

Locally Optimal Fixed-Budget Best Arm Identification in Two-Armed Gaussian Bandits with Unknown Variances

Anonymous authors
Paper under double-blind review

Abstract

We address the problem of best arm identification (BAI) with a fixed budget for two-armed Gaussian bandits. In BAI, given multiple arms, we aim to find the best arm, an arm with the highest expected reward, through an adaptive experiment. Kaufmann et al. (2016) develops a lower bound for the probability of misidentifying the best arm. They also propose a strategy, assuming that the variances of rewards are known, and show that it is asymptotically optimal in the sense that its probability of misidentification matches the lower bound as the budget approaches infinity. However, an asymptotically optimal strategy is unknown when the variances are unknown. For this open issue, we propose a strategy that estimates variances during an adaptive experiment and draws arms with a ratio of the estimated standard deviations. We refer to this strategy as the *Neyman Allocation (NA)-Augmented Inverse Probability weighting (AIPW)* strategy. We then demonstrate that this strategy is asymptotically optimal by showing that its probability of misidentification matches the lower bound when the budget approaches infinity, and the gap between the expected rewards of two arms approaches zero (*small-gap regime*). Our results suggest that under the worst-case scenario characterized by the small-gap regime, our strategy, which employs estimated variance, is asymptotically optimal even when the variances are unknown.

1 Introduction

This study investigates the problem of *best arm identification (BAI) with a fixed budget* in stochastic two-armed Gaussian bandits. In this problem, we consider an adaptive experiment with a fixed number of rounds, called a *budget*. At each round, we can draw an arm and observe the reward. The goal of the problem is to identify the best arm with the highest expected reward at the end of the experiment (Bubeck et al., 2009; Audibert et al., 2010).

Formally, we consider the following adaptive experiment with two arms and Gaussian rewards. There are two arms 1 and 2, and an arm $a \in \{1, 2\}$ has an \mathbb{R} -valued Gaussian reward $Y_a \sim \mathcal{N}(\mu_a, \sigma_a^2)$ with the mean $\mu_a \in (-\infty, \infty)$ and the variance $\sigma_a^2 \in (0, \infty)$, where μ_a and σ_a^2 are constants independent of T . Given fixed (σ_1^2, σ_2^2) , let

$$\mathcal{P}^G := \mathcal{P}_{(\sigma_1^2, \sigma_2^2)}^G := \left\{ P = \left(\mathcal{N}(\mu_1, \sigma_1^2), \mathcal{N}(\mu_2, \sigma_2^2) \right) : \mu_1 \text{ and } \mu_2 \text{ are finite, } \mu_1 \neq \mu_2 \right\}$$

be a set of distributions generating the data, which is referred to as the *Gaussian bandit models*, where $P \in \mathcal{P}$ is a pair of distributions that generate (Y_1, Y_2) , and $\mathcal{N}(\mu, \sigma^2)$ is a Gaussian distribution with a mean μ and a variance σ^2 . For an instance P , the best arm $a^*(P) \in \{1, 2\}$ is defined as $a^*(P) = \arg \max_{a \in \{1, 2\}} \mu_a$, which is assumed to exist uniquely.

In the adaptive experiment, we consider a strategy to identify the best arm. A fixed budget T is given. For each round $t \in [T] := \{1, 2, \dots, T\}$, let $(Y_{1,t}, Y_{2,t})$ be an independent and identically distributed (i.i.d.) copy of (Y_1, Y_2) generated from $P^* = \left(\mathcal{N}(\mu_1^*, \sigma_1^2), \mathcal{N}(\mu_2^*, \sigma_2^2) \right) \in \mathcal{P}^G$. At each round t , we draw arm $A_t \in \{1, 2\}$ and observe a reward $Y_t = \sum_{a \in \{1, 2\}} \mathbb{1}[A_t = a] Y_{a,t}$. At the end of the experiment (after round T), we

recommend an estimated best arm $\hat{a}_T \in \{1, 2\}$. During an experiment, we follow a *strategy* that determines which arm to draw and which arm to recommend as the best arm. The performance of strategies is evaluated by a minimal probability of misidentification $\mathbb{P}_P(\hat{a}_T \neq a^*(P))$, where \mathbb{P}_P is the probability law under P .

Background. In fixed-budget BAI, it has been an important question of interest to investigate the probability of misidentification $\mathbb{P}_P(\hat{a}_T \neq a^*(P))$ in the limit $T \rightarrow \infty$. For the interest, a typical approach is to derive an upper and lower bound of the probability separately and specify its value.

For a lower bound of the probability of misidentification, Kaufmann et al. (2016) develops a general theory for deriving lower bounds of the probability. Their theory applies the change-of-measure argument, which has been employed in various problems (van der Vaart, 1998), including studies for regret minimization (Lai & Robbins, 1985). Their lower bound is general and can be applied to a wide range of settings, such as the fixed confidence setting (Garivier & Kaufmann, 2016) as well as the fixed budget setting.

In contrast, an upper bound of the misidentification probability has not been fully clarified. A typical way to derive upper bounds is to construct a specific strategy and evaluate its misidentification probability. Kaufmann et al. (2016) develops a strategy under a setting in which the variance (σ_1^2, σ_2^2) of the reward is known and shows its misidentification probability corresponds to the lower bound. However, this strategy is not available under the usual setting with unknown variance. Based on these situations, the current results are insufficient to establish an upper bound for the misidentification probability when the variances are unknown.

Based on the situation above, our interest is in strategies for identifying misidentification probabilities in the adaptive experimental setting described above. Specifically, we need a strategy such that an upper bound on its misidentification probability aligns with the lower bound proposed in Kaufmann et al. (2016). Further, this strategy must be valid when the variance is unknown.

Our approach and contribution. In this study, we develop a strategy whose probability of misidentification aligns with the lower bound under an additional setting. To accomplish this, we develop the *Neyman allocation-augmented inverse probability weight* (NA-AIPW) strategy. Then, we show that the probability of misidentification aligns with the lower bound under a *small-gap regime*. The details of each are described below.

The NA-AIPW strategy consists of a sampling rule using the Neyman allocation (NA) and a recommendation rule using the augmented inverse probability weighting (AIPW) estimator. NA is a method of sampling arms using a ratio of the root of the variance of rewards, as utilized in Neyman (1934); Kaufmann et al. (2016). The NA-AIPW strategy samples the arms by estimating this variance during the adaptive experiment. At the end of the experiment, the NA-AIPW strategy recommends an arm with the highest expected reward estimated by using the AIPW estimator, which is an unbiased estimator with a small asymptotic variance.

The small-gap regime considers a situation $\mu_1 - \mu_2 \rightarrow 0$ as $T \rightarrow \infty$. Although this additional setting slightly simplifies the problem with BAI, the problem is still sufficiently complicated since the small gap makes it difficult to identify the best arm. This setting has been utilized in BAI with fixed confidence, such as the analysis of lil'UCB (Jamieson et al., 2014). In the realm of statistical testing, such an evaluation framework is known as the local Bahadur efficiency (Bahadur, 1960; Wieand, 1976; Akritas & Kourouklis, 1988; He & Shao, 1996). From a technical perspective, the small-gap regime is a situation where we can ignore the estimation error of the variances compared to the difficulty of identifying the best arm. Since the error of the estimation of the variance is relatively negligible in the small-gap setting, we can show that the misidentification probability of the NA-AIPW strategy matches the lower bound.

We summarize the backgrounds and our contributions. In BAI with two-armed Gaussian rewards and a fixed budget, a strategy has been needed in which its misidentification probability achieves the lower bound derived by Kaufmann et al. (2016). Although Kaufmann et al. (2016) demonstrates an asymptotically optimal strategy that satisfies the requirement with known variances, it remains an unresolved issue to find a strategy whose upper bound matches their derived lower bound when variances are unknown. For this issue, this study proposes the NA-AIPW strategy whose probability of misidentification matches the lower bound under the small-gap regime.

Organization. This study is organized as follows. First, in Section 2, we review the lower bound of Kaufmann et al. (2016). Then, in Section 3, we propose our NA-AIPW strategy. In Section 4, we show that the misidentification probability of the strategy asymptotically corresponds to the lower bound by Kaufmann et al. (2016) under the small-gap setting. We show the proof in Section 5, where we also provide a novel concentration inequality based on the Chernoff bound. In Appendix A, we discuss the difficulty of this problem. In Appendix B, we introduce related work and remaining problems, which includes an extension of our small-gap setting to a setting with multi-armed bandits and non-Gaussian rewards.

Notation. Let \mathcal{F}_t be the sigma-algebra generated by all observations up to round t . We define a truncation operator: for a variable $v \in \mathbb{R}$ and constants $c_1 < c_2$, $\text{thre}(v; c_1, c_2) := \min\{\max\{v, c_1\}, c_2\}$.

2 Lower Bound of Probability of Misidentification

As a preparation, we introduce a lower bound for the probability of misidentification in BAI with a fixed budget. We call a strategy is *consistent*, if for any $P \in \mathcal{P}^G$, $\mathbb{P}_P(\hat{a}_T \neq a^*(P)) \rightarrow 0$ as $T \rightarrow \infty$. To evaluate the performance of strategies For any $P \in \mathcal{P}^G$, we focus on the following metric for $\mathbb{P}_P(\hat{a}_T \neq a^*(P))$ used in many studies, such as Kaufmann et al. (2016):

$$-\frac{1}{T} \log \mathbb{P}_P(\hat{a}_T \neq a^*(P)).$$

Note that the upper bound (resp. lower bound) of this term works as a lower bound (resp. upper bound) of the probability of misidentification $\mathbb{P}_P(\hat{a}_T \neq a^*(P))$ since $x \mapsto -\log x$ is a strictly decreasing function.

For two-armed Gaussian bandits, Kaufmann et al. (2016) presents the following lower bounds.

Proposition 2.1 (Theorem 12 in Kaufmann et al. (2016)). *For any $P^* \in \mathcal{P}^G$, if $\{(Y_{1,t}, Y_{2,t})\}_{t \in [T]}$ is generated from P^* , any consistent strategy satisfies*

$$\limsup_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*}(\hat{a}_T \neq a^*(P^*)) \leq \frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2} := \text{LowerBound}(\Delta).$$

From the statement, there are some important aspects of this lower bound: (i) The term $\Delta = \mu_1^* - \mu_2^*$, which referred to a *gap*, appears in the numerator and the magnitude of the error is described by the gap. (ii) The variances (σ_1^2, σ_2^2) appear in the denominator, which plays an important role.

It has been discussed to find a strategy in which the upper bound of its probability of misidentification coincides with this lower bound in Proposition 2.1. Although Kaufmann et al. (2016) develops a strategy that satisfies the requirement, it needs to sample arms with some probability depending on the known variances (σ_1^2, σ_2^2) . To the best of our knowledge, if the variances are unknown and need to be estimated during adaptive experiments, no one has found the desired strategy.

3 The NA-AIPW Strategy

In this section, we define our strategy. Formally, a strategy gives a pair $((A_t)_{t \in [T]}, \hat{a}_T)$, where (i) $(A_t)_{t \in [T]} \in \{1, 2\}^T$ is a sequence of arms generated by a sampling rule that determines which arm A_t is chosen in each t based on \mathcal{F}_{t-1} , and (ii) $\hat{a}_T \in \{1, 2\}$ is a recommended arm by a recommendation rule based on \mathcal{F}_T . Our proposed NA-AIPW strategy consists of (i) a sampling rule with the Neyman Allocation (NA) (Neyman, 1923), and (ii) a recommendation rule using the Augmented Inverse Probability Weighting (AIPW) estimator (Robins et al., 1994; Bang & Robins, 2005). Based on these rules, we refer to this strategy as the NA-AIPW strategy¹.

¹Similar strategies are often used in the context of the average treatment effect estimation by an adaptive experiment (van der Laan, 2008; Kato et al., 2020).

3.1 Target Allocation Ratio

As preparation, we introduce the notion of a target allocation ratio, which will be used for the sampling rule. We define target allocation ratios $w_1^*, w_2^* \in (0, 1)$ as

$$w_1^* := \frac{\sigma_1}{\sigma_1 + \sigma_2}, \quad \text{and} \quad w_2^* := 1 - w_1^*.$$

A sampling rule following this target allocation ratio is known as the Neyman allocation rule (Neyman, 1934). Glynn & Juneja (2004) and Kaufmann et al. (2016) also propose this allocation. This target allocation ratio is characterized by the variances (standard deviations); therefore, the target allocation ratio is unknown when the variances are unknown. Therefore, to use this ratio, we need to estimate it from observations. Let $\overline{C}_\mu > 0$ and $\overline{C}_{\sigma^2} > 0$ be some large positive constants predetermined us. They are technical terms for the proof, and we just set them as sufficiently large values such that $\overline{C}_\mu = \overline{C}_{\sigma^2} = 10^{10000000000}$. For the condition, see Theorem 4.1.

3.2 Sampling Rule with Neyman Allocation (NA)

We present the sampling rule with the NA. At each round $t \in [T]$, our sampling rule randomly draws an arm $a \in \{1, 2\}$ with a probability identical to an estimated version of the target allocation ratio w_a^* . To estimate the target allocation ratio w_a^* , we estimate the variances during the adaptive experiment. For $a \in \{1, 2\}$, let $\{\hat{\sigma}_a\}_{t \in [T]}$ and $\{\hat{w}_{a,t}\}$ be sequences of estimators of σ_a, μ_a and w_a^* , that will be defined below.

We use the rounds $t = 1$ and $t = 2$ for initialization. Specifically, we draw the arm 1 at round $t = 1$ and the arm 2 at round $t = 2$, and also set $\hat{w}_{a,1} = \hat{w}_{a,2} = 1/2$ for $a \in \{1, 2\}$.

At the round $t \geq 3$, we estimate the target allocation ratio (variances) w_a^* for $a \in \{1, 2\}$ using past observations \mathcal{F}_{t-1} . For each $t \geq 3$, we first define an estimator of the expected reward μ_a as

$$\tilde{\mu}_{a,t} := \frac{1}{\sum_{s=1}^{t-1} \mathbb{1}[A_s = a]} \sum_{s=1}^{t-1} \mathbb{1}[A_s = a] Y_{a,s}.$$

Also, we define a second moment estimator $\tilde{\zeta}_{a,t} := (\sum_{s=1}^{t-1} \mathbb{1}[A_s = a])^{-1} \sum_{s=1}^{t-1} \mathbb{1}[A_s = a] Y_{a,s}^2$, and a root of variance estimator $\tilde{\sigma}_{a,t} = \{\tilde{\zeta}_{a,t} - (\tilde{\mu}_{a,t})^2\}^{1/2}$. Then, we define the estimator $\hat{\sigma}_{a,t} = \text{thre}(\tilde{\sigma}_{a,t}; 1/\overline{C}_{\sigma^2}, \overline{C}_{\sigma^2})$. Note that this truncation is introduced for a technical purpose to draw each arm infinitely many times as $T \rightarrow \infty$ and avoid the estimators of μ_1^* and μ_2^* , defined below, diverging to infinity. We just use a sufficiently small value for $\overline{C}_{\sigma^2} > 0$. Also, we define the estimator $\hat{w}_{1,t}$ and $\hat{w}_{2,t}$ as

$$\hat{w}_{1,t} := \frac{\hat{\sigma}_{1,t}}{\hat{\sigma}_{1,t} + \hat{\sigma}_{2,t}}, \quad \text{and} \quad \hat{w}_{2,t} := 1 - \hat{w}_{1,t}. \quad (1)$$

In each round $t \geq 3$, we draw arm $A_t = 1$ with probability $\hat{w}_{1,t}$ and $A_t = 2$ with probability $\hat{w}_{2,t}$.

We note the possibility of increasing the number of initialization rounds, although our strategy utilizes only the first two rounds for this purpose. The additional rounds of initialization serve to stabilize the sampling rule in practical applications, akin to the concept of forced sampling (Garivier & Kaufmann, 2016). We can change the number of initialization rounds if the condition $\hat{w}_{a,t} \xrightarrow{\text{a.s.}} w_a^*$ is satisfied as $t \rightarrow \infty$ for every $a \in \{1, 2\}$, which is crucial for our theoretical analysis. For instance, instead of using $\hat{w}_{1,t}$ directly, an alternative arm-drawing probability could be defined as $\tilde{w}_{1,t} = \alpha_t/2 + (1 - \alpha_t)\hat{w}_{1,t}$, assuming $\alpha_t \in [0, 1]$ and converges to zero as t approaches infinity (here, we define $\tilde{w}_{2,t} = 1 - \tilde{w}_{1,t}$). Moreover, the number of initialization rounds can be made dependent upon the number of arms without impacting the theoretical outcomes.

3.3 Recommendation Rule with the Augmented Inverse Probability Weighting (AIPW) Estimator

We present our recommendation rule. In the recommendation phase after round T , we estimate μ_a^* for each $a \in \{1, 2\}$ and recommend an arm with the bigger estimated expected reward. With a truncated version

Algorithm 1 NA-AIPW Strategy**Parameter:** Positive constants \overline{C}_μ and \overline{C}_{σ^2} .**Initialization:**At $t = 1$, sample $A_t = 1$; at $t = 2$, sample $A_t = 2$. For $a \in \{1, 2\}$, set $\widehat{w}_{a,1} = \widehat{w}_{a,2} = 0.5$.**for** $t = 3$ to T **do**Construct $\widehat{w}_{a,t}$ following equation 1.Draw $A_t = 1$ with probability $\widehat{w}_{1,t}$ and $A_t = 2$ with probability $\widehat{w}_{2,t} = 1 - \widehat{w}_{1,t}$.Observe Y_t .**end for**Construct $\widehat{\mu}_{a,T}^{\text{AIPW}}$ for $a \in \{1, 2\}$. following equation 2.Recommend \widehat{a}_T following equation 3.

of the estimated expected reward $\widehat{\mu}_{a,t} := \text{thre}(\widetilde{\mu}_{a,t}, -\overline{C}_\mu, \overline{C}_\mu)$, we define the *augmented inverse probability weighting* (AIPW) estimator of μ_a^* for each $a \in \{1, 2\}$ as

$$\widehat{\mu}_{a,T}^{\text{AIPW}} := \frac{1}{T} \sum_{t=1}^T \psi_{a,t}, \quad \text{where } \psi_{a,t} := \frac{\mathbb{1}[A_t = a](Y_{a,t} - \widehat{\mu}_{a,t})}{\widehat{w}_{a,t}} + \widehat{\mu}_{a,t}. \quad (2)$$

At the end of the experiment (after the round $t = T$), we recommend \widehat{a}_T as

$$\widehat{a}_T := \begin{cases} 1 & \text{if } \widehat{\mu}_{1,T}^{\text{AIPW}} \geq \widehat{\mu}_{2,T}^{\text{AIPW}}, \\ 2 & \text{otherwise.} \end{cases} \quad (3)$$

We adopt the AIPW estimator for our strategy because it has several advantages. First, the AIPW estimator has the property of semiparametric efficiency, which indicates that it has the smallest asymptotic variance among a certain class (Hahn, 1998). The property is necessary to prove that the strategy using the AIPW estimator is optimal, which means the misidentification probability is small enough to achieve its lower bound. The second reason is more technical; the AIPW estimator simplifies the theoretical analysis (see Section A.3). Specifically, we can decompose an estimation error into a sum of random variables with martingale properties, making it suitable for analysis using the central limit theorem. We cannot use this property for an empirical average. Details will be given in Section 5 and Appendix A.

We provide the pseudo-code for our proposed strategy in Algorithm 1. Note again that we introduce \overline{C}_μ and \overline{C}_{σ^2} for technical purposes to bound the estimators and any large positive value can be used.

4 Misidentification Probability and Asymptotic Optimality

In this section, we show an upper bound of the misspecification probability of the NA-AIPW strategy, which also implies that the strategy is asymptotically optimal. Then, the following theorem holds.

Theorem 4.1 (Upper bound of the NA-AIPW strategy). *If \overline{C}_μ and \overline{C}_{σ^2} are sufficiently large positive values such that $\mu_a \in (-\overline{C}_\mu, \overline{C}_\mu)$ and $\sigma_a \in (1/\overline{C}_{\sigma^2}, \overline{C}_{\sigma^2})$ hold for $a \in \{1, 2\}$ and any $P^* \in \mathcal{P}^G$, then it holds that*

$$\lim_{\Delta \rightarrow 0} \liminf_{T \rightarrow \infty} -\frac{1}{\Delta^2 T} \log \mathbb{P}_{P^*} \left(\widehat{\mu}_{a^*(P),T}^{\text{AIPW}} \leq \widehat{\mu}_{a,T}^{\text{AIPW}} \right) \geq \frac{1}{2(\sigma_1 + \sigma_2)^2}.$$

Here, \overline{C}_μ and \overline{C}_{σ^2} are introduced for technical purpose. For example, we can set them as $\overline{C}_\mu = \overline{C}_{\sigma^2} = 10^{10000000000}$ if $\mu_a \in (-\overline{C}_\mu, \overline{C}_\mu)$ and $\sigma_a \in (1/\overline{C}_{\sigma^2}, \overline{C}_{\sigma^2})$ are satisfied.² Our result can be applied to more general distributions not only for Gaussian distributions, under some regularity conditions, such as the

²The constant \overline{C}_μ is introduced to guarantee the boundedness of the estimators, and it is sufficient to use a sufficiently large value for it. The constant \overline{C}_{σ^2} is used to draw each arm infinitely many times as $T \rightarrow \infty$ and avoid the estimators of the means μ_1^* and μ_2^* diverging to infinity, and it is sufficient to use a sufficiently large value for \overline{C}_μ and \overline{C}_{σ^2} .

existence of the second moment of $Y_{a,t}$. However, we focus on Gaussian distributions since our goal is to develop a strategy whose upper bound aligns with the lower bound in Kaufmann et al. (2016). Lower bounds for more general distributions have been still explored by the related work (Wang et al., 2023a).

The lower bound of $-\frac{1}{T} \log \mathbb{P}_{P^*}(\hat{a}_T^{\text{AIPW}} \neq a^*(P^*))$ implies the upper bound of $\mathbb{P}_{P^*}(\hat{a}_T^{\text{AIPW}} \neq a^*(P^*))$. This theorem implies us to evaluate the probability of misidentification up to the constant term, even when it is exponentially small, as $\Delta \rightarrow 0$. In the statement, we suppressed the dependency of P^* on Δ .

Theorem 4.1 directly implies the asymptotic optimality of the NA-AIPW strategy. As $\Delta \rightarrow 0$, the upper bound matches the lower bound in Proposition 2.1. This asymptotic optimality result suggests that the estimation error of the target allocation ratio (variances) w^* is negligible when $\Delta \rightarrow 0$. This is because the estimation error is insignificant compared to the challenges of identifying the best arm due to the small gap.

The statement in Theorem 4.1 has different mathematical representations that have the same meanings, while the statement in Theorem 4.1 is compatible with the main result in Shin et al. (2018). For example, if we write the definition of \lim , we can state the result as follows: for all $\epsilon > 0$, there exists $t(\epsilon) > 0$ such that for all $T > t(\epsilon)$ and for all $\epsilon' > 0$, there exists $\delta_T(\epsilon') > 0$ such that for all $0 < \Delta < \delta_T(\epsilon')$, we have

$$-\frac{1}{T} \log \mathbb{P}_{P^*}(\hat{\mu}_{1,T}^{\text{AIPW}} \leq \hat{\mu}_{2,T}^{\text{AIPW}}) \geq \frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2} - (\epsilon + \epsilon') \Delta^2.$$

By using $\text{LowerBound}(\Delta) = \frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2}$, we can also denote the result as

$$\lim_{\Delta \rightarrow 0} \liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*}(\hat{\mu}_{a^*(P),T}^{\text{AIPW}} \leq \hat{\mu}_{a,T}^{\text{AIPW}}) / \text{LowerBound}(\Delta) \geq 1, \quad (4)$$

which directly implies that the lower and upper bounds match exactly asymptotically.

We conjecture that even if we replace the AIPW estimator with the sample average estimator, defined as $\tilde{\mu}_{a,t} = (\sum_{s=1}^{t-1} \mathbb{1}[A_s = a])^{-1} \sum_{s=1}^{t-1} \mathbb{1}[A_s = a] Y_{a,s}$ in Section 3.2, the upper bound of the strategy still matches the lower bound under the small-gap regime. However, the proof is an open issue. Hirano et al. (2003) and Hahn et al. (2011) show that the sample average estimator $\tilde{\mu}_{a,t}$ and the AIPW estimator have the same asymptotic variance (or asymptotic distribution). To show the result, we need to employ empirical process arguments. One of the problems in extending the result to analysis for BAI is that their result focuses on the asymptotic distribution, not the tail probability. Therefore, to show the asymptotic optimality of the strategy with the sample average in the sense of the probability of misidentification, we need to modify the result in Hirano et al. (2003) and Hahn et al. (2011) to analyze the tail probability. See also Appendix A.

Finite sample properties of the AIPW estimator have been exploited by Kato et al. (2020) and Cook et al. (2023). Kato et al. (2020) develops uniformly valid confidence intervals for the AIPW estimator using the Law of the Iterated Logarithms (LIL), inspired by Balsubramani & Ramdas (2016) and Jamieson et al. (2014). However, Howard et al. (2021) points out that the concentration inequalities in Balsubramani & Ramdas (2016) are not tight. Cook et al. (2023) refines and generalizes the results of Kato et al. (2020) by using the arguments in Howard et al. (2021) and removes dependency on the unknown parameters.

5 Proof of Theorem 4.1

To show Theorem 4.1, we derive the upper bound of $\mathbb{P}_{P^*}(\hat{\mu}_{a^*(P^*),T}^{\text{AIPW}} \leq \hat{\mu}_{b,T}^{\text{AIPW}})$ for $b \in \{1, 2\} \setminus \{a^*(P^*)\}$, which is equivalent to $\mathbb{P}_{P^*}(\hat{a}_T^{\text{AIPW}} \neq a^*(P^*))$. Without loss of generality, we assume that $a^*(P^*) = 1$ and $b = 2$. Let us define $V := \frac{\sigma_1^2}{w_1^*} + \frac{\sigma_2^2}{w_2^*} = (\sigma_1 + \sigma_2)^2$ and

$$\Psi_t := \frac{\psi_{1,t} - \psi_{2,t} - \Delta}{\sqrt{V}}.$$

Therefore, in the following parts, we aim to derive the upper bound of $\mathbb{P}_{P^*}(\hat{\mu}_{a^*(P^*),T}^{\text{AIPW}} \leq \hat{\mu}_{b,T}^{\text{AIPW}}) = \mathbb{P}_{P^*}(\hat{\mu}_{1,T}^{\text{AIPW}} \leq \hat{\mu}_{2,T}^{\text{AIPW}}) = \mathbb{P}_{P^*}(\sum_{t=1}^T \Psi_t \leq -\frac{T\Delta}{\sqrt{V}})$. Let \mathbb{E}_P be the expectation under $P \in \mathcal{P}^G$. We de-

rive the upper bound using the Chernoff bound. This proof is partially inspired by techniques in Hadad et al. (2021), and Kato et al. (2020).

First, because there exists a constant $C > 0$ independent of T such that $\widehat{w}_{a,t} > C$ by construction, the following lemma holds.

Lemma 5.1. *For any $P^* \in \mathcal{P}^G$ and all $a \in \{1, 2\}$, $\widehat{\mu}_{a,t} \xrightarrow{\text{a.s.}} \mu_a^*$ and $\widehat{\sigma}_a^2 \xrightarrow{\text{a.s.}} \sigma_a^2$.*

Furthermore, from $\widehat{\sigma}_a^2 \xrightarrow{\text{a.s.}} \sigma_a^2$ and continuous mapping theorem, for any $P^* \in \mathcal{P}^G$ and all $a \in \{1, 2\}$, $\widehat{w}_{a,t} \xrightarrow{\text{a.s.}} w_{a,t}^*$.

Step 1: Sequence $\{\Psi_t\}_{t=1}^T$ is a martingale difference sequence (MDS)

We prove that $\{\Psi_t\}_{t=1}^T$ is an MDS; that is, $\mathbb{E}_{P^*} [\Psi_t | \mathcal{F}_{t-1}] = 0$. Although this fact is well-known in the literature of causal inference (van der Laan, 2008; Hadad et al., 2021; Kato et al., 2020), we show the proof for the sake of completeness.

Lemma 5.2. *For any $P^* \in \mathcal{P}^G$, $\{\Psi_t\}_{t=1}^T$ is an MDS.*

Proof. For each $t \in [T]$, it holds that

$$\begin{aligned} \sqrt{V} \mathbb{E}_{P^*} [\Psi_t | \mathcal{F}_{t-1}] &= \mathbb{E}_{P^*} \left[\frac{\mathbb{1}[A_t = 1](Y_{1,t} - \widehat{\mu}_{1,t})}{\widehat{w}_{1,t}} + \widehat{\mu}_{1,t} | \mathcal{F}_{t-1} \right] - \mathbb{E}_{P^*} \left[\frac{\mathbb{1}[A_t = 2](Y_{2,t} - \widehat{\mu}_{2,t})}{\widehat{w}_{2,t}} + \widehat{\mu}_{2,t} | \mathcal{F}_{t-1} \right] - \Delta \\ &= \frac{\widehat{w}_{1,t}(\mu_1^* - \widehat{\mu}_{1,t})}{\widehat{w}_{1,t}} + \widehat{\mu}_{1,t} - \frac{\widehat{w}_{2,t}(\mu_2^* - \widehat{\mu}_{2,t})}{\widehat{w}_{2,t}} + \widehat{\mu}_{2,t} - \Delta = \{(\mu_1^* - \mu_2^*) - (\mu_1^* - \mu_2^*)\} = 0. \end{aligned}$$

□

Step 2: Evaluation by using the Chernoff Bound with Martingales

By applying the Chernoff bound, for any $v < 0$ and any $\lambda < 0$, it holds that

$$\mathbb{P}_{P^*} \left(\frac{1}{T} \sum_{t=1}^T \Psi_t \leq v \right) \leq \mathbb{E}_{P^*} \left[\exp \left(\lambda \sum_{t=1}^T \Psi_t \right) \right] \exp(-T\lambda v).$$

From the Chernoff bound and a property of an MDS, we have

$$\mathbb{E}_{P^*} \left[\exp \left(\lambda \sum_{t=1}^T \Psi_t \right) \right] = \mathbb{E}_{P^*} \left[\prod_{t=1}^T \mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}] \right] = \mathbb{E}_{P^*} \left[\exp \left(\sum_{t=1}^T \log \mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}] \right) \right].$$

Then, from the Taylor expansion around $\lambda = 0$, for each $t \in \mathbb{N}$, for all $\epsilon > 0$, there exists $\ell_t(\epsilon) > 0$ such that for all $0 < \lambda < \ell_t(\epsilon)$, we have

$$\log \mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}] \leq \frac{\lambda^2}{2} \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] + \epsilon \lambda^2. \quad (5)$$

This is given as follows. Since there exists $\bar{\lambda} > 0$ such that $\mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}]$ exists for any $\lambda \in (0, \bar{\lambda})$, and the second moment $\mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}]$ also exists, the Taylor expansion yields the following (for the details, see the textbook such as page 75 in Bulmer (1967) and Theorem 5.19 in Apostol (1974)): for each $t \in \mathbb{N}$, for all $\epsilon > 0$, there exists $\ell_t(\epsilon) > 0$ such that for all $0 < \lambda < \ell_t(\epsilon)$, it holds that $\mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}] \leq 1 + \sum_{k=1}^2 \frac{\lambda^k}{k!} \mathbb{E}_{P^*} [\Psi_t^k | \mathcal{F}_{t-1}] + \epsilon \lambda^2$. Note that the existence of $\mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}]$ comes from the following in Ψ_t : (i) $Y_{a,t}$ is a Gaussian random variable, (ii) $\widehat{\mu}_{a,t}$ and $\widehat{w}_{a,t}$ are bounded random variables by our truncation, and (iii) the lower bound of \widehat{w} is given by $1/\overline{C}\sigma^2$.

By using the Taylor expansion again, we approximate $\log(1+z)$ around $z=0$ as $\log(1+z) = z - z^2/2 + z^3/3 - \dots$. Therefore, for each $t \in \mathbb{N}$, for all $\epsilon > 0$, there exists $\ell_t(\epsilon) > 0$ such that for all $0 < \lambda < \ell_t(\epsilon)$, it holds that

$$\log \mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}] \leq \left\{ \lambda \mathbb{E}_{P^*} [\Psi_t | \mathcal{F}_{t-1}] + \frac{\lambda^2}{2} \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] + \epsilon \lambda^2 \right\} - \frac{1}{2} \{ \lambda \mathbb{E}_{P^*} [\Psi_t | \mathcal{F}_{t-1}] + \epsilon \lambda \}^2$$

Therefore, for each $t \in [T]$, for all $\epsilon > 0$, there exists $\ell_t(\epsilon) > 0$ such that for all $0 < \lambda < \ell_t(\epsilon)$, it holds that

$$\log \mathbb{E}_{P^*} [\exp(\lambda \Psi_t) | \mathcal{F}_{t-1}] \leq \frac{\lambda^2}{2} \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] + \epsilon \lambda^2 - \epsilon^2 \lambda^2.$$

Here, we used $\mathbb{E}_{P^*} [\Psi_t | \mathcal{F}_{t-1}] = 0$. Thus, the equation 5 holds.

Step 3: Convergence of the Second Moment

We next show $\mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1 \xrightarrow{\text{a.s.}} 0$. This result is a direct consequence of Lemma 5.1.

Lemma 5.3. *For any $P^* \in \mathcal{P}^G$, we have $\mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1 \xrightarrow{\text{a.s.}} 0$ as $t \rightarrow \infty$.*

Proof. We have

$$\begin{aligned} V \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] &= \mathbb{E}_{P^*} \left[(\psi_{1,t} - \psi_{2,t} - \Delta)^2 | \mathcal{F}_{t-1} \right] \\ &= \mathbb{E}_{P^*} \left[\left(\frac{\mathbb{1}[A_t = 1](Y_{1,t} - \hat{\mu}_{1,t})}{\hat{w}_{1,t}} - \frac{\mathbb{1}[A_t = 2](Y_{2,t} - \hat{\mu}_{2,t})}{\hat{w}_{2,t}} \right)^2 \right. \\ &\quad \left. + 2 \left(\frac{\mathbb{1}[A_t = a^*(P^*)](Y_{1,t} - \hat{\mu}_{1,t})}{\hat{w}_{1,t}} - \frac{\mathbb{1}[A_t = a](Y_{2,t} - \hat{\mu}_{2,t})}{\hat{w}_{2,t}} \right) \times (\hat{\mu}_{1,t} - \hat{\mu}_{2,t} - (\mu_1^* - \mu_2^*)) \right. \\ &\quad \left. + (\hat{\mu}_{1,t} - \hat{\mu}_{2,t} - (\mu_1^* - \mu_2^*))^2 | \mathcal{F}_{t-1} \right] \\ &= \mathbb{E}_{P^*} \left[\frac{\mathbb{1}[A_t = 1](Y_{1,t} - \hat{\mu}_{1,t})^2}{(\hat{w}_{1,t})^2} + \frac{\mathbb{1}[A_t = 2](Y_{2,t} - \hat{\mu}_{2,t})^2}{(\hat{w}_{2,t})^2} \right. \\ &\quad \left. + 2 \left(\frac{\mathbb{1}[A_t = 1](Y_{1,t} - \hat{\mu}_{1,t})}{\hat{w}_{1,t}} - \frac{\mathbb{1}[A_t = a](Y_{2,t} - \hat{\mu}_{2,t})}{\hat{w}_{2,t}} \right) (\hat{\mu}_{1,t} - \hat{\mu}_{2,t} - (\mu_1^* - \mu_2^*)) \right. \\ &\quad \left. + ((\hat{\mu}_{1,t} - \hat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 | \mathcal{F}_{t-1} \right] \\ &= \sum_{a \in \{1,2\}} \mathbb{E}_{P^*} \left[\frac{(Y_{a,t} - \hat{\mu}_{a,t})^2}{(\hat{w}_{a,t})^2} | \mathcal{F}_{t-1} \right] - \mathbb{E}_{P^*} \left[((\hat{\mu}_{1,t} - \hat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 | \mathcal{F}_{t-1} \right]. \end{aligned}$$

Here, for $a \in \{1,2\}$, the followings hold:

$$\begin{aligned} \mathbb{E}_{P^*} \left[\frac{\mathbb{1}[A_t = a](Y_{a,t} - \hat{\mu}_{a,t})^2}{(\hat{w}_{a,t})^2} | \mathcal{F}_{t-1} \right] &= \mathbb{E}_{P^*} \left[\frac{(Y_{a,t} - \hat{\mu}_{a,t})^2}{\hat{w}_{a,t}} | \mathcal{F}_{t-1} \right] = \frac{\mathbb{E}_{P^*} [(Y_{a,t})^2] - 2\mu_a^* \hat{\mu}_{a,t} + (\hat{\mu}_{a,t})^2}{\hat{w}_{a,t}} \\ &= \frac{\mathbb{E}_{P^*} [(Y_{a,t})^2] - (\mu_a^*)^2 + (\mu_a^* - \hat{\mu}_{a,t})^2}{\hat{w}_{a,t}} = \frac{\sigma_a^2 + (\mu_a^* - \hat{\mu}_{a,t})^2}{\hat{w}_{a,t}}, \end{aligned}$$

and

$$\mathbb{E}_{P^*} \left[\frac{\mathbb{1}[A_t = a](Y_{a,t} - \hat{\mu}_{2,t})^2}{(\hat{w}_{1,t})^2} - \frac{\mathbb{1}[A_t = a](Y_{a,t} - \hat{\mu}_{2,t})^2}{(\hat{w}_{2,t})^2} | \mathcal{F}_{t-1} \right] = 0,$$

where we used $\mathbb{E}_{P^*}[(Y_{a,t})^2|x] - (\mu_a^*)^2 = \sigma_a^2$. Therefore, the following holds:

$$\begin{aligned} & \mathbb{E}_{P^*} \left[\frac{(Y_{1,t} - \widehat{\mu}_{1,t})^2}{\widehat{w}_{1,t}} \middle| \mathcal{F}_{t-1} \right] + \mathbb{E}_{P^*} \left[\frac{(Y_{2,t} - \widehat{\mu}_{2,t})^2}{\widehat{w}_{2,t}} \middle| \mathcal{F}_{t-1} \right] - \mathbb{E}_{P^*} \left[((\widehat{\mu}_{1,t} - \widehat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 \middle| \mathcal{F}_{t-1} \right] \\ &= \mathbb{E}_{P^*} \left[\frac{\sigma_1^2 + (\mu_1^* - \widehat{\mu}_{1,t})^2}{\widehat{w}_{1,t}} \right] + \mathbb{E}_{P^*} \left[\frac{\sigma_2^2 + (\mu_2^* - \widehat{\mu}_{2,t})^2}{\widehat{w}_{2,t}} \right] - \mathbb{E}_{P^*} \left[((\widehat{\mu}_{1,t} - \widehat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 \right]. \end{aligned}$$

Because $\widehat{\mu}_{a,t} \xrightarrow{\text{a.s.}} \mu_a^*$ and $\widehat{w}_{a,t} \xrightarrow{\text{a.s.}} w_a^*$, we have

$$\begin{aligned} & \lim_{t \rightarrow \infty} \left| \left(\frac{\sigma_1^2 + (\mu_1^* - \widehat{\mu}_{1,t})^2}{\widehat{w}_{1,t}} \right) + \left(\frac{\sigma_2^2 + (\mu_2^* - \widehat{\mu}_{2,t})^2}{\widehat{w}_{2,t}} \right) - ((\widehat{\mu}_{1,t} - \widehat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 \right. \\ & \qquad \qquad \qquad \left. - \left(\frac{\sigma_1^2}{w_1^*} + \frac{\sigma_a^2}{w_2^*} + ((\mu_1^* - \mu_2^*) - (\mu_1^* - \mu_2^*))^2 \right) \right| \\ & \leq \lim_{t \rightarrow \infty} \left| \frac{\sigma_1^2}{\widehat{w}_{1,t}} - \frac{\sigma_1^2}{w_1^*} \right| + \lim_{t \rightarrow \infty} \left| \frac{\sigma_2^2}{\widehat{w}_{2,t}} - \frac{\sigma_2^2}{w_2^*} \right| + \lim_{t \rightarrow \infty} \frac{(\mu_1^* - \widehat{\mu}_{1,t})^2}{\widehat{w}_{1,t}} + \lim_{t \rightarrow \infty} \frac{(\mu_2^* - \widehat{\mu}_{2,t})^2}{\widehat{w}_{2,t}} \\ & \qquad \qquad \qquad + \lim_{t \rightarrow \infty} |((\widehat{\mu}_{1,t} - \widehat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 - ((\mu_1^* - \mu_2^*) - (\mu_1^* - \mu_2^*))^2| = 0, \end{aligned}$$

with probability 1. Therefore, from Lebesgue's dominated convergence theorem, we obtain

$$\begin{aligned} & V\mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - V \\ &= \mathbb{E}_{P^*} \left[\frac{\sigma_1^2 + (\mu_1^* - \widehat{\mu}_{1,t})^2}{\widehat{w}_{1,t}} \middle| \mathcal{F}_{t-1} \right] + \mathbb{E}_{P^*} \left[\frac{\sigma_2^2 + (\mu_2^* - \widehat{\mu}_{2,t})^2}{\widehat{w}_{2,t}} \middle| \mathcal{F}_{t-1} \right] \\ & \quad - \mathbb{E}_{P^*} \left[((\widehat{\mu}_{1,t} - \widehat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 \middle| \mathcal{F}_{t-1} \right] - \mathbb{E}_{P^*} \left[\frac{\sigma_1^2}{w_1^*} + \frac{\sigma_2^2}{w_2^*} + ((\mu_1^* - \mu_2^*) - (\mu_1^* - \mu_2^*))^2 \middle| \mathcal{F}_{t-1} \right] \\ & \xrightarrow{\text{a.s.}} 0. \end{aligned}$$

□

This lemma immediately yields the following lemma.

Lemma 5.4. *For any $P^* \in \mathcal{P}^G$ and any $\epsilon > 0$, there exists $t(\epsilon) > 0$ such that for all $T > t(\epsilon)$, we have $\frac{1}{T} \sum_{t=1}^T |\mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1| < \epsilon$ with probability one.*

This result is a variant of the Cesàro lemma for a case with almost sure convergence. For completeness, we show the proof, which is based on the proof of Lemma 10 in Hadad et al. (2021).

Proof. Let u_t be $u_t = \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1$. Note that $\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1 = \frac{1}{T} \sum_{t=1}^T u_t$.

From the proof of Lemma 5.3, we can find that u_t is a bounded random variable. Recall that

$$\begin{aligned} & V\mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] \\ &= \mathbb{E}_{P^*} \left[\frac{\sigma_1^2 + (\mu_1^* - \widehat{\mu}_{1,t})^2}{\widehat{w}_{1,t}} \middle| \mathcal{F}_{t-1} \right] + \mathbb{E}_{P^*} \left[\frac{\sigma_2^2 + (\mu_2^* - \widehat{\mu}_{2,t})^2}{\widehat{w}_{2,t}} \middle| \mathcal{F}_{t-1} \right] - \mathbb{E}_{P^*} \left[((\widehat{\mu}_{1,t} - \widehat{\mu}_{2,t}) - (\mu_1^* - \mu_2^*))^2 \middle| \mathcal{F}_{t-1} \right]. \end{aligned}$$

We assumed that $(\mu_1^*, \mu_2^*, \widehat{\mu}_{1,t}, \widehat{\mu}_{2,t}, \widehat{w}_{1,t}, \widehat{w}_{2,t})$ are all bounded random variables. Let C be a constant independent of T such that $|u_t| < C$ for all $t \in \mathbb{N}$.

Almost-sure convergence of u_t to zero as $t \rightarrow \infty$ implies that for all $\epsilon > 0$, there exists $t(\epsilon)$ such that $|u_t| < \epsilon$ for all $t \geq t(\epsilon)$ with probability one. Let $\mathcal{E}(\epsilon)$ denote the event in which this happens; that is, $\mathcal{E}(\epsilon) = \{|u_t| < \epsilon \ \forall t \geq t(\epsilon)\}$. Under this event, for $T > t(\epsilon)$, the following holds:

$$\frac{1}{T} \sum_{t=1}^T |u_t| \leq \frac{1}{T} \sum_{t=1}^{t(\epsilon)} C + \frac{1}{T} \sum_{t=t(\epsilon)+1}^T \epsilon = \frac{1}{T} t(\epsilon) C + \epsilon,$$

where $\frac{1}{T}t(\epsilon)C \rightarrow 0$ as $T \rightarrow \infty$.

Therefore, for all $\epsilon > 0$, there exists $t(\epsilon) > 0$ such that for all $T > t(\epsilon)$, $\frac{1}{T} \sum_{t=1}^T |u_t| < \epsilon$ holds with probability one. \square

Step 4: Tail Bound with the Approximated Second Moment

Let $v = \lambda$. Then, we have

$$\mathbb{P}_{P^*} \left(\sum_{t=1}^T \Psi_t \leq Tv \right) \leq \mathbb{E}_{P^*} \left[\exp \left(-\frac{T\lambda^2}{2} + \frac{\lambda^2}{2} \left\{ \sum_{t=1}^T \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1 \right\} + T\epsilon\lambda^2 \right) \right].$$

From Lemma 5.4, for all $\epsilon > 0$, there exists $t(\epsilon) > 0$ such that for all $T > t(\epsilon)$ and for all $\epsilon' > 0$, there exists $\underline{\ell}_T(\epsilon') = \min\{\ell_1(\epsilon'), \ell_2(\epsilon'), \dots, \ell_T(\epsilon')\} > 0$ such that for all $0 < \lambda < \underline{\ell}_T(\epsilon')$, we have³

$$\begin{aligned} & \mathbb{E}_{P^*} \left[\exp \left(-\frac{T\lambda^2}{2} + \frac{\lambda^2}{2} \left\{ \sum_{t=1}^T \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1 \right\} + T\epsilon'\lambda^2 \right) \right] \\ &= \exp \left(-\frac{T\lambda^2}{2} + T\epsilon'\lambda^2 \right) \mathbb{E}_{P^*} \left[\exp \left(\frac{\lambda^2}{2} \left\{ \sum_{t=1}^T \mathbb{E}_{P^*} [\Psi_t^2 | \mathcal{F}_{t-1}] - 1 \right\} \right) \right] \\ &\leq \exp \left(-\frac{T\lambda^2}{2} + T\epsilon'\lambda^2 \right) \exp \left(\frac{\lambda^2}{2} T\epsilon \right) = \exp \left(-\frac{T\lambda^2}{2} + T\epsilon'\lambda^2 + \frac{\lambda^2}{2} T\epsilon \right). \end{aligned}$$

Step 5: Final Step of the Proof of Theorem 4.1

For all $\epsilon > 0$, there exists $t(\epsilon) > 0$ such that for all $T > t(\epsilon)$ and for all $\epsilon' > 0$, there exists $\underline{\ell}_T(\epsilon') > 0$ such that for all $0 < \lambda < \underline{\ell}_T(\epsilon')$, we have

$$-\frac{1}{T} \log \mathbb{P}_{P^*} \left(\widehat{\mu}_{1,T}^{\text{AIPW}} \leq \widehat{\mu}_{2,T}^{\text{AIPW}} \right) \geq - \left\{ -\frac{\lambda^2}{2} + \epsilon'\lambda^2 + \frac{\lambda^2}{2}\epsilon \right\} = \frac{\lambda^2}{2} - \epsilon'\lambda^2 - \frac{\lambda^2}{2}\epsilon.$$

Let $\lambda = -\frac{\Delta}{\sqrt{V}}$. Then, we obtain $\lim_{\Delta \rightarrow 0} \liminf_{T \rightarrow \infty} -\frac{1}{\Delta^2 T} \log \mathbb{P}_{P^*} \left(\widehat{\mu}_{a^*(P),T}^{\text{AIPW}} \leq \widehat{\mu}_{a,T}^{\text{AIPW}} \right) \geq \frac{1}{2V}$. Thus, the proof is complete.

6 Simulation Studies

This section provides simulation studies to investigate the empirical performance of the NA-AIPW strategy. For comparison, we also investigate the performances of the NA-IPW and NA-SA strategies defined in Section A.3. Furthermore, we also conduct simulation studies of the ‘oracle’ strategy with the known variances, denoted by Oracle, and the uniform strategy that draws an equal number of arms, denoted by Uniform. We recommend an arm with the highest sample average in the Oracle and Uniform strategies.⁴

Throughout the experiment, we set arm 1 as the best arm and $\bar{C}_\mu = \bar{C}_{\sigma^2} = 1000$. We conduct experiments with the three settings.

In the first experiment, we set $\mu_1^* = 1.00$ and choose μ_2^* from the set $\{0.80, 0.85, 0.90, 0.95, 0.99\}$. The variances (σ_1, σ_2) are selected with a probability of 1/2 from either $(1, v_2)$ or $(v_2, 1)$, where v_2 is chosen from 5, 10, 20, 50. We continue the strategies until $T = 10000$ and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9900, 10000\}$.⁵ We conduct 1000 independent trials for each choice of parameters and plot the results in Figures 1 and 2.

³Note that for each $t \in \mathbb{N}$ and for all $\epsilon' > 0$, there exists $\ell(\epsilon') > 0$ such that we can apply the approximation for all $0 < \lambda < \ell(\epsilon')$.

⁴The oracle strategy is the one proposed by Glynn & Juneja (2004) and Kaufmann et al. (2016). The Uniform strategy with the recommendation rule is referred to as the Uniform-Empirical Best Arm (EBA) strategy by Bubeck et al. (2011).

⁵This means we report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9900, 10000\}$.

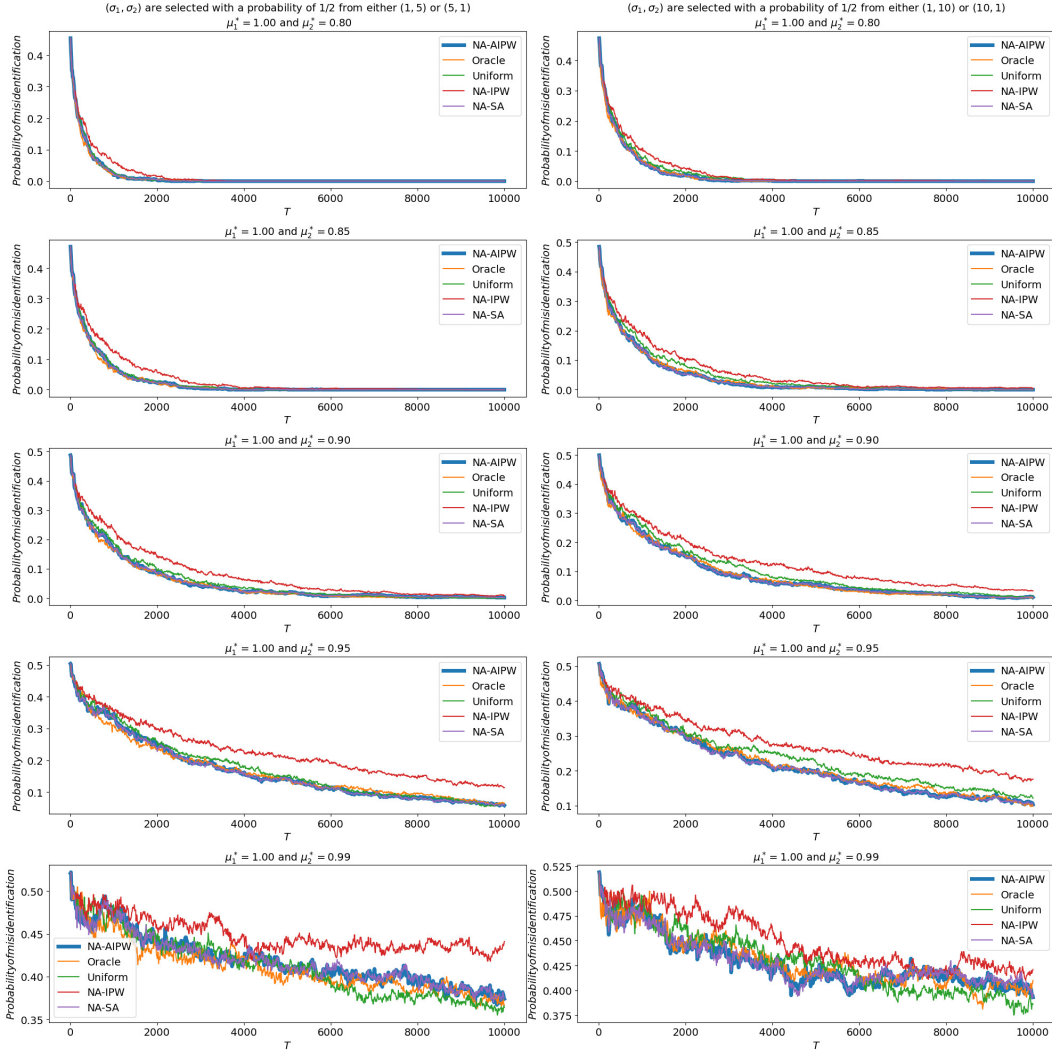


Figure 1: The results are under the first setting. We set $\mu_1^* = 1.00$ and choose μ_2^* from the set $\{0.80, 0.85, 0.90, 0.95, 0.99\}$. The variances (σ_1, σ_2) are selected with a probability of 1/2 from either $(1, v_2)$ or $(v_2, 1)$, where v_2 is chosen from 5, 10. We conduct 1,000 independent trials and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9900, 10000\}$.

In the second experiment, the variances (σ_1, σ_2) are selected with a probability of 1/2 from either $(5, v_2)$ or $(v_2, 5)$, where v_2 is chosen from 5, 10, 20, 50. The other settings are the same as the first experiment. The results are shown in Figures 3 and 4 in Appendix C.

In the third experiment, we set $\mu_1^* = 10.00$ and choose μ_2^* from the set $\{9.80, 9.85, 9.90, 9.95, 9.99\}$. The other settings are the same as the first experiment. The results are shown in Figures 5 and 6 in Appendix C.

Our theoretical results imply that the probability of misidentifications of the NA-AIPW and Oracle strategies approach the same as $\Delta \rightarrow 0$. We can confirm the phenomenon. In the results, the Oracle strategy is a bit better than the NA-AIPW strategy when Δ is large. However, the gap approaches zero as $\Delta \rightarrow 0$. Note that when Δ is large, the convergence of the probability of misidentification is very fast, so the gap is still not so large even if Δ is large because both the probability of misidentifications of the NA-AIPW and Oracle strategies converge to zero very fast.

We can also find that the performance improvement of the NA-AIPW strategy from the Uniform strategy is large as the difference of variances is large. For example, in the second experiment, when $(\sigma_1, \sigma_2) =$

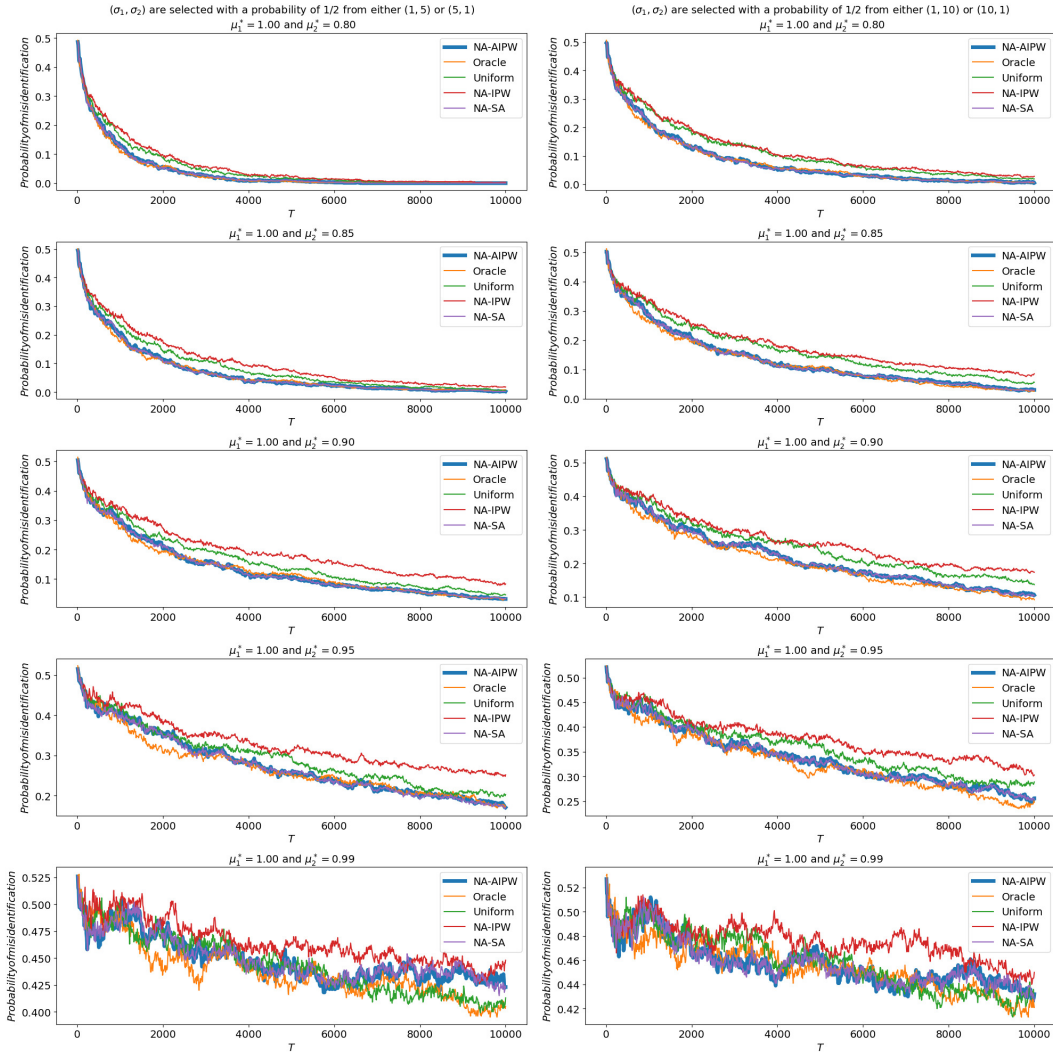


Figure 2: The results are under the first setting. We set $\mu_1^* = 1.00$ and choose μ_2^* from the set $\{0.80, 0.85, 0.90, 0.95, 0.99\}$. The variances (σ_1, σ_2) are selected with a probability of 1/2 from either $(1, v_2)$ or $(v_2, 1)$, where v_2 is chosen from 20, 50. We conduct 1,000 independent trials and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9,900, 10,000\}$.

(5.5), there is no improvement by using the NA-AIPW strategy from the Uniform strategy because the NA allocation also leads us to draw each arm with equal ratio.

The difference between the NA-AIPW and NA-IPW strategies becomes large as the mean outcome of each arm becomes large. We can find that in the third setting, the NA-IPW strategy behaves badly since $\mu_1^* = 10.00$ and μ_2^* is chosen from $\{9.80, 9.85, 9.90, 9.95, 9.99\}$, while $\mu_1^* = 1.00$ and μ_2^* is chosen from $\{0.80, 0.85, 0.90, 0.95, 0.99\}$ in the first and second settings.

7 Conclusion

This study investigated fixed-budget BAI for two-armed Gaussian bandits with unknown variances. We first reviewed the lower bound shown by Kaufmann et al. (2016). Then, we proposed the NA-AIPW strategy and found that its probability of misidentification matches the lower bound when the budget approaches infinity and the gap between the expected rewards of the two arms approaches zero. We referred to this setting as the small-gap regime and the optimality as the local asymptotic optimality. Although there are several remaining open questions, our result provides insight into long-standing open problems in BAI.

References

- Karun Adusumilli. Neyman allocation is minimax optimal for best arm identification with two arms, 2022. arXiv:2204.05527.
- Dohyun Ahn, Dongwook Shin, and Assaf Zeevi. Online ordinal optimization under model misspecification, 2021. URL <https://api.semanticscholar.org/CorpusID:235389954>. SSRN.
- Michael G. Akritas and Stavros Kourouklis. Local bahadur efficiency of score tests. *Journal of Statistical Planning and Inference*, 19(2):187–199, 1988.
- Tom M Apostol. *Mathematical analysis; 2nd ed.* Addison-Wesley series in mathematics. Addison-Wesley, 1974.
- Kaito Ariu, Masahiro Kato, Junpei Komiyama, Kenichiro McAlinn, and Chao Qin. Policy choice and best arm identification: Asymptotic analysis of exploration sampling, 2021. arXiv:2109.08229.
- Timothy B. Armstrong. Asymptotic efficiency bounds for a class of experimental designs, 2022. arXiv:2205.02726.
- Alexia Atsidakou, Sumeet Katariya, Sujay Sanghavi, and Branislav Kveton. Bayesian fixed-budget best-arm identification, 2023. arXiv:2211.08572.
- Jean-Yves Audibert, Sébastien Bubeck, and Remi Munos. Best arm identification in multi-armed bandits. In *Conference on Learning Theory*, pp. 41–53, 2010.
- R. R. Bahadur. Stochastic Comparison of Tests. *The Annals of Mathematical Statistics*, 31(2):276 – 295, 1960.
- Akshay Balsubramani and Aaditya Ramdas. Sequential nonparametric testing with the law of the iterated logarithm. In Alexander T. Ihler and Dominik Janzing (eds.), *Conference on Uncertainty in Artificial Intelligence*, 2016.
- Heejung Bang and James M. Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.
- Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems. In *Algorithmic Learning Theory*, pp. 23–37. Springer Berlin Heidelberg, 2009.
- Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 2011.
- Michael George Bulmer. *Principles of statistics*. M.I.T. Press, 2. ed. edition, 1967.
- Alexandra Carpentier and Andrea Locatelli. Tight (lower) bounds for the fixed budget best arm identification bandit problem. In *COLT*, 2016.
- Chun-Hung Chen, Jianwu Lin, Enver Yücesan, and Stephen E. Chick. Simulation budget allocation for further enhancing the efficiency of ordinal optimization. *Discrete Event Dynamic Systems*, 10(3):251–270, 2000.
- Thomas Cook, Alan Mishler, and Aaditya Ramdas. Semiparametric efficient inference in adaptive experiments. In *NeurIPS 2023 Workshop on Adaptive Experimental Design and Active Learning in the Real World*, 2023. URL <https://openreview.net/forum?id=xfj5jjp0aL>. arXiv:2311.18274.
- Rémy Degenne. On the existence of a complexity in fixed budget bandit identification. In *Conference on Learning Theory*, volume 195, pp. 1131–1154. PMLR, 2023.
- Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In *Conference on Learning Theory*, 2016.

- Peter Glynn and Sandeep Juneja. A large deviations perspective on ordinal optimization. In *Proceedings of the 2004 Winter Simulation Conference*, volume 1. IEEE, 2004.
- Vitor Hadad, David A. Hirshberg, Ruohan Zhan, Stefan Wager, and Susan Athey. Confidence intervals for policy evaluation in adaptive experiments. *Proceedings of the National Academy of Sciences*, 118(15), 2021.
- Jinyong Hahn. On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, 66(2):315–331, 1998.
- Jinyong Hahn, Keisuke Hirano, and Dean Karlan. Adaptive experimental design using the propensity score. *Journal of Business and Economic Statistics*, 2011.
- Xuming He and Qi-man Shao. Bahadur efficiency and robustness of studentized score tests. *Annals of the Institute of Statistical Mathematics*, 48(2):295–314, 1996.
- Keisuke Hirano, Guido Imbens, and Geert Ridder. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 2003.
- Steven R. Howard, Aaditya Ramdas, Jon D. McAuliffe, and Jasjeet S. Sekhon. Time-uniform, nonparametric, nonasymptotic confidence sequences. *Annals of Statistics*, 2021.
- Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lil' ucb : An optimal exploration algorithm for multi-armed bandits. In *Conference on Learning Theory*, 2014.
- Marc Jourdan, Degenne Rémy, and Kaufmann Emilie. Dealing with unknown variances in best-arm identification. In *Proceedings of The 34th International Conference on Algorithmic Learning Theory*, volume 201, pp. 776–849, 2023.
- Maximilian Kasy and Anja Sautmann. Adaptive treatment assignment in experiments for policy choice. *Econometrica*, 89(1):113–132, 2021.
- Masahiro Kato. Worst-case optimal multi-armed gaussian best arm identification with a fixed budget, 2023. arXiv:2310.19788.
- Masahiro Kato, Takuya Ishihara, Junya Honda, and Yusuke Narita. Efficient adaptive experimental design for average treatment effect estimation, 2020. arXiv:2002.05308.
- Masahiro Kato, Masaaki Imaizumi, Takuya Ishihara, and Toru Kitagawa. Asymptotically minimax optimal fixed-budget best arm identification for expected simple regret minimization, 2023a. arXiv:2302.02988.
- Masahiro Kato, Masaaki Imaizumi, Takuya Ishihara, and Toru Kitagawa. Fixed-budget hypothesis best arm identification: On the information loss in experimental design. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023b.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17(1):1–42, 2016.
- Edward H. Kennedy. Semiparametric theory and empirical processes in causal inference, 2016. arXiv:1510.04740.
- Junpei Komiyama, Taira Tsuchiya, and Junya Honda. Minimax optimal algorithms for fixed-budget best arm identification. In *Advances in Neural Information Processing Systems*, 2022.
- Junpei Komiyama, Kaito Ariu, Masahiro Kato, and Chao Qin. Rate-optimal bayesian simple regret in best arm identification. *Mathematics of Operations Research*, 2023.
- T.L Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 1985.

- Jerzy Neyman. Sur les applications de la theorie des probabilités aux expériences agricoles: Essai des principes. *Statistical Science*, 5:463–472, 1923.
- Jerzy Neyman. On the two different aspects of the representative method: the method of stratified sampling and the method of purposive selection. *Journal of the Royal Statistical Society*, 97:123–150, 1934.
- Chao Qin. Open problem: Optimal best arm identification with fixed-budget. In *Conference on Learning Theory*, 2022.
- James M. Robins, Andrea Rotnitzky, and Lue Ping Zhao. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89(427):846–866, 1994.
- Daniel Russo. Simple bayesian algorithms for best-arm identification. *Operations Research*, 68(6):1625–1647, 2020.
- Dongwook Shin, Mark Broadie, and Assaf Zeevi. Tractable sampling strategies for ordinal optimization. *Operations Research*, 66(6):1693–1712, 2018.
- Max Tabord-Meehan. Stratification Trees for Adaptive Randomisation in Randomised Controlled Trials. *The Review of Economic Studies*, 90(5):2646–2673, 2022.
- Anastasios Tsiatis. *Semiparametric Theory and Missing Data*. Springer Series in Statistics. Springer New York, 2007.
- Mark J. van der Laan. The construction and analysis of adaptive group sequential designs, 2008. URL <https://biostats.bepress.com/ucbbiostat/paper232>.
- A.W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.
- Po-An Wang, Kaito Ariu, and Alexandre Proutiere. On uniformly optimal algorithms for best arm identification in two-armed bandits with fixed budget, 2023a. arXiv:2308.12000.
- Po-An Wang, Ruo-Chun Tzeng, and Alexandre Proutiere. Best arm identification with fixed budget: A large deviation perspective. In *Advances in Neural Information Processing Systems*, 2023b.
- Harry S. Wieand. A Condition Under Which the Pitman and Bahadur Approaches to Efficiency Coincide. *The Annals of Statistics*, 4(5):1003 – 1011, 1976.
- Jinglong Zhao. Adaptive neyman allocation, 2023. arXiv:2309.08808.

A Discussion

In this section, we discuss related topics.

A.1 Neyman Allocation with Unknown Variances

For two-armed Gaussian bandits with known variances, Chen et al. (2000), Glynn & Juneja (2004), and Kaufmann et al. (2016) conclude that sampling each arm with a proportion of the standard deviation is optimal, which corresponds to the Neyman allocation Neyman (1934).

The Neyman allocation with unknown variances has been long studied in various fields. van der Laan (2008) and Hahn et al. (2011) develop algorithms for estimating the gap parameter Δ itself in an adaptive experiment with the Neyman allocation. They estimate the variances and show their algorithms' optimalities under the framework of semiparametric efficiency, which closely connects to the Gaussian approximation of estimators using the central limit theorem. Although they show their optimality under the framework, they do not investigate the asymptotic optimality in the large-deviation framework. Tabord-Meehan (2022), Kato et al. (2020), and Zhao (2023) also attempt to address related problems.

Jourdan et al. (2023) examines BAI with unknown variances in a fixed-confidence setting. Beyond the difference in settings (we focus on fixed-budget BAI), the methods of deriving lower bounds differ between our approach and theirs. They determine the lower bound while incorporating the assumption that the variances are unknown. Moreover, under a large-gap regime (Δ is fixed), they confirm a discrepancy between the lower bounds when variances are known versus unknown. Specifically, they consider alternative hypotheses related to both variances and means. In contrast, the lower bounds presented by Kaufmann et al. (2016) and ourselves are based on alternative hypotheses with fixed variances. While Jourdan et al. (2023) suggests that the upper bounds of strategies with unknown variances cannot align with the lower bound when variances are known, our findings indicate a match under the small-gap regime.

A.2 Necessity of the Small-Gap Regime

First, we discuss the necessity of the small-gap regime.

Estimation error of the variances. The most critical reason we employ the small-gap regime is that the estimation error of the variances cannot be ignored in evaluating the probability of misidentification. To clarify this point, we review the probability of misidentification when we know the variances.

Probability of misidentification of the Kaufmann et al. (2016)'s strategy with known variances.

Kaufmann et al. (2016) proposes drawing arm a in $\frac{\sigma_a^2}{\sigma_1 + \sigma_2} T = w_a^* T$ rounds (for simplicity, we deal with $w_a^* T$ as an integer). Without loss of generality, we consider draw arm 1 in the first $w_1^* T$ rounds and draw arm 2 in the following $T - w_1^* T = w_2^* T$ rounds. Then, they estimate the best arm as $\hat{a}_T^{\text{KCG}} := \arg \max_{a \in \{1, 2\}} \hat{\mu}_{a, T}^{\text{SA}}$, where $\hat{\mu}_{a, T}^{\text{SA}}$ is the sample average defined as

$$\hat{\mu}_{a, T}^{\text{SA}} := \frac{1}{\sum_{t=1}^T \mathbb{1}[A_t^{\text{KCG}} = a]} \sum_{t=1}^T \mathbb{1}[A_t^{\text{KCG}} = a] Y_t,$$

where A_t^{KCG} denotes an arm drawn by the Kaufmann et al. (2016)'s strategy. For the strategy, they show that its probability of misidentification is given as

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*} (\hat{a}_T^{\text{KCG}} \neq a^*(P^*)) \geq \frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2}.$$

Note that this upper bound comes from the upper bound of

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*} (\hat{a}_T^{\text{KCG}} \neq a^*(P^*)) = \liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*} (\hat{\mu}_{1, T}^{\text{SA}} - \hat{\mu}_{2, T}^{\text{SA}} \leq 0),$$

in case where arm 1 is the best arm ($a^*(P^*) = 1$). In contrast, in Theorem 4.1, we show that our strategy's upper bound is $\lim_{\Delta \rightarrow 0} \liminf_{T \rightarrow \infty} -\frac{1}{\Delta^2 T} \log \mathbb{P}_{P^*}(\hat{a}_T^{\text{IPW}} \neq a^*(P^*)) \geq \frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2}$. The difference between our upper bounds and theirs is the existence of a term, which vanishes as $\Delta \rightarrow 0$ with a speed faster than Δ^2 . This difference comes from the estimation error of the variances.

Intuitive explanation about the influence of the variance estimation. To understand the variance estimation, we rewrite the sample average as

$$\hat{\mu}_{a,T}^{\text{SA}} = \frac{1}{T} \sum_{t=1}^T \frac{1}{\frac{1}{T} \sum_{t=1}^T \mathbb{1}[A_t^{\text{KCG}} = a]} \mathbb{1}[A_t^{\text{KCG}} = a] Y_t = \frac{1}{T} \sum_{t=1}^T \frac{1}{w_a^*} \mathbb{1}[A_t^{\text{KCG}} = a] Y_t,$$

where we used $\sum_{t=1}^T \mathbb{1}[A_t^{\text{KCG}} = a] = w_a^* T$. Here, we consider a strategy that estimates w_a^* by estimating the variances. Let \tilde{w}_a be some estimator of w_a^* . Then, we design a strategy that draws arm a in $\tilde{w}_a T$ rounds in some way. In that case, the sample average roughly becomes

$$\tilde{\mu}_{a,T}^{\text{SA}} = \frac{1}{T} \sum_{t=1}^T \frac{1}{\frac{1}{T} \sum_{t=1}^T \mathbb{1}[\tilde{A}_t = a]} \mathbb{1}[\tilde{A}_t = a] Y_t = \frac{1}{T} \sum_{t=1}^T \frac{1}{\tilde{w}_a} \mathbb{1}[\tilde{A}_t = a] Y_t,$$

where \tilde{A}_t denotes an arm drawn by some strategy that draws arm a in $\tilde{w}_a T$ rounds. Then, if we recommend arm $\tilde{a}_T := \arg \max_{a \in \{1,2\}} \tilde{\mu}_{a,T}^{\text{SA}}$ as the best arm, we evaluate

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*}(\tilde{a}_T \neq a^*(P^*)) = \liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*}(\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}} \leq 0),$$

From Markov's inequality, for any $\lambda > 0$, we have

$$\log \mathbb{P}_{P^*}(\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}} \leq 0) \leq \mathbb{E} [\exp(T\lambda(\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}))].$$

To obtain the same upper bound as that in the Kaufmann et al. (2016)'s strategy, we consider the following decomposition:

$$\mathbb{E} [\exp(T\lambda(\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}))] = \mathbb{E} [\exp(T\lambda(\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}})) \exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))].$$

Suppose that the following holds in some way:

$$\begin{aligned} & \mathbb{E} [\exp(T\lambda(\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}})) \exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))] \\ &= \mathbb{E} [\exp(T\lambda(\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}))] \mathbb{E} [\exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))]. \end{aligned}$$

Note that this decomposition does not generally hold, but we assume it since it makes it easy to understand the variance estimation problem. Under the assumption, it holds that

$$\begin{aligned} & -\frac{1}{T} \log \mathbb{P}_{P^*}(\tilde{a}_T \neq a^*(P^*)) \\ & \geq -\frac{1}{T} \log \mathbb{E} [\exp(T\lambda(\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}))] - \frac{1}{T} \log \mathbb{E} [\exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))] \\ & \geq \frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2} - \frac{1}{T} \log \mathbb{E} [\exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))]. \end{aligned}$$

This inequality implies that to obtain the same upper bound as that of Kaufmann et al. (2016)'s strategy, we need to bound

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{E} [\exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))]$$

with an arbitrage rate of convergence; more exactly, we need to show that for any $\varepsilon > 0$,

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{E} [\exp(T\lambda(\{\tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}}\} - \{\hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}}\}))] \geq -\varepsilon$$

holds. However, it is impossible to achieve that convergence rate with commonly known theorems about convergence. Therefore, we introduced the small-gap regime, which evaluates the term as follows: for all $\epsilon > 0$, there exists $\delta(\epsilon) > 0$ such that for all $0 < \Delta < \delta(\epsilon)$, we have

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{E} \left[\exp \left(T\lambda \left(\{ \tilde{\mu}_{1,T}^{\text{SA}} - \tilde{\mu}_{2,T}^{\text{SA}} \} - \{ \hat{\mu}_{1,T}^{\text{SA}} - \hat{\mu}_{2,T}^{\text{SA}} \} \right) \right) \right] = -\epsilon \Delta^2.$$

Note that this argument is not rigorous and is simplified for explanation.

A.3 The AIPW, IPW, and Sample Average Estimators

A key component of our analysis is the AIPW estimator, which comprises an MDS and boasts minimum asymptotic variance. By using the properties of an MDS, we tackle the dependence among observations. The upper bound can also be applied to the Inverse Probability Weighting (IPW) estimator, but in this case, the upper bound may not coincide with the lower bound. This discrepancy occurs because the AIPW estimator's asymptotic variance is smaller than the IPW estimator's. The minimum variance property of the AIPW estimator stems from the efficient influence function (Hahn, 1998; Tsiatis, 2007).

We conjecture that the asymptotic optimality of strategies employing the naive sample average estimator in the recommendation rule can be demonstrated, although we do not prove it in this study. This is because Hahn et al. (2011) shows that, using the CLT, the AIPW and sample average estimators have the same asymptotic distribution. However, due to the inability to utilize MDS properties and the presence of sample dependency, the analysis becomes challenging when we derive a corresponding result for a large deviation (exponential rate of the probability of misidentification).

For the reader's reference, we detail the problems related to the IPW estimator and the sample average estimator.

The NA-IPW strategy. We consider the following strategy. In the NA-AIPW strategy, instead of the AIPW estimator, we use the following IPW estimator to estimate the means:

$$\hat{\mu}_{a,T}^{\text{IPW}} := \frac{1}{T} \sum_{t=1}^T \psi_{a,t}^{\text{IPW}}, \quad \text{where } \psi_{a,t}^{\text{IPW}} := \frac{\mathbb{1}[A_t = a] Y_{a,t}}{\hat{w}_{a,t}}.$$

At the end of the experiment (after the round $t = T$), we recommend \hat{a}_T^{IPW} as

$$\hat{a}_T^{\text{IPW}} := \begin{cases} 1 & \text{if } \hat{\mu}_{1,T}^{\text{IPW}} \geq \hat{\mu}_{2,T}^{\text{IPW}}, \\ 2 & \text{otherwise.} \end{cases}$$

We refer to this strategy as the NA-IPW strategy, whose probability of misidentification of this strategy is given as follows.

Theorem A.1 (Upper bound of the NA-IPW strategy). *for all $\epsilon > 0$, there exist $t(\epsilon) > 0$ and $\delta(\epsilon) > 0$ such that for all $T > t(\epsilon)$ and for any $P^* \in \{P \in \mathcal{P}^G \mid 0 < \Delta = \mu_1 - \mu_2 < \delta(\epsilon)\}$, if $\{(Y_{1,t}, Y_{2,t})\}_{t \in [T]}$ is generated from P^* , and \bar{C}_μ and \bar{C}_{σ^2} are sufficiently large positive values such that $\mu_a \in (-\bar{C}_\mu, \bar{C}_\mu)$ and $\sigma_a \in (1/\bar{C}_{\sigma^2}, \bar{C}_{\sigma^2})$ hold for $a \in \{1, 2\}$, then it holds that*

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*} \left(\hat{a}_T^{\text{IPW}} \neq a^*(P^*) \right) \geq \frac{\Delta^2}{2(\sigma_1 + \sigma_2) \left(\frac{\zeta_1^*}{\sigma_1} + \frac{\zeta_2^*}{\sigma_2} \right)} - \epsilon \Delta^2,$$

where $\zeta_a^* := \mathbb{E}_{P^*} [Y_{a,t}^2]$.

Proof. Let us define $V^{\text{IPW}} := \frac{\zeta_1^*}{w_1^*} + \frac{\zeta_2^*}{w_2^*} - \Delta^2$ and $\Psi^{\text{IPW}} := \left\{ \psi_{1,t}^{\text{IPW}} - \psi_{2,t}^{\text{IPW}} - \Delta \right\} / V^{\text{IPW}}$. Then, we have

$$\mathbb{E}_{P^*} \left[\left(\Psi_t^{\text{IPW}} \right)^2 \mid \mathcal{F}_{t-1} \right] = \mathbb{E}_{P^*} \left[\left(\psi_{1,t}^{\text{IPW}} - \psi_{2,t}^{\text{IPW}} - \Delta \right)^2 \mid \mathcal{F}_{t-1} \right]$$

$$\begin{aligned}
&= \mathbb{E}_{P^*} \left[\left(\frac{\mathbb{1}[A_t = 1]Y_{1,t}}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t = 2]Y_{2,t}}{\widehat{w}_{2,t}} - \Delta \right)^2 \middle| \mathcal{F}_{t-1} \right] \\
&= \mathbb{E}_{P^*} \left[\left(\frac{\mathbb{1}[A_t = 1]Y_{1,t}}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t = 2]Y_{2,t}}{\widehat{w}_{2,t}} \right)^2 \middle| \mathcal{F}_{t-1} \right] - \Delta^2 \\
&= \mathbb{E}_{P^*} \left[\frac{Y_{1,t}^2}{\widehat{w}_{1,t}} - \frac{Y_{2,t}^2}{\widehat{w}_{2,t}} \middle| \mathcal{F}_{t-1} \right] - \Delta^2 \rightarrow V^{\text{IPW}}
\end{aligned}$$

as $T \rightarrow \infty$. By replacing V and Ψ in the proof of Theorem 4.1 with V^{IPW} and Ψ^{IPW} , for all $\epsilon > 0$, there exist $t(\epsilon) > 0$ and $\delta(\epsilon) > 0$ such that for all $T > t(\epsilon)$ and for any $P^* \in \left\{ P \in \mathcal{P}^G \mid 0 < \Delta = \mu_1 - \mu_2 < \delta(\epsilon) \right\}$, we have

$$\liminf_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{P^*} (\widehat{a}_T^{\text{IPW}} \neq a^*(P^*)) \geq \frac{\Delta^2}{2 \left(\frac{\zeta_1^*}{w_1^*} + \frac{\zeta_2^*}{w_2^*} - \Delta^2 \right)} - \epsilon \Delta^2.$$

Then, for all $\epsilon > 0$, there exist $t(\epsilon) > 0$ and $\delta(\epsilon) > 0$ such that for all $T > t(\epsilon)$ and for any $P^* \in \left\{ P \in \mathcal{P}^G \mid 0 < \Delta = \mu_1 - \mu_2 < \delta(\epsilon) \right\}$, we have $\frac{\Delta^2}{2 \left(\frac{\zeta_1^*}{w_1^*} + \frac{\zeta_2^*}{w_2^*} \right)} - \epsilon \Delta^2$. The proof is complete. \square

Note that $2(\sigma_1 + \sigma_2) \left(\frac{\zeta_1^*}{\sigma_1} + \frac{\zeta_2^*}{\sigma_2} \right) \geq 2(\sigma_1 + \sigma_2)^2$ since $\left(\frac{\zeta_1^*}{\sigma_1} + \frac{\zeta_2^*}{\sigma_2} \right) \geq \left(\frac{\zeta_1^* - (\mu_1^*)^2}{\sigma_1} + \frac{\zeta_2^* - (\mu_2^*)^2}{\sigma_2} \right) = (\sigma_1 + \sigma_2)$. Therefore, the upper bound of probability of misidentification of the NA-IPW strategy is larger than that of the NA-AIPW strategy (Note that the inequality is flipped due to $-\frac{1}{T} \log \mathbb{P}_{P^*}$; that is, $\frac{\Delta^2}{2(\sigma_1 + \sigma_2)^2} \geq \frac{\Delta^2}{2(\sigma_1 + \sigma_2) \left(\frac{\zeta_1^*}{\sigma_1} + \frac{\zeta_2^*}{\sigma_2} \right)}$ implies that the upper bound of the probability of misidentification of the NA-AIPW strategy is smaller than that of the NA-IPW strategy). In the case of the evaluation using the CLT, similar results have been known in existing studies, such as Hirano et al. (2003) and Kato et al. (2020).

The NA-SA strategy. Next, we consider the following strategy. In the NA-AIPW strategy, instead of the AIPW estimator, we use the following sample average estimator:

$$\widehat{\mu}_{a,T}^{\text{SA}} := \frac{1}{\sum_{t=1}^T \mathbb{1}[A_t = a]} \sum_{t=1}^T \mathbb{1}[A_t = a] Y_t = \frac{1}{T} \sum_{t=1}^T \psi_{a,t}^{\text{SA}}, \quad \text{where } \psi_{a,t}^{\text{SA}} := \frac{\mathbb{1}[A_t = a] Y_t}{\frac{1}{T} \sum_{t=1}^T \mathbb{1}[A_t = a]}.$$

At the end of the experiment (after the round $t = T$), we recommend $\widehat{a}_T^{\text{IPW}}$ as

$$\widehat{a}_T^{\text{SA}} := \begin{cases} 1 & \text{if } \widehat{\mu}_{1,T}^{\text{SA}} \geq \widehat{\mu}_{2,T}^{\text{SA}}, \\ 2 & \text{otherwise.} \end{cases}$$

We refer to this strategy as the NA-SA strategy. Evaluation of the probability of misidentification of this strategy is not easy since we cannot employ a martingale property, which has been used in the analysis of the NA-AIPW strategy and the NA-IPW strategy. In order to derive its upper bound, we need to evaluate $\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}}$. Here, note that

$$\begin{aligned}
\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}} &= \widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}} - \left\{ \frac{\mathbb{1}[A_t = 1](Y_{1,t} - \mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t = 2](Y_{2,t} - \mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\} \\
&\quad + \left\{ \frac{\mathbb{1}[A_t = 1](Y_{1,t} - \mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t = 2](Y_{2,t} - \mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\}
\end{aligned}$$

holds, and the variance of $\Psi_t^* := \left\{ \frac{\mathbb{1}[A_t=1](Y_{1,t}-\mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t=2](Y_{2,t}-\mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\}$ is V in the proof of Theorem 4.1, and $\{\Psi_t^*\}_{t=1}^T$ consists of an MDS. Therefore, if $\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}} - \left\{ \frac{\mathbb{1}[A_t=1](Y_{1,t}-\mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t=2](Y_{2,t}-\mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\} = 0$, we can directly apply the proof of Theorem 4.1 to obtain the same upper bound in Theorem 4.1. However, $\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}} - \left\{ \frac{\mathbb{1}[A_t=1](Y_{1,t}-\mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t=2](Y_{2,t}-\mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\}$ is not zero and remains as a bias term, and it is known that its evaluation requires several techniques. For example, to show $\sqrt{T}(\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}} - \Delta) \xrightarrow{d} \mathcal{N}(0, V)$, Hahn et al. (2011) bounds $\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}} - \left\{ \frac{\mathbb{1}[A_t=1](Y_{1,t}-\mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t=2](Y_{2,t}-\mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\}$ using the property of the stochastic equicontinuity, which is based on the arguments in Hirano et al. (2003). This problem is related to the use of the Donsker condition in semiparametric analysis, as explained in Kennedy (2016). We may show the upper bound of the NA-SA strategy by using a similar approach used in Hahn et al. (2011), but there are two issues. First, it is unknown what condition corresponds to the stochastic equicontinuity in the setting of BAI, where the samples are dependent. Second, it is unclear whether we can directly apply the stochastic equicontinuity or similar properties to show the large deviation upper bound since such conditions have been used for the central limit evaluation. Therefore, although the findings of Hirano et al. (2003) and Hahn et al. (2011) may aid in resolving this issue, it is an open issue how we use it. Note that this issue caused by the bias of $\widehat{\mu}_{1,T}^{\text{SA}} - \widehat{\mu}_{2,T}^{\text{SA}}$, which is non-zero. Also note that in contrast, the bias of $\widehat{\mu}_{1,T}^{\text{AIPW}} - \widehat{\mu}_{2,T}^{\text{AIPW}}$ is zero due to the properties of an MDS.

A.4 The Tracking Strategy

In fixed-confidence BAI, the tracking strategy is popular, as used in Garivier & Kaufmann (2016). In the existing studies of the Neyman allocation, such a strategy has been used. For example, Hahn et al. (2011) splits the whole samples into two groups. In the first stage, we uniformly randomly draw each arm and estimate w_a^* . In the second stage, for the estimators \widehat{w}_a of w_a^* we draw each arm so that $\frac{1}{T} \sum_{t=1}^T \mathbb{1}[A_t = a] = \widehat{w}_a$ holds. Then, Hahn et al. (2011) estimates $\Delta = \mu_1^* - \mu_2^*$ using the sample average estimator. This strategy is quite similar to that in Garivier & Kaufmann (2016), since it draws arms to track the ratio of w_a^* .

However, the strategy of Hahn et al. (2011) makes analyzing upper bounds difficult. As we explained in Section A.3, in our analysis, the unbiasedness of the AIPW estimator plays an important role. In contrast, if we use the tracking strategy, we cannot employ the property of $\mathbb{E}_{P^*} \left[\frac{\mathbb{1}[A_t=a]}{\widehat{w}_t} \mid \mathcal{F}_{t-1} \right] = 1$. Note that the NA-AIPW strategy draws arm A_t with probability \widehat{w}_t , but the tracking strategy draws arm A_t more complicatedly, under which we cannot use the martingale property.

As well as we explained in Section A.3, the bias term makes the analysis significantly difficult. According to the existing studies, we need to use some techniques for the analysis, such as the Donsker condition (Hirano et al., 2003; Hahn et al., 2011). Existing studies have proposed using the AIPW estimator to avoid this issue, as shown in van der Laan (2008) and Kato et al. (2020).

Thus, although we acknowledge the possibility of using the tracking strategy, the analysis requires some sophisticated techniques. We expect that existing studies such as (Hirano et al., 2003) and Hahn et al. (2011) will help the analysis, but it is still an open issue. Note that even in the tracking strategy, existing strategy such as (Hirano et al., 2003) and Hahn et al. (2011) bypass the evaluation of the AIPW-type estimators in the analysis. This proof procedure is also related to the semiparametric efficiency bound (Hahn, 1998), under which the semiparametric efficient score is given as $\Psi_t^* = \left\{ \frac{\mathbb{1}[A_t=1](Y_{1,t}-\mu_1^*)}{\widehat{w}_{1,t}} - \frac{\mathbb{1}[A_t=2](Y_{2,t}-\mu_2^*)}{\widehat{w}_{2,t}} - \Delta \right\}$.

B Related Work

This section presents related works.

B.1 On the Asymptotic Optimality in Fixed-Budget BAI

There is a long debate on the optimal strategies for fixed-budget BAI. Glynn & Juneja (2004) develops their strategies by using the large deviation principles. However, while they justify their strategies using the large deviation principles, they do not provide lower bounds for strategies. Therefore, there remains a question about whether their strategies are truly asymptotically optimal.

Kaufmann et al. (2016) establishes distribution-dependent lower bounds for BAI with fixed confidence and budget, utilizing change-of-measure arguments. According to their results, we can confirm that for two-armed Gaussian bandits, the strategy of Glynn & Juneja (2004) is optimal.

However, Kaufmann et al. (2016) leaves lower bounds for multi-armed fixed-budget BAI as an open issue. Based on the arguments of Glynn & Juneja (2004) and Russo (2020), Kasy & Sautmann (2021) attempts to derive an asymptotically optimal strategy, but their attempt does not succeed. As pointed out by Ariu et al. (2021), without additional assumptions, there exists an instance P^* whose lower bound is larger than that of Kaufmann et al. (2016). This result is based on another lower bound discovered by Carpentier & Locatelli (2016). These arguments are summarized by Qin (2022).

To address this issue, Kato et al. (2023b) and Degenne (2023) consider a restriction such that sampling rules do not depend on P^* . Under this restriction, we can show the asymptotic optimality of the strategy provided by Glynn & Juneja (2004), which requires full knowledge about P^* and is practically infeasible.

Komiyama et al. (2022) and Atsidakou et al. (2023) discuss asymptotically optimal strategies from minimax and Bayesian perspectives, respectively, where the leading factor ignoring some constant terms of lower and upper bounds match, unlike our optimality up to constant terms. This open issue is further explored by Komiyama et al. (2022), Wang et al. (2023a), Wang et al. (2023b), and Kato (2023).

Note that in the fixed confidence BAI setting, Garivier & Kaufmann (2016) proposes a strategy with an upper bound matching the derived lower bound. However, in the fixed-budget BAI, it remains unclear whether a strategy with an upper bound matching Kaufmann et al. (2016)’s lower bound exists.

Alternative lower bounds have been proposed by Audibert et al. (2010), Bubeck et al. (2011), Komiyama et al. (2023) and Kato et al. (2023a) for the expected simple regret minimization, which is another performance measure different from the probability of misidentification.

Some research employs local asymptotics to examine the asymptotic optimality of the Neyman allocation rule in this context, such as Armstrong (2022) and Adusumilli (2022).

Ordinal optimization in the operations research community is another related field (Ahn et al., 2021; Chen et al., 2000).

B.2 Extension to BAI in Multi-Armed Bandit (MAB) Problems

In contrast to two-armed bandit problems and BAI with fixed confidence, lower bounds for MAB problems remain unknown. One primary reason is the reversal of KL divergence. Kato et al. (2023b), Degenne (2023), Kato (2023) consider strategies that use sampling rules that are (asymptotically) invariant for any $P^* \in \mathcal{P}^G$. Such a class of strategies is sometimes called *static* in the sense that it cannot estimate parameters during an adaptive experiment to avoid the dependency on P^* . However, if we consider Gaussian bandit models, sampling strategies that are invariant for P^* do not imply non-adaptive (static) strategies because we can still adaptively estimate the variances during an adaptive experiment (the variances are assumed to be the same for any P^*).

C Experimental Results

We present experimental results in Section 6.

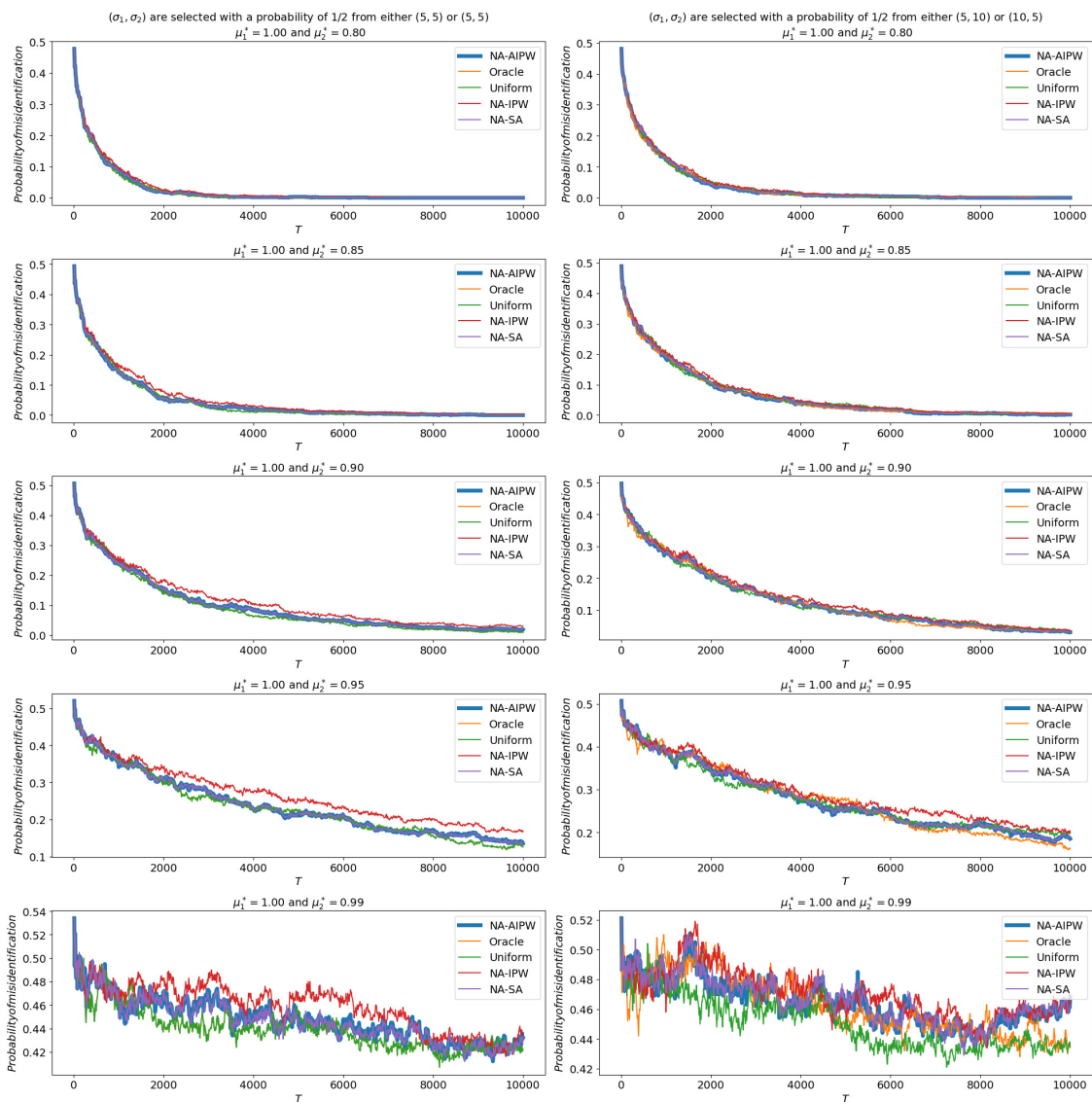


Figure 3: The results are under the first setting. We set $\mu_1^* = 1.00$ and choose μ_2^* from the set $\{0.80, 0.85, 0.90, 0.95, 0.99\}$. The variances (σ_1, σ_2) are selected with a probability of $1/2$ from either $(5, v_2)$ or $(v_2, 5)$, where v_2 is chosen from $5, 10$. We conduct 1,000 independent trials and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9,900, 10,000\}$.

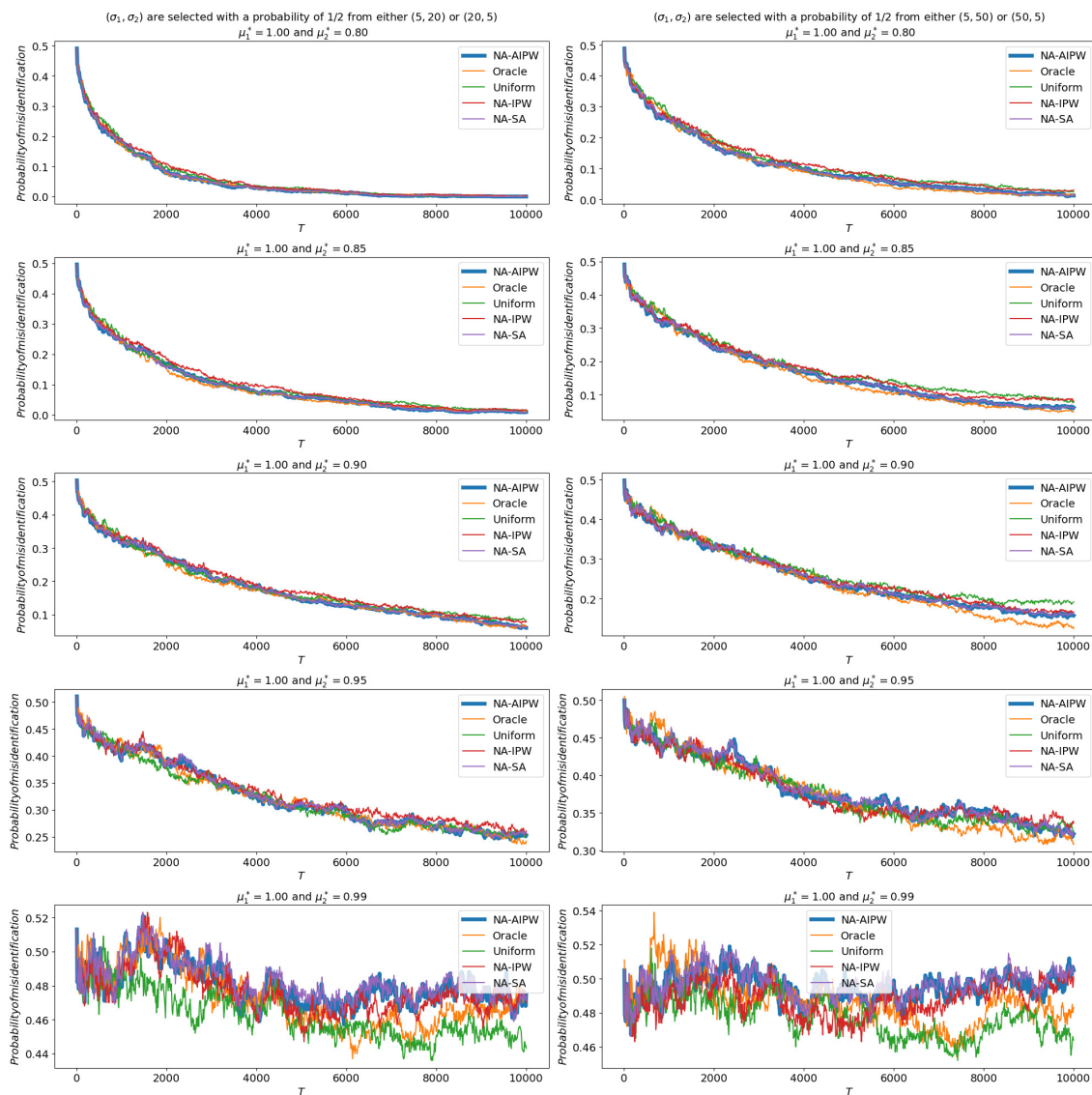


Figure 4: The results are under the first setting. We set $\mu_1^* = 1.00$ and choose μ_2^* from the set $\{0.80, 0.85, 0.90, 0.95, 0.99\}$. The variances (σ_1, σ_2) are selected with a probability of $1/2$ from either $(5, v_2)$ or $(v_2, 5)$, where v_2 is chosen from $20, 50$. We conduct 1,000 independent trials and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9,900, 10,000\}$.

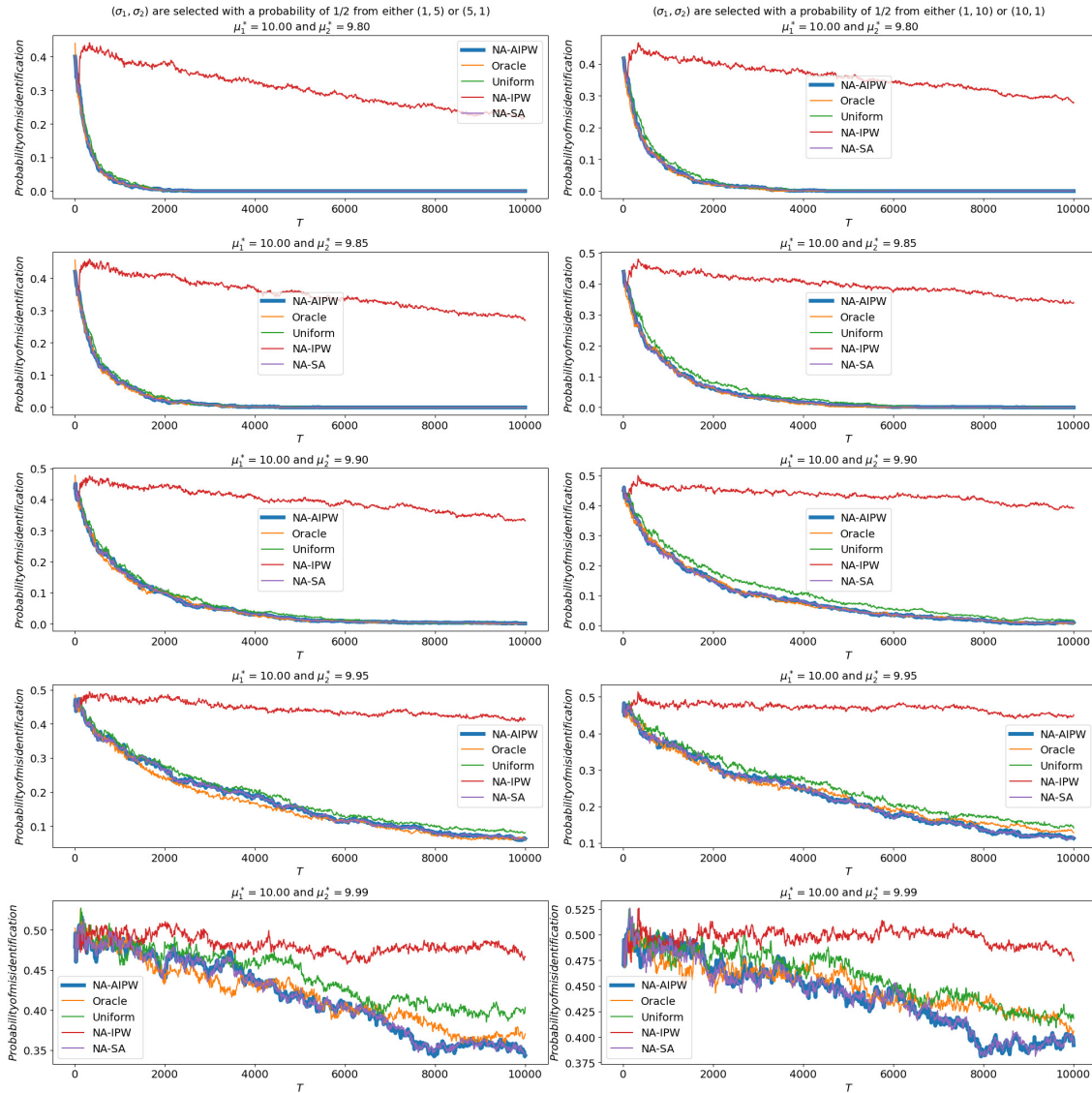


Figure 5: The results are under the first setting. We set $\mu_1^* = 10.00$ and choose μ_2^* from the set $\{9.80, 9.85, 9.90, 9.95, 9.99\}$. The variances (σ_1, σ_2) are selected with a probability of $1/2$ from either $(1, v_2)$ or $(v_2, 1)$, where v_2 is chosen from $5, 10$. We conduct 1,000 independent trials and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9,900, 10,000\}$.

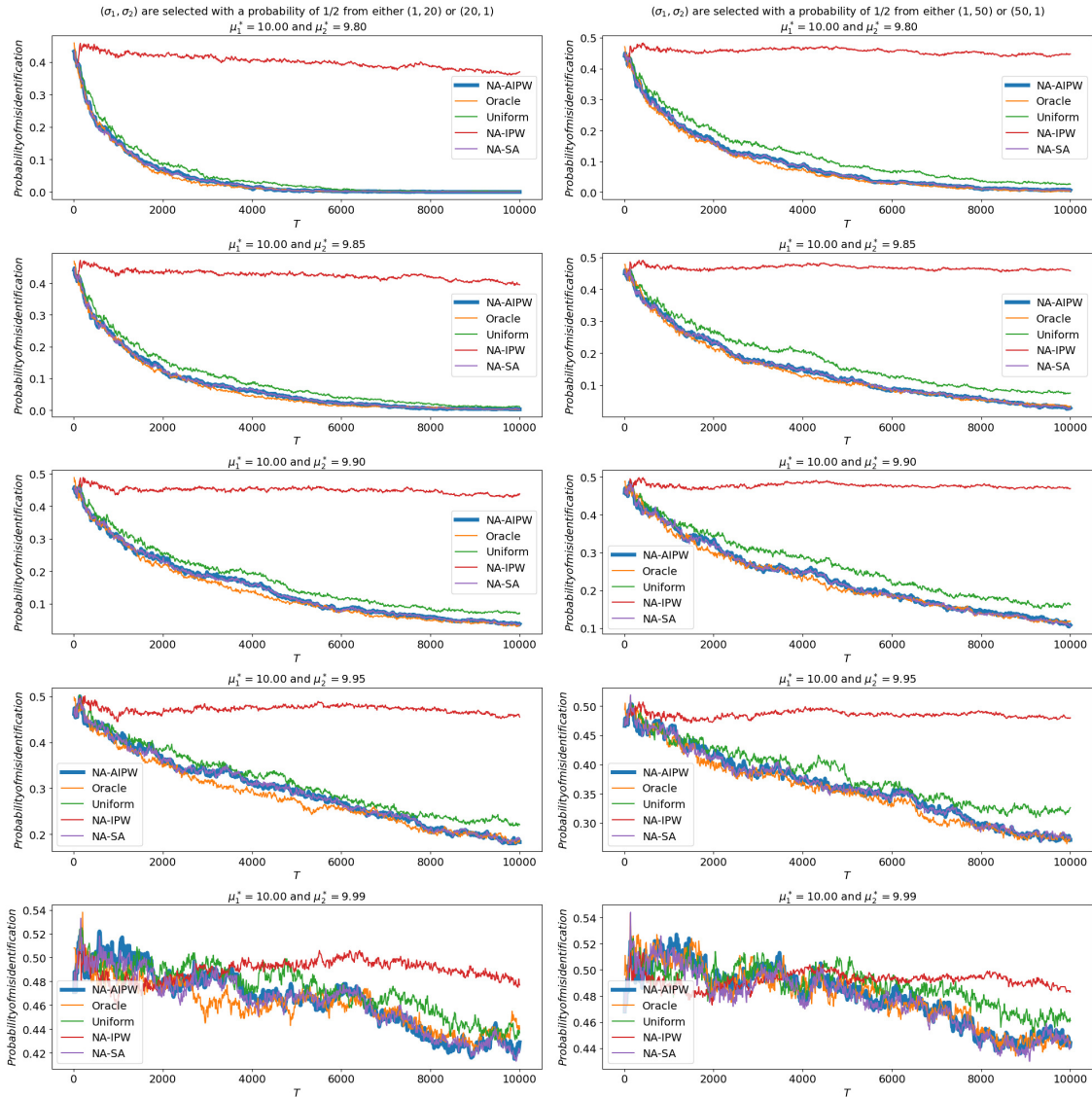


Figure 6: The results are under the first setting. We set $\mu_1^* = 10.00$ and choose μ_2^* from the set $\{9.80, 9.85, 9.90, 9.95, 9.99\}$. The variances (σ_1, σ_2) are selected with a probability of $1/2$ from either $(1, v_2)$ or $(v_2, 1)$, where v_2 is chosen from $20, 50$. We conduct $1,000$ independent trials and report the empirical probability of misidentification at $T \in \{100, 200, 300, \dots, 9,900, 10,000\}$.