ(A_T^\circ) \geq \lfloor \frac{T}{2} \rfloor (> N)$. We prove the claim for the following two cases:

Case 1: If $n^\circ(A^{\text{TS}}) \geq N$, consider an allocation $\mathbf{n}_{1:K}^\dagger$ that is a deviation from the given optimal allocation $\mathbf{n}_{1:K}^\circ$ such that arm A^{TS} gets one pull whereas arm A_T° gets one less pull: i.e., $n^\dagger(A^{\text{TS}}) = n^\circ(A^{\text{TS}}) + 1$, $n^\dagger(A_T^\circ) = n^\circ(A_T^\circ) - 1$, and $n^\dagger(a) = n^\circ(a)$ for any $a \notin \{A^{\text{TS}}, A_T^\circ\}$. The change in the total payoff from this deviation is

$$\sum_{a=1}^K \sum_{i=1}^{n^\dagger(a)} \hat{\mu}_{a,i-1} - \sum_{a=1}^K \sum_{i=1}^{n^\circ(a)} \hat{\mu}_{a,i-1} = \hat{\mu}_{A^{\text{TS}},n^\circ(A^{\text{TS}})} - \hat{\mu}_{A_T^\circ,n^\circ(A_T^\circ)-1} \geq \frac{\Delta}{2} > 0, \quad (117)$$

where the inequality follows from (115) and that $n^\circ(A^{\text{TS}}) \geq N$ and $n^\circ(A_T^\circ) \geq N$. The allocation $\mathbf{n}_{1:K}^\dagger$ is strictly better than $\mathbf{n}_{1:K}^\circ$, which contradicts the assumption that $\mathbf{n}_{1:K}^\circ$ is an optimal allocation.

Case 2: If $n^\circ(A^{\text{TS}}) < N$, consider an allocation $\mathbf{n}_{1:K}^\dagger$ that is a deviation from the given optimal allocation $\mathbf{n}_{1:K}^\circ$ such that arm A_T° gets no more than N pulls whereas arm A^{TS} gets the remains: i.e.,

$$n^\dagger(a) \triangleq \begin{cases} n^\circ(A^{\text{TS}}) + (n^\circ(A_T^\circ) - N) & \text{if } a = A^{\text{TS}}, \\ N & \text{if } a = A_T^\circ, \\ n^\circ(a) & \text{if } a \notin \{A^{\text{TS}}, A_T^\circ\}. \end{cases} \quad (118)$$

By making this the deviation, the total payoff should increase by

$$\sum_{a=1}^K \sum_{i=1}^{n^\dagger(a)} \hat{\mu}_{a,i-1} - \sum_{a=1}^K \sum_{i=1}^{n^\circ(a)} \hat{\mu}_{a,i-1} \quad (119)$$

$$= \sum_{i=n^\circ(A^{\text{TS}})+1}^{n^\circ(A^{\text{TS}})+(n^\circ(A_T^\circ)-N)} \hat{\mu}_{A^{\text{TS}},i-1} - \sum_{i=N+1}^{n^\circ(A_T^\circ)} \hat{\mu}_{A_T^\circ,i-1} \quad (120)$$

$$\geq -(N - n^\circ(A^{\text{TS}})) \cdot 2M + \sum_{i=N+1}^{n^\circ(A_T^\circ)} \hat{\mu}_{A^{\text{TS}},i-1} - \sum_{i=N+1}^{n^\circ(A_T^\circ)} \hat{\mu}_{A_T^\circ,i-1} \quad (121)$$

$$\geq -(N - n^\circ(A^{\text{TS}})) \cdot 2M + (n^\circ(A_T^\circ) - N) \cdot \frac{\Delta}{2} \quad (122)$$

$$\geq (n^\circ(A_T^\circ) - N) \cdot \frac{\Delta}{2} - 2NM. \quad (123)$$

Since $T \geq 2N + \frac{8MN}{\Delta} + 2$ and $n^\circ(A_T^\circ) \geq \lfloor \frac{T}{2} \rfloor$, the last term is strictly positive, which is a contradiction.

We've shown that for almost all $\tilde{\omega}$, when T is large enough, the optimal allocation $\mathbf{n}_{1:K}^\circ$ must allocate more than a half of the pulls on arm $A^{\text{TS}} = \operatorname{argmax}_a \mu_a(\tilde{\theta}_a)$. This concludes the proof.

D.3. Proof of Theorem 2

D.3.1. Proof of “ $W^{\text{TS}}(T, \mathbf{y}) \geq W^{\text{IRS.FH}}(T, \mathbf{y})$ ”

Proof. It immediately follows from Jensen's inequality: since $\max(\dots)$ is a convex function,

$$W^{\text{TS}}(T, \mathbf{y}) = T \times \mathbb{E}_{\mathbf{y}} \left[\max_a \mu_a(\theta_a) \right] \geq T \times \mathbb{E}_{\mathbf{y}} \left[\max_a \mathbb{E}(\mu_a(\theta_a) | \hat{\boldsymbol{\mu}}_{T-1}) \right] = W^{\text{IRS.FH}}(T, \mathbf{y}). \quad (124)$$

■

D.3.2. Proof of “ $W^{\text{IRS.FH}}(T, \mathbf{y}) \geq W^{\text{IRS.V-ZERO}}(T, \mathbf{y})$ ”

Lemma 1 (Variant of Jensen's inequality). *Suppose that $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is an **increasing** (deterministic) function. Then, for any real-valued random variable X such that $\mathbb{E}|X| < \infty$,*

$$\mathbb{E}[\max\{X + \varphi(X), 0\}] \geq \mathbb{E}[\max\{\mathbb{E}(X) + \varphi(X), 0\}]. \quad (125)$$

Proof. Define $\mu \triangleq \mathbb{E}(X)$ and $f_x(t) \triangleq \max\{t + \varphi(x), 0\}$. Since $f_x(\cdot)$ is a convex function for each $x \in \mathbb{R}$,

$$f_x(t) \geq f_x(\mu) + (t - \mu) \cdot f'_x(\mu) = \max\{\mu + \varphi(x), 0\} + (t - \mu) \cdot \mathbf{1}\{\mu + \varphi(x) \geq 0\}, \quad \forall t, \quad \forall x. \quad (126)$$

By setting $t = x$, we get

$$\max\{x + \varphi(x), 0\} = f_x(x) \geq \max\{\mu + \varphi(x), 0\} + (x - \mu) \cdot \mathbf{1}\{\mu + \varphi(x) \geq 0\}, \quad \forall x. \quad (127)$$

Note that, since $\mathbf{1}\{\mu + \varphi(x) \geq 0\}$ is increasing in x , (i) for any $x \geq \mu$, $(x - \mu) \geq 0$ and $\mathbf{1}\{\mu + \varphi(x)\} \geq \mathbf{1}\{\mu + \varphi(\mu)\}$, and (ii) for any $x < \mu$, $(x - \mu) < 0$ and $\mathbf{1}\{\mu + \varphi(x)\} \leq \mathbf{1}\{\mu + \varphi(\mu)\}$. Therefore,

$$(x - \mu) \cdot \mathbf{1}\{\mu + \varphi(x) \geq 0\} \geq (x - \mu) \cdot \mathbf{1}\{\mu + \varphi(\mu) \geq 0\}, \quad \forall x \in \mathbb{R}. \quad (128)$$

Combining this with (127), we get

$$\max\{x + \varphi(x), 0\} \geq \max\{\mu + \varphi(x), 0\} + (x - \mu) \cdot \mathbf{1}\{\mu + \varphi(\mu) \geq 0\}, \quad \forall x \in \mathbb{R}. \quad (129)$$

For random variable X , by taking expectation, we get

$$\mathbb{E}[\max\{X + \varphi(X), 0\}] \geq \mathbb{E}[\max\{\mu + \varphi(X), 0\} + (X - \mu) \cdot \mathbf{1}\{\mu + \varphi(\mu) \geq 0\}] \quad (130)$$

$$\geq \mathbb{E}[\max\{\mu + \varphi(X), 0\}] + \mathbb{E}(X - \mu) \cdot \mathbf{1}\{\mu + \varphi(\mu) \geq 0\} \quad (131)$$

$$= \mathbb{E}[\max\{\mu + \varphi(X), 0\}]. \quad (132)$$

■

Corollary 1. *On a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let $\varphi(x, \omega) : \mathbb{R} \times \Omega \rightarrow \mathbb{R}$ be a function such that (i) the mapping $x \mapsto \varphi(x, \omega)$ is **increasing** for each $\omega \in \Omega$ and (ii) for some sub- σ -field $\mathcal{H} \subseteq \mathcal{F}$, the mapping $\omega \mapsto \varphi(x, \omega)$ is \mathcal{H} -measurable for each $x \in \mathbb{R}$ (i.e., $\varphi(\cdot, \omega)$ is a deterministic function conditioned on \mathcal{H}). Then*

$$\mathbb{E}[\max\{X(\omega) + \varphi(X(\omega), \omega), 0\}] \geq \mathbb{E}[\max\{\mathbb{E}(X|\mathcal{H})(\omega) + \varphi(X(\omega), \omega), 0\}]. \quad (133)$$

Proof. Define

$$\mu(\omega) \triangleq \mathbb{E}(X|\mathcal{H})(\omega), \quad I(\omega) \triangleq \mathbf{1}\{\mu(\omega) + \varphi(\mu(\omega), \omega) \geq 0\}. \quad (134)$$

By (129), we have

$$\max\{x + \varphi(x, \omega), 0\} \geq \max\{\mu(\omega) + \varphi(x, \omega), 0\} + (x - \mu(\omega)) \cdot I(\omega), \quad \forall x \in \mathbb{R}, \quad \text{for each } \omega \in \Omega. \quad (135)$$

Since $\mu(\omega)$ and $I(\omega)$ are \mathcal{H} -measurable,

$$\mathbb{E}[\max\{X(\omega) + \varphi(X(\omega), \omega), 0\}] \geq \mathbb{E}[\max\{\mu(\omega) + \varphi(X(\omega), \omega), 0\} + (X(\omega) - \mu(\omega)) \cdot I(\omega)] \quad (136)$$

$$= \mathbb{E}[\mathbb{E}(\max\{\mu(\omega) + \varphi(X(\omega), \omega), 0\} + (X(\omega) - \mu(\omega)) \cdot I(\omega) | \mathcal{H})] \quad (137)$$

$$= \mathbb{E}[\max\{\mu(\omega) + \varphi(X(\omega), \omega), 0\}] + \mathbb{E}[\mathbb{E}((X(\omega) - \mu(\omega)) \cdot I(\omega) | \mathcal{H})] \quad (138)$$

$$= \mathbb{E}[\max\{\mathbb{E}(X|\mathcal{H})(\omega) + \varphi(X(\omega), \omega), 0\}] \quad (139)$$

$$+ \mathbb{E} \left[\underbrace{(\mathbb{E}(X|\mathcal{H})(\omega) - \mu(\omega)) \cdot I(\omega)}_{=0} \right] \quad (140)$$

$$= \mathbb{E}[\max\{\mathbb{E}(X|\mathcal{H})(\omega) + \varphi(X(\omega), \omega), 0\}]. \quad (141)$$

■

Corollary 2. *On a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let (C_0, \dots, C_T) be \mathcal{H} -measurable real-valued ran-*

dom variables for some sub- σ -field $\mathcal{H} \subseteq \mathcal{F}$ (i.e., C_i 's are constants conditioned on \mathcal{H}). Then

$$\mathbb{E} \left[\max_{0 \leq i \leq T} \left\{ (i-n)^+ \times X + C_i \right\} \right] \geq \mathbb{E} \left[\max_{0 \leq i \leq T} \left\{ \mathbb{E}(X | \mathcal{H}) \cdot \mathbf{1}\{i \geq n+1\} + (i-n-1)^+ \times X + C_i \right\} \right] \quad (142)$$

for any $n = 0, 1, \dots, T$.

Proof. When $n = T$, both sides become $\mathbb{E}[\max_{0 \leq i \leq T} \{C_i\}]$, which makes the claim true. Fix $n < T$ and define

$$\varphi(x, \omega) \triangleq \max_{n+1 \leq i \leq T} \{(i-n-1) \times x + C_i(\omega)\} - \max_{0 \leq i \leq n} \{C_i(\omega)\}. \quad (143)$$

Note that $\varphi(x, \omega)$ satisfies the conditions in Corollary [1](#). By Corollary [1](#),

$$\mathbb{E} \left[\max_{0 \leq i \leq T} \{(i-n)^+ \times X + C_i\} \right] \quad (144)$$

$$= \mathbb{E} \left[\max \left\{ \max_{n+1 \leq i \leq T} \{(i-n) \times X + C_i\}, \max_{0 \leq i \leq n} C_i \right\} \right] \quad (145)$$

$$= \mathbb{E} \left[\max \left\{ X + \max_{n+1 \leq i \leq T} \{(i-n-1) \times X + C_i\}, \max_{0 \leq i \leq n} C_i \right\} \right] \quad (146)$$

$$= \mathbb{E} \left[\max \left\{ X(\omega) + \underbrace{\max_{n+1 \leq i \leq T} \{(i-n-1) \times X(\omega) + C_i(\omega)\} - \max_{0 \leq i \leq n} C_i(\omega)}_{=\varphi(X(\omega), \omega)}, 0 \right\} + \max_{0 \leq i \leq n} C_i(\omega) \right] \quad (147)$$

$$\geq \mathbb{E} \left[\max \left\{ \mathbb{E}(X | \mathcal{H})(\omega) + \max_{n+1 \leq i \leq T} \{(i-n-1) \times X(\omega) + C_i(\omega)\} - \max_{0 \leq i \leq n} C_i(\omega), 0 \right\} + \max_{0 \leq i \leq n} C_i(\omega) \right] \quad (148)$$

$$= \mathbb{E} \left[\max \left\{ \max_{n+1 \leq i \leq T} \{\mathbb{E}(X | \mathcal{H}) + (i-n-1) \times X + C_i\}, \max_{0 \leq i \leq n} C_i \right\} \right] \quad (149)$$

$$= \mathbb{E} \left[\max_{0 \leq i \leq T} \left\{ \mathbb{E}(X | \mathcal{H}) \cdot \mathbf{1}\{i \geq n+1\} + (i-n-1)^+ \times X + C_i \right\} \right]. \quad (150)$$

■

Proof of “ $W^{\text{IRS.FH}}(T, \mathbf{y}) \geq W^{\text{IRS.V-ZERO}}(T, \mathbf{y})$.” Define

$$N_T \triangleq \left\{ \mathbf{n}_{1:K} \in \mathbb{N}_0^K : \sum_{a=1}^K n_a = T \right\} \quad \text{and} \quad S_a(n_a) \triangleq \sum_{i=1}^{n_a} \hat{\mu}_{a,i-1}. \quad (151)$$

What we want to show is

$$W^{\text{IRS.FH}} \equiv \mathbb{E} \left[T \times \max_a \{\hat{\mu}_{a,T-1}\} \right] = \mathbb{E} \left[\max_{\mathbf{n}_{1:K} \in N_T} \left\{ \sum_{a=1}^K n_a \times \hat{\mu}_{a,T-1} \right\} \right] \quad (152)$$

$$\geq \mathbb{E} \left[\max_{\mathbf{n}_{1:K} \in N_T} \left\{ \sum_{a=1}^K S_a(n_a) \right\} \right] \equiv W^{\text{IRS.V-ZERO}}. \quad (153)$$

Further define

$$U_{k,n} \triangleq \mathbb{E} \left[\max_{\mathbf{n}_{1:K} \in N_T} \left\{ \left(\sum_{a=1}^{k-1} S_a(n_a) \right) + (S_k(n_k \wedge n) + (n_k - n)^+ \times \hat{\mu}_{a,T-1}) + \left(\sum_{a=k+1}^K n_a \times \hat{\mu}_{a,T-1} \right) \right\} \right], \quad (154)$$

where $a \wedge b \triangleq \min(a, b)$. Observe that $W^{\text{IRS.FH}} = U_{1,0}$, $W^{\text{IRS.V-ZERO}} = U_{K,T}$, and $U_{k+1,0} = U_{k,T}$. Therefore, it suffices to show that

$$U_{k,n} \geq U_{k,n+1}, \quad \forall k = 1, \dots, K, \quad \forall n = 0, \dots, T-1. \quad (155)$$

Fix k and n . Define a sub- σ -field

$$\mathcal{H} \triangleq \sigma(\{R_{a,s}\}_{a=k,1 \leq s \leq n} \cup \{R_{a,s}\}_{a \neq k, 1 \leq s \leq T-1}). \quad (156)$$

For each $i = 0, \dots, T$, define

$$C_i \triangleq \max \left\{ \left(\sum_{a=1}^{k-1} S_a(n_a) \right) + S_k(i \wedge n) + \left(\sum_{a=k+1}^K n_a \times \hat{\mu}_{a,T-1} \right) : \mathbf{n}_{1:K} \in N_T, n_k = i \right\}. \quad (157)$$

Note that C_i 's are \mathcal{H} -measurable and

$$U_{k,n} = \mathbb{E} \left[\max_{0 \leq i \leq T} \{(i - n)^+ \times \hat{\mu}_{k,T-1} + C_i\} \right]. \quad (158)$$

With $X \triangleq \hat{\mu}_{a,T-1}$,

$$U_{k,n} = \mathbb{E} \left[\max_{0 \leq i \leq T} \{(i - n)^+ \times X + C_i\} \right] \quad (159)$$

$$\stackrel{\text{Corollary 2}}{\geq} \mathbb{E} \left[\max_{0 \leq i \leq T} \{\mathbb{E}(X | \mathcal{H}) \cdot \mathbf{1}\{i \geq n+1\} + (i - n - 1)^+ \times X + C_i\} \right] \quad (160)$$

$$\stackrel{(a)}{=} \mathbb{E} \left[\max_{0 \leq i \leq T} \{\hat{\mu}_{k,n} \cdot \mathbf{1}\{i \geq n+1\} + (i - n - 1)^+ \times \hat{\mu}_{a,T-1} + C_i\} \right] \quad (161)$$

$$\stackrel{(b)}{=} U_{k,n+1}. \quad (162)$$

Equation (a) holds since $\mathbb{E}(X | \mathcal{H}) = \mathbb{E}(\hat{\mu}_{k,T-1} | \mathcal{H}) = \mathbb{E}(\hat{\mu}_{k,T-1} | R_{k,1}, \dots, R_{k,n}) = \hat{\mu}_{a,n}$, and equation (b) holds since $S_k(i \wedge n) + \hat{\mu}_{k,n} \cdot \mathbf{1}\{i \geq n+1\} = \sum_{s=1}^n \hat{\mu}_{k,s-1} \cdot \mathbf{1}\{i \geq s\} + \hat{\mu}_{k,n} \cdot \mathbf{1}\{i \geq n+1\} = \sum_{s=1}^{n+1} \hat{\mu}_{k,s-1} \cdot \mathbf{1}\{i \geq s\} = S_k(i \wedge (n+1))$. ■

A note on the proof. One may wonder if the above result can be derived in a simpler way by exploiting the properties of nested filtration (e.g., Proposition 2.3 of [Brown et al., 2010](#)). Unlike the proof of $W^{\text{TS}} \geq W^{\text{IRS.FH}}$, however, the proof of $W^{\text{IRS.FH}} \geq W^{\text{IRS.V-ZERO}}$ does not simply follow from the fact that $\sigma(\hat{\mu}_{T-1})$ is larger than $\sigma(H_{t-1})$.

Consider a Bernoulli MAB with $K = 2$, $T = 2$, and a prior distribution Beta(1,1), and let us

introduce its variation whose reward function is given by $r'_t(\cdot)$ as follows:

$$r'_1(a_1) = r_1(a_1), \quad r'_2(\mathbf{a}_{1:2}) = -\kappa r_2(\mathbf{a}_{1:2}), \quad (163)$$

where $r_t(\cdot)$ is the reward function of the original Bernoulli MAB. When $\kappa > 0$, one can show that

$$W^{\text{IRS.FH}} = \mathbb{E} \left[\max_{\mathbf{a}_{1:T}} \left\{ \sum_{t=1}^T \mathbb{E}(r'_t(\mathbf{a}_{1:t}) | \hat{\boldsymbol{\mu}}_{T-1}) \right\} \right] = \frac{7}{12} - \frac{5}{12} \kappa, \quad (164)$$

$$W^{\text{IRS.V-ZERO}} = \mathbb{E} \left[\max_{\mathbf{a}_{1:T}} \left\{ \sum_{t=1}^T \mathbb{E}(r'_t(\mathbf{a}_{1:t}) | H_{t-1}) \right\} \right] = \frac{1}{2} - \frac{3}{8} \kappa. \quad (165)$$

If κ is large enough, we obtain $W^{\text{IRS.FH}} < W^{\text{IRS.V-ZERO}}$, which is opposite to the above result.

Recall that the additional gain from knowing the future information can be decomposed into two components; the gain from knowing the immediate reward and the gain from knowing the next belief state, where IRS.V-ZERO considers the former component only. When those two components are not aligned as in this example (i.e., a higher r'_1 leads to a worse next belief state), the DM can exploit the penalties if they penalize only for the first component (e.g., when r'_1 is smaller than expected, the DM will get compensated for this difference but she can still earn the larger reward in the next period).

This is also related to the fact that $z_t^{\text{IRS.V-ZERO}}$ does not correspond to zero penalty under the some (partial) information relaxation, but should be understood as an approximation of z_t^{ideal} under the perfect information relaxation. As opposed to TS and IRS.FH, the optimal solution to the IRS.V-ZERO's inner problem may depend on the entire outcome ω . With the terminology of [Brown et al. \(2010\)](#), there is a mismatch between the filtration that generates the penalties and the filtration that characterizes the relaxed policy space.

D.3.3. Proof of “ $W^{\text{TS}}(T, \mathbf{y}) \geq W^{\text{IRS.V-EMAX}}(T, \mathbf{y})$ ”

To show that $W^{\text{TS}} \geq W^{\text{IRS.V-EMAX}}$, we take a completely different approach that utilizes Theorem 4 in [Desai et al. \(2012a\)](#). We here rephrase the definition and the theorem therein using our notation.

Definition 2 (Supersolution). *An approximate value function $\hat{V} : \mathbb{N}_0 \times \mathcal{Y} \rightarrow \mathbb{R}$ is a **supersolution** to the Bellman equation [\(15\)](#) if*

$$\hat{V}(T, \mathbf{y}) \geq \max_{a \in \mathcal{A}} \left\{ \mathbb{E}_{y_a} \left[R_{a,1} + \hat{V}(T-1, \mathcal{U}(\mathbf{y}, R_{a,1}, r)) \right] \right\}, \quad \forall \mathbf{y} \in \mathcal{Y}, \quad \forall T \geq 1, \quad (166)$$

with $\hat{V}(0, \mathbf{y}) = 0$ for all $\mathbf{y} \in \mathcal{Y}$.

Remark 8. If $\widehat{V}(\cdot, \cdot)$ is a supersolution, then for any given ω , T , and \mathbf{y} ,

$$\widehat{V}(T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})) \geq \mathbb{E}_{\mathbf{y}} \left[r_t(\mathbf{a}_{1:t-1} \oplus a, \omega; \mathbf{y}) + \widehat{V}(T-t, \mathbf{y}_t(\mathbf{a}_{1:t-1} \oplus a, \omega; \mathbf{y})) \middle| H_{t-1}(\mathbf{a}_{1:t-1}, \omega) \right], \quad (167)$$

for all $a \in \mathcal{A}$, $\mathbf{a}_{1:t-1} \in \mathcal{A}^{t-1}$ and $t \in \{1, \dots, T\}$.

Lemma 2 (Theorem 4 of [Desai et al. \(2012a\)](#), rephrased). Consider a penalty function \hat{z}_t generated by $\widehat{V}(\cdot, \cdot)$:

$$\begin{aligned} \hat{z}_t(\mathbf{a}_{1:t}, \omega; T, \mathbf{y}) &\triangleq r_t(\mathbf{a}_{1:t}, \omega) - \mathbb{E}_{\mathbf{y}} [r_t(\mathbf{a}_{1:t}, \omega) | H_{t-1}(\mathbf{a}_{1:t-1}, \omega)] \\ &\quad + \widehat{V}(T-t, \mathbf{y}_t(\mathbf{a}_{1:t}, \omega; \mathbf{y})) - \mathbb{E}_{\mathbf{y}} \left[\widehat{V}(T-t, \mathbf{y}_t(\mathbf{a}_{1:t}, \omega; \mathbf{y})) \middle| H_{t-1}(\mathbf{a}_{1:t-1}, \omega) \right]. \end{aligned} \quad (168)$$

If $\widehat{V}(\cdot, \cdot)$ is a supersolution, then the performance bound induced by penalty function \hat{z}_t is tighter than \widehat{V} : i.e.,

$$W^{\hat{z}}(T, \mathbf{y}) \leq \widehat{V}(T, \mathbf{y}). \quad (169)$$

And this holds in a stronger sense: for each outcome ω , the maximal value of the inner problem with respect to ω (denoted by $V_1^{\hat{z}, \text{in}}(\emptyset, \omega; T, \mathbf{y})$ in the proof) is smaller than or equal to $\widehat{V}(T, \mathbf{y})$.

Proof. Let $V_t^{\hat{z}, \text{in}}(\cdot)$ be the DP solution of inner problem [\(*\)](#) for a given penalty \hat{z}_t with respect to a particular outcome ω :

$$V_t^{\hat{z}, \text{in}}(\mathbf{a}_{1:t-1}, \omega; T, \mathbf{y}) = \max_{a \in \mathcal{A}} \left\{ r_t(\mathbf{a}_{1:t-1} \oplus a, \omega) - \hat{z}_t(\mathbf{a}_{1:t-1} \oplus a, \omega; T, \mathbf{y}) + V_{t+1}^{\hat{z}, \text{in}}(\mathbf{a}_{1:t-1} \oplus a, \omega; T, \mathbf{y}) \right\}, \quad (170)$$

with $V_{T+1}^{\hat{z}, \text{in}}(\cdot, \omega; T, \mathbf{y}) = 0$. Then, we have $W^{\hat{z}}(T, \mathbf{y}) = \mathbb{E} [V_1^{\hat{z}, \text{in}}(\emptyset, \omega; T, \mathbf{y})]$. To prove the claim, it suffices to show that, for any given ω ,

$$V_t^{\hat{z}, \text{in}}(\mathbf{a}_{1:t-1}, \omega; T, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})) \leq \widehat{V}(T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})), \quad (171)$$

for all $\mathbf{a}_{1:t-1} \in \mathcal{A}^{t-1}$ and for all $t = 1, \dots, T+1$.

We argue by induction. As a terminal case, when $t = T+1$, the inequality [\(171\)](#) holds trivially since both sides are zero. Fix t and suppose that the inequality [\(171\)](#) holds for $t+1$. Omitting ω

for brevity, we get

$$\widehat{V}(T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1})) - V_t^{\hat{z}, \text{in}}(\mathbf{a}_{1:t-1}; T, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1})) \quad (172)$$

$$= \widehat{V}(T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1})) - \max_{a \in \mathcal{A}} \left\{ r_t(\mathbf{a}_{1:t-1} \oplus a) - \hat{z}_t(\mathbf{a}_{1:t-1} \oplus a; T, \mathbf{y}) + V_{t+1}^{\hat{z}, \text{in}}(\mathbf{a}_{1:t-1} \oplus a; T, \mathbf{y}) \right\} \quad (173)$$

$$= \min_{a \in \mathcal{A}} \left\{ \underbrace{\widehat{V}(T-t, \mathbf{y}_t(\mathbf{a}_{1:t})) - V_{t+1}^{\hat{z}, \text{in}}(\mathbf{a}_{1:t-1} \oplus a; T, \mathbf{y})}_{\geq 0 \text{ } (\because \text{induction hypothesis})} + \underbrace{\widehat{V}(T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1})) - \mathbb{E} \left[r_t(\mathbf{a}_{1:t-1} \oplus a) + \widehat{V}(T-t, \mathbf{y}_t(\mathbf{a}_{1:t-1} \oplus a)) \middle| H_{t-1} \right]}_{\geq 0 \text{ } (\because \text{Remark 8})} \right\} \quad (174)$$

$$\geq 0. \quad (175)$$

■

Proof of “ $W^{\text{TS}}(T, \mathbf{y}) \geq W^{\text{IRS.V-EMAX}}(T, \mathbf{y})$.” Recall that $z_t^{\text{IRS.V-EMAX}}$ is a penalty function generated by W^{TS} . We observe that $W^{\text{TS}}(\cdot, \cdot)$ is a supersolution: for any T and \mathbf{y} ,

$$W^{\text{TS}}(T, \mathbf{y}) = \mathbb{E}_{\mathbf{y}} \left[T \times \max_{a \in \mathcal{A}} \mu_a(\theta_a) \right] \quad (176)$$

$$= \mathbb{E}_{\mathbf{y}} \left[\max_{a \in \mathcal{A}} \mu_a(\theta_a) \right] + W^{\text{TS}}(T-1, \mathbf{y}) \quad (177)$$

$$\geq \max_{a \in \mathcal{A}} \left\{ \mathbb{E}_{y_a} [\mu_a(\theta_a)] + W^{\text{TS}}(T-1, \mathbf{y}) \right\} \quad (178)$$

$$= \max_{a \in \mathcal{A}} \left\{ \mathbb{E}_{\mathbf{y}} \left[R_{a,1} + W^{\text{TS}}(T-1, \mathbf{y}) \right] \right\} \quad (179)$$

$$= \max_{a \in \mathcal{A}} \left\{ \mathbb{E}_{\mathbf{y}} \left[R_{a,1} + W^{\text{TS}}(T-1, \mathcal{U}(\mathbf{y}, a, R_{a,1})) \right] \right\}. \quad (180)$$

The last equality holds since $\mathbb{E} \left[W^{\text{TS}}(T-1, \mathcal{U}(\mathbf{y}, a_1, r_1(a_1, \omega))) \right] = W^{\text{TS}}(T-1, \mathbf{y})$, as argued in (39). By Lemma 2, we have $W^{\text{IRS.V-EMAX}}(T, \mathbf{y}) \leq W^{\text{TS}}(T, \mathbf{y})$ which also holds in a stronger sense. ■

D.4. Proof of Theorem 3

D.4.1. Suboptimality Decomposition

As in §C.1, we define the Q-values of the inner problem given a particular outcome ω , a penalty function $z_t(\cdot)$, a time horizon T , and a prior belief \mathbf{y} .

$$Q_t^{z, \text{in}}(\mathbf{a}_{1:t-1}, a, \omega; T, \mathbf{y}) = r_t(\mathbf{a}_{1:t-1} \oplus a, \omega) - z_t(\mathbf{a}_{1:t-1} \oplus a, \omega; T, \mathbf{y}) + V_{t+1}^{z, \text{in}}(\mathbf{a}_{1:t-1} \oplus a, \omega; T, \mathbf{y}), \quad (181)$$

$$V_t^{z, \text{in}}(\mathbf{a}_{1:t-1}, \omega; T, \mathbf{y}) = \max_{a \in \mathcal{A}} \left\{ Q_t^{z, \text{in}}(\mathbf{a}_{1:t-1}, a, \omega; T, \mathbf{y}) \right\}, \quad (182)$$

with $V_{T+1}^{z,\text{in}}(\cdot, \omega; T, \mathbf{y}) \equiv 0$. Additionally define the total payoff of an action sequence and the hindsight best action under penalties:

$$\mathcal{S}^z(\mathbf{a}_{1:T}, \omega; T, \mathbf{y}) \triangleq \sum_{t=1}^T r_t(\mathbf{a}_{1:t}, \omega) - z_t(\mathbf{a}_{1:t}, \omega; T, \mathbf{y}), \quad (183)$$

$$a_t^{z,*}(\mathbf{a}_{1:t-1}, \omega; T, \mathbf{y}) \triangleq \operatorname{argmax}_{a \in \mathcal{A}} \left\{ Q_t^{z,\text{in}}(\mathbf{a}_{1:t-1}, a, \omega; T, \mathbf{y}) \right\}. \quad (184)$$

We have $V_1^{z,\text{in}}(\emptyset, \omega; T, \mathbf{y}) = \max_{\mathbf{a}_{1:T} \in \mathcal{A}^T} \mathcal{S}^z(\mathbf{a}_{1:T}, \omega; T, \mathbf{y})$.

Proposition 6 (Suboptimality decomposition). *Given a non-anticipating policy $\pi \in \Pi_{\mathbb{F}}$ and a dual-feasible penalty function z_t , the suboptimality gap is the sum of the instantaneous suboptimality of individual actions taken by π along the sample path: i.e.,*

$$W^z(T, \mathbf{y}) - V(\pi, T, \mathbf{y}) = \mathbb{E}_{\mathbf{y}} \left[\max_{\mathbf{a}_{1:T}} \{ \mathcal{S}^z(\mathbf{a}_{1:T}, \omega; T, \mathbf{y}) \} - \mathcal{S}^z(\mathbf{A}_{1:T}^{\pi}, \omega; T, \mathbf{y}) \right] \quad (185)$$

$$= \mathbb{E}_{\mathbf{y}} \left[\sum_{t=1}^T \max_a \left\{ Q_t^{z,\text{in}}(\mathbf{A}_{1:t-1}^{\pi}, a, \omega; T, \mathbf{y}) \right\} - Q_t^{z,\text{in}}(\mathbf{A}_{1:t-1}^{\pi}, A_t^{\pi}, \omega; T, \mathbf{y}) \right], \quad (186)$$

where the expectation is taken with respect to the randomness of outcome ω and the randomness of policy π .

Proof. The first equality immediately follows from the definition of W^z and mean equivalence (Proposition 3). Now fix ω , T , and \mathbf{y} . Consider the (pathwise) suboptimality of the action sequence $\mathbf{A}_{1:T}^{\pi}$ compared to the clairvoyant optimal solution. It can be decomposed into the instantaneous suboptimality incurred by the individual action at each time:

$$\max_{\mathbf{a}_{1:T}} \{ \mathcal{S}^z(\mathbf{a}_{1:T}) \} - \mathcal{S}^z(\mathbf{A}_{1:T}^{\pi}) = \sum_{t=1}^T \max_a \left\{ Q_t^{z,\text{in}}(\mathbf{A}_{1:t-1}^{\pi}, a) \right\} - Q_t^{z,\text{in}}(\mathbf{A}_{1:t-1}^{\pi}, A_t^{\pi}). \quad (187)$$

By taking expectation, we obtain the second equality. ■

The next lemma shows that the instantaneous suboptimality of the first action can be expressed in terms of mean reward metrics for each of the IRS penalty functions.

Lemma 3. *Fix time horizon T , prior belief \mathbf{y} , and the true outcome ω , and hide the dependency on them in notation for $Q_1^{z,\text{in}}(\cdot)$, $a_1^{z,*}(\cdot)$, $\mu_a(\cdot)$ and $\hat{\mu}_{a,n}(\cdot)$. For each of the penalty functions z^{TS} , $z^{\text{IRS.FH}}$, and $z^{\text{IRS.V-ZERO}}$, the instantaneous suboptimality of action $a \in \mathcal{A}$ satisfies the following:*

(1) When $z \equiv z^{TS}$,

$$Q_1^{z,\text{in}}(a_1^{z,*}) - Q_1^{z,\text{in}}(a) = \mu_{a_1^{z,*}} - \mu_a. \quad (188)$$

(2) When $z \equiv z^{\text{IRS.FH}}$,

$$Q_1^{z,\text{in}}(a_1^{z,*}) - Q_1^{z,\text{in}}(a) = \hat{\mu}_{a_1^{z,*}, T-1} - \hat{\mu}_{a, T-1}. \quad (189)$$

(3) When $z \equiv z^{V-ZERO}$,

$$Q_1^{z,\text{in}}(a_1^{z,*}) - Q_1^{z,\text{in}}(a) \leq \max_{0 \leq n \leq T-1} \left\{ \hat{\mu}_{a_1^{z,*},n} \right\} - \hat{\mu}_{a,0}. \quad (190)$$

Proof. (1) When $z \equiv z^{\text{TS}}$, we have

$$Q_1^{z,\text{in}}(a) = \mu_a + (T-1) \times \max_{a'} \mu_{a'}. \quad (191)$$

Since the last term does not depend on action a , the claim follows.

(2) When $z \equiv z^{\text{IRS.FH}}$, we obtain the claim by replacing μ_a with $\hat{\mu}_{a,T-1}$ in the above proof.

(3) When $z \equiv z^{\text{IRS.V-ZERO}}$, recall that the associated inner problem is to find an optimal allocation: i.e.,

$$\max_{\mathbf{n}_{1:K} \in N_T} \left\{ \sum_{a=1}^K \sum_{i=0}^{n_a-1} \hat{\mu}_{a,i} \right\}. \quad (192)$$

Let $\mathbf{n}_{1:K}^*$ be the optimal allocation. Observe that the suboptimality is incurred only when $n_a^* = 0$, it is no worse than $\hat{\mu}_{a^*,n_{a^*}^*} - \hat{\mu}_{a,0}$ (the loss if the payoff when pulling a one more time but pulling $a_1^{z,*}$ one less time). Since $n_{a^*}^* \leq T-1$, the claim follows. \blacksquare

D.4.2. Recursive Structure of IRS Penalty Functions

To describe the recursive structure of Bayesian MAB problems explicitly, we define a shift operator $\mathcal{M}_t : \mathcal{A}^t \times \Omega \rightarrow \Omega$,

$$\mathcal{M}_t(\mathbf{a}_{1:t}, \omega) \triangleq (R_{a,n_a}; \forall n_a > n_t(\mathbf{a}_{1:t}, a), \forall a \in \mathcal{A}). \quad (193)$$

The shifted outcome $\mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega)$ encodes the remaining reward realizations after taking $\mathbf{a}_{1:t-1}$.

Remark 9 (Recursive structure of remaining uncertainties). *Conditioned on $\mathcal{H}_{t-1}(\mathbf{a}_{1:t-1}, \omega)$, the remaining uncertainties are sufficiently described by $\mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})$, i.e.,*

$$\mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega) | \mathcal{H}_{t-1}(\mathbf{a}_{1:t-1}, \omega) \sim \mathcal{I}(\mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})). \quad (194)$$

Remark 10 (Recursive structure of IRS penalties). *Each of penalty functions (22)–(26) has the following form:*

$$z_t(\mathbf{a}_{1:t}, \omega; T, \mathbf{y}) = \varphi^z(\mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega), T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})), \quad (195)$$

for some function $\varphi^z : \Omega \times \mathbb{N} \times \mathcal{Y} \rightarrow \mathbb{R}$, i.e., the penalty at each time is completely determined by the remaining rewards $\mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega)$, the remaining time horizon $T-t+1$, and the prior belief $\mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega)$ at that moment.

Remark 9 immediately follows from Bayes' rule, and Remark 10 can be easily verified. We

observe the recursive structure of the sequential inner problems that the DM solves throughout the decision-making process, which can be characterized by the following property.

Proposition 7 (Generalized posterior sampling). *For each of penalty functions (22)–(26), the IRS policy π is randomized in such a way that it takes an action a with the probability that the action a is indeed the best action $a_t^{z,*}$ at that moment, i.e.,*

$$\mathbb{P}[A_t^\pi = a | \mathcal{F}_{t-1}] = \mathbb{P}[a_t^{z,*}(\mathbf{A}_{1:t-1}^\pi, \omega) = a | \mathcal{F}_{t-1}], \quad \forall a, \quad \forall t. \quad (196)$$

The source of uncertainty in the LHS is the randomness of the policy (embedded in $\tilde{\omega}$) and that in the RHS is the randomness of nature (embedded in ω). Here we assume that the tie-breaking rule in argmax of (184) is identical to the one used when π^z solves the inner problem.

Proof. Observe that the IRS's action A_t^π can be represented as

$$A_t^\pi = a_1^{z,*}(\emptyset, \tilde{\omega}; T - t + 1, \mathbf{y}_{t-1}(\mathbf{A}_{1:t-1}^\pi, \omega; \mathbf{y})), \quad (197)$$

where $\tilde{\omega} \sim \mathcal{I}(\mathbf{y}_{t-1}(\mathbf{A}_{1:t-1}^\pi, \omega; \mathbf{y}))$, i.e., the action that the clairvoyant DM will take in an MAB instance specified by horizon $T - t + 1$, prior belief $\mathbf{y}_{t-1}(\mathbf{A}_{1:t-1}^\pi, \omega; \mathbf{y})$, and the outcome $\tilde{\omega}$. Therefore, it suffices to verify that the inner problem that π solves at time t is identically distributed with the sub-inner problem with respect to ground-truth ω (i.e., the subproblem given the past action sequence $\mathbf{A}_{1:t-1}^\pi$).

Fix time t , past actions $\mathbf{a}_{1:t-1} = \mathbf{A}_{1:t-1}^\pi$, and the true outcome ω . The sub-inner problem determining $a_t^{z,*}(\mathbf{a}_{1:t-1}, \omega)$ is

$$\max_{\mathbf{a}'_{t:T}} \left\{ \sum_{s=t}^T r_s(\mathbf{a}_{1:t-1} \oplus \mathbf{a}'_{t:s}, \omega) - z_s(\mathbf{a}_{1:t-1} \oplus \mathbf{a}'_{t:s}, \omega; T, \mathbf{y}) \right\}. \quad (198)$$

By Remark 10, for any $s \in \{t, \dots, T\}$, the penalty at (inner) time s is given by

$$z_s(\mathbf{a}_{1:t-1} \oplus \mathbf{a}'_{t:s}, \omega; T, \mathbf{y}) \quad (199)$$

$$= \varphi^z(\mathcal{M}_{s-1}(\mathbf{a}_{1:t-1} \oplus \mathbf{a}'_{t:s-1}, \omega), T - s + 1, \mathbf{y}_{s-1}(\mathbf{a}_{1:t-1} \oplus \mathbf{a}'_{t:s-1}, \omega; \mathbf{y})) \quad (200)$$

$$= \varphi^z \left(\begin{array}{c} \mathcal{M}_{s-t}(\mathbf{a}'_{t:s-1}, \mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega)), \\ (T - t + 1) - (s - t), \\ \mathbf{y}_{s-t}(\mathbf{a}'_{t:s-1}, \mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega); \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})) \end{array} \right) \quad (201)$$

$$= z_{s-t+1}(\mathbf{a}'_{t:s}, \mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega); T - t + 1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})). \quad (202)$$

For rewards, similarly, we have $r_s(\mathbf{a}_{1:t-1} \oplus \mathbf{a}'_{t:s}, \omega) = r_{s-t+1}(\mathbf{a}'_{t:s}, \mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega))$. Therefore, the

sub-inner problem (198) is reformulated as

$$\max_{\mathbf{a}'_{t:T}} \left\{ \sum_{s=t}^T r_{s-t+1}(\mathbf{a}'_{t:s}, \mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega)) - z_{s-t+1}(\mathbf{a}'_{t:s}, \mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega); T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})) \right\}. \quad (203)$$

Given the fact that the shifted outcome $\mathcal{M}_{t-1}(\mathbf{a}_{1:t-1}, \omega)$ and the sampled outcome $\tilde{\omega}$ are identically distributed with $\mathcal{I}(\mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y}))$ conditionally on $H_{t-1}(\mathbf{a}_{1:t-1}, \omega)$ (Remark 9), this sub-inner problem follows the same distribution with

$$\max_{\mathbf{a}'_{1:T-t+1}} \left\{ \sum_{s=1}^{T-t+1} r_s(\mathbf{a}'_{1:s}, \tilde{\omega}) - z_s(\mathbf{a}'_{1:s}, \tilde{\omega}, T-t+1, \mathbf{y}_{t-1}(\mathbf{a}_{1:t-1}, \omega; \mathbf{y})) \right\}, \quad (204)$$

which characterizes the IRS's action A_t^π . Therefore, $a_t^{z,*}(\mathbf{A}_{1:t-1}^\pi, \omega)$ is identically distributed with A_t^π conditioned on \mathcal{F}_{t-1} . ■

Remark 11. Utilizing the recursive structure of IRS penalty functions, Lemma 3 can be extended to describe the instantaneous suboptimality of the t^{th} action. Fix true outcome ω and past actions $\mathbf{a}_{1:t-1}$, and hide the dependency on them in notation for $Q_t^{z,\text{in}}(\cdot)$, $a_t^{z,*}(\cdot)$, $n_t(\cdot)$, $\mu_a(\cdot)$ and $\hat{\mu}_{a,n}(\cdot)$.

(1) When $z \equiv z^{\text{TS}}$,

$$Q_t^{z,\text{in}}(a_t^{z,*}) - Q_t^{z,\text{in}}(a) = \mu_{a_t^{z,*}} - \mu_a. \quad (205)$$

(2) When $z \equiv z^{\text{IRS.FH}}$,

$$Q_t^{z,\text{in}}(a_t^{z,*}) - Q_t^{z,\text{in}}(a) = \hat{\mu}_{a_t^{z,*}, n_{t-1}(a_t^{z,*})+T-t} - \hat{\mu}_{a, n_{t-1}(a)+T-t}. \quad (206)$$

(3) When $z \equiv z^{\text{V-ZERO}}$,

$$Q_t^{z,\text{in}}(a_t^{z,*}) - Q_t^{z,\text{in}}(a) \leq \max_{0 \leq n \leq T-t} \left\{ \hat{\mu}_{a_t^{z,*}, n_{t-1}(a_t^{z,*})+n} \right\} - \hat{\mu}_{a, n_{t-1}(a)}. \quad (207)$$

D.4.3. Preliminary Lemmas on MAB with Natural Exponential Family Distributions

We first describe the notion of sub-Gaussian random variable as an effective tool for bounding its tail behavior.

Definition 3 (Sub-Gaussian random variable). A random variable X is σ -sub-Gaussian if

$$\mathbb{E}[\exp(\lambda(X - \mathbb{E}X))] \leq \exp\left(\frac{\sigma\lambda^2}{2}\right), \quad \forall \lambda \in \mathbb{R}, \quad (208)$$

for some $\sigma > 0$.

Lemma 4. Given a random variable X , suppose that there exists $\sigma > 0$ such that

$$\mathbb{P}[X \geq \mathbb{E}X + z\sigma] \leq e^{-z^2/2}, \quad \forall z \geq 0. \quad (209)$$

Then, the following holds:

$$\mathbb{E} \left[(X - (\mathbb{E}X + z\sigma))^+ \right] \leq \frac{\sigma}{z} e^{-z^2/2}, \quad \forall z > 0. \quad (210)$$

Corollary 3. *If a random variable X is σ -sub-Gaussian, it satisfies the condition of Lemma 4 and hence the inequality (210) holds.*

Proof. With $\mu \triangleq \mathbb{E}X$, we have

$$\mathbb{E} \left[(X - (\mu + z\sigma))^+ \right] = \int_{x=\mu+z\sigma}^{\infty} \mathbb{P}[X \geq x] dx = \int_{t=z}^{\infty} \mathbb{P}[X \geq \mu + t\sigma] \sigma dt \leq \sigma \int_{t=z}^{\infty} e^{-t^2/2} dt. \quad (211)$$

Utilizing the tail bound established for the standard normal distribution, we can show that

$$\int_{t=z}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \leq \frac{1}{z} \frac{e^{-z^2/2}}{\sqrt{2\pi}}. \quad (212)$$

By combining these two inequalities, we obtain the desired result.

The corollary simply follows from Markov inequality: for any $z \geq 0$ and $\lambda \geq 0$, we have

$$\mathbb{P}[X \geq \mu + z\sigma] = \mathbb{P} \left[e^{\lambda(X-\mu)} \geq e^{\lambda z\sigma} \right] \leq \frac{\mathbb{E}[e^{\lambda(X-\mu)}]}{e^{\lambda z\sigma}} \leq \exp \left(\frac{\sigma^2 \lambda^2}{2} - \lambda z\sigma \right). \quad (213)$$

By taking $\lambda = \frac{z}{\sigma}$, it follows that $\mathbb{P}[X \geq \mu + z\sigma] \leq e^{-z^2/2}$. ■

We now return to the context of MAB problems and show that the mean reward metrics are sub-Gaussian.

Lemma 5 (Sub-Gaussianity of mean reward metrics). *Consider the setting of Theorem 3, i.e., the reward distribution of arm a is described by an L -smooth log-partition function $A_a(\theta_a)$ and hyperparameters (ξ_a, ν) . Then, the conditional mean reward μ_a is $\sqrt{L/\nu}$ -sub-Gaussian: i.e.,*

$$\mathbb{E}_{(\xi_a, \nu)} [\exp(\lambda(\mu_a - \bar{\mu}_a))] \leq \exp \left(\frac{L\lambda^2}{2\nu} \right), \quad \forall \lambda \in \mathbb{R}, \quad (214)$$

where $\bar{\mu}_a = \mathbb{E}_{(\xi_a, \nu)}[\mu_a] = \frac{\xi_a}{\nu}$ is the prior predictive mean reward (i.e., the unconditional mean reward). Furthermore, the posterior predictive mean reward $\hat{\mu}_{a,n}$ is $\sqrt{\frac{Ln}{\nu(\nu+n)}}$ -sub-Gaussian: i.e.,

$$\mathbb{E}_{(\xi_a, \nu)} [\exp(\lambda(\hat{\mu}_{a,n} - \bar{\mu}_a))] \leq \exp \left(\frac{\lambda^2}{2} \times \frac{Ln}{\nu(\nu+n)} \right), \quad \forall \lambda \in \mathbb{R}. \quad (215)$$

Proof. We first prove that μ_a is $\sqrt{L/\nu}$ -sub-Gaussian. Due to L -smoothness condition, $A_a(\theta_a)$ is

finite valued for all $\theta_a \in \mathbb{R}$. For any $\lambda \in \mathbb{R}$, we have

$$\mathbb{E}_{(\xi_a, \nu)} [\exp(\lambda \mu_a)] \stackrel{(i)}{=} \mathbb{E}_{(\xi_a, \nu)} [\exp(\lambda A'_a(\theta_a))] \quad (216)$$

$$= \int_{-\infty}^{\infty} \exp(\lambda A'_a(\theta_a)) \times f_a(\xi_a, \nu) \exp(\xi_a \theta_a - \nu A_a(\theta_a)) d\theta_a \quad (217)$$

$$= \int_{-\infty}^{\infty} f_a(\xi_a, \nu) \exp\{\xi_a \theta_a - \nu A_a(\theta_a) + \lambda A'_a(\theta_a)\} d\theta_a \quad (218)$$

$$= \int_{-\infty}^{\infty} f_a(\xi_a, \nu) \exp\{\xi_a \theta_a - \nu (A_a(\theta_a) - \lambda/\nu \cdot A'_a(\theta_a))\} d\theta_a \quad (219)$$

$$\stackrel{(ii)}{\leq} \int_{-\infty}^{\infty} f_a(\xi_a, \nu) \exp\left\{\xi_a \theta_a - \nu \left(A_a(\theta_a - \lambda/\nu) - \frac{L\lambda^2}{2\nu^2}\right)\right\} d\theta_a \quad (220)$$

$$= \exp\left(\frac{L\lambda^2}{2\nu}\right) \times \int_{-\infty}^{\infty} f_a(\xi_a, \nu) \exp\{\xi_a \theta_a - \nu A_a(\theta_a - \lambda/\nu)\} d\theta_a \quad (221)$$

$$= \exp\left(\frac{\xi_a \lambda}{\nu} + \frac{L\lambda^2}{2\nu}\right) \times \int_{-\infty}^{\infty} f_a(\xi_a, \nu) \exp\{\xi_a(\theta_a - \lambda/\nu) - \nu A_a(\theta_a - \lambda/\nu)\} d\theta_a \quad (222)$$

$$= \exp\left(\frac{\xi_a \lambda}{\nu} + \frac{L\lambda^2}{2\nu}\right) \times \int_{-\infty}^{\infty} f_a(\xi_a, \nu) \exp\{\xi_a \theta_a - \nu A_a(\theta_a)\} d\theta_a \quad (223)$$

$$= \exp\left(\frac{\xi_a \lambda}{\nu} + \frac{L\lambda^2}{2\nu}\right), \quad (224)$$

where we have utilized that (i) $\mu_a(\theta_a) = A'_a(\theta_a)$ and (ii) $A_a(\theta_a + \delta) \leq A_a(\theta_a) + \delta A'_a(\theta_a) + \frac{L}{2}\delta^2$. Since $\bar{\mu}_a = \xi_a/\nu$, we obtained the desired result.

Next we focus on the posterior predictive mean reward $\hat{\mu}_{a,n}$. Recall that we have

$$\hat{\mu}_{a,n} = \frac{\xi_a + \sum_{i=1}^n R_{a,i}}{\nu + n}. \quad (225)$$

For any $\lambda \in \mathbb{R}$, we have

$$\mathbb{E}_{(\xi_a, \nu)} \left[\exp \left(\lambda \sum_{i=1}^n R_{a,i} \right) \right] = \mathbb{E}_{(\xi_a, \nu)} \left[\mathbb{E} \left\{ \exp \left(\lambda \sum_{i=1}^n R_{a,i} \right) \middle| \theta_a \right\} \right] \quad (226)$$

$$\stackrel{(i)}{=} \mathbb{E}_{(\xi_a, \nu)} [\mathbb{E} \{ \exp(\lambda R_{a,1}) | \theta_a \}^n] \quad (227)$$

$$\stackrel{(ii)}{=} \mathbb{E}_{(\xi_a, \nu)} [\exp \{ A_a(\theta_a + \lambda) - A_a(\theta_a) \}^n] \quad (228)$$

$$\stackrel{(iii)}{\leq} \mathbb{E}_{(\xi_a, \nu)} \left[\exp \left\{ \lambda \cdot A'_a(\theta_a) + \frac{L\lambda^2}{2} \right\}^n \right] \quad (229)$$

$$\stackrel{(iv)}{=} \mathbb{E}_{(\xi_a, \nu)} \left[\exp \left\{ n\lambda \cdot \mu_a + \frac{Ln\lambda^2}{2} \right\} \right] \quad (230)$$

$$= \exp \left(n\lambda\bar{\mu}_a + \frac{Ln\lambda^2}{2} \right) \times \mathbb{E}_{(\xi_a, \nu)} [\exp(n\lambda(\mu_a - \bar{\mu}_a))] \quad (231)$$

$$\stackrel{(v)}{\leq} \exp \left(n\lambda\bar{\mu}_a + \frac{Ln\lambda^2}{2} \right) \times \exp \left(\frac{Ln^2\lambda^2}{2\nu} \right) \quad (232)$$

$$= \exp(n\lambda\bar{\mu}_a) \times \exp \left(\frac{\lambda^2}{2} \times \frac{Ln(\nu+n)}{\nu} \right), \quad (233)$$

where we have utilized that (i) $R_{a,i}$'s are conditionally independent given θ_a , (ii) the moment-generating function of $R_{a,1}$ is given by $\mathbb{E}[\lambda R_a | \theta_a] = \exp(A_a(\theta_a + \lambda) - A_a(\theta_a))$, (iii) $A_a(\cdot)$ is L -smooth, (iv) $A'_a(\theta_a) = \mu_a(\theta_a)$, and (v) μ_a is $\sqrt{L/\nu}$ -sub-Gaussian. Given that $\mathbb{E}[\sum_{i=1}^n R_{a,i}] = n\bar{\mu}_a$, we just have shown that the sum $\sum_{i=1}^n R_{a,i}$ is $\sqrt{\frac{Ln(\nu+n)}{\nu}}$ -sub-Gaussian. Therefore, its scaled version $\frac{\sum_{i=1}^n R_{a,i}}{\nu+n}$ is $\sqrt{\frac{Ln}{\nu(\nu+n)}}$ -sub-Gaussian, and so is $\hat{\mu}_{a,n}$. \blacksquare

Lemma 6. Consider the setting of Theorem [3](#). With $\sigma_n \triangleq \sqrt{\frac{Ln}{\nu(\nu+n)}}$, the following holds:

$$\mathbb{E} \left[\left(\max_{0 \leq i \leq n} \hat{\mu}_{a,i} - (\bar{\mu}_a + z\sigma_n) \right)^+ \right] \leq \frac{\sigma_n}{z} e^{-z^2/2}, \quad \forall z > 0. \quad (234)$$

Proof. Recall that the posterior predictive mean reward process $\{\hat{\mu}_{a,n}\}_{n \geq 0}$ is the martingale with respect to the filtration generated by reward realizations $R_{a,1}, R_{a,2}, \dots$ and whose mean is $\bar{\mu}_a$. Therefore, $\{\exp(\lambda \hat{\mu}_{a,n})\}_{n \geq 0}$ is a positive submartingale for any given $\lambda \geq 0$. By Doob's maximal inequality, we deduce that

$$\mathbb{P} \left[\max_{0 \leq i \leq n} \hat{\mu}_{a,i} \geq \bar{\mu}_a + z\sigma_n \right] = \mathbb{P} \left[\max_{0 \leq i \leq n} \exp(\lambda(\hat{\mu}_{a,i} - \bar{\mu}_a)) \geq \exp(\lambda z\sigma_n) \right] \leq \frac{\mathbb{E}[\exp(\lambda(\hat{\mu}_{a,n} - \bar{\mu}_a))]}{\exp(\lambda z\sigma_n)}. \quad (235)$$

By Lemma 5, since $\hat{\mu}_{a,n}$ is σ_n -sub-Gaussian, we further have

$$\frac{\mathbb{E} [\exp(\lambda(\hat{\mu}_{a,n} - \bar{\mu}_a))]}{\exp(\lambda z \sigma_n)} \leq \frac{\exp\left(\frac{\lambda^2 \sigma_n^2}{2}\right)}{\exp(\lambda z \sigma_n)} = \exp\left(\frac{\lambda^2 \sigma_n^2}{2} - \lambda z \sigma_n\right). \quad (236)$$

Therefore, by taking $\lambda \triangleq \frac{z}{\sigma_n}$, we have $\mathbb{P}[\max_{0 \leq i \leq n} \hat{\mu}_{a,i} \geq \bar{\mu}_a + z \sigma_n] \leq e^{-z^2/2}$, and by invoking Lemma 4, we obtain the claim. \blacksquare

D.4.4. Proof of Theorem 3

Lemma 7. Consider one of the IRS penalty functions z^{TS} , $z^{IRS.FH}$, and $z^{IRS.V-ZERO}$. As discussed in Remark 11, we have

$$Q_t^{z,\text{in}}(\mathbf{a}_{1:t-1}, a_t^{z,*}, \omega) - Q_t^{z,\text{in}}(\mathbf{a}_{1:t-1}, a, \omega) \leq \mu_t^U(\mathbf{a}_{1:t-1}, a_t^{z,*}, \omega) - \mu_t^L(\mathbf{a}_{1:t-1}, a, \omega), \quad (237)$$

for some $\mu_t^U(\mathbf{a}_{1:t-1}, a_t^{z,*}, \omega)$ and $\mu_t^L(\mathbf{a}_{1:t-1}, a, \omega)$, where $a_t^{z,*}$ abbreviates $a_t^{z,*}(\mathbf{a}_{1:t-1}, \omega)$. Suppose that there exists a sequence of confidence intervals $\{(L_t(a), U_t(a))\}_{a \in \mathcal{A}, t \in \mathbb{N}}$ such that $(L_t(\cdot), U_t(\cdot))$ is $\sigma(H_{t-1})$ -measurable, and

$$\mathbb{E}_{\mathbf{y}} \left[\left(\mu_t^U(\mathbf{a}_{1:t-1}, a, \omega) - U_t(a) \right)^+ \middle| H_{t-1}(\mathbf{a}_{1:t-1}, \omega) \right] \leq \frac{C_U}{T}, \quad \forall a, \forall t \quad (238)$$

$$\mathbb{E}_{\mathbf{y}} \left[\left(L_t(a) - \mu_t^L(\mathbf{a}_{1:t-1}, a, \omega) \right)^+ \middle| H_{t-1}(\mathbf{a}_{1:t-1}, \omega) \right] \leq \frac{C_L}{T}, \quad \forall a, \forall t \quad (239)$$

for some constants $C_U > 0$ and $C_L > 0$. Then, for IRS policy π induced by the chosen penalty function, we have

$$W^z(T, \mathbf{y}) - V(\pi, T, \mathbf{y}) \leq C_U + C_L + \sum_{t=1}^T \mathbb{E}[U_t(A_t^\pi) - L_t(A_t^\pi)]. \quad (240)$$

Proof. Let $A_t^* \triangleq a_t^{z,*}(\mathbf{A}_{1:t-1}^\pi, \omega)$, and let $\mathbb{E}_t[\cdot]$ denote $\mathbb{E}[\cdot | \mathcal{F}_{t-1}]$. By Proposition 7 we have

$$\mathbb{E}_t[U_t(A_t^\pi)] = \sum_{a \in \mathcal{A}} U_t(a) \cdot \mathbb{P}_t[A_t^\pi = a] = \sum_{a \in \mathcal{A}} L_t(a) \cdot \mathbb{P}_t[A_t^* = a] = \mathbb{E}_t[U_t(A_t^*)]. \quad (241)$$

Therefore, we have

$$\mathbb{E}_t \left[\mu_t^U(A_t^*) - \mu_t^L(A_t^\pi) \right] \quad (242)$$

$$= \mathbb{E}_t \left[\mu_t^U(A_t^*) - \mu_t^L(A_t^\pi) \right] + \mathbb{E}_t [U_t(A_t^\pi) - U_t(A_t^*)] + \mathbb{E}_t [L_t(A_t^\pi) - L_t(A_t^*)] \quad (243)$$

$$= \mathbb{E}_t \left[\mu_t^U(A_t^*) - U_t(A_t^*) \right] + \mathbb{E}_t \left[L_t(A_t^\pi) - \mu_t^L(A_t^\pi) \right] + \mathbb{E}_t [U_t(A_t^\pi) - L_t(A_t^\pi)] \quad (244)$$

$$\leq \mathbb{E}_t \left[\left(\mu_t^U(A_t^*) - U_t(A_t^*) \right)^+ \right] + \mathbb{E}_t \left[\left(L_t(A_t^\pi) - \mu_t^L(A_t^\pi) \right)^+ \right] + \mathbb{E}_t [U_t(A_t^\pi) - L_t(A_t^\pi)]. \quad (245)$$

We further observe that

$$\mathbb{E}_t \left[\left(\mu_t^U(A_t^*) - U_t(A_t^*) \right)^+ \right] = \sum_{a \in \mathcal{A}} \mathbb{E}_t \left[\left(\mu_t^U(a) - U_t(a) \right)^+ \right] \mathbb{P}_t[A_t^* = a] \leq \frac{C_U}{T} \sum_{a \in \mathcal{A}} \mathbb{P}_t[A_t^* = a] = \frac{C_U}{T}. \quad (246)$$

Similarly, we have $\mathbb{E}_t \left[\left(L_t(A_t^\pi) - \mu_t^L(A_t^\pi) \right)^+ \right] \leq \frac{C_L}{T}$. Combining all these results, we have

$$W(T, \mathbf{y}) - V(\pi, T, \mathbf{y}) \stackrel{\text{Prop 6}}{=} \mathbb{E} \left[\sum_{t=1}^T Q_t^{z, \text{in}}(A_t^*) - Q_t^{z, \text{in}}(A_t^\pi) \right] \quad (247)$$

$$\leq \mathbb{E} \left[\sum_{t=1}^T \mu_t^U(A_t^*) - \mu_t^L(A_t^\pi) \right] \quad (248)$$

$$= \mathbb{E} \left[\sum_{t=1}^T \mathbb{E}_t \left[\mu_t^U(A_t^*) - \mu_t^L(A_t^\pi) \right] \right] \quad (249)$$

$$\leq \mathbb{E} \left[\sum_{t=1}^T \left(\frac{C_U}{T} + \frac{C_L}{T} + \mathbb{E}_t [U_t(A_t^\pi) - L_t(A_t^\pi)] \right) \right] \quad (250)$$

$$\leq C_U + C_L + \sum_{t=1}^T \mathbb{E} [U_t(A_t^\pi) - L_t(A_t^\pi)]. \quad (251)$$

■

We are now ready to prove Theorem 3. To facilitate simpler notation, we define

$$N_{t-1}^\pi(a) \triangleq n_{t-1}(\mathbf{A}_{1:t-1}^\pi, a), \quad \hat{\mu}_t^\pi(a, n) \triangleq \hat{\mu}_{a, N_{t-1}^\pi(a)+n}, \quad (252)$$

which represent, respectively, the number of pulls on arm a prior to time t under policy π , and the posterior predictive mean reward process given the past actions $\mathbf{A}_{1:t-1}^\pi$. Observe that for each $a \in \mathcal{A}$, the process $\{\hat{\mu}_t^\pi(a, n)\}_{n \geq 0}$ is a martingale, as discussed Remark 1.

Further define

$$\Delta_t^\pi(a, n) \triangleq \sqrt{\frac{L}{\nu + N_{t-1}^\pi(a)} \times \frac{n}{\nu + N_{t-1}^\pi(a) + n}}, \quad (253)$$

which is measurable with respect to \mathcal{F}_{t-1} . In the context of Theorem 3, the prior/posterior of arm a at time t is described by the hyperparameters $\left(\xi_a + \sum_{i=1}^{N_{t-1}^\pi(a)} R_{a,i}, \nu + N_{t-1}^\pi(a) \right)$ that converges to μ_a , and therefore Lemma 5 implies that $\hat{\mu}_t^\pi(a, n)$ is $\Delta_t^\pi(a, n)$ -sub-Gaussian *conditioned* on \mathcal{F}_{t-1} .

(1) Suboptimality analysis for TS (60). As discussed in Remark 11, for TS, we have

$$Q_t^{z, \text{in}}(a_t^{z,*}) - Q_t^{z, \text{in}}(a) = \mu_{a_t^{z,*}} - \mu_a = \hat{\mu}_t^\pi(a_t^{z,*}, \infty) - \hat{\mu}_t^\pi(a, \infty). \quad (254)$$

We construct the confidence intervals as follows:

$$U_t(a) \triangleq \hat{\mu}_t^\pi(a, 0) + \sqrt{2 \log T} \times \Delta_t^\pi(a, \infty), \quad L_t(a) \triangleq \hat{\mu}_t^\pi(a, 0) - \sqrt{2 \log T} \times \Delta_t^\pi(a, \infty), \quad (255)$$

where $\Delta_t^\pi(a, \infty) = \lim_{n \rightarrow \infty} \Delta_t^\pi(a, n) = \sqrt{\frac{L}{\nu + N_{t-1}^\pi(a)}}$ so that μ_a is $\Delta_t^\pi(a, \infty)$ -sub-Gaussian conditioned on \mathcal{F}_{t-1} . By Lemma 4, we have

$$\mathbb{E} \left[(\mu_a - U_t(a))^+ \mid \mathcal{F}_{t-1} \right] \leq \frac{\Delta_t^\pi(a, \infty)}{\sqrt{2 \log T}} e^{-\frac{2 \log T}{2}} \leq \frac{\sqrt{L/\nu}}{T}, \quad (256)$$

where we use the fact that $2 \log T \geq 1$ for any $T \geq 2$. Symmetrically, we have $\mathbb{E} \left[(L_t(a) - \mu_a)^+ \mid \mathcal{F}_{t-1} \right] \leq \frac{\sqrt{L/\nu}}{T}$. By Lemma 7, we have

$$W^{\text{TS}}(T, \mathbf{y}) - V(\pi^{\text{TS}}, T, \mathbf{y}) \leq 2\sqrt{L/\nu} + \sum_{t=1}^T \mathbb{E} [U_t(A_t^\pi) - L_t(A_t^\pi)] \quad (257)$$

$$= 2\sqrt{L/\nu} + 2\sqrt{2 \log T} \sum_{t=1}^T \Delta_t^\pi(A_t^\pi, \infty). \quad (258)$$

Further observe that

$$\sum_{t=1}^T \Delta_t^\pi(A_t^\pi, \infty) = \sum_{t=1}^T \sqrt{\frac{L}{\nu + N_{t-1}^\pi(A_t^\pi)}} = \sum_{a \in \mathcal{A}} \sum_{n=0}^{N_T^\pi(a)-1} \frac{\sqrt{L}}{\sqrt{\nu + n}} = \sum_{a \in \mathcal{A}} \left(\frac{\sqrt{L}}{\sqrt{\nu}} + \sum_{n=1}^{N_T^\pi(a)-1} \frac{\sqrt{L}}{\sqrt{\nu + n}} \right) \quad (259)$$

$$\leq \sum_{a \in \mathcal{A}} \left(\frac{\sqrt{L}}{\sqrt{\nu}} + \sum_{n=1}^{N_T^\pi(a)-1} \frac{\sqrt{L}}{\sqrt{n}} \right) \leq \sum_{a \in \mathcal{A}} \left(\frac{\sqrt{L}}{\sqrt{\nu}} + \int_{x=0}^{N_T^\pi(a)} \frac{\sqrt{L}}{\sqrt{x}} dx \right) = \frac{K\sqrt{L}}{\sqrt{\nu}} + 2\sqrt{L} \sum_{a \in \mathcal{A}} \sqrt{N_T^\pi(a)}. \quad (260)$$

By utilizing Cauchy–Schwartz inequality, we deduce that

$$\sum_{a \in \mathcal{A}} \sqrt{N_T^\pi(a)} \leq \sqrt{K \sum_{a \in \mathcal{A}} N_T(a)} = \sqrt{KT}. \quad (261)$$

Combining all these results, we conclude that

$$W^{\text{TS}}(T, \mathbf{y}) - V(\pi^{\text{TS}}, T, \mathbf{y}) \leq 2\sqrt{L} \left[\frac{1}{\sqrt{\nu}} + \sqrt{2 \log T} \left(\frac{K}{\sqrt{\nu}} + 2\sqrt{KT} \right) \right]. \quad (262)$$

(2) Suboptimality analysis for IRS.FH (61). As discussed in Remark 11, for IRS.FH, we have

$$Q_t^{z, \text{in}}(a_t^{z, *}) - Q_t^{z, \text{in}}(a) = \hat{\mu}_t^\pi(a_t^{z, *}, T-t) - \hat{\mu}_t^\pi(a, T-t). \quad (263)$$

We construct the confidence intervals as follows:

$$U_t(a) \triangleq \hat{\mu}_t^\pi(a, 0) + \sqrt{2 \log T} \times \Delta_t^\pi(a, T-t), \quad L_t(a) \triangleq \hat{\mu}_t^\pi(a, 0) + \sqrt{2 \log T} \times \Delta_t^\pi(a, T-t). \quad (264)$$

Given that $\hat{\mu}_t^\pi(a, T-t)$ is $\Delta_t^\pi(a, T-t)$ -sub-Gaussian conditioned on \mathcal{F}_{t-1} , by Lemma 4, we have

$$\mathbb{E} \left[(\hat{\mu}_t^\pi(a, T-t) - U_t(a))^+ \middle| \mathcal{F}_{t-1} \right] \leq \frac{\Delta_t^\pi(a, T-t)}{\sqrt{2 \log T}} e^{-\frac{2 \log T}{2}} \leq \frac{\Delta_t^\pi(a, \infty)}{\sqrt{2 \log T}} e^{-\frac{2 \log T}{2}} \leq \frac{\sqrt{L/\nu}}{T}. \quad (265)$$

Symmetrically, we have $\mathbb{E} \left[(L_t(a) - \hat{\mu}_t^\pi(a, T-t))^+ \middle| \mathcal{F}_{t-1} \right] \leq \frac{\sqrt{L/\nu}}{T}$.

On the other hand, since $N_{t-1}(a) \leq t$ in any case, we have

$$\frac{1}{\nu + N_{t-1}^\pi(a)} \times \frac{T-t}{\nu + N_{t-1}^\pi(a) + T-t} = \frac{1}{\nu + N_{t-1}^\pi(a)} \times \left(1 - \frac{\nu + N_{t-1}^\pi(a)}{\nu + N_{t-1}^\pi(a) + T-t} \right) \quad (266)$$

$$= \frac{1}{\nu + N_{t-1}^\pi(a)} - \frac{1}{\nu + N_{t-1}^\pi(a) + T-t} \quad (267)$$

$$\leq \frac{1}{\nu + N_{t-1}^\pi(a)} - \frac{1}{\nu + T}. \quad (268)$$

Consequently,

$$\sum_{t=1}^T \sqrt{\frac{1}{\nu + N_{t-1}^\pi(a)} - \frac{1}{\nu + T}} = \sum_{a \in \mathcal{A}} \sum_{n=0}^{N_T^\pi(a)-1} \sqrt{\frac{1}{\nu + n} - \frac{1}{\nu + T}} \quad (269)$$

$$= \sum_{a \in \mathcal{A}} \left(\sqrt{\frac{1}{\nu} - \frac{1}{\nu + T}} + \sum_{n=1}^{N_T^\pi(a)-1} \sqrt{\frac{1}{\nu + n} - \frac{1}{\nu + T}} \right) \quad (270)$$

$$\leq \frac{K}{\sqrt{\nu}} + \sum_{a \in \mathcal{A}} \sum_{n=1}^{N_T^\pi(a)-1} \sqrt{\frac{1}{n} - \frac{1}{T}} \quad (271)$$

$$\stackrel{(i)}{\leq} \frac{K}{\sqrt{\nu}} + \sum_{a \in \mathcal{A}} \sum_{n=1}^{N_T^\pi(a)-1} \left(\frac{1}{\sqrt{n}} - \frac{\sqrt{n}}{2T} \right) \quad (272)$$

$$\leq \frac{K}{\sqrt{\nu}} + \sum_{a \in \mathcal{A}} \int_0^{N_T^\pi(a)} \left(\frac{1}{\sqrt{x}} - \frac{\sqrt{x}}{2T} \right) dx \quad (273)$$

$$= \frac{K}{\sqrt{\nu}} + \sum_{a \in \mathcal{A}} \left(2\sqrt{N_T^\pi(a)} - \frac{(N_T^\pi(a))^{3/2}}{2T} \right) \quad (274)$$

$$\stackrel{(ii)}{\leq} \frac{K}{\sqrt{\nu}} + 2\sqrt{KT} - \frac{1}{3}\sqrt{T/K}, \quad (275)$$

where we have utilized that (i) the concavity of $\sqrt{\cdot}$, and (ii) $\min\{\sum_{a=1}^K n_a^{3/2}; \sum_{a=1}^K n_a = T\} = \sum_{a=1}^K (T/K)^{3/2} = \sqrt{T^3/K}$.

Combining all these results, we conclude that

$$W^{\text{IRS.FH}}(T, \mathbf{y}) - V(\pi^{\text{IRS.FH}}, T, \mathbf{y}) \leq 2\sqrt{\frac{L}{\nu}} + 2\sqrt{2\log T} \sum_{t=1}^T \Delta_t^\pi(A_t^\pi, T-t) \quad (276)$$

$$\leq 2\sqrt{L} \left[\frac{1}{\sqrt{\nu}} + \sqrt{2\log T} \left(\frac{K}{\sqrt{\nu}} + 2\sqrt{KT} - \frac{1}{3}\sqrt{T/K} \right) \right]. \quad (277)$$

(3) Suboptimality analysis for IRS.V-ZERO (62). As discussed in Remark 11, for IRS.FH, we have

$$Q_t^{z,\text{in}}(a_t^{z,*}) - Q_t^{z,\text{in}}(a) = \max_{0 \leq n \leq T-t} \{ \hat{\mu}_t^\pi(a_t^{z,*}, n) \} - \hat{\mu}_t^\pi(a, 0). \quad (278)$$

We construct the confidence intervals as follows:

$$U_t(a) \triangleq \hat{\mu}_t^\pi(a, 0) + \sqrt{2\log T} \times \Delta_t^\pi(a, T-t), \quad L_t(a) \triangleq \hat{\mu}_t^\pi(a, 0). \quad (279)$$

By Lemma 6, we have

$$\mathbb{E} \left[\left(\max_{0 \leq n \leq T-t} \hat{\mu}_t^\pi(a, n) - U_t(a) \right)^+ \middle| \mathcal{F}_{t-1} \right] \leq \frac{\Delta_t^\pi(a, T-t)}{\sqrt{2\log T}} e^{-\frac{2\log T}{2}} \leq \frac{\sqrt{L/\nu}}{T}, \quad (280)$$

where

$$\mathbb{E} [\hat{\mu}_t^\pi(a, 0) - L_t(a) | \mathcal{F}_{t-1}] = 0. \quad (281)$$

The rest of the proof is almost identical to the case of IRS.FH:

$$W^{\text{IRS.V-ZERO}}(T, \mathbf{y}) - V(\pi^{\text{IRS.V-ZERO}}, T, \mathbf{y}) \leq \sqrt{\frac{L}{\nu}} + \sum_{t=1}^T \mathbb{E} [U_t(A_t^\pi) - L_t(A_t^\pi)]. \quad (282)$$

$$= \sqrt{\frac{L}{\nu}} + \sqrt{2\log T} \sum_{t=1}^T \Delta_t^\pi(A_t^\pi, T-t) \quad (283)$$

$$\leq \sqrt{L} \left[\frac{1}{\sqrt{\nu}} + \sqrt{2\log T} \left(\frac{K}{\sqrt{\nu}} + 2\sqrt{KT} - \frac{1}{3}\sqrt{T/K} \right) \right]. \quad (284)$$