
Future Information-Directed Sampling for Bayesian Nonstationary Bandits

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Exploration–exploitation is a central trade-off in bandit learning. While classical
2 algorithms such as upper confidence bound methods and Thompson Sampling effectively
3 balance this trade-off in stationary environments, their exploration strategies
4 mainly reduce uncertainty about the current optimal arm, which can be insufficient
5 in nonstationary settings where future optimal arms may differ substantially from
6 current ones. In this paper, we propose Future Information-Directed Sampling
7 (FIDS), a new algorithm for Bayesian nonstationary bandits that explicitly explores
8 to gather information about future optimal arms. We show that FIDS achieves
9 regret comparable to Thompson Sampling up to a small constant factor, while
10 being able to exploit predictive information structures that conventional exploration
11 objectives fail to capture. To address the practical difficulty of posterior inference,
12 we further propose a supervised-learning-based approximation framework that
13 learns the FIDS policy from offline data, and demonstrate its effectiveness on
14 synthetic benchmarks.

15 1 Introduction

16 In this paper, we study multi-armed bandits (MABs) in nonstationary environments. MAB problems
17 have been extensively studied over the past decades as a fundamental sequential decision-making
18 framework [10, 11, 13], primarily under the assumption that the underlying environment is stationary.
19 However, this assumption is often unrealistic in real-world applications such as recommendation
20 systems [4] and online advertising [19], where user preferences and reward distributions may evolve
21 over time. Motivated by these applications, we consider a nonstationary setting in which the
22 environments appearing over time are generated according to a known prior distribution. We refer to
23 this setting as Bayesian nonstationary bandits.

24 A core challenge in nonstationary bandits is balancing the exploration-exploitation tradeoff [13],
25 as the environment changes over time. Nevertheless, it is possible to build an intuition by looking
26 at the literature of stationary MABs: in stationary Bayesian bandit problems, this trade-off can be
27 tackled by using methods like Thompson Sampling (TS) [25], which achieves both strong theoretical
28 guarantees and empirical performance [23]. An information-theoretic analysis of TS is given by [22],
29 who show that the ratio of instantaneous regret to mutual information gained about the optimal arm is
30 bounded uniformly over time. In stationary environments, gathering information about the current
31 optimal arm is a natural exploration objective, since the optimal arm remains informative for future
32 decisions. However, this intuition breaks down in nonstationary environments, where the optimal
33 arm may change over time. As a result, exploration strategies that focus only on reducing uncertainty
34 about the current optimal arm may fail to gather information that is useful for future decision-making.

35 Motivated by this observation, we propose a new exploration objective that explicitly gathers informa-
36 tion about future optimal arms, while retaining instantaneous expected reward maximization as the

37 exploitation objective. To balance this new exploration–exploitation trade-off, we introduce Future
38 Information-Directed Sampling (FIDS). Our main contributions are summarized as follows:

- 39 1. We prove in Theorem 3.2 that FIDS achieves the same regret upper bound as Thompson
40 Sampling in Bayesian nonstationary environments, up to a small constant factor.
- 41 2. Through several analytic examples, we demonstrate that FIDS can exploit predictive in-
42 formation structures that conventional exploration objectives fail to capture, leading to
43 significantly improved performance over baseline methods.
- 44 3. In practice, the prior distribution may be unknown and posterior inference may be com-
45 putationally intractable. To sidestep these difficulties, we propose a supervised-learning
46 framework that trains an approximate FIDS policy from offline data, inspired by DPT [14],
47 which similarly trains a Thompson Sampling policy via supervision. Empirically, the learned
48 FIDS policy outperforms DPT and other baselines by effectively exploiting the information
49 structure of the environment.

50 2 Preliminaries and Problem Formulation

51 In this section we introduce the problem formulation, and provide background knowledge on Thomp-
52 son Sampling and Information Directed Sampling.

53 We consider a finite MAB problem with horizon n and an action set \mathcal{A} with $|\mathcal{A}| = K$ arms. We
54 model the non-stationary environment by a random variable ν drawn from a prior distribution known
55 to the learner. Here $\nu = (\nu_{t,a})_{t \in [n], a \in \mathcal{A}}$ where each $\nu_{t,a}$ specifies the reward distribution of arm a at
56 round t . Consider $R_{t,a}$ to be the reward one will see if they play arm a at round t . We assume that
57 conditional on $\nu_{t,a}$, the reward $R_{t,a}$ is sampled independently as $R_{t,a} \sim \nu_{t,a}$. As a special case, if
58 for every $a \in \mathcal{A}$, we have $\nu_{t,a} = \nu_{s,a}$ for all $t, s \in [n]$ almost surely, then the model reduces to the
59 standard stationary Bayesian bandit setting. We also make the following assumption on the reward
60 distributions.

61 **Assumption 2.1.** *There exists a constant $\sigma \in \mathbb{R}$, such that for all $a \in \mathcal{A}$, conditioned on \mathcal{F}_t ,
62 $R_{t,a} - \mathbb{E}[R_{t,a} | \mathcal{F}_t]$ is σ sub-Gaussian. And we assume σ is known to the learner.*

63 In each round t , the learner interacts with the environment by selecting an action A_t according to a
64 policy π_t , which maps the history $\mathcal{F}_{t-1} = (A_1, R_1, \dots, A_{t-1}, R_{t-1})$ to a distribution over \mathcal{A} , and
65 observes reward $R_t \triangleq R_{t,A_t}$. Let $\pi = (\pi_t)_{t=1}^n$. We evaluate π using the following Bayesian regret:

$$\text{Regret}(\pi, n) = \mathbb{E} \left[\sum_{t=1}^n (R_{t,A_t^*} - R_{t,A_t}) \right], \quad (1)$$

66 where $A_t^* \in \arg \max_{a \in \mathcal{A}} \mu(\nu_{t,a})$. And $\mu(\nu_{t,a})$ denotes the mean of the distribution $\nu_{t,a}$.

67 **Notation.** We introduce the following shorthand notations that will be used throughout the analysis.
68 For $i, j \in [n]$ with $i \leq j$, let $A_{a,b}^* \triangleq [A_i^*, \dots, A_j^*]$. Let $\mathbb{P}_t(\cdot) \triangleq \mathbb{P}(\cdot | \mathcal{F}_{t-1})$ denote the posterior
69 measure given the history up to round $t - 1$. For random variables X and Y , we denote $\mathbb{E}_t[X]$,
70 $H_t(X)$, $H_t(X | Y)$, and $I_t(X; Y)$ to be the expectation, entropy, conditional entropy, and mutual
71 information under the posterior \mathbb{P}_t , respectively. Formal definitions of these information-theoretic
72 quantities can be found in Appendix B. Note that $\mathbb{E}_t[X]$, $H_t(X)$, $H_t(X | Y)$, and $I_t(X; Y)$ are all
73 random variables, since they are functions of the history \mathcal{F}_{t-1} which is a random variable.

74 2.1 Preliminaries about TS and IDS

75 We now recall a key result from the information-theoretic analysis of TS by [22], which underpins
76 both the design of FIDS and its regret analysis.

77 **Proposition 2.1.** *(Corollary 1 in [22]) Denote the Thompson sampling policy by $\pi^{TS} = (\pi_t^{TS})_{t=1}^n$.
78 Suppose that conditioned on \mathcal{F}_t , $R_{t,a} - \mathbb{E}[R_{t,a} | \mathcal{F}_t]$ is σ sub-Gaussian, where σ is a constant
79 number in \mathbb{R} . Then,*

$$\frac{\sum_{a \in \mathcal{A}} \pi_t^{TS}(a) \mathbb{E}_t [R_{t,A_t^*} - R_{t,a}]}{\sqrt{\sum_{a \in \mathcal{A}} \pi_t^{TS}(a) I_t(A_t^*; R_{t,a})}} \leq \sigma \sqrt{2K}, \quad \forall t \in [n], \text{ a.s.} \quad (2)$$

80 The left-hand side of Equation (2) is known as the *information ratio*. Using this bound, [22] show
 81 that the Bayesian regret of TS is at most $\sigma\sqrt{2nKH(A^*)}$, where $A^* \triangleq A_1^* = \dots = A_n^*$ since the
 82 environment is stationary.

83 In addition, IDS [21] selects a policy π_t^{IDS} that explicitly minimizes the information ratio at each
 84 round t , i.e.,

$$\pi_t^{\text{IDS}} \in \operatorname{argmin}_{\pi_t \in \Delta(\mathcal{A})} \frac{\sum_{a \in \mathcal{A}} \pi_t(a) \mathbb{E}_t [R_{t,A_t^*} - R_{t,a}]}{\sqrt{\sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_t^*; R_{t,a})}}, \quad (3)$$

85 IDS enjoy the same worst-case regret bound as TS, but can achieve significantly improved perfor-
 86 mance in certain scenarios (see [21] for more details).

87 Using the same proof technique, [18] analyze TS in the nonstationary setting of Section 2 and obtain
 88 a regret bound of $O(\sigma\sqrt{2nKH([A_1^*, \dots, A_n^*])})$. This subsumes the stationary result: when the
 89 environment is stationary, $A_1^* = \dots = A_n^* = A^*$ almost surely, so $H([A_1^*, \dots, A_n^*]) = H(A^*)$.

90 3 FIDS Algorithm

91 In this section, we develop FIDS by addressing a key limitation of TS and IDS in nonstationary
 92 environments. In stationary settings, both algorithms achieve low regret by bounding the information
 93 ratio

$$\frac{\sum_{a \in \mathcal{A}} \pi_t(a), \mathbb{E}_t [R_{t,A_t^*} - R_{t,a}]}{\sqrt{\sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_t^*; R_{t,a})}},$$

94 which controls the trade-off between instantaneous regret and information gained about the optimal
 95 arm A_t^* . This works because A_t^* is fixed, so reducing uncertainty about it directly improves future
 96 decisions. In nonstationary environments, however, A_t^* may vary with t , so information gained about
 97 the current optimum need not reduce uncertainty about future optimal arms. FIDS addresses this by
 98 replacing the per-round information term with one that targets the sequence of optimal arms.

99 To address the exploration-exploration problem in a meaningful way, we introduce *Forward-looking*
 100 *Information-Directed Sampling* (FIDS), which balances instantaneous regret against information
 101 gained about future optimal arms. Under Assumption 2.1, FIDS is defined by the policy sequence
 102 $\pi^{\text{FIDS}} = (\pi_t^{\text{FIDS}})_{t=1}^n$, where π_t^{FIDS} is any policy that solves the following optimization problem

$$\pi_t^{\text{FIDS}} \in \operatorname{argmin}_{\pi_t \in \Delta(\mathcal{A})} \sum_{a \in \mathcal{A}} \pi_t(a) \mathbb{E}_t [R_{t,A_t^*} - R_{t,a}] - \sigma \sqrt{2K \cdot \sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_{t+1:n}^*; R_{t,a})} \quad (4)$$

103 Note that solving Equation (4) is equivalent to solving the following Equation (5)

$$\pi_t^{\text{FIDS}} \in \operatorname{argmax}_{\pi_t \in \Delta(\mathcal{A})} \sum_{a \in \mathcal{A}} \pi_t(a) \mathbb{E}_t [R_{t,a}] + \sigma \sqrt{2K \cdot \sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_{t+1:n}^*; R_{t,a})}, \quad (5)$$

104 as $\sum_{a \in \mathcal{A}} \pi_t(a) \mathbb{E}_t [R_{t,A_t^*}] = \mathbb{E}_t [R_{t,A_t^*}]$ is a constant with respect to π_t .

105 A useful property of Equation (5) is that the optimization in Equation (5) always admits an optimal
 106 policy supported on at most two arms, as shown in the following lemma (proof deferred to the
 107 appendix in section E.1).

108 **Lemma 3.1.** *Given $\mathbb{E}_t [R_{t,a}]$ and $I_t([A_{t+1}^*, \dots, A_n^*]; R_{t,a})$, Equation (5) is maximized at some π^*
 109 which has at most 2 non-zero elements.*

110 Intuitively, the objective in (5) is concave in π_t , so at any optimum every arm in the support must have
 111 the same marginal contribution. This equal-gradient condition forces a linear relationship between
 112 the reward and information values of the supported arms, and in the proof we show that any optimal
 113 mixture can be replicated by a two-point mixture of the arms in the support.

114 Before proceeding with the regret bound of FIDS, we note a subtlety in applying IDS to the nonsta-
 115 tionary setting.

116 **Remark 3.1.** In nonstationary environments, $\sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_t^*; R_{t,a})$ can equal zero: even if the
 117 current optimal arm A_t^* is known with certainty, the future optimal arms A_{t+1}^*, \dots, A_n^* may remain
 118 unknown, so the problem is still nontrivial. Because the standard ratio form of IDS is undefined
 119 when this information term vanishes, we adopt the following additive variant for nonstationary
 120 environments:

$$\pi_t^{IDS(ns)} \in \operatorname{argmax}_{\pi_t \in \Delta(\mathcal{A})} \sum_{a \in \mathcal{A}} \pi_t(a) \mathbb{E}_t[R_{t,a}] + \sigma \sqrt{2K \cdot \sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_t^*; R_{t,a})} \quad (6)$$

121 All experiments involving IDS use this formulation.

122 We now show that FIDS achieves a regret bound comparable to that of TS in nonstationary environ-
 123 ments.

124 **Theorem 3.2.** Under Assumption 2.1,

$$\operatorname{Regret}(\pi^{\text{FIDS}}, n) \leq 2\sigma \sqrt{2K \cdot n \cdot H([A_1^*, \dots, A_n^*])} \quad (7)$$

125 This bound matches the TS regret bound of [18] up to a constant factor of 2. Despite this worst-case
 126 similarity, FIDS can exploit additional structure in the environment; in Section 3.1, we show that it
 127 achieves strictly better performance in several settings. Here is the proof of the theorem.

128 *Proof.* Here we outline the proof with several supporting lemmas. The proof of lemmas
 129 can be found in Section C. For ease of the notation, in this section, we use $\Delta_t(\pi_t) :=$
 130 $\sum_{a \in \mathcal{A}} \pi_t(a) \mathbb{E}_t[R_{t,A_t^*} - R_{t,a}]$, $\tilde{g}_t(\pi_t) := \sum_{a \in \mathcal{A}} \pi_t(a) I_t([A_{t+1}^*, \dots, A_n^*]; R_{t,a})$ and $g_t(\pi_t) :=$
 131 $\sum_{a \in \mathcal{A}} \pi_t(a) I_t(A_t^*; R_{t,a})$. When an action A_t is sampled from π_t , we slightly abuse the nota-
 132 tion by letting $\Delta_t(A_t) := \mathbb{E}_t[R_{t,A_t^*} - R_{t,A_t}]$, $\tilde{g}_t(A_t) := I_t([A_{t+1}^*, \dots, A_n^*]; (A_t, R_{t,A_t}))$ and
 133 $g_t(A_t) := I_t(A_t^*; (A_t, R_{t,A_t}))$, as Lemma C.1 shows that $\Delta_t(A_t) = \Delta_t(\pi_t)$, $\tilde{g}_t(A_t) = \tilde{g}_t(\pi_t)$,
 134 and $g_t(A_t) = g_t(\pi_t)$.

135 The first part of the proof connects π^{FIDS} with π^{TS} . By the Tower Rule, $\operatorname{Regret}(\pi^{\text{FIDS}}, n) =$
 136 $\mathbb{E} \left[\sum_{t=1}^n \Delta_t(\pi_t^{\text{FIDS}}) \right]$, then by adding and subtracting $\sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})}$, we get

$$\begin{aligned} \operatorname{Regret}(\pi^{\text{FIDS}}, n) &= \mathbb{E} \left[\sum_{t=1}^n \Delta_t(\pi_t^{\text{FIDS}}) - \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})} + \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})} \right] \\ &\stackrel{(a)}{\leq} \mathbb{E} \left[\sum_{t=1}^n \Delta_t(\pi_t^{\text{TS}}) - \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{TS}})} + \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})} \right] \\ &\stackrel{(b)}{=} \mathbb{E} \left[\sum_{t=1}^n \Delta_t(\pi_t^{\text{TS}}) - \sigma \sqrt{2K \cdot g_t(\pi_t^{\text{TS}})} - \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{TS}})} \right. \\ &\quad \left. + \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})} + \sigma \sqrt{2K \cdot g_t(\pi_t^{\text{TS}})} \right] \\ &\stackrel{(c)}{\leq} \mathbb{E} \left[\sum_{t=1}^n \sigma \sqrt{2K \cdot g_t(\pi_t^{\text{TS}})} - \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{TS}})} + \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})} \right] \end{aligned}$$

137 where (a) uses the fact that π_t^{FIDS} is the minimizer of Equation (4); (b) adds and subtracts
 138 $\sigma \sqrt{2K \cdot g_t(\pi_t^{\text{TS}})}$; (c) holds because of Proposition 2.1. Then by Lemma C.2,

$$\mathbb{E} \left[\sum_{t=1}^n \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{FIDS}})} \right] \leq \sigma \sqrt{2n \cdot K \cdot H(A_{1:n}^*)} \quad (8)$$

139 by Lemma C.4,

$$\mathbb{E} \left[\sum_{t=1}^n \sigma \sqrt{2K \cdot g_t(\pi_t^{\text{TS}})} - \sigma \sqrt{2K \cdot \tilde{g}_t(\pi_t^{\text{TS}})} \right] \leq \sigma \sqrt{2n \cdot K \cdot H(A_{1:n}^*)} \quad (9)$$

140 Combining Equation (8) and (9) proves the theorem. \square

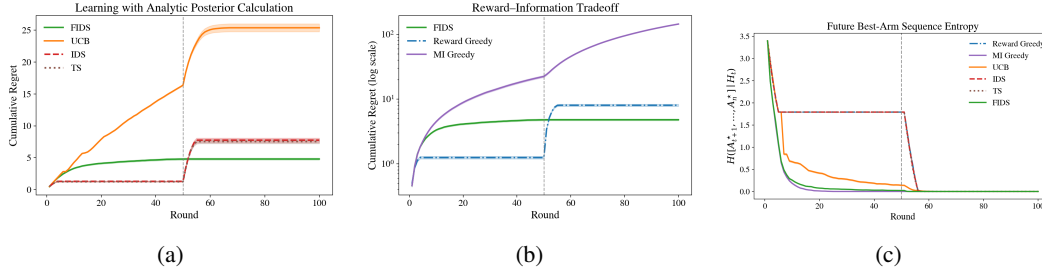


Figure 1: (a) Cumulative regret of FIDS, IDS and TS. (b) Cumulative regret (log scale) of FIDS, Reward Greedy (only optimize for the first reward term of Equation (5)), and MI only (only optimize for the second mutual information term of Equation (5)). (c) The entropy of $[A_{t+1}^*, \dots, A_n^*]$ conditioned on the history.

141 In the next subsection, we illustrate settings in which TS and IDS fail to exploit the information
 142 structure of the environment, while FIDS does not.

143 3.1 Illustrative Examples: When Does Information Structure Matter?

144 We present two examples that illustrate the role of information structure. The first is an i.i.d.
 145 environment in which the optimal arm is redrawn independently each round, so observations carry
 146 no information about future optimal arms. The second features a change point where information
 147 about the future environment is encoded in a currently suboptimal arm; we complement the analytical
 148 results for this last example with numerical experiments.

149 **Example 3.1** (Memoryless environment). *At each round t , the environment ν_t is sampled i.i.d.*
 150 *with $K = 2$, where $(\mu_{t,1}, \mu_{t,2}) = (0.4, 0.2)$ with probability $1/2$ and $(\mu_{t,1}, \mu_{t,2}) = (0.3, 0.4)$ with*
 151 *probability $1/2$. Conditional on ν_t , rewards satisfy $R_{t,a} \sim \text{Bernoulli}(\mu_{t,a})$ for $a \in \{1, 2\}$.*

152 In this setting, exploration is futile: no observation at round t can reduce uncertainty about future
 153 environments. As shown in Proposition E.1, FIDS recognizes this and always selects arm 1, the arm
 154 with the highest expected instantaneous reward, achieving cumulative reward $0.35n$. In contrast,
 155 TS and IDS allocate probability to unnecessary exploration, achieving only $0.325n$ and $0.3294n$
 156 respectively.

157 Example 3.1 serves as an environment where no future information can be gathered. We next consider
 158 a setting where information about future optimal arms is encoded in a currently suboptimal arm, and
 159 FIDS is able to exploit this cross-environment structure.

160 **Example 3.2** (Magic arm environment). *We consider a non-stationary 6-arm bandit with horizon*
 161 *$n = 100$ and a single change point at $t = 50$. The environment is **Env 1** for $t \leq 50$ and **Env 2** for*
 162 *$t > 50$, and is stationary within each segment.*

163 *Env 1. Let $Z_1 \sim \text{Unif}(\{1, \dots, 5\})$. For $a \in \{1, \dots, 5\}$, $\mu_a^1 = 0.9$ if $a = Z_1$ and $\mu_a^1 = 0.4$*
 164 *otherwise, with $R_{t,a} \sim \mathcal{N}(\mu_a^1, 0.1^2)$. For $a = 6$, $\mu_6^1 = Z_2/7$ and $R_{t,6} \sim \mathcal{N}(\mu_6^1, 0.1^2)$.*

165 *Env 2. Let $Z_2 \sim \text{Unif}(\{1, \dots, 6\})$, independent of Z_1 . For all $a \in \mathcal{A}$, $\mu_a^2 = 3.0$ if $a = Z_2$ and*
 166 *$\mu_a^2 = 0.5$ otherwise, with $R_{t,a} \sim \mathcal{N}(\mu_a^2, 0.5^2)$.*

167 Note that the information about Env 2 is encoded in the mean reward of arm 6 in Env 1. In particular,
 168 if $\mu_{1,6}$ is known exactly, then the means of all arms in Env 2 are immediately determined. FIDS is
 169 able to exploit this structure through its objective of gathering information about future optimal arms.
 170 In contrast, TS and IDS are designed to gather information only about the currently optimal arm.
 171 Since arm 6 is never optimal in Env 1, these methods fail to actively explore it and therefore cannot
 172 effectively exploit the information it provides about Env 2.

173 **Numerical results on the magic arm environment.** We evaluate FIDS, UCB [2], IDS and TS on
 174 the magic arm environment. The results are shown in Figure 1. In Figure 1a, compared with the
 175 other methods, FIDS spends more rounds identifying the optimal arm in Env 1, since it allocates
 176 part of its sampling probability to the informative but suboptimal arm 6. However, FIDS almost
 177 immediately identifies the optimal arm in Env 2, resulting in nearly zero regret during the second

Algorithm 1 Infer π^{FIDS} from Dataset \mathcal{D} with prediction window length k : Training and Deployment

```
// Training phase - supervised learning with offline dataset  $\mathcal{D}$ 
Construct training samples  $\{z_t^\xi\}_{t=1}^n$  from each data instance  $\xi \in \mathcal{D}$ ; initialize the model  $M_\theta$  with parameter  $\theta$ .
for  $t = 1$  to  $N_{\text{epochs}}$  do
  | Calculate loss (10) and do back-propagation to update  $\theta$ 
end
Output the trained model  $M_\theta$ 
// Deployment phase
for  $t = 1$  to  $n$  do
  | For each action  $a$  compute  $\hat{\mathbb{E}}_t(R_{t,a})$  and  $\hat{I}_t(A_{t+1:t+k}^*; R_{t,a})$  using Eqs. (11)-(12)
  | Solve
      
$$(i^*, j^*, \alpha^*) = \underset{i, j \in \mathcal{A}, \alpha \in [0, 1]}{\operatorname{argmax}} \left\{ \alpha \hat{\mathbb{E}}_t(R_{t,i}) + (1 - \alpha) \hat{\mathbb{E}}_t(R_{t,j}) \right. \\ \left. + \sigma \sqrt{2\sigma K \left( \alpha \hat{I}_t(A_{t+1:t+k}^*; R_{t,i}) + (1 - \alpha) \hat{I}_t(A_{t+1:t+k}^*; R_{t,j}) \right)} \right\}.$$

  | Output  $\pi_t^{\text{FIDS}}(i^*) = \alpha^*$ ,  $\pi_t^{\text{FIDS}}(j^*) = 1 - \alpha^*$ , and  $\pi_t^{\text{FIDS}}(a) = 0$  for all  $a \notin \{i^*, j^*\}$ .
end
```

178 half of the horizon. This phenomenon is also reflected in Figure 1c: by pulling arm 6 in *Env 1*, FIDS
179 rapidly reduces the entropy of the optimal arms in future rounds, rather than only the entropy of the
180 optimal arm in the current environment.

181 In Figure 1b, we additionally compare FIDS with two ablated policies: a “pure exploitation” policy
182 that optimizes only the reward term in Equation (5), and a “pure exploration” policy that opti-
183 mizes only the mutual information term. FIDS outperforms both policies by effectively balancing
184 exploration and exploitation.

185 4 Inferring FIDS policy from Data via Supervised Learning

186 Computing π^{FIDS} exactly is often impractical: the prior over environments may be unknown, and even
187 when it is known, evaluating Equation (5) requires posterior inference that is generally intractable. A
188 further difficulty is that the mutual information $I_t(A_{t+1:n}^*; R_{t,a})$ involves the full sequence of future
189 optimal arms, whose support grows exponentially with the horizon.

190 To address these challenges, we propose a supervised-learning approach (Algorithm 1). Given an
191 offline dataset of environment instances, we train models to estimate the two quantities needed by
192 FIDS: the expected reward $\mathbb{E}_t[R_{t,a}]$ and a *windowed* approximation $I_t(A_{t+1:t+k}^*; R_{t,a})$, where k is
193 a fixed lookahead horizon that controls the trade-off between fidelity and computational cost. At
194 deployment, these estimates are plugged into Equation (5) to compute π^{FIDS} . We describe the training
195 and deployment phases in detail below.

196 4.1 On the Training Phase

197 Algorithm 1 takes an offline dataset \mathcal{D} as input. Each data instance $\xi \in \mathcal{D}$ is generated by first sampling
198 an environment ν from the prior, and then collecting: (i) the full reward table of ν , $\{r_{t,a}^\nu\}_{t \in [n], a \in \mathcal{A}}$.
199 (ii) the best-arm sequence of ν , $(a_1^{*\nu}, \dots, a_n^{*\nu})$, and (iii) an action sequence (a_1, \dots, a_n) generated
200 by some behavior policy.

201 For each $\xi \in \mathcal{D}$ and $t \in [n]$, we construct one training sample: $z_t^\xi := \left(h_t^\xi, a_{t:t+k}^*, \{r_{t,a}^\xi\}_{a \in \mathcal{A}} \right)$,
202 where $h_t^\xi := (a_1, r_1^\nu, \dots, a_{t-1}, r_{t-1}^\nu)$ denotes the observed history up to round $t - 1$, and $r_s^\nu := r_{s, a_s}^\nu$
203 for $s \in [t]$. $a_{t:t+k}^{*\nu} := (a_t^{*\nu}, \dots, a_{t+k}^{*\nu})$ denotes the future best-arm sequence.

204 Our goal is to estimate $\mathbb{E}_t[R_{t,a}]$ and $I_t(A_{t+1:t+k}^*; R_{t,a})$. To this end, we approximate the following
205 conditional distributions: $\mathbb{P}_t(R_{t,a} \in \cdot)$ for all $a \in \mathcal{A}$, $\mathbb{P}_t(A_{t+1:t+k}^* \in \cdot)$, $\mathbb{P}_t(A_{t:t+k-1}^* \in \cdot)$. In practice,
206 we model the reward posterior by a Gaussian, and the best-arm posteriors by categorical distributions.
207 To learn these models, we use a sequential architecture, specifically GPT-2 [20]. The architecture has

208 three heads, and outputs:

$$M_{\theta}^{\text{rew}}(\cdot | h_t^{\xi}, a), \quad M_{\theta}^{\text{fut}+}(\cdot | h_t^{\xi}), \quad M_{\theta}^{\text{fut}}(\cdot | h_t^{\xi}),$$

209 which model the conditional distributions of $R_{t,a}$, $A_{t+1:t+k}^*$, and $A_{t:t+k-1}^*$, respectively. Then, to
 210 train the reward posterior, we simply use negative log-likelihood, while the best-arm posteriors are
 211 trained via cross-entropy loss. The overall training objective is

$$\mathcal{L}_{\theta} = \sum_{\xi \in \mathcal{D}} \sum_{t=1}^n \left[\sum_{a \in \mathcal{A}} \text{NLL}(M_{\theta}^{\text{rew}}(\cdot | h_t^{\xi}, a), r_{t,a}^{\xi}) + \text{CE}(M_{\theta}^{\text{fut}+}(\cdot | h_t^{\xi}), a_{t+1:t+k}^*) \right. \\ \left. + \text{CE}(M_{\theta}^{\text{fut}}(\cdot | h_t^{\xi}), a_{t:t+k-1}^*) \right]. \quad (10)$$

212 4.2 On the Deployment Phase

213 At deployment, we use the inferred posteriors to compute FIDS, which relies on $\mathbb{E}_t [R_{t,a}]$ and
 214 $I_t(A_{t+1:t+k}^*; R_{t,a})$. Given a history h_t , we first estimate the conditional reward mean by the mean
 215 of the learned reward posterior:

$$\widehat{\mathbb{E}}_t [R_{t,a}] = \mathbb{E}_{\tilde{R}_{t,a} \sim M_{\theta}^{\text{rew}}(\cdot | h_t, a)} [\tilde{R}_{t,a}]. \quad (11)$$

216 Since $M_{\theta}^{\text{rew}}(\cdot | h_t, a)$ is Gaussian, this is simply the predicted Gaussian mean. We next estimate the
 217 conditional mutual information via an entropy decomposition

$$I_t(A_{t+1:t+k}^*; R_{t,a}) = H_t(A_{t+1:t+k}^*) - H_t(A_{t+1:t+k}^* | R_{t,a}).$$

218 The first term is estimated directly from the future head: $\widehat{H}_t(A_{t+1:t+k}^*) = H(M_{\theta}^{\text{fut}+}(\cdot | h_t))$. For
 219 the conditional entropy term, we draw L Monte Carlo samples: $\tilde{r}_{t,a}^{(1)}, \dots, \tilde{r}_{t,a}^{(L)} \sim M_{\theta}^{\text{rew}}(\cdot | h_t, a)$.

220 For each sampled reward $\tilde{r}_{t,a}^{(\ell)}$, we form the hypothetical updated history $\tilde{h}_{t+1}^{(\ell,a)} = (h_t, a, \tilde{r}_{t,a}^{(\ell)})$. Since
 221 $M_{\theta}^{\text{fut}}(\cdot | \tilde{h}_{t+1}^{(\ell,a)})$ models the posterior distribution of $A_{t+1:t+k}^*$ after observing action a and reward
 222 $\tilde{r}_{t,a}^{(\ell)}$, we estimate

$$\widehat{H}_t(A_{t+1:t+k}^* | R_{t,a}) = \frac{1}{L} \sum_{\ell=1}^L H(M_{\theta}^{\text{fut}}(\cdot | \tilde{h}_{t+1}^{(\ell,a)})).$$

223 Therefore, the plug-in Monte Carlo estimator of the conditional mutual information is

$$\widehat{I}_t(A_{t+1:t+k}^*; R_{t,a}) = H(M_{\theta}^{\text{fut}+}(\cdot | h_t)) - \frac{1}{L} \sum_{\ell=1}^L H(M_{\theta}^{\text{fut}}(\cdot | \tilde{h}_{t+1}^{(\ell,a)})). \quad (12)$$

224 We choose $L = 16$ for all the experiments in Section 5.

225 **Solving Equation (5).** With the learned $\mathbb{E}_t [R_{t,a}]$ and $I_t([A_{t+1}^*, \dots, A_{t+k}^*]; R_{t,a})$, FIDS policy
 226 can be obtained by just plugging in the learned two quantities into Eq. (5) and solve the corresponding
 227 optimization problem. With Lemma 3.1, for the plugged-in $\widehat{\mathbb{E}}_t(R_{t,a})$ and $\widehat{I}_t(A_{t+1:n}^*; R_{t,a})$, one can
 228 solve Eq. (5) by first enumerating all pairs of arms and using line search to find the optimal weights
 229 assigned on the arms. Details are presented in Alg. 1.

230 5 Experiments

231 We evaluate the FIDS policy learned via Algorithm 1 on two environment classes for which exact
 232 posterior inference is intractable: the *One-Step Predictive* environment and the *Pair-Revealing*
 233 environment. In both environments, information about future optimal arms is encoded in the rewards
 234 of a currently suboptimal arm: the first reveals the next-round optimum exactly, while the second only
 235 narrows the future best arm to a pair. Further details are provided in the next subsections. Throughout,
 236 we set $\sigma = 0.7$ in Equation (5), which is a valid sub-Gaussian parameter by Lemma E.2.

237 **Comparison.** We compare the following methods: (i) FIDS with prediction window length $k = 5$;
 238 (ii) FIDS with $k = 3, 1$ (used only in the second Pair-Revealing Environment); (iii) DPT; and (iv)
 239 IDS. We train IDS using the same procedure as Algorithm 1. The only difference is that for IDS we
 240 estimate $I_t(A_t^*; R_{t,a})$ instead of $I_t(A_{t+1:t+k}^*; R_{t,a})$.

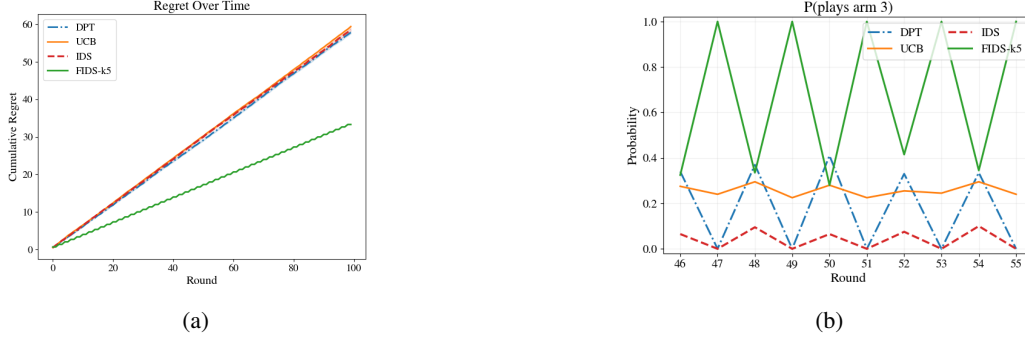


Figure 2: (a) Cumulative Regret of DPT, IDS, UCB and FIDS-k5 in One Step Predictive Environment. (b) The probability of choosing arm 3 in each round.

241 5.1 One-Step Predictive Environment

242 **Environment Description** (Next round best arm encoding). We consider a nonstationary 3-arm
 243 bandit with horizon $n = 100$, divided into 50 independent windows of length 2. For each window $j \in$
 244 $\{1, \dots, 50\}$, corresponding to rounds $2j-1$ and $2j$, we independently sample $m_j \sim \text{Unif}(\{1, 2, 3\})$.
 245 At round $2j-1$, we sample $Z_j \sim \text{Unif}(\{1, 2\})$. Arms 1 and 2 have rewards $\mu_{2j-1,a} = \mathbf{1}[a = Z_j]$,
 246 while arm 3 has reward $\mu_{2j-1,3} = 1/(2m_j)$. At round $2j$, the best arm is determined by m_j :
 247 $\mu_{2j,a} = \mathbf{1}[a = m_j]$ for all $a \in \{1, 2, 3\}$. All rewards are deterministic.

248 The key feature is that arm 3’s reward in odd rounds encodes the identity of the best arm in the
 249 following even round, creating a one-step-ahead information channel that a forward-looking policy
 250 can exploit.

251 **Results.** Figures 2a and 2b summarize the main results (additional plots in Figures 5a–5b in the
 252 appendix). Although the windows are independent, each contains a one-step information channel:
 253 pulling arm 3 in the first round reveals the optimal arm in the second round. The learned FIDS policy
 254 recovers exactly this behavior: it selects arm 3 in the first round of each window and the true optimal
 255 arm in the second (Figures 5a, 2b, and 5b). In contrast, TS and IDS never learn to exploit arm 3,
 256 since their exploration targets only the current optimal arm.

257 5.2 Pair-Revealing Environment

258 **Environment Description.** Unlike the previous environment, where the informative arm directly
 259 identifies the future optimum, here it only reveals the top-2 arm pair in the future environment. We
 260 consider a nonstationary 6-arm bandit with horizon $n = 100$ and a single change point at $t = 50$.
 261 The environment is Env 1 for $t \leq 50$ and Env 2 for $t > 50$, stationary within each segment.

- 262 • In Env 1, let $Z_1 \sim \text{Unif}(\{1, \dots, 5\})$. For $a \in \{1, \dots, 5\}$, $\mu_a^1 = 1$ if $a = Z_1$ and $\mu_a^1 = U_a$
 263 otherwise, where $U_a \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 1])$, with $R_{t,a} \sim \mathcal{N}(\mu_a^1, 0.5^2)$. Arm 6 has mean $\mu_6^1 =$
 264 $m/30$, where $m \sim \text{Unif}(\{0, \dots, 14\})$ indexes one of the $\binom{6}{2} = 15$ unordered arm pairs;

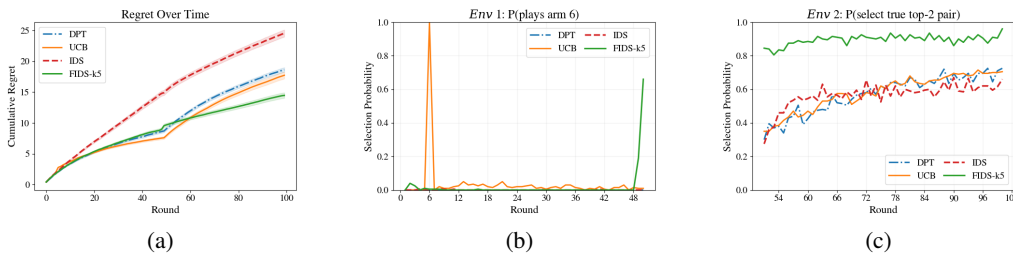


Figure 3: Evaluation Results for the Pair-Revealing Environment with $\sigma_6^2 = 0.0$. (a) Cumulative regret. (b) The probability of play arm 6 at each round. (c) The probability of selecting an arm that belongs to the true top-2 set of the underlying environment.

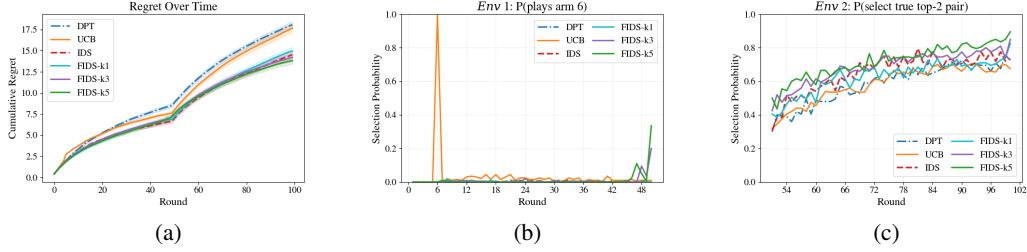


Figure 4: Evaluation Results for the Pair-Revealing Environment with $\sigma_6^2 = 0.05^2$. (a) Cumulative regret. (b) The probability of play arm 6 at each round. (c) The probability of selecting an arm that belongs to the true top-2 set of the underlying environment.

265 we denote this pair by $\mathcal{S}(m)$. We consider two noise settings for arm 6: $\sigma_6^2 = 0$ (exact
 266 observation) and $\sigma_6^2 = 0.05^2$ (noisy observation).

267 • In Env 2, we sample $V_1, \dots, V_6 \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 1])$ and assign the two largest values to the
 268 arms in $\mathcal{S}(m)$, with the remaining values randomly permuted among the other arms, yielding
 269 means $\{\mu_a^2\}$. Rewards satisfy $R_{t,a} \sim \mathcal{N}(\mu_a^2, 0.5^2)$.

270 The information structure is that arm 6 in Env 1, while never optimal, encodes which pair of arms
 271 will dominate in Env 2. The two noise settings let us examine how the precision of this signal affects
 272 the benefit of forward-looking exploration.

273 **Results.** We evaluate under two noise settings for arm 6.

274 • ($\sigma_6^2 = 0$). We set the prediction window to $k = 5$. Figure 3a shows that FIDS matches
 275 the baselines in Env 1, while intentionally pulling arm 6 (Figure 3b) to gather information
 276 about Env 2. Upon entering Env 2, FIDS identifies the candidate top-2 arms and focuses
 277 exploration on them rather than the full arm set.

278 • ($\sigma_6^2 = 0.05^2$). When arm 6’s reward is noisy, recovering μ_6^1 requires multiple pulls, making
 279 a longer prediction window more valuable. We compare FIDS with $k = 5$ (FIDS-K5) and
 280 a shorter window. As shown in Figures 4b and 4c, pulling arm 6 earlier in Env 1 allows
 281 FIDS-K5 to infer the top-2 pair more accurately, leading to faster adaptation after the change
 282 point.

283 Across both environments, the learned FIDS policy consistently exploits informative but suboptimal
 284 arms that TS and IDS ignore, confirming that Algorithm 1 successfully recovers forward-looking
 285 exploration behavior from data.

286 6 Conclusion

287 We introduced FIDS, an algorithm for Bayesian nonstationary bandits that balances instantaneous
 288 reward against information gained about future optimal arms. FIDS matches the regret guarantee
 289 of Thompson Sampling, yet analytical examples show it can significantly outperform both TS and
 290 IDS by exploiting predictive structure across environments. To handle settings where exact posterior
 291 inference is intractable, we proposed a supervised-learning framework that approximates the FIDS
 292 policy from offline data. Experiments confirm that the learned policy consistently outperforms
 293 supervised-learning baselines with alternative exploration objectives.

294 **Limitations and Future Work.** One limitation of Algorithm 1 is that training requires access to the
 295 full reward table $\{R_{t,a}\}_{t \in [n], a \in \mathcal{A}}$ in the offline dataset, which may be unrealistic in many practical
 296 settings. In addition, estimating $I_t(A_{t+1:n}^*; R_{t,a})$ during deployment requires sampling multiple $\tilde{R}_{t,a}$
 297 values and performing Monte Carlo estimation, which can become computationally expensive when a
 298 large number of samples is needed. Developing more scalable implementations of the FIDS principle
 299 therefore remains an important direction for future work.

References

- [1] Yasin Abbasi-Yadkori, András György, and Nevena Lazić. A new look at dynamic regret for non-stationary stochastic bandits. *Journal of Machine Learning Research*, 24(288):1–37, 2023.
- [2] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- [3] Peter Auer, Pratik Gajane, and Ronald Ortner. Adaptively tracking the best bandit arm with an unknown number of distribution changes. In *Conference on learning theory*, pages 138–158. PMLR, 2019.
- [4] Giulia Clerici, Pierre Laforgue, and Nicolo Cesa-Bianchi. Linear bandits with memory: from rotting to rising. *arXiv preprint arXiv:2302.08345*, 2023.
- [5] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- [6] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. R^2 : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- [7] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- [8] Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for non-stationary bandit problems. *arXiv preprint arXiv:0805.3415*, 2008.
- [9] Neha Gupta, Ole-Christoffer Granmo, and Ashok Agrawala. Thompson sampling for dynamic multi-armed bandits. In *2011 10th International Conference on Machine Learning and Applications and Workshops*, volume 1, pages 484–489. IEEE, 2011.
- [10] T. L. Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, March 1985. ISSN 0196-8858. doi: 10.1016/0196-8858(85)90002-8. URL <https://www.sciencedirect.com/science/article/pii/0196885885900028>.
- [11] Tze Leung Lai. Adaptive Treatment Allocation and the Multi-Armed Bandit Problem. *The Annals of Statistics*, 15(3):1091 – 1114, 1987. doi: 10.1214/aos/1176350495. URL <https://doi.org/10.1214/aos/1176350495>.
- [12] Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforcement learning with algorithm distillation. *arXiv preprint arXiv:2210.14215*, 2022.
- [13] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [14] Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma Brunskill. Supervised pretraining can learn in-context reinforcement learning. *Advances in Neural Information Processing Systems*, 36:43057–43083, 2023.
- [15] Yueyang Liu, Benjamin Van Roy, and Kuang Xu. Nonstationary bandit learning via predictive sampling. In *International Conference on Artificial Intelligence and Statistics*, pages 6215–6244. PMLR, 2023.
- [16] Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob Foerster, Satinder Singh, and Feryal Behbahani. Structured state space models for in-context reinforcement learning. *Advances in Neural Information Processing Systems*, 36:47016–47031, 2023.
- [17] Joseph Mellor and Jonathan Shapiro. Thompson sampling in switching environments with bayesian online change detection. In *Artificial intelligence and statistics*, pages 442–450. PMLR, 2013.
- [18] Seungki Min and Daniel Russo. An information-theoretic analysis of nonstationary bandit learning. In *International Conference on Machine Learning*, pages 24831–24849. PMLR, 2023.

- 347 [19] Sandeep Pandey, Deepak Agarwal, Deepayan Chakrabarti, and Vanja Josifovski. Bandits
348 for taxonomies: A model-based approach. In *Proceedings of the 2007 SIAM international*
349 *conference on data mining*, pages 216–227. SIAM, 2007.
- 350 [20] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
351 Language models are unsupervised multitask learners.
- 352 [21] Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling.
353 *Advances in neural information processing systems*, 27, 2014.
- 354 [22] Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling.
355 *Journal of Machine Learning Research*, 17(68):1–30, 2016.
- 356 [23] Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial
357 on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- 358 [24] Joe Suk and Samory Kpotufe. Tracking most significant arm switches in bandits. In *Conference*
359 *on Learning Theory*, pages 2160–2182. PMLR, 2022.
- 360 [25] William R Thompson. On the likelihood that one unknown probability exceeds another in view
361 of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 362 [26] Francesco Trovo, Stefano Paladino, Marcello Restelli, and Nicola Gatti. Sliding-window
363 thompson sampling for non-stationary settings. *Journal of Artificial Intelligence Research*, 68:
364 311–364, 2020.

365 **A More Related Work**

366 **Nonstationary Bandit Learning.** Predictive Sampling [15] considers the same Bayesian nonsta-
 367 tionary setting and shares a similar intuition with our work, namely that exploration strategies should
 368 account for future information in nonstationary environments. However, our work adopts a stronger
 369 regret notion in the theoretical analysis. As a result, our regret guarantees directly imply guarantees
 370 under their setting, while the converse does not hold.

371 In addition, several variants of Thompson Sampling have been proposed for nonstationary environ-
 372 ments (e.g., [26, 9, 17]). However, these methods typically rely on specific structural assumptions
 373 about the environment and do not provide meaningful guarantees in general Bayesian nonstationary
 374 settings.

375 Nonstationary bandits have also been extensively studied in the frequentist setting, e.g., [8, 3]. More
 376 recent results in [24, 1] establish regret bounds of order $O(\sqrt{Ln})$, where L denotes the number of
 377 changes of the optimal arm. Interestingly, these bounds share a similar flavor with our result, as both
 378 depend on the temporal variation of the optimal arm sequence.

379 **Learning Decision-Making Algorithms from Data.** The proposed supervised-learning-based
 380 framework for approximating the FIDS policy is primarily inspired by [14], which introduces the
 381 Decision-Pretrained Transformer (DPT), a model that theoretically learns the Thompson Sampling
 382 policy. More broadly, the idea of learning sequential decision-making algorithms from data is also
 383 related in spirit to the meta-RL literature [12, 6, 16, 7].

384 **B Information Theoretic Preliminaries**

385 In this section, X is limited to discrete random variables and Y, Z can be general random variables,
 386 under measure \mathbb{P} .

387 **Lemma B.1.** (*Mutual Information and Entropy, Theorem 2.4.1 in [5]*)

$$H(X) = - \sum_x \mathbb{P}(X = x) \log(\mathbb{P}(X = x))$$

388

$$I(X; Y) = H(X) - H(X|Y)$$

389 **Lemma B.2.** (*Chain Rule for Entropy, Theorem 2.5.1 in [5]*)

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1})$$

390 **Lemma B.3.** (*Chain Rule for Mutual Information, Theorem 2.5.2 in [5]*)

$$I([X_1, X_2, \dots, X_n]; Y) = \sum_{i=1}^n I(X_i; Y | X_1, \dots, X_{i-1})$$

391 A direct application of the above lemma gives the following Corollary,

Corollary B.4.

$$I(X; Y) - I(Z; Y) = I(X; Y|Z) - I(Z; Y|X)$$

392 *Proof.* The proof is done by simply noticing that $I(X_1; Y) + I(X_2; Y|X_1) = I(X_2; Y) +$
 393 $I(X_1; Y|Z_2) = I([X_1, X_2]; Y)$ by Lemma B.3. \square

Corollary B.5.

$$H(X|Y, Z) \leq H(X|Y)$$

Proof.

$$H(X|Y, Z) \leq H(X|Y) \Leftrightarrow H(X) - I(X; [Y, Z]) \leq H(X) - I(X; Y) \Leftrightarrow I(X; [Y, Z]) \geq I(X; Y)$$

394 which is true by Lemma B.3. \square

395 **Lemma B.6.** (*Non-negativity of Mutual information and Entropy*)

$$H(X) \geq 0; H(X|Y) \geq 0; I(X; Y) \geq 0; I(X; Y|Z) \geq 0$$

396 **Lemma B.7.** *Let $a, b \geq 0$, $\sqrt{a} - \sqrt{b} \leq \sqrt{|a - b|}$*

397 *Proof.* Clearly we only need to consider the case where $a \geq b$. In such a case,

$$\sqrt{a} - \sqrt{b} \leq \sqrt{|a - b|} \Leftrightarrow a + b - 2\sqrt{ab} \leq a - b \Leftrightarrow a \geq b$$

398

□

399 C Supporting Lemmas in the Proof of Theorem 3.2

Lemma C.1.

$$\begin{aligned} \mathbb{E}_t [R_{t,A_t^*} - R_{t,A_t}] &= \sum_{a \in [K]} \mathbb{P}_t(A_t = a) \mathbb{E}_t [R_{t,A_t^*} - R_{t,a} \mid A_t = a] \\ &\stackrel{(a)}{=} \sum_{a \in [K]} \pi_t(a) \mathbb{E}_t [R_{t,A_t^*} - R_{t,a}] \end{aligned} \quad (13)$$

400 (a) uses the fact that conditioned on \mathcal{F}_{t-1} , A_t is jointly independent of R_{t,A_t^*} and $R_{t,a}$.

$$\begin{aligned} I_t(A_t^*; (A_t, R_{t,A_t})) &\stackrel{(a)}{=} I_t(A_t^*; A_t) + I_t(A_t^*; R_{t,A_t} \mid A_t) \\ &\stackrel{(b)}{=} I_t(A_t^*; R_{t,A_t} \mid A_t) \\ &= \sum_{a \in [K]} \mathbb{P}_t(A_t = a) I_t(A_t^*; R_{t,a} \mid A_t = a) \\ &\stackrel{(c)}{=} \sum_{a \in [K]} \mathbb{P}_t(A_t = a) I_t(A_t^*; R_{t,a}) \\ &= \sum_{a \in [K]} \pi_t(a) I_t(A_t^*; R_{t,a}) \end{aligned} \quad (14)$$

401 (a) uses Lemma B.3; (b) uses the fact conditioned on \mathcal{F}_{t-1} , A_t is independent of A_t^* , and the mutual
402 information between two independent variables is 0; (c) uses the fact that conditioned on \mathcal{F}_{t-1} , A_t is
403 jointly independent of A_t^* and $R_{t,a}$.

404 Similarly,

$$I_t([A_{t+1}^*, \dots, A_n^*]; (A_t, R_{t,A_t})) = \sum_{a \in [K]} \pi_t(a) I_t([A_{t+1}^*, \dots, A_n^*]; R_{t,a}) \quad (15)$$

Lemma C.2.

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^n \sigma \sqrt{2K \cdot I_t(A_{t+1:n}^*; (A_t, R_{t,A_t}))} \right] &\stackrel{(a)}{\leq} \sum_{t=1}^n \sigma \sqrt{2K \cdot \mathbb{E} [I_t(A_{t+1:n}^*; (A_t, R_{t,A_t}))]} \\ &= \sum_{t=1}^n \sigma \sqrt{2K \cdot I(A_{t+1:n}^*; (A_t, R_{t,A_t}) | \mathcal{F}_{t-1})} \\ &\stackrel{(b)}{\leq} \sigma \sqrt{2n \cdot K \cdot \sum_{t=1}^n I(A_{t+1:n}^*; (A_t, R_{t,A_t}) | \mathcal{F}_{t-1})} \\ &\stackrel{(c)}{\leq} \sigma \sqrt{2n \cdot K \cdot \sum_{t=1}^n I(A_{t+1:n}^*; (A_t, R_{t,A_t}) | \mathcal{F}_{t-1})} \\ &\stackrel{(d)}{=} \sigma \sqrt{2n \cdot K \cdot I(A_{1:n}^*; \mathcal{F}_n)} \\ &\stackrel{(e)}{\leq} \sigma \sqrt{2n \cdot K \cdot H(A_{1:n}^*)} \end{aligned} \quad (16)$$

405 (a) uses Jensen's Inequality; (b) uses Cauchy-Schwarz inequality; (c) uses Lemma B.3 and Lemma
406 B.6; (d) uses Lemma B.3; (e) uses Lemma B.1 and B.6.

Lemma C.3.

$$\sigma\sqrt{2K \cdot g_t(\pi_t^{TS})} - \sigma\sqrt{2K \cdot \tilde{g}_t(\pi_t^{TS})} \leq \sigma\sqrt{2K} \cdot \sqrt{I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t}) | A_{t+1:n}^*)}$$

407 *Proof.* When $I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t})) \leq I_t(A_{t+1:n}^*; (\tilde{A}_t, R_{t, \tilde{A}_t}))$. We have $\text{term1}(t) \leq 0$ directly. When
 408 $I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t})) > I_t(A_{t+1:n}^*; (\tilde{A}_t, R_{t, \tilde{A}_t}))$. Then

$$\begin{aligned} \text{term1}(t) &\stackrel{(a)}{\leq} \sigma\sqrt{2K} \cdot \sqrt{I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t})) - I_t(A_{t+1:n}^*; (\tilde{A}_t, R_{t, \tilde{A}_t}))} \\ &\stackrel{(b)}{=} \sigma\sqrt{2K} \cdot \sqrt{I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t}) | A_{t+1:n}^*) - I_t(A_{t+1:n}^*; (\tilde{A}_t, R_{t, \tilde{A}_t}) | A_t^*)} \\ &\leq \sigma\sqrt{2K} \cdot \sqrt{I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t}) | A_{t+1:n}^*)} \end{aligned}$$

409 (a) holds because of Lemma B.7; (b) holds because of Lemma B.4.

410 We conclude the proof by combining the two cases. \square

Lemma C.4.

$$\begin{aligned} \sigma\sqrt{2K} \cdot \mathbb{E} \left[\sum_{t=1}^n \sqrt{g_t(\pi_t^{TS})} - \sqrt{\tilde{g}_t(\pi_t^{TS})} \right] &\stackrel{(a)}{\leq} \sigma\sqrt{2K} \cdot \mathbb{E} \left[\sum_{t=1}^n \sqrt{I_t(A_t^*; (\tilde{A}_t, R_{t, \tilde{A}_t}) | A_{t+1:n}^*)} \right] \\ &\leq \sigma\sqrt{2K} \cdot \mathbb{E} \left[\sum_{t=1}^n \sqrt{H_t(A_t^* | A_{t+1:n}^*)} \right] \\ &\leq \sigma\sqrt{2K} \cdot \sum_{t=1}^n \sqrt{\mathbb{E} [H_t(A_t^* | [A_{t+1}^*, \dots, A_n^*])]} \\ &= \sigma\sqrt{2K} \cdot \sum_{t=1}^n \sqrt{H(A_t^* | ([A_{t+1}^*, \dots, A_n^*], \mathcal{F}_{t-1}))} \\ &\stackrel{(b)}{\leq} \sigma\sqrt{2K} \cdot \sum_{t=1}^n \sqrt{H(A_t^* | [A_{t+1}^*, \dots, A_n^*])} \\ &\leq \sigma\sqrt{2K} \cdot \sqrt{n \cdot \sum_{t=1}^n H(A_t^* | [A_{t+1}^*, \dots, A_n^*])} \\ &\stackrel{(c)}{=} \sigma\sqrt{2K \cdot n \cdot H([A_1^*, \dots, A_n^*])} \end{aligned} \tag{17}$$

411 (a) uses Lemma C.3; (b) uses Corollary B.5; (c) uses Lemma B.2.

D Numerical Results

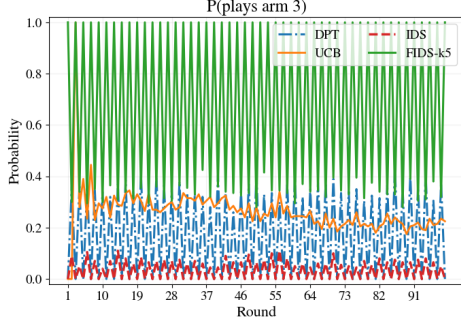
413 Here we report additional numerical results for the one-step predictive environment in Sec. 5.1.

E Other Results

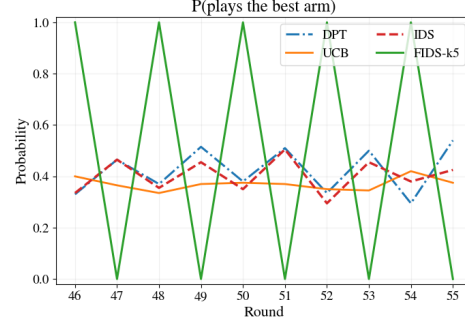
E.1 Proof of Lemma 3.1

416 *Proof.* Denote $E_{t,a} := \mathbb{E}_t [R_{t,a}]$, $g_{t,a} := I_t([A_{t+1}^*, \dots, A_n^*]; R_{t,a})$, and $E_t := [E_{t,1}, \dots, E_{t,K}]^T$,
 417 $g_t := [g_{t,1}, \dots, g_{t,K}]^T$. Our objective function can then be written as $\rho(\pi_t) := \pi_t \cdot E_t + \sqrt{\Gamma} \cdot \pi_t \cdot g_t$,
 418 where $\Gamma = 2\sigma^2 K$. Consider some fixed maximizer π^* of Equation (5). The partial derivative of ρ
 419 with respect to $\pi_{t,a}$ at π^* is $\frac{\partial \rho}{\partial \pi_{t,a}}(\pi^*) := E_{t,a} + \frac{\Gamma g_{t,a}}{2\sqrt{\Gamma} g_t \cdot \pi^*}$.

420 Denote $d^* := \max_{a \in \mathcal{A}} \frac{\partial \rho}{\partial \pi_{t,a}}(\pi^*)$. The important observation here is that for any a such that $\pi_a^* > 0$,
 421 $\frac{\partial \rho}{\partial \pi_{t,a}}(\pi^*) = d^*$, because otherwise one must be able to construct a new valid policy by transferring



(a) The probability of choosing arm 3



(b) The probability of choosing the true best arm

422 some probability from a to other arms with larger partial derivative and that will only increase the
 423 objective value. Therefore, for all $a \in B_t := \{a \in \mathcal{A} : \pi_a^* > 0\}$,

$$E_{t,a} + \frac{\Gamma g_{t,a}}{2\sqrt{\Gamma g_t \cdot \pi^*}} = d^* \quad (18)$$

424 Reorder the set B_t such that $E_{t,a_1} \geq \dots \geq E_{t,a_{|B_t|}}$. Note that there always exists a $\beta \in [0, 1]$
 425 such that $\sum_{a \in B_t} \pi_a^* E_{t,a} = \beta E_{t,a_1} + (1 - \beta) E_{t,a_{|B_t|}}$, which, by Equation (18), also implies that
 426 $\sum_{a \in B_t} \pi_a^* g_{t,a} = \beta g_{t,a_1} + (1 - \beta) g_{t,a_{|B_t|}}$. Therefore, we can construct a new policy that only assigns
 427 possibilities on the two arms a_1 and $a_{|B_t|}$, and this policy has the same objective value as π^* and thus
 428 optimal. \square

429 E.2 Analytic Solutions

430 Depending on which round we are currently in, the mutual information term
 431 $I_t([A_{t+1}^*, \dots, A_n^*]; R_{t,a})$ is equal to $I_t([Z_1, Z_2]; R_{t,a})$, if we are in the first half, or $I_t(Z_2; R_{t,a})$
 432 if we are in the second half. By independency, $I_t([Z_1, Z_2]; R_{t,a}) = I_t(Z_1; R_{t,a}) + I_t(Z_2; R_{t,a})$.

433 In the first half, by independency, when $a \neq 6$ $I_t([Z_1, Z_2]; R_{t,a}) = I_t(Z_1; R_{t,a})$, while $a = 6$
 434 $I_t([Z_1, Z_2]; R_{t,a}) = I_t(Z_2; R_{t,6})$. The main problem is to calculate the posterior distribution.
 435 Essentially, we need to maintain $q_t(Z_1)$ and $w_t(Z_2)$.

436 if $a_t = 6$, $q_{t+1}(Z_1) = q_t(Z_1)$, and the posterior update rule for $w_{t+1}(Z_2)$ is

$$w_{t+1}(Z_2) \propto w_t(Z_2) \mathcal{N}\left(r_t; \frac{Z_2}{7}, 0.1^2\right)$$

437 While for $a_t \neq 6$, $w_{t+1}(Z_2) = w_t(Z_2)$, and the posterior update rule for $q_{t+1}(Z_1)$ is

$$q_{t+1}(Z_1) \propto q_t(Z_1) \mathcal{N}(r_t; \mu_z, 0.1^2)$$

438 This equation can be calculated analytically because z is finite. The posterior reward distribution can
 439 also be calculated by

$$\mathbb{P}_t(R_{t,a} = r) = \sum_z \mathbb{P}_t(R_{t,a} = r | Z_1 = z) \mathbb{P}_t(Z_1 = z)$$

440 In the second half, the posterior update rule for $w_{t+1}(Z_2)$ is

$$w_{t+1}(Z_2) \propto w_t(Z_2) \mathcal{N}(r_t; \mu_z, 0.5^2)$$

441 E.3 Other Lemmas

442 **Proposition E.1.** *In Example 3.1, FIDS is the optimal policy and collects $0.35 \cdot n$ cumulative reward,*
 443 *while TS and IDS collect $0.325 \cdot n$ and $0.3294 \cdot n$ cumulative reward respectively.*

444 *Proof.* Notice that $\mathbb{E}_t [R_{t,1}] = 0.35, \mathbb{E}_t [R_{t,2}] = 0.3$ for every t . Thus, the optimal policy is to choose
 445 arm 1 w.p. 1 at every t . FIDS will actually do that by noticing that $I_t ([A_{t+1}^*, \dots, A_n^*]; R_{t,a}) = 0$
 446 for every t and a . And this will give cumulative reward $0.35 \cdot n$.

447 TS will play each arm with equal probability at each round, and the cumulative reward is $0.5 \cdot (0.35 +$
 448 $0.3) \cdot n = 0.325 \cdot n$.

449 We calculate the mutual information terms $I_t (A_t^*; (R_{t,1}))$ and $I_t (A_t^*; (R_{t,2}))$ explicitly as shown
 450 below

$$451 I_t (A_t^*; (R_{t,1})) = \mathbb{P}(A_t^* = 1) \text{KL}(\mathbb{P}(R_{t,1} \in \cdot | A_t^* = 1) || \mathbb{P}(R_{t,1} \in \cdot)) + \mathbb{P}(A_t^* = 2) \text{KL}(\mathbb{P}(R_{t,1} \in$$

 452 $\cdot | A_t^* = 2) || \mathbb{P}(R_{t,1} \in \cdot)) \approx 0.0055$

$$453 I_t (A_t^*; (R_{t,2})) = \mathbb{P}(A_t^* = 1) \text{KL}(\mathbb{P}(R_{t,2} \in \cdot | A_t^* = 1) || \mathbb{P}(R_{t,2} \in \cdot)) + \mathbb{P}(A_t^* = 2) \text{KL}(\mathbb{P}(R_{t,2} \in$$

 454 $\cdot | A_t^* = 2) || \mathbb{P}(R_{t,2} \in \cdot)) \approx 0.0242$

455 And then we calculate the IDS policy by solving Equation (3) - $\pi_t^{\text{IDS}}(1) \approx 0.588$ and $\pi_t^{\text{IDS}}(2) \approx 0.412$.
 456 The cumulative reward will be $(0.588 \cdot 0.35 + 0.412 \cdot 0.3) \cdot n = 0.3294 \cdot n$. \square

457 **Lemma E.2.** Fix constants $\alpha, \beta, \sigma \in \mathbb{R}$, for the environment with $\mu_{t,a} \in [\alpha, \beta]$ and $R_{t,a} \sim$
 458 $\mathcal{N}(\mu_{t,a}, \sigma_{t,a}^2)$ where $\sigma_{t,a}$ is a deterministic non-negative constant no larger than σ . Then, con-
 459 ditioned on \mathcal{F}_t ,

$$R_{t,a} - \mathbb{E}[R_{t,a} | \mathcal{F}_t]$$

460 is $\sqrt{((b-a)/2)^2 + \sigma^2}$ sub-Gaussian.

461 *Proof.* Note that $R_{t,a} - \mathbb{E}[R_{t,a} | \mathcal{F}_t] = \mu_{t,a} + \epsilon - \mathbb{E}[\mu_{t,a} | \mathcal{F}_t]$. Therefore,

$$\begin{aligned} \mathbb{E}[\exp(\lambda(R_{t,a} - \mathbb{E}[R_{t,a} | \mathcal{F}_t])) | \mathcal{F}_t] &= \mathbb{E}[\exp(\lambda(\mu_{t,a} + \epsilon - \mathbb{E}[\mu_{t,a} | \mathcal{F}_t])) | \mathcal{F}_t] \\ &= \mathbb{E}[\exp(\lambda(\mu_{t,a} - \mathbb{E}[\mu_{t,a} | \mathcal{F}_t])) \cdot \exp(\lambda\epsilon) | \mathcal{F}_t] \\ &= \mathbb{E}[\exp(\lambda(\mu_{t,a} - \mathbb{E}[\mu_{t,a} | \mathcal{F}_t])) | \mathcal{F}_t] \cdot \mathbb{E}[\exp(\lambda\epsilon)] \\ &\leq \exp(\lambda^2(\beta - \alpha)^2/8) \cdot \exp(\lambda^2\sigma^2/2) \\ &= \exp\left(\frac{\lambda^2((\beta - \alpha)^2/4 + \sigma^2)}{2}\right) \end{aligned}$$

462 \square