
On Adaptivity and Confounding in Contextual Bandit Experiments

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Multi-armed bandit algorithms minimize experimentation costs required
2 to converge on optimal behavior. They do so by rapidly adapting
3 experimentation effort away from poorly performing actions as feedback is
4 observed. But this desirable feature makes them sensitive to confounding.
5 We highlight, for instance, that popular bandit algorithms cannot address
6 the problem of identifying the best action when day-of-week effects may
7 confound inferences. In response, this paper formulates a general model of
8 contextual bandit experiments with nonstationary contexts, which act as
9 the confounders for inferences and can be also viewed as the distribution
10 shifts in the earlier periods of the experiments. In addition, this general
11 model allows the target distribution or population distribution that is
12 used to determine the best action to be different from the empirical
13 distribution over the contexts observed during the experiments. The
14 paper proposes *deconfounded Thompson sampling*, which makes simple, but
15 critical, modifications to the way Thompson sampling is usually applied.
16 Theoretical guarantees suggest the algorithm strikes a delicate balance
17 between adaptivity and robustness to confounding and distribution shifts.
18 It attains asymptotic lower bounds on the number of samples required to
19 confidently identify the best action — suggesting optimal adaptivity — but
20 also satisfies strong performance guarantees in the presence of day-of-week
21 effects and delayed observations — suggesting unusual robustness.

22 1 Paper Summary

23 Multi-armed bandit algorithms are designed to adapt their experimentation rapidly as
24 evidence is gathered. By quickly shifting measurements away from less promising actions
25 or ‘arms’, they focus measurement effort where it is most useful. This desirable feature
26 can make these same algorithms brittle in the face of delayed observations or confounding
27 factors. We highlight this challenge through an example of a week-long experiment where
28 observations are influenced by specific day-of-week effects.

29 **Example 1** (Day-of-week effects). *In any period $t \in [T] := \{1, \dots, T\}$, the decision-maker
30 observes context $X_t \in [7]$, selects arm $I_t \in [k]$, and observes a noisy reward R_t that reflects the
31 performance of the chosen arm in the current context. For concreteness, one might imagine that a
32 period corresponds to a customer visiting an online retailer, the context indicates the day of the week,
33 an arm indicates the price set for a particular product, and the reward is the resulting revenue. The
34 context at time t is $X_t = \lceil t/m \rceil$, meaning the first m periods are Sunday, the next m are Monday*

35 and so on. Assume that $R_t = \theta_{I_t, X_t} + W_t$ where $W_t \mid \theta, I_t \sim N(0, 1)$ is independent Gaussian noise
 36 and $\theta \in \mathbb{R}^{7k}$ is an unknown parameter vector that encodes the day and arm specific mean rewards.
 37 By intelligently adapting measurement effort, the decision-maker hopes to identify the arm

$$I^*(\theta) = \arg \max_{i \in [k]} \frac{\theta_{i,1} + \dots + \theta_{i,7}}{7} \quad (1)$$

38 that maximizes expected revenue if employed throughout the entire week. The goal is to learn a single
 39 price and not a sequence of seven prices to charge on separate days of the week. Such predictable price
 40 variations might, for instance, lead to unintended strategic customer behavior if implemented across
 41 many future weeks.

42 The decision-maker begins with prior belief under which $\theta \sim N(\mu, \Sigma)$. This might, for instance, arise
 43 from a latent variable model where $\theta_{i,x} = \theta_{i,x}^{\text{idio}} + \theta_i^{\text{arm}} + \theta_x^{\text{day}}$ is determined by an effect $\theta_{i,x}^{\text{idio}}$ that is
 44 idiosyncratic to a specific arm and day, an effect θ_i^{arm} associated with the chosen arm, and a shared
 45 day-of-week effect θ_x^{day} . Placing an independent normal prior on the idiosyncratic, arm-specific, and
 46 day-specific effects induces a structured covariance matrix Σ . When the idiosyncratic terms have
 47 large variance, the decision-maker must guard against almost arbitrary non-stationary patterns. If
 48 these are believed to have smaller magnitude, the decision-maker may be able rule out some very poor
 49 arms early in the experiment.

50 Day-of-week effects are a standard concern when practitioners run A/B tests [Kohavi et al.,
 51 2020], so it is concerning that popular bandit algorithms like Thompson sampling and upper
 52 confidence bound [Lattimore and Szepesvári, 2020] fail in this example. The issue is that
 53 these algorithms either risk confounding by ignoring contextual information or aim to find
 54 the best action for every specific context, which is often not the experimenter’s goal (see
 55 the discussion of coarse segmentation below). In light of this discussion, it may not be
 56 surprising that many real life experiments implement uniformly random arm selection. This
 57 experimental design is highly robust, but it is inefficient when there are many arms and
 58 some can be quickly identified as inferior. The adaptivity of multi-armed bandit algorithms
 59 is important in those settings.

60 We propose *deconfounded Thompson sampling* (DTS). This method involves simple but critical
 61 modifications to Thompson sampling — an algorithm that is widely used in industry
 62 in academia. Our results suggest that DTS strikes a delicate balance: it is aggressive in
 63 shifting measurement effort away from alternatives that appear inferior while being robust
 64 to observed confounders like the day-of-week effects in Example 1.

65 1.1 A Model of Contextual Bandit Experiments

66 We formulate a general model of contextual bandit experiments that encompasses Example
 67 1 as a special case. The model captures the defining features of Example 1 including:

- 68 • *Coarse segmentation*: The ultimate decision-rule (1) pools together all seven contexts
 69 into a single segment over which decisions are held constant. If we think of the
 70 customers as belonging to one of seven different groups, then specifying a price for
 71 each day would be the most granular segmentation and (1) specifies the coarsest.
 72 In settings where the context contains all information available about a customer,
 73 coarse segmentation reduces data requirements, reduces the risk of bias, and avoids
 74 complex strategic incentives that occur when customers’ interactions affect their
 75 future service. In practice, products, public policies, and health interventions are
 76 often designed to serve a segment of the population (e.g. rural millennials) without
 77 being specialized to each individual.
- 78 • *Nonstationary confounders*: The experimenter needs to account for day-of-week
 79 effects in order to correctly infer which arm is best. They may need to model
 80 granular contextual information when performing inference even if they want to
 81 employ a coarse segmentation for the decisions they implement. The nonstationary

82 contexts act as the confounders for inferences and can be also viewed as the
83 distribution shifts in the earlier periods of the experiments.

84 • *Pure exploration*: In the lingo of the multi-armed bandit literature, what we have
85 described is a “pure-exploration” problem [Bubeck et al., 2009]. In common bandit
86 formulations, experimentation continues indefinitely but is costly only if suboptimal
87 action is selected. In our formulation, one hopes to quickly stop the experimentation
88 process and commit to given strategy for selecting actions going forward. This is
89 natural in settings where the process of experimentation is inherently costly, as it is
90 in clinical trials or many public policy experiments. Even in internet experiments,
91 the dominant workflow involves running a finite length experiment to validate or
92 select among alternatives. After an option is selected, engineering resources might
93 be invested toward productionizing it. One salient feature of the model is that the
94 target distribution or population distribution that is used to determine the best
95 action can be different from the empirical distribution over the contexts observed
96 during the experiments.

97 We formulate a general model that combines these features. A decision-maker experiments
98 across a sequence of periods. In each, they observe a context vector, select from a finite set
99 of possible actions, and observe a reward whose probability distribution depends on the
100 chosen context and action. After the experimentation process stops, the decision-maker
101 commits to a given strategy for selecting actions going forward. Specifically, they pick
102 among a class of candidate policies, each of which is a rule that prescribes an action for
103 every context. Restricting the class of candidate policies enforces coarse segmentation. The
104 decision-maker’s choice is judged by how it performs on average under contexts drawn
105 from a population distribution, effectively capturing how that policy will perform when
106 employed throughout an extremely large number of remaining periods. We assume the
107 population distribution is known, which would be essential if the contexts observed during
108 the experiment are not representative of the distribution anticipated in the future. More
109 generally, web companies typically have rich historical data on their users and should not
110 try to estimate this population’s attributes separately in each experiment they run.

111 1.2 Failure of Popular Bandit Algorithms

112 The two most popular approaches to (stochastic) multi-armed bandit problems are upper
113 confidence bound (UCB) and Thompson sampling (TS) algorithms [see e.g. Slivkins et al.,
114 2019, Lattimore and Szepesvári, 2020]. UCB selects the arm with the highest UCB on its
115 mean reward. TS is a randomized strategy under which the probability of sampling an arm
116 is “matched” to the posterior probability that arm is optimal. Both algorithm have been
117 applied to a variety of complex and interesting online decision-making problems.

118 We show that neither TS nor UCB, as usually applied, can address Example 1. In problems
119 with contexts, TS and UCB aim to select an action that could plausibly maximize the
120 expected reward earned in the current context. UCB does this by forming a confidence
121 bound on each arm’s performance under the current context and TS performs probability
122 matching with respect to the optimal arm in the current context. These strategies do not
123 gather sufficient information about arms that are suboptimal on the current day but might
124 be optimal throughout the week. They also could waste measurement effort on arms
125 that appear almost certain to offer suboptimal average performance throughout the week.
126 Heuristic versions of TS or UCB that disregard contextual information when performing
127 inference would risk confounding due to un-modeled day-of-week effects.

128 A potential adaptation of UCB to Example 1 would form UCBs on the weeklong average
129 reward in (1). We show this may sample only a single arm on a given day because UCBs do
130 not diminish until later days are observed. As a result, the data it collects cannot be used to
131 identify the best arm in (1), regardless of the length of the problem’s time horizon.

132 1.3 Deconfounded Thompson Sampling

133 Our proposed algorithm makes two modifications to Thompson sampling as it is usually
134 defined in contextual bandit problems. The first makes the algorithm suitable for learning
135 about a target policy with coarse segmentation. In the setting of Example 1, rather than
136 perform probability matching with respect to the best action for the current day, it performs
137 probability matching with respect to the arm with best performance throughout the week as
138 in (1). More generally, the proposed algorithm performs probability matching with respect
139 to the action prescribed at the current context by the target policy in the policy class. This
140 idea limits exploration to important distinctions between the candidate policies. The second
141 modification makes the algorithm suitable for pure-exploration problems by adapting the
142 top-two sampling strategy of Russo [2020]. This modification explores suboptimal arms
143 more aggressively by running Thompson sampling until two distinct actions are drawn
144 and then randomly picking among those “top-two”. We call this algorithm deconfounded
145 Thompson sampling (DTS). Unlike standard TS, it can control for confounding factors
146 without segmenting its decisions on the basis of those confounders.

147 1.4 Theoretical Results

148 It is difficult to give a single theoretical analysis that illuminates all the issues that are
149 relevant in practice. Instead, we focus on a single algorithm and prove three distinct results
150 that stress different capabilities. All results study simple regret [Bubeck et al., 2009], which
151 measures the shortfall in the expected future per-period reward earned by the decision-
152 maker’s selected policy relative to the best the best policy in the policy class. We elaborate
153 on the results below:

- 154 1. *Robustness to delay and confounding*: Our first result removes the assumption that
155 contexts are drawn i.i.d. For analytical tractability, we assume a Gaussian linear
156 model governs reward observations and focus on a best-arm learning problem,
157 where the goal is to identify the best fixed arm to employ in the future. Example
158 1 serves as a special case. We study the expected simple regret incurred by DTS,
159 conditioned on an arbitrary sequence of contexts. We provide a bound that depends
160 only on the information contained in the contexts and is completely independent
161 of the order in which they arrive, demonstrating robustness to non-stationary
162 confounders that are modeled by the algorithm. This result also allows for an
163 arbitrary delay in observing reward realizations.
- 164 2. *Adapting optimally to the problem instance*: Our next result fixes some arbitrary
165 parameter vector and studies expected simple regret conditioned on this vector
166 being the true draw from nature. This can be thought of as a “frequentist” bound,
167 whereas the previous two were “Bayesian.” This section again imposes the
168 assumption that contexts are drawn i.i.d. and, for analytical tractability, again
169 assumes a Gaussian linear model governs reward observations and focuses on a
170 best-arm learning problem. A fundamental lower bound shows how the expected
171 sample size of an adaptive experiment must grow in order to guarantee some
172 vanishing level of simple regret. The sampling requirements are milder for problem
173 instances where some arms are far from optimal and can be effectively discarded
174 with few samples. We prove that DTS meets attains this asymptotic lower bound.
175 In this sense it optimally adapts its experimentation to the problem instance.

176 It may not be difficult to design an algorithm that attains one of the results above. It is
177 remarkable, however, that these distinct properties are satisfied simultaneously by one
178 simple heuristic algorithm. Attaining both simultaneously seems to require a delicate
179 balance between robustness and adaptivity.

180 **References**

- 181 Susan Athey and Stefan Wager. Policy learning with observational data. *Econometrica*, 89(1):133–161,
182 2021.
- 183 Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems.
184 In *International conference on Algorithmic learning theory*, pages 23–37. Springer, 2009.
- 185 Herman Chernoff et al. Sequential design of experiments. *Annals of Mathematical Statistics*, 30(3):
186 755–770, 1959.
- 187 Peter Glynn and Sandeep Juneja. A large deviations perspective on ordinal optimization. In *Proceedings*
188 *of the 2004 Winter Simulation Conference, 2004.*, volume 1. IEEE, 2004.
- 189 Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification
190 in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- 191 Ron Kohavi, Diane Tang, and Ya Xu. *Trustworthy online controlled experiments: A practical guide to a/b*
192 *testing*. Cambridge University Press, 2020.
- 193 Tor Lattimore and Csaba Szepesvári. The end of optimism? an asymptotic analysis of finite-armed
194 linear bandits. In *Artificial Intelligence and Statistics*, pages 728–737. PMLR, 2017.
- 195 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 196 Chao Qin, Diego Klabjan, and Daniel Russo. Improving the expected improvement algorithm.
197 *Advances in Neural Information Processing Systems*, 2017:5382–5392, 2017.
- 198 Daniel Russo. Simple bayesian algorithms for best-arm identification. *Operations Research*, 2020.
- 199 Aleksandrs Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine*
200 *Learning*, 12(1-2):1–286, 2019.

201 **A Problem Formulation**

202 After running an experiment, a decision-maker must select among k arms. The performance
 203 of an arm depends on the context in which it is employed. Each context is represented by a
 204 d dimensional feature vector and the set of possible contexts is denoted by \mathcal{X} . For each arm
 205 $i \in [k] := \{1, \dots, k\}$, there is an uncertain arm specific parameter $\theta^{(i)}$, which we model as
 206 a draw $\theta^{(i)} \sim N(\mu_{1,i}, \Sigma_{1,i})$ from a multi-variate Gaussian prior. We let $\theta = (\theta^{(1)}, \dots, \theta^{(k)})$
 207 denote the concatenation of the vectors. A linear function $\mu(\theta, i, x) = \langle \theta^{(i)}, x \rangle$ determines
 208 the performance of arm i in context $x \in \mathcal{X}$.

209 We assume the decision-maker has access to a probability distribution w over contexts that
 210 encodes the frequency with which they expect contexts to occur in the future. We call this
 211 either the *target distribution* or the *population distribution*, where the latter suggests that w
 212 denotes the characteristics of a population of individuals. If employed across a large number
 213 of future periods, arm i would generate average reward

$$\mu(\theta, i, w) := \sum_{x \in \mathcal{X}} w(x) \langle \theta^{(i)}, x \rangle = \langle \theta^{(i)}, X_{\text{pop}} \rangle \quad \text{where} \quad X_{\text{pop}} := \sum_{x \in \mathcal{X}} w(x)x. \quad (2)$$

214 In Example 1, X_{pop} is the vector $(1/7, \dots, 1/7)$ and $\mu(\theta, i, w)$ is the average that appears in
 215 Equation (1). If the decision-maker knew θ , the optimal arm to employ in the future would
 216 be $I^* = I^*(\theta) \in \arg \max_{i \in [k]} \mu(\theta, i, w)$.

217 For technical or notational convenience, we make several additional assumptions. First, we
 218 assume \mathcal{X} is finite (though possibly enormous), which allows us later to analyze a lower
 219 bound on performance that is expressed through a finite dimensional optimization problem.
 220 Second, we assume that the arm-specific parameters $\theta^{(i)}$ are drawn independently across
 221 arms, allowing us to track beliefs separately across arms in the analysis. Assume also that
 222 the prior covariance matrix $\Sigma_{1,i}$ is the same for each arm i and is positive definite. We denote
 223 this by Σ_1 .

224 **Sequential learning.** The decision-maker can reduce uncertainty about θ through
 225 experimentation. In each period, $t \in \mathbb{N}$, they select an arm $I_t \in [k]$ in some context $X_t \in \mathcal{X}$
 226 and observe a real valued reward signal $R_t = \langle \theta^{(I_t)}, X_t \rangle + W_t$, where $W_t \mid \theta, X_t \sim N(0, \sigma^2)$ is
 227 Gaussian noise drawn independently across time. Rewards are observed after a lag of $L \geq 1$
 228 periods. The information available when choosing I_t is the history $H_t = (X_{1:t}, I_{1:t-1}, R_{1:t-L})$.
 229 Formally, the action I_t must be chosen as a function of H_t and some random seed ζ_t that is
 230 independent of all else. We assume the context sequence $(X_t)_{t \in \mathbb{N}}$ is independent of θ , so
 231 that the decision-maker cannot passively learn the impact of their actions by observing the
 232 contexts.

233 The distribution of $\theta^{(i)}$ conditioned on H_t is multivariate Gaussian with covariance and
 234 mean given by $\Sigma_{t,i} = \Sigma_{1,i}$ and $\mu_{t,i} = \mu_{1,i}$ for $t \leq L$ and

$$\Sigma_{t,i} = \left(\Sigma_1^{-1} + \sigma^{-2} \sum_{\ell=1}^{t-L} \mathbb{1}\{I_\ell = i\} X_\ell X_\ell^\top \right)^{-1} \quad \mu_{t,i} = \Sigma_{t,i} \left(\Sigma_1^{-1} \mu_{1,i} + \sum_{\ell=1}^{t-L} \mathbb{1}\{I_\ell = i\} X_\ell R_\ell \right). \quad (3)$$

235 for $t > L$. Posterior beliefs about θ induce posterior beliefs about I^* . We set $\alpha_{t,i} = \mathbb{P}(I^* = i \mid$
 236 $H_t)$ for any period $t \in \mathbb{N}$ and arm $i \in [k]$. Since $\mu(\theta, i, w)$ is a linear function of $\theta^{(i)}$, it also
 237 has a Gaussian posterior. We write $\mu(\theta, i, w) \mid H_t \sim N(m_{t,i}, s_{t,i}^2)$ where

$$s_{t,i}^2 = X_{\text{pop}}^\top \Sigma_{t,i} X_{\text{pop}} \quad m_{t,i} = \langle X_{\text{pop}}, \mu_{t,i} \rangle. \quad (4)$$

238 Notice that the Latin alphabet is used for the posterior parameters of the scalar quantity
 239 $\mu(\theta, i, w)$ and the Greek alphabet is used for the posterior parameters of the vector $\theta^{(i)}$.

240 **Performance measures.** Let $H_T^+ = (X_{1:T}, I_{1:T}, R_{1:T})$ denote all information generated by a
 241 T -period experiment, including the delayed reward outcomes. The non-negative random
 242 variable

$$\Delta_T = \mu(\theta, I^*, w) - \mu(\theta, \hat{I}_T^+, w)$$

243 measures the shortfall in future performance caused by selecting the greedy decision at
 244 time T by $\hat{I}_T^+ \in \arg \max_{i \in [k]} \mathbb{E} [\mu(\theta, i, w) \mid H_T^+]$ with only the incomplete information about
 245 θ accrued after T measurements. We call Δ_T the *simple regret* at time T , after Bubeck et al.
 246 [2009]. Having in mind policy decisions where $\mu(\theta, i, x)$ denotes the utility generated for
 247 an individual with features x , Athey and Wager [2021] call this the *utilitarian regret*. Notice
 248 that the decision \hat{I}_T^+ can be made using the full results of the experiment H_T^+ while a
 249 measurement decision I_t must be made in real-time based on partial information H_t .

250 The goal in the problem, informally, is to experiment intelligently so that simple regret is
 251 small after using as few measurements as possible. This objective is can be formalized in
 252 several ways. We focus on two ways of studying performance that allow for clear analytical
 253 insight into specific properties of deconfounded Thompson sampling:

- 254 1. (Fixed budget and Bayesian) In Section D, we study the expected simple regret
 255 $\mathbb{E} [\Delta_T \mid X_{1:T} = x_{1:T}]$ at some finite time T , conditioned on the sequence of realized
 256 contexts $X_{1:T} := (X_1, \dots, X_T)$ taking on some specific value. This expectation
 257 integrates over most randomness in the problem, including over the prior
 258 distribution, and emphasizes dependence on the observed contexts and their order.
 259 Our goal in this section is to show that deconfounded TS satisfies an important
 260 robustness property other adaptive algorithms do not: roughly speaking, we have a
 261 result of the form $\mathbb{E} [\Delta_T \mid X_{1:T} = x_{1:T}] \leq \tilde{O}(\sqrt{k/T})$, where the big- O hides a natural
 262 dependence on the second moment $\frac{1}{T} \sum_{t=1}^T x_t x_t^\top$ but has no dependence on context
 263 order.
- 264 2. (Adaptive stopping and frequentist) In Section E, we aim to verify that the algorithm
 265 adapts its measurement effort optimally, in an appropriate sense, as it learns
 266 about the true problem instance. To do this, we study performance conditional
 267 on the unknown parameter θ but integrate over the distribution of the contexts
 268 (X_1, X_2, \dots) , which we assume to drawn i.i.d. Following the style of result in Russo
 269 [2020], Glynn and Juneja [2004] or Kaufmann et al. [2016], we would hope to show
 270 that $\mathbb{E} [\Delta_T \mid \theta]$ goes to zero at an exponential rate T grows, and that the problem-
 271 dependent exponent is in appropriate sense the best-possible among adaptive
 272 algorithms. This is called the “fixed-budget” formulation in the literature on the
 273 best-arm identification literature, because there is a hard constraint on the number
 274 of samples (i.e T) that can be collected that must be satisfied with probability one.
 275 Unfortunately, the sharp asymptotic limits in that setting are poorly understood,
 276 even in problems without contexts. We instead look at formulations in which there
 277 is a soft-constraint on the number of measurements. There we study performance
 278 at a adaptively chosen stopping time τ , which essentially, stops once the posterior
 279 expectation of simple regret is small. We study the combined cost $\mathbb{E} [c\tau + \Delta_\tau \mid \theta]$
 280 as $c \rightarrow 0$, measuring the *expected* number of samples required to deliver vanishing
 281 simple regret. This formulation follows classic work of Chernoff et al. [1959]; very
 282 similar results follow if one instead imposes a constraint on the simple regret or
 283 the probability of incorrect selection, which is called the “fixed-confidence” setting.
 284 Kaufmann et al. [2016].

285 B Deconfounded Thompson Sampling

286 Deconfounded Thompson sampling (DTS) can be defined succinctly. At each time period
 287 $t \in \mathbb{N}$, it selects an arm to measure through the following procedure:

288 *Continue sampling from α_t until two distinct arms are chosen.*
 289 *Flip a (biased) coin to select among these two.*

290 Recall that $\alpha_t \in \mathbb{R}^k$ is defined by $\alpha_{t,i} = \mathbb{P}(I^* = i \mid H_t)$. We explain below how to efficiently
 291 sample from this distribution. Throughout the paper, we take $\beta_t \in (0, 1]$ to be the probability
 292 the first sample from α_t is played. By default, we recommend an unbiased coin ($\beta_t = 1/2$)
 293 but this is discussed further below.

294 DTS can be understood as making two modifications to Thompson sampling in contextual
 295 bandits:

296 1. *Changing the learning target:* Thompson sampling for contextual bandits usually
 297 samples an action according to the probability it maximizes the mean reward
 298 in the current context. In particular, one sets $\mathbb{P}(I_t = i \mid H_t) = \mathbb{P}(i =$
 299 $\arg \max_{i \in [k]} \mu(\theta, i, X_t) \mid H_t)$. DTS is instead based on sampling from the posterior
 300 distribution of the arm I^* , which is the arm that maximizes the average reward
 301 in the target population rather than in the current context. Defining α_t carefully
 302 controls for confounders while directing exploration toward learning about the
 303 target arm of interest.

304 2. *Resampling:* Consider a problem without contexts. Then standard TS draws I_t from
 305 α_t , without resampling. This algorithm is designed to maximize the reward earned
 306 throughout the experiment, implicitly imagining that the experimentation process
 307 never ends. But it performs poorly if there is an interest also in being able to rapidly
 308 stop and commit confidently to a decision. To understand the issue, imagine that
 309 $\alpha_{t,1} = .95$, so the algorithm believes there is a 95% chance that arm 1 is optimal.
 310 Then TS plays arm 1 in roughly 19/20 periods, making it very slow to gather
 311 information about alternatives. TS would be very slow to reach 99% confidence as
 312 result and this is exacerbated if even higher confidence is desired.

313 To overcome this issue, Russo [2020] suggests a “top-two sampling” version of TS,
 314 which continues drawing arms from TS until two distinct options are drawn and
 315 then flips a biased coin to select among these two. To understand the resampling
 316 step, imagine that $\alpha_{t,1} \rightarrow 1$ as $t \rightarrow \infty$. In this limit, the first sample from α_t
 317 is nearly always arm 1 and this is played with probability β_t . Otherwise, an
 318 arm is chosen by resampling, and the chance of picking arm $j > 1$ is roughly
 319 $\mathbb{P}(I_t = j \mid I_t \neq 1) \sim \frac{\alpha_{t,j}}{1 - \alpha_{t,1}} = \mathbb{P}(I^* = j \mid I^* \neq 1)$. Resampling shifts $1 - \beta_t$ fraction
 320 of measurement effort away from arm 1 and assigns it to the strongest challengers.
 321 In particular, a challenger is sampled according to its conditional probability of
 322 being optimal.

323 By default in this paper, we have in mind that DTS is implemented with a fair coin ($\beta_t = 1/2$).
 324 Fixing a higher bias might be helpful to a practitioner. This would focus more measurement
 325 effort on the most promising arm, providing more confidence about the rewards it generates
 326 and reducing the expected regret incurred during the experiment. On the other hand, a
 327 longer experiment might be required to reach confidence about the best arm if a high bias is
 328 used. We discuss in Section E how the bias might be tuned adaptively as data is observed to
 329 maximize certain asymptotic performance measures.

330 **Notable features of DTS.** Before proceeding, it is worth highlighting a few important
 331 features of DTS. First, let us draw a contrast with another popular strategy, UCB algorithms.
 332 These are based on the principle of *optimism in the face of uncertainty*. The decision-maker
 333 responds to uncertainty by playing whichever action is best in the best plausible model
 334 given current information. Notice that DTS, by default, *randomizes in the face of uncertainty*.
 335 Indeed, with a symmetric prior, one would have $\alpha_{1,1} = \dots = \alpha_{1,k} = 1/k$ and so the initial
 336 arm I_1 is sampled uniformly at random. As information is gathered, beliefs are updated
 337 and the decision-maker becomes less likely to sample inferior arms. The algorithm’s
 338 randomization gives it a chance of sampling all plausibly optimal arms in all contexts. This
 339 appears to be critical to some of its robustness properties.

340 Another striking feature of the algorithm is that decisions at time t do not depend on the
 341 context at time t . That decisions are *context independent* in this way could offer substantial

342 practical benefits. Even if contexts are logged, enormous engineering resources might be
 343 required to develop a system that observes contexts and responds in real time. For instance,
 344 assessing X_t could easily require querying several different datasets containing the current
 345 user’s interaction history and then applying a trained machine learning algorithm that
 346 generates a compact feature vector from this history. With a context independent algorithm,
 347 this could be done without substantial latency requirements.

348 **Efficient computation.** Following conventional implementation of Thompson sampling,
 349 a generic approach sampling from α_t , is to sample a parameter vector $\tilde{\theta}$ from the
 350 posterior distribution of θ and then to find the arm $\arg \max_{i \in [k]} \mu(\tilde{\theta}, i, w)$ that is best under
 351 this sample. The structure of Gaussian linear belief models allows for an even cleaner
 352 implementation of DTS. Because the population average reward of arm i , $\mu(\theta, i, w)$, has a
 353 Gaussian posterior with posterior parameters given in (4), one can directly perform inference
 354 on the population average rewards.

355 The pseudocode below almost perfectly mirrors top-two TS in problems without contexts,
 356 except that the posterior parameters $(m_{t+1,i}, s_{t+1,i}^2)$ are updated in a manner that controls
 357 for observed confounders, reflects the target population of contexts, and may be affected
 358 by delayed observations. By default, we imagine $\beta_t = 1/2$, but the pseudocode allows for
 359 adaptive tuning of the coin’s bias.

360 A possible concern is that it might take an enormous number of samples until the top-two
 361 arms differ (i.e. until $I_t^{(1)} \neq I_t^{(2)}$). However, each fresh sample has chance $1 - \alpha_{t, I_t^{(1)}}$ of
 362 generating a different arm, so this while-loop is expected to require many iterations only if
 363 the the posterior has already concentrated on a single action. In that case, it makes sense
 364 to terminate the experiment. When the posterior concentrates, there are also a variety of
 365 asymptotic approximations that could be used to calculate selection probabilities and avoid
 366 repeated sampling.

Algorithm 1: DTS allocation rule in Gaussian best-arm learning

Input prior parameters $(\mu_{1,i}, \Sigma_{1,i})_{i \in [k]}$, population weights X_{pop} and noise variance σ^2 .

for $t = 1, 2, \dots$ **do**

Sample $v_i \sim N(m_{t,i}, s_{t,i}^2)$ for $i \in [k]$ and set $I_t^{(1)} = \arg \max_{i \in [k]} v_i$;

do

Sample $v_i \sim N(m_{t,i}, s_{t,i}^2)$ for $i \in [k]$ and set $I_t^{(2)} = \arg \max_{i \in [k]} v_i$;

367 **while** $I_t^{(1)} = I_t^{(2)}$;

Flip coin $C_t \in \{0, 1\}$ with bias $\mathbb{P}(C_t = 1) = \beta_t$;

Play arm $I_t = I_t^{(1)} C_t + I_t^{(2)} (1 - C_t)$;

Gather delayed observation $o = (I_{t-L}, X_{t-L}, R_{t-L})$.;

Calculate posterior parameters $m_{t+1,i}, s_{t+1,i}^2$ for $i \in [k]$ according to (4) to reflect o ;

Calculate new tuning parameter β_{t+1} if using adaptive tuning;

end

368 **C Failure of Alternative Bandit Algorithms**

369 This section provides examples showing that alternative bandit algorithms can fail for
 370 simple examples within the scope of our problem formulation. Most interesting, perhaps, is
 371 that a deconfounded UCB algorithm cannot address a simplified version of the example with
 372 day-of-week effects described in the introduction. Past theory on TS highlights connections
 373 to UCB, so any theory of DTS in that example will need to push well beyond current
 374 understanding. We also show the failure of a context unaware algorithm and a usual
 375 contextual bandit algorithm.

376 What does it mean that these algorithms ‘fail’? We show formally that simple regret is
 377 bounded as $\mathbb{E} [\Delta_T] \geq c$, where c is some absolute numerical constant that does not depend
 378 on T . Regardless of the time dedicated to experimentation, the data the collected by these
 379 algorithms is inadequate and cannot be used to make near-optimal decisions. The examples
 380 we describe are meant to give insight into what can go wrong with alternative algorithms
 381 and the subtleties of designing an algorithm like DTS. They are purposefully simplistic.

382 C.1 Deconfounded UCB

383 Consider the following simplification of Example 1. Here, there are two contexts instead of
 384 seven and we restrict to the case of two actions.

385 **Example 2** (Day of week effects). *The context set is $\mathcal{X} = \{1, 2\}$ and there are $k = 2$ arms. The*
 386 *reward at time t is $R_t = \theta_{X_t}^{(I_t)} + W_t$ where each $\theta_x^{(i)}$ is independent and Gaussian and $W_t \sim N(0, \sigma^2)$*
 387 *is i.i.d Gaussian noise. Observations are not subject to delay (i.e $L = 1$). The the goal is to identify*
 388 *the best arm under equal context weights w :*

$$I^* = \arg \max_{i \in [2]} \frac{\theta_1^{(i)} + \theta_2^{(i)}}{2}.$$

389 *The context sequence is non-random, with $X_t = 1$ for $t \leq \lfloor T/2 \rfloor$, $X_t = 2$ for $t > \lfloor T/2 \rfloor$.*

390 Consider a UCB analogue of our Thompson sampling based algorithm. Reflecting that the
 391 true goal is to select an arm with strong performance throughout the week, not on a specific
 392 day, it plays the arm with the highest UCB on its average performance throughout the week:
 393

$$I_t \in \arg \max_{i \in [k]} m_{t,i} + z \cdot s_{t,i} \quad \text{for all } t \in \mathbb{N}. \quad (5)$$

394 where $m_{t,i}$ and $s_{t,i}$ are defined in (4) and $z > 0$ is a tuning parameter. When $z = 1.645$, the
 395 term $m_{t,i} + z \cdot s_{t,i}$ measures the 95% quantile of the posterior distribution. Like DTS, this can
 396 be thought of as a *deconfounded* UCB, which still selects the arm with the highest upside but
 397 accounts for observed confounders when performing inference.

398 The next result shows formally that deconfounded UCB fails to collect adequate data,
 399 regardless of the length of the time horizon. The issue is that the UCB in (5) is sometimes
 400 higher for action 2 for each of the first $T/2$ periods. Action 1 is then never sampled in
 401 context 1, so learning is incomplete. This holds true regardless of how z is set and holds for
 402 time dependent tuning parameters. The issue is that, unlike common bandit settings, UCBs
 403 do not diminish when actions are repeatedly sampled in a single context.

404 **Lemma 1.** *Consider Example 2. Suppose that the components of the vector $\theta = (\theta_x^{(i)})_{i \in [2], x \in [2]}$ are*
 405 *independent with $\theta_x^{(1)} \sim N(0, 1)$ and $\theta_x^{(2)} \sim N(0, 2)$ for $x \in \{1, 2\}$, and $\sigma^2 = 0$. If (5) holds, then*
 406 *there is an absolute numerical constant $c > 0$ such that $\mathbb{E} [\Delta_T] \geq c$ for any $T \in \mathbb{N}$.*

407 C.2 Context Unaware Algorithms

408 Our next example highlights the risk of confounding for an algorithm that does not model
 409 day-of-week effects when performing inference. We set $\tilde{s}_{t,i}^2 = \left(1 + \sigma^{-2} \sum_{\ell=1}^{t-1} \mathbb{1}(I_\ell = i)\right)^{-1}$
 410 and $\tilde{m}_{t,i} = \tilde{s}_{t,i}^2 \left(\sum_{\ell=1}^{t-1} \mathbb{1}(I_\ell = i) R_\ell\right)$. We define these expressions for $\sigma^2 = 0$ by taking the
 411 limit as $\sigma^2 \downarrow 0$. In particular, we set $\tilde{s}_{t,i}^2 = 0$ if arm i has been played previously and $\tilde{m}_{t,i}$
 412 to be 0 if arm i was never played previously and to be the empirical average reward
 413 otherwise. These are the posterior updating equations if $\theta_1^{(i)} \sim N(0, 1)$ and the algorithm
 414 (incorrectly) ignores day of week effects and assumes $\theta_2^{(i)} = \theta_1^{(i)}$ almost surely.

415 Based on this, define context unaware Thompson sampling. It chooses an arm at time t
 416 according to

$$I_t = \arg \max_{i \in [2]} v_{t,i} \quad \text{where} \quad v_{t,i} \mid H_t \sim N(\tilde{m}_{t,i}, \tilde{s}_{t,i}^2). \quad (6)$$

417 In the above equation $v_{t,1}$ and $v_{t,2}$ are sampled independently. The next lemma formalizes
 418 that this algorithm risks confounding. The same result applies to a context unaware form
 419 UCB, which forms UCBs based on $\tilde{m}_{t,i}$ and $\tilde{s}_{t,i}^2$. A context-unaware top-two TS algorithm
 420 fails in a similar way in problems with more than two actions.

421 **Lemma 2** (Failure of context unaware Thompson sampling). *Consider Example 2. Suppose*
 422 *the components of the vector $\theta = (\theta_x^{(i)})_{i \in [2], x \in [2]}$ are independent with $\theta_x^{(1)} \sim N(0, 1)$ and $\theta_x^{(2)} \sim$*
 423 *$N(0, 2)$ for $x \in \{1, 2\}$, and $\sigma^2 = 0$. If (6) holds then there exists an absolute numerical constant*
 424 *$c > 0$ such that for all $T \in \mathbb{N}$, $\mathbb{E} [\Delta_T] \geq c$.*

425 C.3 Contextual Bandit Algorithms

426 The goal in our formulation is to select among a very restricted set of decision-rules: those
 427 that choose a common action, irrespective of context. Experimentation should be tailored to
 428 this objective. Here, we give insight into potential failures when the exploration algorithm
 429 is designed with a different learning target in mind. Consider the following example. There
 430 are three actions, and the decision-maker would like to identify the best action to employ on
 431 average, across all contexts. Imagine that the context set describes two customer segments.
 432 Action 1 appeals to one segment, but is highly unappealing to the other. For action 2,
 433 the situation is reversed. Action 3 is not ideal for either segment, but is also not disliked
 434 by either. When personalization is inappropriate or costly, action 3 may be the preferred
 435 communal option.

436 The next example does not align with our formulation, because we take the prior distribution
 437 to be non-Gaussian. Similar issues can arise with a Gaussian prior, but its unbounded nature
 438 always allows for a nonzero— even if very small — chance that the mainstream action is better
 439 even for a specific segment.

440 **Example 3** (A mainstream action). *Consider a problem with $k = 3$ arms and 2 contexts given*
 441 *as $\mathcal{X} = \{1, 2\}$. The population distribution w is uniform over \mathcal{X} and $(X_t)_{t \in \mathbb{N}}$ are drawn i.i.d*
 442 *from w . The components of the parameter vector $\theta = (\theta_0, \theta_1, \theta_2)$ are drawn independently with*
 443 *$\theta_0 \sim \text{Uniform}([0, 1])$ and $\theta_x \sim \text{Uniform}(\{1, 2\})$ for $x \in [2]$. Rewards are noiseless, with*
 444 *$R_t = \mu(\theta, I_t, X_t)$. Observations are not subject to delay (i.e. $L = 1$). Action 3's performance is*
 445 *insensitive to the context, and it always generates mean-reward $\mu(\theta, 3, x) = \theta_0$. Actions 1 and 2*
 446 *generate mean rewards in context $x \in \mathcal{X}$ given by*

$$\mu(\theta, 1, x) = 1/2 + (1/2)\mathbb{1}(\theta_x = 1), \quad \mu(\theta, 2, x) = 1/2 + (1/2)\mathbb{1}(\theta_x = 2).$$

447 The next lemma formalizes that contextual Thompson sampling, which selects an action
 448 according to the posterior probability it is the optimal action for the current context, has
 449 simple regret that does not vanish even as the horizon grows. The same result applies to
 450 appropriate contextual versions of UCB. The simple reason is that action 3 is never sampled,
 451 because it does not maximize the reward in either context. This means no information about
 452 θ_0 is gathered and the decision-maker cannot determine whether action 3 is the best arm
 453 to select. If the goal is to identify the best policy within a restricted class, the exploration
 454 algorithm needs to be designed so that it gathers the right information for this task. The
 455 proof follows from this argument and is omitted for brevity.

456 **Lemma 3.** (Failure of contextual Thompson sampling) *Suppose that $\mathbb{P}(I_t = i \mid H_t) =$*
 457 *$\mathbb{P}(I^*(\theta; X_t) = i \mid H_t)$ for each $i \in [k]$. Under Example 3, there is an absolute numerical constant*
 458 *$c > 0$ such that for all $T \in \mathbb{N}$, $\mathbb{E} [\Delta_T] \geq c$.*

459 **D Result 1: Robustness to Delay and Confounding**

460 In this paper, we provide the first of two guarantees for DTS. The focus here is on
 461 assurances of robustness. We do this by establishing generic bounds on simple regret
 462 that essentially mirror regret guarantees satisfied when actions are selected uniformly at
 463 random. The challenge is to show that the adaptivity of DTS does not make the algorithm’s
 464 performance brittle, in contrast to the algorithms described in Section C. In the next section,
 465 we complement this study of robustness with a study of the adaptivity benefits of DTS.

466 **D.1 Performance Guarantee**

467 Because we do not require contexts to be i.i.d, there is no guarantee that the observed context
 468 sequence provides the information required to select the best-arm. We measure this through
 469 the quantity

$$V(X_{1:T}) = X_{\text{pop}}^\top \left(\Sigma_1^{-1} + \sigma^{-2} \sum_{t=1}^T X_t X_t^\top \right)^{-1} X_{\text{pop}}. \quad (7)$$

470 The matrix $\left(\Sigma_1^{-1} + \sigma^{-2} \sum_{t=1}^T X_t X_t^\top \right)^{-1}$ appearing in (7) would be the posterior covariance
 471 matrix of $\theta^{(i)}$ at the end of the experimentation horizon if that arm were played in every
 472 period. We similarly think of $V(X_{1:T})$ as the posterior variance $\text{Var}(\mu(\theta, i, w) \mid H_T^+)$ of the
 473 population effect of arm i if we observed the reward it generated in every period of the
 474 experiment. Notice that what makes the day-of-week effects in Example 2 challenging is *the*
 475 *order* in which contexts arrive. But observing a single arm throughout the entire experiment
 476 would be informative, and so $V(X_{1:T})$ would be small if T were large.

477 If arms were selected uniformly at random, we might expect the posterior variance of each
 478 one to scale roughly as $k \cdot V(X_{1:T})$, reflecting that information is divided equally across the
 479 arms. The next result establishes a simple regret bound for DTS that scales as $\sqrt{k \cdot V(X_{1:T})}$.
 480 One can think of this result as indicating a robustness property: the algorithm can cope with
 481 arbitrary context order and delayed reward observations, offering a guarantee matching
 482 what we would attain under a uniform allocation even when the context order and delay
 483 are severe. Of course, DTS is actually a highly adaptive algorithm, so it is subtle to show it
 484 satisfies this kind of robustness property and avoids the pitfalls described in Section C.

485 For random variables X and Y , let $\mathbb{H}(X)$ and $\mathbb{H}(X|Y)$ denote the Shannon entropy and
 486 conditional Shannon entropy of X .

487 **Proposition 1.** *Suppose that $\|X_t\|_2 \leq 1$ almost surely for $t \in \mathbb{N}$. If DTS applied with tuning*
 488 *parameters satisfying $\inf_{t \in \mathbb{N}} \beta_t \geq 1/2$ almost surely, then for any $T \in \mathbb{N}$,*

$$\mathbb{E} [\Delta_T \mid X_{1:T}] \leq \sqrt{2\iota \cdot k \cdot \mathbb{H}(I^* \mid H_T^+) \cdot V(X_{1:T})}$$

489 *where $\iota = \max \left\{ 9 \log \left(d \lambda_{\max}(\Sigma_1) \left[\lambda_{\max}(\Sigma_1^{-1}) + T \right] \right) \cdot \lambda_{\max}(\Sigma_1), 9 \right\}$.*

490 Under a natural condition that ensures the context sequence contains sufficient information
 491 about the population distribution, the next corollary of Proposition 1 gives a simple-
 492 regret bound that scales as $\tilde{O}(\sqrt{k/T})$. Notice that this result is nearly-independent of
 493 the dimension of the linear model d . If $X_t \sim w$, then $\mathbb{E} [X_t X_t^\top] = X_{\text{pop}} X_{\text{pop}}^\top + \text{Cov}(X_t)$. In
 494 this sense, if context vectors have high variance in every direction, the bound $\frac{1}{T} \sum_{t=1}^T x_t x_t^\top \succeq$
 495 $X_{\text{pop}} X_{\text{pop}}^\top$ may underestimate the information they provide and make this corollary
 496 conservative.

497 **Corollary 1.** Under the conditions of Proposition 1, for any sequence $x_{1:T} \in \mathcal{X}^T$, with
 498 $\frac{1}{T} \sum_{t=1}^T x_t x_t^\top \succeq X_{\text{pop}} X_{\text{pop}}^\top$,

$$\mathbb{E} [\Delta_T \mid X_{1:T} = x_{1:T}] \leq \sigma \sqrt{\frac{2\iota \cdot k \cdot \mathbb{H}(I^* \mid H_T^+)}{T}} \leq \sigma \sqrt{\frac{2\iota \cdot k \cdot \log(k)}{T}}$$

499 where ι is given in Proposition 1.

500 E Result 2: Adaptivity and Asymptotic Optimality

501 Like most popular multi-armed bandit algorithms, DTS allocates measurement effort
 502 adaptively. As time proceeds, it learns about the quality of different policies or
 503 arms. By shifting most measurements away from clearly inferior alternatives, it focuses
 504 experimentation effort where it is most useful. The previous section showed, although
 505 adaptivity makes other natural algorithms brittle in the face of nonstationary confounders,
 506 DTS has certain robustness guarantees. This section aims to formalize that DTS also adapts
 507 its measurement effort very effectively and, in a sense, *optimally* in a meaningful special case
 508 of our formulation.

509 E.1 Asymptotic Optimality Notion

510 We assess how effectively the algorithm uses its limited measurements, essentially, by
 511 understanding the rate at which simple regret decays as measurements are gathered. Among
 512 the several natural ways of studying this, we focus on one that allows for a sharp and
 513 enlightening asymptotic theory. We build on asymptotic limits of sequentially designed
 514 experiments that have been understood since classic work of Chernoff et al. [1959]. We
 515 allow the decision-maker to decide adaptively when to stop collecting measurements. The
 516 total cost incurred is $c\tau + \Delta_\tau$, where τ denotes the chosen stopping time, $c > 0$ is a cost
 517 per-period of experimentation, and Δ_τ is the simple-regret of the final decision.

518 We study the expected cost incurred under problem instance θ_0 , given by

$$\mathbb{E} [c\tau + \Delta_\tau \mid \theta = \theta_0]. \quad (8)$$

519 In this section, we focus on the parameter class

$$\Theta \triangleq \left\{ \theta \in \mathbb{R}^{dk} : \arg \max_{i \in [k]} \mu(\theta, i, w) \text{ is unique} \right\}.$$

520 In other words, each parameter in Θ corresponds to a problem instance with a unique best
 521 arm under the population distribution.

522 Sharp results can be established through asymptotic analysis as c tends to zero. This is a
 523 regime where the cost of gathering one more observation is negligible relative to the cost
 524 committing to a sub-optimal final decision. It arises naturally if one imagines the final
 525 decision will later be implemented for a very large number of periods. We establish a
 526 kind of uniform optimal guarantee, roughly showing that DTS minimizes (8) to first-order
 527 asymptotically for every specific instance θ_0 . This is only possible under an algorithm that
 528 tailors its experimentation optimally to θ_0 as information is gathered.

529 This theory requires an appropriate stopping rule is used. Attaining the exact optimal
 530 constant also requires tuning the β_t parameter as information about θ_0 is acquired. We
 531 discuss how this can be done with low computational cost and also discuss robustness with
 532 some non-adaptive choices of β .

533 Our result directly builds on previous analyses that have established similar results for
 534 top-two sampling rules [Russo, 2020, Qin et al., 2017]. It may be surprising, however, that
 535 these results extend to a contextual setting with linear models, given that more complex

536 exploration rules are often required to attain asymptotic optimality results in problems with
 537 parametric dependencies [See e.g. Lattimore and Szepesvari, 2017].

538 E.2 Notation

539 Recall the definitions of $m_{t,1} = \langle X_{\text{pop}}, \mathbb{E} [\theta^{(i)} | H_t] \rangle$ and $s_{t,i}^2 = \text{Var} (\mu(\theta, i, w) | H_t)$ given
 540 in . Since the analysis in this section is conditioned on θ , it is helpful to develop analogous
 541 notation for these quantities when an improper prior is used. Were an improper prior is
 542 used, $\mu_{t,i}$ and $\sigma_{t,i}^2$ would have the formulas:

$$\hat{m}_{t,i} = X_{\text{pop}}^\top \left[\sum_{\ell=1}^t \mathbb{1}(I_\ell = i) X_\ell X_\ell^\top \right]^{-1} \sum_{\ell=1}^t \mathbb{1}(I_\ell = i) X_\ell R_\ell$$

$$\hat{s}_{t,i}^2 = X_{\text{pop}}^\top \left[\sigma^{-2} \sum_{\ell=1}^t \mathbb{1}(I_\ell = i) X_\ell X_\ell^\top \right]^{-1} X_{\text{pop}}.$$

543 Note that $\hat{m}_{t,i}$ is simply the inner product of X_{pop} with the least-squares estimate for $\theta^{(i)}$. If
 544 the chosen arms $\{I_\ell\}$ were fixed in advance rather than selected adaptively, then $\hat{s}_{t,i}^2$ would
 545 be the formula for the sampling variance of $\hat{m}_{t,i}$.

546 It will be important to measure the strength of evidence that one arm outperforms another in
 547 the population. For this purpose, consider the natural test of the null hypothesis $\mu(\theta, i, w) \neq$
 548 $\mu(\theta, j, w)$. The classic test would be based on the z-score for the difference in means,

$$Z_{t,i,j} := \frac{m_{t,i} - \mu_{t,j}}{\sqrt{s_{t,i}^2 + s_{t,j}^2}}. \quad (9)$$

549 Each $Z_{t,i,j}$ follows a normal distribution with unit variance when I_1, \dots, I_{t-1} are chosen
 550 non-adaptively.

551 E.3 Lower Bound

552 Let $\mathcal{S} = \{v \in \mathbb{R}_+^k : \sum_{i=1}^k v_i = 1\}$ denote the $k - 1$ dimensional probability simplex. Define
 553 the complexity measure $\Gamma(\theta)$ by

$$\Gamma(\theta)^{-1} = \sup_{p: \mathcal{X} \rightarrow \mathcal{S}} \min_{i \neq I^*} \frac{1}{2\sigma^2} \frac{(\mu(\theta, I^*, w) - \mu(\theta, i, w))^2}{X_{\text{pop}}^\top \left(\mathbb{E} [p(X_1, I^*) X_1 X_1^\top]^{-1} + \mathbb{E} [p(X_1, i) X_1 X_1^\top]^{-1} \right) X_{\text{pop}}} \quad (10)$$

554 where p is a stochastic kernel, which is associates any $x \in \mathcal{X}$ with an element $p(x, \cdot) \in \mathcal{S}$.
 555 In this optimization problem, we imagine the experimenter the action at time t by sampling
 556 from $p(\cdot | X_t)$. The problem (10) seeks a measurement rule p that maximizes the growth
 557 rate of the *minimal* z-score $\min_{j \neq I^*} Z_{t,I^*,j}$. One can then think of $\Gamma(\theta)^{-1}$ as determining a
 558 fundamental limit on the rate at which an experimenter can gather evidence against all
 559 alternative arms. A peculiar feature of this complexity term is that actually optimizing over
 560 p as (10) prescribes would requiring knowing θ , which is circular as uncertainty about θ is
 561 point of experimenting in the first place. Nevertheless this complexity measure serves to
 562 produce a valid lower bound, as evidenced by the next proposition. The lower bound here
 563 applies ideas that has been known since Chernoff et al. [1959], but our proof specifically
 564 applies inequalities of Kaufmann et al. [2016].

565 **Proposition 2.** *If*

$$\mathbb{E} [c\tau + \Delta_\tau | \theta = \theta_0] \leq O(c \log(1/c)) \quad \text{for all } \theta_0 \in \Theta,$$

566 *as* $c \rightarrow 0$, *then*

$$\mathbb{E} [c\tau + \Delta_\tau | \theta = \theta_0] \geq \Gamma(\theta_0)[c + o(1)] \log(1/c) \quad \text{for all } \theta_0 \in \Theta. \quad (11)$$

567 The idea of this lower bound is that any algorithm that outperforms (11) on some instance
568 must attain a loss an *an order of magnitude* larger on some other instance. The result is shown,
569 essentially, by establishing that any algorithm with uniformly vanishing simple regret —
570 meaning $\mathbb{E}[\Delta_\tau \mid \theta = \theta_0] = o(1)$ for all θ_0 — must gather an expected number of samples
571 that scales as $\mathbb{E}[\tau \mid \theta = \theta_0] \geq \Gamma(\theta_0)(\log(1/c) + o(1))$.

572 E.4 Optimality of Context Independent Sampling frequencies

573 The lower bound above turns out to be tight. It is matched by adaptive algorithms that
574 learn as information is gathered to adjust their measurement proportions rapidly enough
575 toward proportions that attain the maximum in (10). As such, the form of the solution is of
576 particular importance. Here we show a striking simplification. The maximal information
577 rate in (10) can be attained by context independent allocation, which samples each arm with a
578 probability that is independent of context. One we reduce a context-independent allocations,
579 it is easy to characterize the solution in terms of the first-order necessary conditions of
580 optimality. Equations (12) and (13) are known for problems without contexts [Glynn and
581 Juneja, 2004].

582 **Lemma 4** (Optimality of context independent sampling frequencies). *Suppose \mathcal{X} is finite.*
583 *There exists a vector $p^* = p^*(\theta) \in \mathcal{S}$ such that the rule given by $p(x, i) = p_i^*$ for all $x \in \mathcal{X}$ attains*
584 *the supremum in (10). The vector p^* is the unique solution to the k nonlinear equations:*

$$\frac{\mu(\theta, I^*, w) - \mu(\theta, i, w)}{\sqrt{(p_{I^*}^*)^{-1} + (p_i^*)^{-1}}} = \frac{\mu(\theta, I^*, w) - \mu(\theta, j, w)}{\sqrt{(p_{I^*}^*)^{-1} + (p_j^*)^{-1}}} \quad \forall i, j \neq I^* \quad (12)$$

$$p_{I^*}^* = \sqrt{\sum_{i \neq I^*} (p_i^*)^2} \quad (13)$$

585 Then Equation (10) becomes

$$\Gamma(\theta)^{-1} = \frac{1}{2\|X_{\text{pop}}\|_A} \frac{(\mu(\theta, I^*, w) - \mu(\theta, i, w))^2}{(p_{I^*}^*)^{-1} + (p_i^*)^{-1}} \quad \forall i \neq I^*$$

586 where $A = \sigma^2 (\mathbb{E}[X_1 X_1^\top])^{-1}$.

587 We refer to equation (12) as imposing *information balance*. It essentially ensures that the
588 z-scores $Z_{t, I^*, j}$ grow at an equal rate for arms $j \neq I^*$, balancing the evidence against each
589 suboptimal arm.

590 E.5 Adaptive Tuning

591 We will show that DTS automatically gathers information in a manner that satisfies an
592 information balance property like Equation (12). By shifting measurement effort away
593 from clearly inferior arms and toward those that could more plausibly be the best arm, the
594 algorithm automatically balances the rate of information acquisition. The precise fraction of
595 measurement effort that (13) suggests should be assigned to the optimal arm is not satisfied

596 automatically, however. In order to do that, the tuning parameter β_t needs to be adjusted
 597 properly.

Algorithm 2: Adaptive Tuning Algorithm

Input posterior means of expected reward $(m_{t,i})_{i \in [k]}$.

if \hat{I}_t is not unique then

 | Set $\beta_t = \beta_{t-1}$

end

else

 Obtain the unique optimal solution $x \in \mathcal{S}$ of the empirical version of Equations (12)
 and (13) with $(\mu(\theta, i, w))_{i \in [k]}$ and I^* replaced by $(m_{t,i})_{i \in [k]}$ and \hat{I}_t , respectively:

598

$$\frac{m_{t,\hat{I}_t} - m_{t,i}}{\sqrt{x_{\hat{I}_t}^{-1} + x_i^{-1}}} = \frac{m_{t,\hat{I}_t} - m_{t,j}}{\sqrt{x_{\hat{I}_t}^{-1} + x_j^{-1}}} \quad \forall i, j \neq \hat{I}_t \quad (14)$$

$$x_{\hat{I}_t} = \sqrt{\sum_{i \neq \hat{I}_t} x_i^2} \quad (15)$$

 Set $\beta_t = x_{\hat{I}_t}$

end

599 **Efficient Implementation of the Tuning Algorithm** For each $i \in [k]$, we define $\Delta_{t,i} \triangleq$
 600 $m_{t,\hat{I}_t} - m_{t,i}$. Equation (14) implies there exists y such that

$$\frac{1 + x_{\hat{I}_t} x_i^{-1}}{\Delta_{t,i}^2} = y, \quad \forall i \neq \hat{I}_t.$$

601 Clearly $y > \max_{i \neq \hat{I}_t} \Delta_{t,i}^{-2}$ and

$$\frac{x_{\hat{I}_t}}{x_i} = \Delta_{t,i}^2 y - 1, \quad \forall i \neq \hat{I}_t. \quad (16)$$

602 Together with Equation (15), Equation (16) implies

$$\sum_{i \neq \hat{I}_t} \left(\Delta_{t,i}^2 y - 1 \right)^{-2} = 1.$$

603 We can solve this fixed-point equation for y using, for example, bisection search or Newton's
 604 method. Notice that if Newton's method is used, one may wish to save the value of y solved
 605 in the previous time period, which provides an effective initial point for finding an updated
 606 value of y . Finally $\sum_{i \in [k]} x_i = 1$ and Equation (16) imply

$$x_{\hat{I}_t} = \frac{1}{1 + \sum_{i \neq \hat{I}_t} \left(\Delta_{t,i}^2 y - 1 \right)^{-1}},$$

607 which is the value assigned to β_t .

608 **E.6 DTS Attains the Lower Bound**

609 We now show that when β_n is tuned as suggested in the previous section, and an appropriate
 610 stopping rule is employed, DTS matches the fundamental lower bound in (11).

611 We consider the empirical selection rule

$$\hat{I}_t \in \arg \max_i \hat{\mu}_{t,i} \quad (17)$$

612 that selects the arm with highest performance under a least-squares estimate. Similar results
 613 can be developed if the Bayes selection rule were used instead, which essentially uses
 614 ridge-regression rather than least-squares.

615 Developing stopping rules is itself an area of active research. We do not try to advance that
 616 literature, and instead focus on a very simple candidate that is sufficient for the results we
 617 wish to prove. Recall that the z-score $Z_{t,\hat{I}_t,j}$ measures the strength of evidence that arm \hat{I}_t
 618 outperforms arm j in the population. The stopping rule

$$\tau = \inf \left\{ t \in \mathbb{N} : \min_{j \neq \hat{I}_t} Z_{t,\hat{I}_t,j} \geq \gamma_t \right\} \quad \text{where} \quad \gamma_t = \Phi^{-1} \left(1 - \frac{c}{t^2 k} \right), \quad (18)$$

619 stops at the first time all z-scores exceed a threshold. The threshold was picked to ensure a
 620 probability of incorrect selection less than c . The specific choice of γ_t is based on a Bonferroni
 621 correction to account for multiple hypothesis testing and could likely be reduced through
 622 more granular analysis.

623 The next proposition gives two upper bounds.

624 **Proposition 3.** *Under the selection rule (17), the stopping rule (18), and allocation rule DTS with*
 625 *β_t defined by Algorithm 2, for any $\theta_0 \in \Theta$,*

$$\mathbb{E} [c\tau + \Delta_\tau \mid \theta = \theta_0] \leq \Gamma(\theta_0)[c + o(1)] \log(1/c) \quad \text{as } c \rightarrow 0.$$

626 *If instead the allocation rule is DTS with fixed $\beta = 1/2$, then for any $\theta_0 \in \Theta$,*

$$\mathbb{E} [c\tau + \Delta_\tau \mid \theta = \theta_0] \leq 2\Gamma(\theta_0)[c + o(1)] \log(1/c) \quad \text{as } c \rightarrow 0.$$

627 This shows that DTS with adaptively tuned $\{\beta_t\}$ attains the exact optimal constant defined in
 628 Equation (10), which matches the lower bound in (11). In addition, DTS with non-adaptive
 629 choice of $\beta = 1/2$ achieves near-optimal statistical guarantee while reduces computational
 630 cost.