
Bandit Learning with Implicit Feedback

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Implicit feedback, such as user clicks, although abundant in online information
2 service systems, does not provide substantial evidence on users' evaluation of
3 system's output. Such incomplete supervision inevitably misleads model estima-
4 tion, especially in a bandit learning setting where the feedback is acquired on the
5 fly. In this work, we study a contextual bandit problem with implicit feedback
6 by modeling the feedback as a composition of user result examination and rele-
7 vance judgment. Since users' examination behavior is unobserved, we introduce
8 latent variables to model it. We perform Thompson sampling on top of variational
9 Bayesian inference for arm selection and model update. Rigorous upper regret
10 bound analysis of the proposed algorithm proves its feasibility of learning from
11 implicit feedback; and extensive empirical evaluations on click logs collected from
12 a major MOOC platform further demonstrate its learning effectiveness in practice.

13 1 Introduction

14 Contextual bandit algorithms [3, 18, 17] provide modern information service systems an effective
15 solution to adaptively find good mappings between available items and users. This family of
16 algorithms sequentially select items to serve users using side information about user and item, while
17 adapting their selection strategies based on the immediate user feedback to maximize users' long-term
18 satisfaction. They have been popularly deployed in practical systems for content recommendation
19 [18, 4, 24] and display advertising [5, 20].

20 However, the most dominant form of user feedback in such systems is implicit feedback, such as
21 clicks, which is known to be biased and incomplete about users' evaluation of system's output
22 [15, 10]. For example, a user skips a recommended item might not be because he/she does not like
23 the item, but he/she just does not examine that display position, i.e., position bias [12]. Unfortunately,
24 a common practice in contextual bandit applications simply treats no click as a form of negative
25 feedback [18, 23, 5]. This introduces inconsistency to model update, since the skipped items might
26 not be truly irrelevant, and it inevitably leads to suboptimal outputs of bandit algorithms over time.

27 In this work, we focus on learning contextual bandits with user click feedback, and model such
28 implicit feedback as a composition of user result examination and relevance judgment. Examination
29 hypothesis [7], which is a fundamental assumption in click modeling, postulates that a user clicks on
30 a system's returned result *if and only if* that result has been examined by the user and it is relevant to
31 the user's information need at the moment. Because a user's examination behavior is unobserved,
32 we model it as a latent variable, and realize the examination hypothesis in a probabilistic model.
33 We define the conditional probabilities of result examination and relevance judgment via logistic
34 functions over the corresponding contextual features. To perform model update, we take a variational
35 Bayesian approach to develop closed form approximation to the posterior distribution of model
36 parameters on the fly. This approximation also paves the way for an efficient Thompson sampling
37 strategy for arm selection in bandit learning. Our rigorous finite time analysis proves that, despite
38 the increased complexity in parameter estimation introduced by the latent variable, our Thompson

39 sampling policy based on the variational lower bound is guaranteed to achieve a sub-linear Bayesian
40 regret with a high probability. We also prove that when one fails to model result examination in click
41 feedback, a linearly increasing regret is possible, as the model cannot differentiate examination driven
42 skip from relevance driven skip in negative feedback.

43 We tested the algorithm in a major Massive Open Online Course (MOOC) platform for personalized
44 education. To personalize students’ learning experience in this platform, we recommend quiz-like
45 questions in a form of banners on top of the lecture videos when students are watching the videos.
46 The algorithm needs to decide where in a video to display which question to a target student. If the
47 student feels the displayed question is helpful for him/her to understand the video lecture, he/she
48 could click on the banner to read the answer and more related online content about the question.
49 Therefore, our goal is to maximize the click through rate (CTR) on the selected questions. There are
50 several properties of this application that amplifies the bias and incompleteness of click feedback.
51 First, based on the consideration of user experience, to minimize the risk of annoying any student,
52 the displayed time of a banner is limited to a few seconds. Second, as this feature is newly introduced
53 to the platform, many users might not realize that they can click on the question to read more related
54 content about it. Hence, no click on a question does not necessarily indicate its irrelevance. We
55 tested the algorithm in this application in a four-month period, where a total of 500 questions are
56 manually compiled for the algorithm to select over 503 videos within more than 200 thousands
57 student video-watching sessions. Based on the unbiased offline evaluation policy [19], our algorithm
58 achieved a 7.4% CTR lift compared to standard contextual bandits [18, 8] which do not model users’
59 examination behavior. By sharing user examination model across videos, the improvement is further
60 boosted to over 8.8%.

61 2 Related Works

62 As having been extensively studied in click modeling of user search results [6], various factors affect
63 users’ click decisions, and among them result examination plays a central role [12, 7]. Unfortunately,
64 most applications of bandit algorithms simply treat user clicks as feedback for model update [18, 23,
65 5, 24], where no click on a selected result is considered as negative feedback. This inevitably leads to
66 inaccurate model update and sub-optimal arm selection.

67 There is a line of research that develops click model based bandit algorithms for learning to rank
68 problems. For example, by assuming that skipped documents are less attractive than later clicked
69 ones in a ranked list, Kveton et al. [16] develop a cascading bandit model to learn from both clicks
70 and skips in search results. To enable learning from multiple clicks in the same result ranking list,
71 they adopt the dependent click model [9] to infer user satisfaction after a sequence of clicks [13],
72 and later further extend to broader types of click models [25]. However, such algorithms aim at
73 estimating the best ranking of results in a per-query basis, without specifying any ranking function
74 nor generalizing to unseen queries. This limits their application scenario in practice.

75 Another line of related research is bandit learning with latent variables. Maillard and Mannor studied
76 the problem of latent bandit [21], which assumes reward distributions are clustered and the clusters
77 are determined by some latent variables. They only studied the problem in a context-free setting, and
78 a very weak performance guarantee is provided when the reward distribution is unknown in those
79 clusters. Kawale et al. developed a Thompson sampling scheme for online matrix-factorization [14].
80 Latent features are extracted via an online low-rank matrix completion based on samples selected
81 from Thompson sampling on the fly. Due to the ad-hoc combination of factorization method and
82 bandit method, little theoretical analysis is provided. Wang et al. studied the problem of latent feature
83 learning for contextual bandits [23]. They extended arms’ context vectors with latent features under a
84 linear reward structure, and applied the upper confidence bound principle over coordinate descent to
85 iteratively estimate the hidden features and model parameters. The linear reward structure prohibits it
86 from recognizing the nonlinear dependency between result examination and relevance judgment in
87 click feedback.

88 3 Problem Setup

89 We consider a contextual bandit problem with finite, but possibly large, number of arms \mathcal{A} . At each
90 trial $t = 1, \dots, T$, the learner observes a subset of candidate arms \mathcal{A}_t with $\mathcal{A}_t \subset \mathcal{A}$, where each arm

91 a is associated with a context vector \mathbf{x}^a summarizing the side information about the arm. Once an
 92 arm $a_t \in \mathcal{A}_t$ is chosen according to some policy π , corresponding implicit binary feedback C_{a_t} ,
 93 e.g., user click, will be given to the learner as the reward. The learner’s goal is to adjust its arm
 94 selection strategy to maximize its cumulative reward over time. What makes this problem unique and
 95 challenging is that C_{a_t} does not truly reflect users’ evaluation of the selected arm a_t . Based on the
 96 examination hypothesis, when $C_{a_t} = 1$, the choice a_t must be relevant to the user’s information need
 97 at time t ; but when $C_{a_t} = 0$, a_t might be relevant but the user just does not examine it. Unfortunately,
 98 the unobserved result examination condition is unobserved to the learner.

99 We model user examination through a binary latent variable E_{a_t} and assume that the context vector
 100 \mathbf{x}_t^a of arm a can be decomposed into $(\mathbf{x}_{C,t}^a, \mathbf{x}_{E,t}^a)$, where the dimension of $\mathbf{x}_{C,t}^a$ and $\mathbf{x}_{E,t}^a$ are d_C
 101 and d_E respectively. Accordingly, users’ result examination and relevance judgment decisions are
 102 assumed to be governed by a conjecture of $(\mathbf{x}_{C,t}^a, \mathbf{x}_{E,t}^a)$ and the bandit parameter $\boldsymbol{\theta}^* = (\boldsymbol{\theta}_C^*, \boldsymbol{\theta}_E^*)$. In
 103 the rest of this paper, when no ambiguity is introduced, we drop the index a to simplify the notations.
 104 As a result, we make the following generative assumption about an observed click C_t on arm a_t ,

$$\begin{aligned}\mathbb{P}(C_t = 1 | E_t = 0, \mathbf{x}_{C,t}) &= 0 \\ \mathbb{P}(C_t = 1 | E_t = 1, \mathbf{x}_{C,t}) &= \rho(\mathbf{x}_{C,t}^\top \boldsymbol{\theta}_C^*) \\ \mathbb{P}(E_t = 1 | \mathbf{x}_{E,t}) &= \rho(\mathbf{x}_{E,t}^\top \boldsymbol{\theta}_E^*)\end{aligned}$$

105 where $\rho(x) = \frac{1}{1+e^{-x}}$. Based on this assumption, we have $\mathbb{E}[C_t | \mathbf{x}_t] = \rho(\mathbf{x}_{C,t}^\top \boldsymbol{\theta}_C^*) \rho(\mathbf{x}_{E,t}^\top \boldsymbol{\theta}_E^*)$. As a
 106 result, the observed click feedback C_t is a sample from this generative process.

107 4 Algorithm

108 The learner needs to estimate the bandit parameters $\boldsymbol{\theta}_C^*$ and $\boldsymbol{\theta}_E^*$ based on its interactively obtained
 109 click feedback $\{C_i\}_{i=1}^t$ over time. Ideally, this estimation can be obtained by maximizing the data
 110 likelihood. However, the inclusion of examination as a latent variable to model implicit feedback in
 111 our bandit setting poses serious challenges to both parameter estimation and arm selection policy
 112 design. Neither conventional least square estimator nor maximum likelihood estimator can be
 113 easily obtained, let alone computational efficiency, due to the non-convexity of the corresponding
 114 optimization problem. Even worse, the two popular bandit learning paradigms, upper confidence
 115 bound principle [1] and Thompson sampling [2], both demand an accurate estimation of bandit
 116 parameters and its uncertainty. In this section, we present an elegant solution to tackle these
 117 two challenges, which makes use of variational Bayesian inference technique to learn parameters
 118 approximately on the fly, as well as to bridge parameter estimation and policy design.

119 4.1 Variational Bayesian for parameter estimation

120 To complete the generative process defined in Section 3, we further assume $\boldsymbol{\theta}_C$ and $\boldsymbol{\theta}_E$ follow
 121 Gaussian distribution $N(\hat{\boldsymbol{\theta}}_C, \boldsymbol{\Sigma}_C)$ and $N(\hat{\boldsymbol{\theta}}_E, \boldsymbol{\Sigma}_E)$ respectively. We are interested in developing
 122 a closed form approximation to their posteriors, when a newly obtained observation $(\mathbf{x}_C, \mathbf{x}_E, C)$
 123 becomes available. By applying Bayes’ rule in the log space, we have,

$$\begin{aligned}\log \mathbb{P}(\boldsymbol{\theta}_C, \boldsymbol{\theta}_E | \mathbf{x}_C, \mathbf{x}_E, C) &= \log \mathbb{P}(C | \boldsymbol{\theta}_C, \boldsymbol{\theta}_E, \mathbf{x}_C, \mathbf{x}_E) + \log \mathbb{P}(\boldsymbol{\theta}_C, \boldsymbol{\theta}_E) + \log \text{const} \\ &= C \log \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) + (1 - C) \log (1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) \\ &\quad - \frac{1}{2} (\boldsymbol{\theta}_C - \hat{\boldsymbol{\theta}}_C)^\top \boldsymbol{\Sigma}_C^{-1} (\boldsymbol{\theta}_C - \hat{\boldsymbol{\theta}}_C) - \frac{1}{2} (\boldsymbol{\theta}_E - \hat{\boldsymbol{\theta}}_E)^\top \boldsymbol{\Sigma}_E^{-1} (\boldsymbol{\theta}_E - \hat{\boldsymbol{\theta}}_E) + \log \text{const}\end{aligned}$$

124 The key idea is to develop a variational lower bound in the quadratic form of $\boldsymbol{\theta}_C$ and $\boldsymbol{\theta}_E$ for the
 125 log-likelihood function $\log \mathbb{P}(C | \boldsymbol{\theta}_C, \boldsymbol{\theta}_E, \mathbf{x}_C, \mathbf{x}_E)$. Because of the convexity of $\log \rho(x) - \frac{x}{2}$ with
 126 respect to x^2 (See Appendix B.1) and a Jensen’s inequality of $\log x$ (See Appendix B.2), the lower
 127 bound in required form is achievable. When $C = 1$, by Equation 15 in Appendix B.3, we have,

$$l_{C=1}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}) := \log (\rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) \geq g(\mathbf{x}_C^\top \boldsymbol{\theta}, \xi_C) + g(\mathbf{x}_E^\top \boldsymbol{\theta}, \xi_E) \quad (1)$$

128 where $g(x, \xi) := \frac{x}{2} - \frac{\xi}{2} + \log \rho(\xi) - \lambda(\xi)(x^2 - \xi^2)$, $\lambda(\xi) = \frac{\tanh \frac{\xi}{2}}{4\xi}$, $x, \xi \in \mathcal{R}$. More specifically,
 129 $g(x, \xi)$ is a polynomial of degree 2 with respect to x . When $C = 0$, by Equation 16 in Appendix B.3,

Algorithm 1 Thompson sampling for E-C Bandit

- 1: Initiate $\Sigma_C = \lambda I, \Sigma_E = \lambda I, \hat{\theta}_C = \theta_{C,0}, \hat{\theta}_E = \theta_{E,0}$.
 - 2: **for** $k = 0, 1, 2, \dots$ **do**
 - 3: Observe the available arm set $\mathcal{A}_k \subset \mathcal{A}$ and its corresponding context set $\mathcal{X}_k := \{(\mathbf{x}_C^a, \mathbf{x}_E^a) : a \in \mathcal{A}_k\}$.
 - 4: Randomly sample $\tilde{\theta}_C \sim N(\hat{\theta}_C, \Sigma_C), \tilde{\theta}_E \sim N(\hat{\theta}_E, \Sigma_E)$.
 - 5: Select:

$$a_k = \arg \max_{a \in \mathcal{A}_k} \rho((\mathbf{x}_C^a)^\top \tilde{\theta}_C) \rho((\mathbf{x}_E^a)^\top \tilde{\theta}_E)$$
 - 6: Play the selected arm a_k and Observe the reward C_k .
 - 7: Update $\Sigma_C, \hat{\theta}_C, \Sigma_E, \hat{\theta}_E$ according to Eq (3), (4), (5), (6) respectively.
 - 8: **end for**
-

130 we have,

$$l_{C=0}(\mathbf{x}_C, \mathbf{x}_E, \theta) := \log(1 - \rho(\mathbf{x}_C^\top \theta_C) \rho(\mathbf{x}_E^\top \theta_E)) \quad (2)$$

$$\geq H(q) + qg(-\mathbf{x}_C^\top \theta, \xi_C) + qg(\mathbf{x}_E^\top \theta, \xi_{E,1}) + (1-q)g(-\mathbf{x}_E^\top \theta, \xi_{E,2})$$

131 where $H(q) := -q \log q - (1-q) \log(1-q)$. Once the lower bound in the quadratic form is
 132 established, we can use a Gaussian distribution to approximate our target posterior, whose mean and
 133 covariance matrix are determined by the following equations:

$$\Sigma_{C,\text{post}} = \Sigma_C + 2q^{1-C} \lambda(\xi_C) \mathbf{x}_C \mathbf{x}_C^\top \quad (3)$$

$$\hat{\theta}_{C,\text{post}} = \Sigma_{C,\text{post}}^{-1} (\Sigma_C \hat{\theta}_C + \frac{1}{2} (-q)^{1-C} \mathbf{x}_C) \quad (4)$$

$$\Sigma_{E,\text{post}} = \Sigma_E + 2\lambda(\xi_E) \mathbf{x}_E \mathbf{x}_E^\top \quad (5)$$

$$\hat{\theta}_{E,\text{post}} = \Sigma_{E,\text{post}}^{-1} (\Sigma_E \hat{\theta}_E + \frac{1}{2} (2q-1)^{1-C} \mathbf{x}_E) \quad (6)$$

134 where the subscript ‘‘post’’ denotes the parameters in the Gaussian distributions that approximate the
 135 posteriors. Consecutive observations can be incorporated into the approximated posterior recursively.
 136 There is one thing left to decide, i.e., the choice of variational parameters (ξ_C, ξ_E, q) . A reasonable
 137 criterion is to choose the values such that likelihood on the observation is maximized. Similar to the
 138 choice made by [11], we choose the closed form update formulas of those variational parameters as:

$$\xi_C = \sqrt{\mathbf{E}_{\theta_C} [\mathbf{x}_C^\top \theta_C]^2}$$

$$\xi_E = \sqrt{\mathbf{E}_{\theta_E} [\mathbf{x}_E^\top \theta_E]^2}$$

$$q = \frac{\exp(g(\mathbf{x}_C^\top \theta_C, \xi_C) + g(\mathbf{x}_E^\top \theta_E, \xi_E) - g(-\mathbf{x}_E^\top \theta_E, \xi_E))}{1 + \exp(g(\mathbf{x}_C^\top \theta_C, \xi_C) + g(\mathbf{x}_E^\top \theta_E, \xi_E) - g(-\mathbf{x}_E^\top \theta_E, \xi_E))}$$

139 where all the expectations are taken under the approximated posteriors. Empirically, we found the
 140 iterative update of the approximated posterior and the variational parameters converge quite rapidly,
 141 such that it usually only needs several rounds of iterations in our experiments.

142 4.2 Thompson sampling with approximated lower bound

143 Thompson sampling, also known as probability matching, is widely used in bandit learning to balance
 144 exploration and exploitation, and shows great empirical performance [5]. Thompson sampling
 145 requires a distribution of the model parameters to sample from. In a standard Thompson sampling
 146 [2], one is required to sample from the true posterior of model parameters. But as logistic regression
 147 does not have a conjugate prior, the model defined in our problem does not have an exact posterior.
 148 We decide to sample from the approximated posterior. Later we will demonstrate this is a very tight
 149 posterior approximation in our experiment. Once the sampling of $(\tilde{\theta}_C, \tilde{\theta}_E)$ is complete, we can select
 150 the corresponding arm $a_k \in \mathcal{A}_k$ which maximizes $\rho(\mathbf{x}_C^\top \tilde{\theta}_C) \rho(\mathbf{x}_E^\top \tilde{\theta}_E)$. We name the resulting bandit
 151 algorithm as examination-click bandit, or E-C Bandit in short, and summarize it in Algorithm 1.

152 5 Regret Analysis

153 Define $f_{\theta}(\mathbf{x}) := \mathbb{E}[C|\mathbf{x}, \theta] = \rho(\mathbf{x}_C^T \theta_C) \rho(\mathbf{x}_E^T \theta_E)$. The accumulated regret of a bandit policy π up to
154 time T is formally defined as,

$$\text{Regret}(T, \pi, \theta^*) = \sum_{t=1}^T \max_{a \in \mathcal{A}_t} f_{\theta^*}(\mathbf{x}^a) - f_{\theta^*}(\mathbf{x}^{a_t})$$

155 where $\mathbf{x}^{a_t} := (\mathbf{x}_C^{a_t}, \mathbf{x}_E^{a_t})$ is the context vector of the arm $a_t \in \mathcal{A}_t$ selected by the policy π at time t
156 based on the history $\mathcal{H}_t := \{(\mathcal{A}_i, \mathbf{x}_i, C_i)\}_{i=1}^{t-1}$. The Bayesian regret is defined by $\mathbb{E}[\text{Regret}(T, \pi, \theta^*)]$,
157 where the expectation is taken with respect to the prior distribution over θ^* , and can be written as,

$$\text{BayesRegret}(T, \pi) = \sum_{t=1}^T \mathbb{E}[\max_{a \in \mathcal{A}_t} f_{\theta^*}(\mathbf{x}^a) - f_{\theta^*}(\mathbf{x}^{a_t})]$$

158 Our algorithm, which is based on a maximum likelihood estimator, is equivalent to an estimator that
159 minimizes a log-loss with binary random variables. In this section, we will first bound the aggregate
160 empirical discrepancy of the log-loss estimator used in our model in Proposition 1, which prepares for
161 the upper bound of the generic Bayesian regret under log-loss estimators with Thompson sampling
162 policy in Theorem 1. Based on this generic Bayesian regret bound under a log-loss estimator, we
163 study the upper bound of Bayesian regret for our proposed E-C Bandit. Due to space limit, we
164 provide all the detailed proofs in the Appendix.

165 To further simplify our notations, we use f for f_{θ} , which is the reward function based on a model's
166 estimated bandit parameter θ , and f_k for $f_{\theta}(\mathbf{x}^{a_k})$. We use f^* for f_{θ^*} , which is the reward function
167 based on the ground-truth bandit parameter, and f_k^* for $f_{\theta^*}(\mathbf{x}^{a_k})$. We assume that f^* lies in a known
168 function space \mathcal{F} , where any $f \in \mathcal{F}$ is a function mapping from the arm set \mathcal{A} to the range $(0, 1)$.
169 Define the log-loss estimator by $\hat{f}_t^{\text{LOGLOSS}} \in \arg \max_{f \in \mathcal{F}} L_{2,t}(f)$ where $L_{2,t}(f)$ is the aggregate
170 log-loss written as $\sum_{k=1}^{t-1} l_k(f)$ where $l_k(f) = -(C_k \log f_k + (1 - C_k) \log(1 - f_k))$. We have the
171 following lemma,

172 **Proposition 1.** Denote the aggregate empirical discrepancy $\sum_{k=1}^t (f_k - f_k^*)^2$ by $\|f - f^*\|_{E,t}$. For
173 all $\delta > 0$ and $\alpha > 0$, if $\mathcal{F}_t = \{f \in \mathcal{F} : \|f - \hat{f}_t^{\text{LOGLOSS}}\|_{E,t} \leq \sqrt{\beta_t^*(\mathcal{F}, \delta, \alpha)}\}$ for all $t \in N$,

$$\mathbb{P}(f^* \in \cap_{t=1}^{\infty} \mathcal{F}_t) > 1 - \delta, \quad (7)$$

174 where $\beta_t^*(\mathcal{F}, \delta, \alpha)$ is an appropriately constructed confidence parameter and is defined as
175 $\beta_t^*(\mathcal{F}, \delta, \alpha) := \frac{2}{\lambda_0} \log(N(\mathcal{F}, \alpha, \|\cdot\|_{\infty})/\delta) + 2\alpha\eta_t$, where $N(\mathcal{F}, \alpha, \|\cdot\|_{\infty})$ denotes the alpha-covering
176 number of \mathcal{F} , $\lambda_0 = \frac{3}{(\frac{1}{m_f} + \frac{1}{1-M_f})^2}$ and $\eta_t = (4M_f + \frac{1}{\min\{m_f, 1-M_f\}})t$, in which $m_f, M_f \in \mathcal{R}$ are
177 upper and lower bounds of f such that $0 < m_f \leq f \leq M_f < 1$ for any $f \in \mathcal{F}$.

178 **Remark 1.** The proof is provided in Appendix C. Here we discuss two important issues re-
179 lated to our later proof about E-C Bandit's regret. First, the precise optimization of $\hat{f}_t^{\text{LOGLOSS}} \in$
180 $\arg \max_{f \in \mathcal{F}} L_{2,t}(f)$ could be hard to achieve in some instances of \mathcal{F} , for example, when \mathcal{F} is a
181 set of non-convex functions. Nevertheless, we can always resort to approximation method to solve
182 the optimization problem as long as the approximation error can be bounded. Indeed, in our E-C
183 Bandit, we resort to variational inference to estimate $\hat{f}_t^{\text{LOGLOSS}}$ on the fly and find it works quite well
184 in practice. Second, when $f^* \notin \mathcal{F}$, this corresponds to the problem of mis-specification of model. In
185 this situation, the regret bound could be very poor, and the real regret could be linear with respect to
186 time. To show this clearly in our case, in Appendix F we constructed a situation in which the regret is
187 inevitably linear if one fails to model the examination condition in click feedback and simply treats
188 no click as negative feedback.

189 With Proposition 1, we have the following Theorem which bounds the Bayesian regret of the
190 Thompson Sampling strategy under a log-loss estimator.

191 **Theorem 1.** For all $T \in N$, $\alpha > 0$ and $\delta < \frac{1}{2T}$, if π^{TS} denotes the policy derived from the log-loss
192 estimator and a Thompson sampling strategy along the time steps, then

$$\text{BayesRegret}(T, \pi^{TS}) \leq 1 + (\dim_E^A(\mathcal{F}, T^{-1}) + 1)C + 4\sqrt{\dim_E^A(\mathcal{F}, T^{-1})\beta_T^*(\mathcal{F}, \alpha, \delta)T} \quad (8)$$

193 where $C = \sup_{f \in \mathcal{F}} \{\sup f\}$, $\dim_E^A(\mathcal{F}, T^{-1})$ is the eluder dimension (see Definition 3 in Russo and
194 Van Roy [22]) of \mathcal{F} with respect to \mathcal{A} .

195 **Remark 2.** We can choose $C = 1$ in our click feedback case since $f \in (0, 1)$. C is kept in the
 196 theorem to show the same form compared to the Proposition 8 in Russo and Van Roy [22]. In fact,
 197 the proof is almost the same once we have Proposition 1. Hence, we omit the proof in our paper.

198 Now we turn to provide an upper regret bound of our E-C Bandit, based on the above generic Bayesian
 199 regret analysis under a log-loss estimator. We add the following assumptions which are standard in
 200 the literature of contextual bandits.

201 **Assumption 1.** The optimal θ^* lies in $\mathcal{B}_s := \{\theta \in \mathcal{R}^d : \|\theta\|_2 \leq s\}$, and s is known as a prior.

202 **Assumption 2.** The norm of context vectors are bounded by x , i.e., $(\mathbf{x}_C, \mathbf{x}_E) \in \mathcal{B}_x$, where $\mathcal{B}_x :=$
 203 $\{\mathbf{x} \in \mathcal{R}^d : \|\mathbf{x}\|_2 \leq x\}$ and x is known as a prior.

204 Based on these two assumptions, it is straightforward to verify that $\rho(\mathbf{x}_C^\top \theta_C)$, $\rho(\mathbf{x}_E^\top \theta_E)$ and $f_\theta(\mathbf{x})$
 205 are bounded. Let $M_\rho = \max_{\theta \in \mathcal{B}_s, \mathbf{x} \in \mathcal{B}_x} \rho(\mathbf{x}^\top \theta)$ and $m_\rho = \min_{\theta \in \mathcal{B}_s, \mathbf{x} \in \mathcal{B}_x} \rho(\mathbf{x}^\top \theta)$. Hence, $0 <$
 206 $m_\rho \leq M_\rho < 1$. Similarly, denoting the maximum of $f_\theta(\mathbf{x})$ by M_f and the minimum of it by m_f , we
 207 have $0 < m_f \leq M_f < 1$. Once the arm set is restricted to a finite cardinality, $\dim_{\mathcal{A}}(\mathcal{F}, T^{-1}) \leq |\mathcal{A}|$
 208 by Appendix C.1. in Russo and Van Roy [22]. Choosing the function class as that in our E-C Bandit,
 209 i.e., $\mathcal{F} = \{f : \mathcal{B}_x \rightarrow \mathcal{R} \mid f = \rho(\mathbf{x}_C^\top \theta_C) \rho(\mathbf{x}_E^\top \theta_E), \theta \in \mathcal{B}_s\}$, by Lemma 8 (See Appendix 8 for its
 210 proof), we have $N(\mathcal{F}, \alpha, \|\cdot\|_\infty) = (\gamma/\alpha)^d$ where $\gamma = 2M_\rho k_\rho x$ (k_ρ is the Lipschitz constant of ρ ,
 211 see Lemma 4). Hence, choosing $\alpha = 1/t^2$ and $\delta = 1/t$ leads to

$$\beta_t^*(\mathcal{F}, 1/t, 1/t^2) = \frac{2}{\lambda_0} d \log(\gamma t^3) + \frac{1}{t} (4M_f + \frac{1}{m_f}). \quad (9)$$

212 Therefore, the upper bound of Bayesian regret of our proposed E-C Bandit takes the following form,

$$\text{BayesRegret}(T, \pi^{TS}) = O(|\mathcal{A}| + \sqrt{d|\mathcal{A}|T \log T}). \quad (10)$$

213 When $T \gg |\mathcal{A}|$ and $T \gg d$, which is the typical case in practice, it becomes $O(\sqrt{T \log T})$.

214 6 Experiment

215 We perform empirical comparisons in simulation and on the student click logs collected from our
 216 MOOC platform to verify the effectiveness of our proposed algorithm, against those that fail to model
 217 the examination condition and directly use click as feedback.

218 6.1 Algorithms for comparison

219 We list the models used for empirical comparisons below. We briefly explain why we choose them
 220 as baseline models and how we adjust them in our scenario. **Logistic Bandit.** This model has been
 221 extensively used for advertisement CTR optimization. In [5, 19], the authors model user clicks
 222 by a regularized logistic regression model over observed context features and makes decisions by
 223 Thompson sampling. In particular, no click is treated as negative feedback. We also used the Laplace
 224 approximation and Gaussian prior presented in [5] to update the model parameters on the fly. We also
 225 want to highlight that despite of mountains of works focusing on generalized linear bandits, most
 226 of them are not truly online algorithms, because the estimation of parameters at each iteration has
 227 to involve all historical observations iteratively, which incurs a time complexity at least $O(t)$ (e.g.,
 228 Filippi et al. [8] require exact optimum of logistic regression on all historical observations at each
 229 round). **E-C Bandit.** This is the algorithm we present in Algorithm 1. We highlight here that in the
 230 experiment on real-world data, the manual separation of examination feature \mathbf{x}_E and click feature
 231 \mathbf{x}_C in the context vector \mathbf{x} offers a principled approach to incorporate one's domain knowledge about
 232 what affects user examination and what affects user result relevance into the model. We explain in
 233 detail what features are chosen for which in Appendix G. Thanks to the approximation achieved by
 234 Bayesian variational inference presented in section 4, truly online model update is feasible in this
 235 algorithm. This provides both computational and storage efficiency. **hLinUCB.** This is the bandit
 236 algorithm proposed by Wang et al. [23]. It is related to our model in a sense that both models learn
 237 hidden features. In particular, hLinUCB extends the linear contextual bandit by inclusion of hidden
 238 features and operates under a UCB-like strategy. However, it still treats click as direct feedback, but
 239 aims at learning more expressive features to describe the observed clicks.

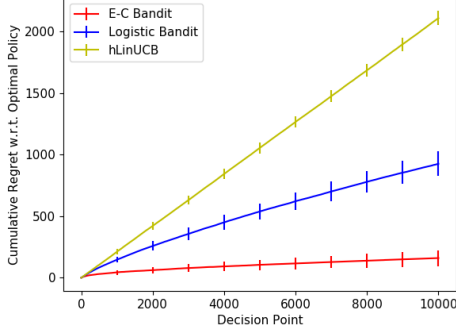


Figure 1: Cumulative regret over 100 simulations.

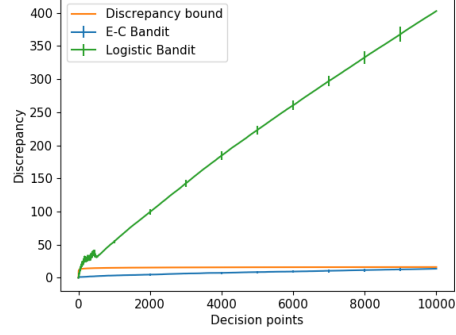


Figure 2: Discrepancy bound given by Proposition 1

240 6.2 Experiments on simulations

241 First we demonstrate the effectiveness of our algorithms by experiment with simulated data. The
 242 experiment is set as follows. The context vector’s dimension d_C and d_E are set to 5, and thus
 243 $d = d_C + d_E = 10$. We set $|\mathcal{A}| = 100$, each of which is associated with a unique context
 244 vector $(\mathbf{x}_C, \mathbf{x}_E)$. The ground-truth parameter $(\boldsymbol{\theta}_C^*, \boldsymbol{\theta}_E^*)$ and the specific value of $(\mathbf{x}_C, \mathbf{x}_E)$ are all
 245 randomly sampled from the unit ball $\mathcal{B} = \{\mathbf{x} \in \mathcal{R}^d : \|\mathbf{x}\|_2 \leq 1\}$. Since $(\boldsymbol{\theta}_C^*, \boldsymbol{\theta}_E^*)$ and $(\mathbf{x}_C, \mathbf{x}_E)$
 246 are both sampled from \mathcal{B} , m_f and M_f can be obtained by taking the minimum and maximum of
 247 $\rho(\mathbf{x}_C^T \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^T \boldsymbol{\theta}_E)$ on \mathcal{B} , respectively, i.e., $m_f = \frac{1}{(1+e)^2}$ and $M_f = \frac{1}{(1+e^{-1})^2}$.

248 At each time t , an available arm set $\tilde{\mathcal{A}}_t$ is randomly sampled from \mathcal{A} such that $|\tilde{\mathcal{A}}_t| = 10$, i.e., each
 249 time we offer 10 randomly selected arms for the algorithm to choose from. Any algorithm selects
 250 an arm from \mathcal{A}_t and observes the corresponding reward C_t^{alg} generated by the Bernoulli distribution
 251 $B(\rho(\mathbf{x}_{C,t}^T \boldsymbol{\theta}_C^*) \rho(\mathbf{x}_{E,t}^T \boldsymbol{\theta}_E^*))$. The actual regret for this algorithm at time t is defined by the received
 252 click, i.e., $\text{regret}(t) = C_{a_t^*} - C_{a_t}$, where a_t^* is the optimum arm to be chosen based on ground-truth
 253 bandit parameters $(\boldsymbol{\theta}_C^*, \boldsymbol{\theta}_E^*)$.

254 We repeat the experiment 100 times using the same simulation setting, each containing 10,000
 255 iterations. The average cumulative regret over 100 runs and the corresponding standard deviation
 256 (plotted per thousand decision point) are illustrated in Figure 1. One can clearly notice that while
 257 the Logistic Bandit suffers from a linear regret with respect to time t , as it mistakenly treats no
 258 click as negative feedback, our E-C Bandit achieves a fast converging sub-linear regret. The result
 259 that hLinUCB performs the worst is expected, since it assumes a linear relation between click and
 260 context feature vectors. We further investigate how the aggregate empirical discrepancy of E-C bandit
 261 and Logistic Bandit increase with respect to time. Figure 2 illustrates that the aggregate empirical
 262 discrepancy of E-C Bandit is well bounded by the upper bound given by Proposition 1, while the
 263 Logistic Bandit’s aggregate empirical discrepancy increase linearly. This directly explains their
 264 accumulative regret in this experiment comparison.

265 6.3 Experiments on MOOC video watching data

266 The MOOC data we used for evaluation is collected from a single course in a 4-month period. The
 267 course has 503 lecture videos in total. About 500 high-quality quiz-like questions has been manually
 268 crafted and each video is assigned with a subset of them based on human-judged relatedness and
 269 fitness. On average a video has 7 questions, each of which is associated with six possible displaying
 270 positions within the video, leading to a total of 56 arms in average.

271 We picked one video as an example to analyze the students’ click behavior. 11 arms are randomly
 272 picked and projected to a plane according to their relative distance in Figure 3. The number in
 273 the parenthesis is the empirical CTR of the corresponding arm. It can be clearly seen that while
 274 arm b and arm j have the same empirical CTR, the arms between them, such as are e or d , have
 275 lower CTR. Logistic Bandit is never able to capture this non-monotonicity relation, since its reward
 276 prediction increases monotonically with respect to a linear predictor. We construct a more general
 277 case in Appendix F to illustrate the scenario that failing to model examination condition would

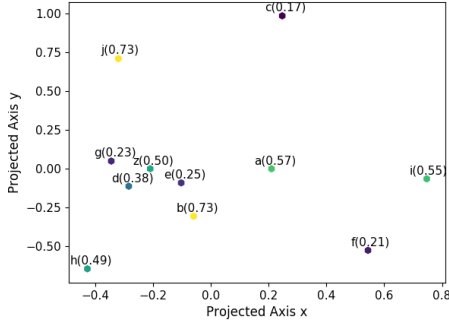


Figure 3: 11 arms’ feature vectors projected onto a plane, such that the relative distances between points are kept. The number in the parenthesis is the arm’s relative CTR.

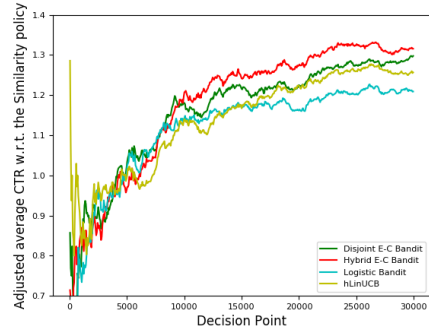


Figure 4: Performance comparison on MOOC videos’ data

278 lead to a linear regret. In reality, arm b and arm j are two different questions displayed at the same
 279 position in the video, while arm b and arm e are the same question displayed at different positions.
 280 This phenomenon strongly suggests a bias in users’ implicit feedback, which again justifies our
 281 decomposition of examination and relevance in click feedback.

282 We followed [19] to develop our online data collection policy in our MOOC platform to prepare
 283 our evaluation data set. In particular, any related questions with respect to a video will have an
 284 equal probability to be selected and displayed at all positions in this video. We name this policy
 285 as *Similarity*. We create an instance of a bandit model for each video to learn its own optimal
 286 question placing policy. Further, we follow [18] to create a hybrid model of our E-C Bandit across
 287 multiple videos. We make the examination component in E-C Bandit, i.e., θ_E , shared by all videos.
 288 Our underlying assumption is that while how the relevance quality of each question to the host
 289 video affects student clicks may vary across videos, the students’ examination decision should be
 290 mostly affected by the platform or their own watching habits, e.g., our user interface design. (See
 291 Appendix G for detailed explanations of our examination feature choice.) We denote the non-hybrid
 292 algorithm by disjoint E-C Bandit and the hybrid one by Hybrid E-C Bandit. It is worth noting that
 293 there is no straightforward way to construct hybrid algorithm for Logistic Bandit, as it only has one
 294 set of parameter for modeling click. And hLinUCB is inherently hybrid since different videos share
 295 information through the learned hidden features, i.e., a collaborative filtering scheme.

296 Li et al. [19] proposed a method to calculate a near-unbiased estimate of the average CTR by any
 297 bandit algorithm based on the collected history data, so that offline evaluation and performance
 298 comparison are possible. We take their offline evaluation protocol here and report the estimated
 299 average CTR in Figure 4. To avoid disclosing any proprietary information about the platform, all
 300 algorithms’ CTRs are normalized by that from the *Similarity* policy. As shown in the figure,
 301 independently estimating E-C Bandits across videos achieves an average 29.9% increase in CTR over
 302 the *Similarity* baseline, while the joint estimation further boosts it to 31.6%. These two versions of
 303 E-C Bandit consistently outperform other two baseline bandits, i.e., hLinUCB and Logistic Bandit.
 304 This clearly suggests the necessity of modeling examination condition in user clicks for improving
 305 the online recommendation performance.

306 7 Conclusion

307 Motivated by the examination hypothesis in user click modeling, in this paper we developed E-C
 308 Bandit, which differentiates result examination and content relevance in user clicks and actively
 309 learns from such implicit feedback. We developed an efficient and effective learning algorithm
 310 based on variational inference and demonstrated its effectiveness on both simulated and real-world
 311 datasets. We proved that despite the complexity of underlying reward generation assumption and the
 312 resulting parameter estimation procedure, the proposed learning algorithm enjoys a sub-linear regret
 313 bound. Currently we only studied click feedback on single items; it is important for us to study it in a
 314 more general setting, e.g., a list of ranked items, where sequential result examination and relevance
 315 judgment introduce richer inter-dependency.

316 References

- 317 [1] Yasin Abbasi-yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear
318 stochastic bandits. In *NIPS*, pages 2312–2320. 2011.
- 319 [2] Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear
320 payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- 321 [3] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*
322 *Learning Research*, 3:397–422, 2002. URL [http://www.jmlr.org/papers/v3/auer02a.](http://www.jmlr.org/papers/v3/auer02a.html)
323 [html](http://www.jmlr.org/papers/v3/auer02a.html).
- 324 [4] Djallel Bouneffouf, Amel Bouzeghoub, and Alda Lopes Gançarski. A contextual-bandit
325 algorithm for mobile context-aware recommender system. In *Neural Information Processing*,
326 pages 324–331. 2012.
- 327 [5] Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. In *Advances*
328 *in neural information processing systems*, pages 2249–2257, 2011.
- 329 [6] Aleksandr Chuklin, Ilya Markov, and Maarten de Rijke. Click models for web search. *Synthesis*
330 *Lectures on Information Concepts, Retrieval, and Services*, 7(3):1–115, 2015.
- 331 [7] Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. An experimental comparison of
332 click position-bias models. In *Proceedings of the 2008 international conference on web search*
333 *and data mining*, pages 87–94. ACM, 2008.
- 334 [8] Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The
335 generalized linear case. In *Advances in Neural Information Processing Systems*, pages 586–594,
336 2010.
- 337 [9] Fan Guo, Chao Liu, and Yi Min Wang. Efficient multiple-click models in web search. In
338 *Proceedings of the Second ACM International Conference on WSDM*, pages 124–131. ACM,
339 2009.
- 340 [10] Yifan Hu, Yehuda Koren, and Chris Volinsky. Collaborative filtering for implicit feedback
341 datasets. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on*, pages
342 263–272. Ieee, 2008.
- 343 [11] Tommi S Jaakkola and Michael I Jordan. Bayesian parameter estimation via variational methods.
344 *Statistics and Computing*, 10(1):25–37, 2000.
- 345 [12] Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, and Geri Gay. Accurately
346 interpreting clickthrough data as implicit feedback. In *ACM SIGIR Forum*, volume 51, pages
347 4–11. Acm, 2017.
- 348 [13] Sumeet Katariya, Branislav Kveton, Csaba Szepesvari, and Zheng Wen. Dcm bandits: Learning
349 to rank with multiple clicks. In *International Conference on Machine Learning*, pages 1215–
350 1224, 2016.
- 351 [14] Jaya Kawale, Hung H Bui, Branislav Kveton, Long Tran-Thanh, and Sanjay Chawla. Efficient
352 thompson sampling for online matrix-factorization recommendation. In *NIPS*, pages 1297–1305,
353 2015.
- 354 [15] Diane Kelly and Jaime Teevan. Implicit feedback for inferring user preference: a bibliography.
355 In *Acm Sigir Forum*, volume 37, pages 18–28. ACM, 2003.
- 356 [16] Branislav Kveton, Csaba Szepesvari, Zheng Wen, and Azin Ashkan. Cascading bandits:
357 Learning to rank in the cascade model. In *International Conference on Machine Learning*,
358 pages 767–776, 2015.
- 359 [17] John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side
360 information. In *NIPS*, pages 817–824, 2008.

- 361 [18] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
362 personalized news article recommendation. In *Proceedings of the 19th international conference*
363 *on World wide web*, pages 661–670. ACM, 2010.
- 364 [19] Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of
365 contextual-bandit-based news article recommendation algorithms. In *Proceedings of the fourth*
366 *ACM international conference on Web search and data mining*, pages 297–306. ACM, 2011.
- 367 [20] Wei Li, Xuerui Wang, Ruofei Zhang, Ying Cui, Jianchang Mao, and Rong Jin. Exploitation
368 and exploration in a performance based contextual advertising system. In *Proceedings of 16th*
369 *SIGKDD*, pages 27–36. ACM, 2010.
- 370 [21] Odalric-Ambrym Maillard and Shie Mannor. Latent bandits. In *International Conference on*
371 *Machine Learning*, pages 136–144, 2014.
- 372 [22] Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics*
373 *of Operations Research*, 39(4):1221–1243, 2014.
- 374 [23] Huazheng Wang, Qingyun Wu, and Hongning Wang. Learning hidden features for contextual
375 bandits. In *Proceedings of the 25th ACM International on Conference on Information and*
376 *Knowledge Management*, pages 1633–1642. ACM, 2016.
- 377 [24] Yisong Yue and Carlos Guestrin. Linear submodular bandits and their application to diversified
378 retrieval. In *NIPS*, pages 2483–2491, 2011.
- 379 [25] Masrour Zoghi, Tomas Tunys, Mohammad Ghavamzadeh, Branislav Kveton, Csaba Szepesvari,
380 and Zheng Wen. Online learning to rank in stochastic click models. In *International Conference*
381 *on Machine Learning*, pages 4199–4208, 2017.

382 **Appendix A Preliminaries**

383 In this section, we present some basic definitions and inequalities for later use.

384 First, let $p, q \in (0, 1)$, $H(p) := -p \log p - (1 - p) \log(1 - p)$ and $KL(p|q) := p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$ denote the entropy and the KL-divergence function of Bernoulli variable respectively.
 386 Note that for later derivations' convenience, the log here is with respect to e instead of 2, which is slightly different from a standard definition.
 387

388 The first inequality is Pinsker's inequality in Bernoulli case.

389 **Lemma 1** (Pinsker's inequality). $KL(p|q) - 2(p - q)^2 \geq 0$.

390 *Proof.* Let $c(q) := KL(p|q) - 2(p - q)^2$. We have $c'(q) = (p - q)(\frac{1}{q(1-q)} - 4)$. Since $q \in (0, 1)$,
 391 $q(1 - q) \leq 1/4$ by the arithmetic average inequality. Thus, when $q < p$, $c'(q) < 0$; when $q > p$,
 392 $c'(q) > 0$. Hence, $c(q) \geq c(p) = 0$. \square

393 The next inequality provides a Lipschitz condition-like equality for log function defined on bounded
 394 set.

395 **Lemma 2.** Given $a, b \in (0, 1)$, for any $x, y \in [a, b]$, $\frac{1}{b}|x - y| \leq |\log x - \log y| \leq \frac{1}{a}|x - y|$.

Proof. According to the mean value theorem, for any $x, y \in [a, b]$, there exists a $\eta \in [x, y]$ such that

$$\log x - \log y = \frac{1}{\eta}(x - y).$$

396 The inequality follows by taking absolute value on two side and the fact that η is bounded in $[a, b]$. \square

397 The next is Hoeffding's lemma.

Lemma 3 (Hoeffding's lemma). Let X be any real-valued random variable such that $E(X) = 0$ and $a \leq X \leq b$ almost surely. Then, for all $\lambda \in \mathcal{R}$,

$$\mathbb{E}[e^{\lambda X}] \leq \lambda^2(b - a)^2/8.$$

398

399 We omit its proof here, since its proof should be available on any standard textbook about the theory
 400 of probability. The last inequality deals with the Lipschitz condition of the logistic function.

401 **Lemma 4.** For the logistic function $\rho(x) = \frac{1}{1+e^{-x}}$, we have:

$$|\rho(x) - \rho(y)| < k_\rho|x - y|$$

402 where $k_\rho = 1/4$.

Proof. According to the mean value theorem, there exists a $\eta \in (x, y)$ such that

$$\rho(x) - \rho(y) = \rho(\eta)(1 - \rho(\eta))(x - y)$$

403 where we have used $\dot{\rho}(x) = \rho(x)(1 - \rho(x))$. Since $\rho(\eta) \in (0, 1)$, we have $\rho(\eta)(1 - \rho(\eta)) \leq 1/4$.
 404 Thus,

$$|\rho(x) - \rho(y)| = |\rho(\eta)(1 - \rho(\eta))| \times |(x - y)| \leq k_\rho|x - y|.$$

405 \square

406 **Appendix B Variational lower bound for E-C Bandit**

407 The variational lower bound in the quadratic form of our E-C bandit is constructed from two basic
 408 variational lower bounds. They are the lower bound to a logistic regression and the variational lower
 409 bound to a log-sum function. In this section, we first provide these two basic lower bounds in Section
 410 B.1 and Section B.2 respectively. We end this section by providing the desired lower bound for our
 411 E-C bandit in Section B.3.

412 **B.1 Variational lower bound for the log-logistic function**

413 Jaakkola and Jordan [11] provide for the log-logistic function a variational lower bound taking the
 414 form of a 2-degree polynomial based on the observation that $\log \rho(x) - \frac{x}{2}$ is convex with respect to
 415 x^2 . We report their result here. Defining the related functions as follows:

$$\begin{aligned}\rho(x) &:= \frac{1}{1 + e^{-x}} (x \in \mathcal{R}) \\ \lambda(\xi) &:= \frac{\tanh \frac{\xi}{2}}{4\xi} \\ g(x, \xi) &:= \frac{x}{2} - \frac{\xi}{2} + \log \rho(\xi) - \lambda(\xi)(x^2 - \xi^2) (x, \xi \in \mathcal{R})\end{aligned}$$

416 we have the variational lower bound:

$$\log(\rho(x)) \geq g(x, \xi) \quad (11)$$

417 When $\xi^2 = x^2$, the lower bound is exact. Note that $1 - \rho(x) = \rho(-x)$. Thus:

$$\log(1 - \rho(x)) \geq g(-x, \xi) \quad (12)$$

418 The equation holds when $\xi^2 = x^2$.

419 **B.2 Variational lower bound for a log-sum function**

420 We are interested in developing a bound for the following log-sum function, which is exactly the
 421 likelihood function of a single sample with reward $C = 0$ in E-C Bandit:

$$\log \left((1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) + 1 - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \right). \quad (13)$$

422 Since $\log(x)$ is convex, we have the following inequality holds for any $q \in (0, 1)$:

$$\begin{aligned}\log \left((1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) + 1 - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \right) \\ \geq H(q) + q \log \left((1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \right) + (1 - q) \log \left(1 - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \right)\end{aligned} \quad (14)$$

423 where $H(q) = -q \log q - (1 - q) \log(1 - q)$. It can be easily verified using Jensen's inequality for a
 424 convex function.

425 The equality holds whenever $q = \frac{\rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)}{1 - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)}$, which always lies in the domain of q ,
 426 i.e., $(0, 1)$, since $\rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \in (0, 1)$ and $\rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) \in (0, 1)$.

427 **B.3 Variational lower bound for the log-likelihood of E-C Bandit**

428 We provide the variational lower bound of a quadratic form for the log-likelihood of a single
 429 observation in E-C Bandit in this section. When $C = 1$, the log-likelihood of a single sample is:

$$l_{C=1}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E) = \log(\rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E))$$

430 Since $\log(x)$ is additive, given Equation 11, we have:

$$\begin{aligned}l_{C=1}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E) &= \log(\rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)) + \log(\rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) \\ &\geq g(\mathbf{x}_C^\top \boldsymbol{\theta}_C, \xi_C) + g(\mathbf{x}_E^\top \boldsymbol{\theta}_E, \xi_E) \\ &=: \tilde{l}_{C=1}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E, \xi_C, \xi_E)\end{aligned} \quad (15)$$

431 where $=:$ means "denoted by". When $C = 0$, the log-likelihood of a single sample is a log-sum
 432 function, which is the first lower-bounded by Equation 14. Then each log-logistic function is
 433 lower-bounded by Equation 11 and Equation 12. The derivation is as follows:

$$\begin{aligned}l_{C=0}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E) &= \log(1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) \\ &= \log((1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E) + 1 - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) \\ &\geq H(q) + q \log((1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) + (1 - q) \log(1 - \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)) \\ &\geq H(q) + qg(-\mathbf{x}_C^\top \boldsymbol{\theta}_C, \xi_C) + qg(\mathbf{x}_E^\top \boldsymbol{\theta}_E, \xi_{E,1}) + (1 - q)g(-\mathbf{x}_E^\top \boldsymbol{\theta}_E, \xi_{E,2})\end{aligned} \quad (16)$$

434 where $\xi_C, \xi_{E,1}, \xi_{E,2}$ are the variational parameter of $\log(1 - \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C)), \log(\rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E)), \log(1 -$
435 $\rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E))$ respectively. Note that the equality holds when $\xi_{E,1}^2 = (\mathbf{x}_E^\top \boldsymbol{\theta}_E)^2$ and $\xi_{E,2}^2 = (-\mathbf{x}_E^\top \boldsymbol{\theta}_E)^2$.
436 We can always impose $\xi_{E,1} = \xi_{E,2}$ without harming the attainability of equality in the lower-bound.
437 Thus, we can get a simplified lower-bound:

$$\begin{aligned} l_{C=0}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E) &\geq H(q) + qg(-\mathbf{x}_C^\top \boldsymbol{\theta}_C, \xi_C) + qg(\mathbf{x}_E^\top \boldsymbol{\theta}_E, \xi_E) + (1-q)g(-\mathbf{x}_E^\top \boldsymbol{\theta}_E, \xi_E) \\ &=: \tilde{l}_{C=0}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E, \xi_C, \xi_E, q) \end{aligned} \quad (17)$$

438 We can unify Equation 15 and Equation 17 as:

$$\tilde{l}_C(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E) := C\tilde{l}_{C=1}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E) + (1-C)\tilde{l}_{C=0}(\mathbf{x}_C, \mathbf{x}_E, \boldsymbol{\theta}_C, \boldsymbol{\theta}_E, \xi_C, \xi_E, q)$$

439 which is a quadratic form of $(\boldsymbol{\theta}_C, \boldsymbol{\theta}_E)$ as desired (Plug in the definition of $g(x)$ to see it clearly).

440

441 Appendix C Proof of confidence bound: Proposition 1

442 Proposition 1 is established by a union lower bound of aggregate empirical discrepancy of any
443 function $f \in \mathcal{F}$ (Lemma 5) with respect to the optimal f^* and a discretization error bound based on
444 an α -covering of the function space \mathcal{F} (Lemma 6). In this section, we first present the proof of these
445 two lemmas first, and end this section by a proof of Proposition 1 based on these two lemmas.

446 C.1 Bound of log-loss estimator by aggregate empirical discrepancy

447 We lower-bound the log-loss of any function $f \in \mathcal{F}$ in terms of the empirical log-loss of the true
448 function f^* and the aggregate empirical discrepancy $\|f - f^*\|_{E,t} := \sum_{i=1}^t (f_k - f_k^*)^2$. The result is
449 stated as the following lemma.

450 **Lemma 5.** For all $\delta > 0$ and $f \in \mathcal{F}$, with probability at least $1 - \delta$,

$$\mathbb{P}\left(L_{2,t+1}(f) \geq L_{2,t+1}(f^*) + \frac{1}{2}\|f - f^*\|_{E,t} - \frac{1}{\lambda_0} \log \frac{1}{\delta}\right) \geq 1 - \delta \quad (18)$$

451 holds simultaneously for all natural number $t \in N$, where $\lambda_0 := \frac{3}{(\frac{1}{m_f} + \frac{1}{1-M_f})^2}$.

452 We first recall some definitions and properties that are helpful here. Let $\mathcal{H}_t^x = \mathcal{H}_t \cup \{\mathcal{A}_t, \mathbf{x}_t\}$ be a
453 new filtration. First, $C_k \in \{0, 1\}$ is a binary variable whose expectation conditioned on \mathcal{H}_k^x is f_k^* .
454 Hence, the random variable $\zeta_k := C_k - f_k^*$ is sub-Gaussian according to Lemma 7 such that:

$$\begin{aligned} \mathbb{E}[\zeta_k | \mathcal{H}_k^x] &= 0 \\ \mathbb{E}[e^{\lambda \zeta_k} | \mathcal{H}_k^x] &\leq \frac{\lambda^2}{2} \end{aligned}$$

455 The log-loss of a single sample is defined as $l_k(f) = -(C_k \log f_k + (1 - C_k) \log(1 - f_k))$. The
456 KL-divergence of two binary variables whose expectation are p and q respectively is defined by
457 $KL(p|q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$, and satisfied Pinsker's Inequality(Lemma 1):

$$KL(p|q) \geq 2(p-q)^2$$

458 *Proof.* We add the following definition for simplicity of notation:

$$\begin{aligned} h_k &:= \log f_k - \log f_k^* + \log(1 - f_k^*) - \log(1 - f_k) \\ Z_k &:= l_k(f^*) - l_k(f) \\ \Phi_k(\lambda) &:= \log \mathbb{E}[e^{\lambda(Z_k - E(Z_k))} | \mathcal{H}_k^x]. \end{aligned}$$

459 By definition, $\sum_{i=1}^t Z_k = L_{2,t+1}(f^*) - L_{2,t+1}(f)$. Plugging in the equalities $C_k = \zeta_k + f_k^*$ and
460 $l_k(f) = -(C_k \log f_k + (1 - C_k) \log(1 - f_k))$, we can get that:

$$Z_k = h_k \zeta_k - KL(f_k^* | f_k). \quad (19)$$

461 Therefore, $\mathbb{E}[Z_k | \mathcal{H}_k^x] = -KL(f_k^* | f_k)$ and $\Phi_k(\lambda) = \log \mathbb{E}[e^{\lambda h_k \zeta_k} | \mathcal{H}_k^x]$. Hence, by Martingale
 462 Exponential Inequality (See Appendix B.1. in Russo and Van Roy [22]), for all $x \geq 0, \lambda \geq 0$,

$$\mathbb{P}\left(\lambda \sum_{i=1}^t Z_k \leq x - \lambda \sum_{i=1}^t KL(f_k^* | f_k) + \sum_{i=1}^t \log \mathbb{E}[e^{\lambda h_k \zeta_k} | \mathcal{H}_k^x]\right) \geq 1 - e^{-x}. \quad (20)$$

463 Since $KL(p|q) \geq 2(p-q)^2$ according to Lemma 1, and $\mathbb{E}[e^{\lambda \zeta_k} | \mathcal{H}_k^x] \leq \frac{\lambda^2}{2}$ according to Lemma 7,
 464 we have,

$$\mathbb{P}\left(\lambda \sum_{i=1}^t Z_k \leq x - 2\lambda \sum_{i=1}^t (f_k - f_k^*)^2 + \sum_{i=1}^t \lambda^2 h_k^2 / 2\right) \geq 1 - e^{-x}. \quad (21)$$

465 Again, according to Lemma 2, $|h_k| \leq |\log f_k - \log f_k^*| + |\log(1 - f_k^*) - \log(1 - f_k)| \leq (\frac{1}{m_f} +$
 466 $\frac{1}{1 - M_f})|f_k - f_k^*|$. Hence,

$$\mathbb{P}\left(\lambda \sum_{i=1}^t Z_k \leq x - 2\lambda \sum_{i=1}^t (f_k - f_k^*)^2 + \sum_{i=1}^t \lambda^2 \left(\frac{1}{m_f} + \frac{1}{1 - M_f}\right)^2 (f_k - f_k^*)^2 / 2\right) \geq 1 - e^{-x}. \quad (22)$$

467 Plugging in the definition of aggregate empirical discrepancy and Z_t , and rearranging the terms we
 468 get,

$$\mathbb{P}\left(L_{2,t+1}(f^*) - L_{2,t+1}(f) \leq \frac{x}{\lambda} + \left(\left(\frac{1}{m_f} + \frac{1}{1 - M_f}\right)^2 \frac{\lambda}{2} - 2\right) \|f - f^*\|_{E,t}\right) \geq 1 - e^{-x}. \quad (23)$$

469 Choosing $\lambda = \lambda_0 := \frac{3}{\left(\frac{1}{m_f} + \frac{1}{1 - M_f}\right)^2}$ and $x := \log \frac{1}{\delta}$ results in

$$\mathbb{P}\left(L_{2,t+1}(f) \geq L_{2,t+1}(f^*) + \frac{1}{2} \|f - f^*\|_{E,t} - \frac{1}{\lambda_0} \log \frac{1}{\delta}\right) \geq 1 - \delta. \quad (24)$$

470 Note that $\lambda_0 > 0$ thus satisfies the assumption of Martingale Exponential Inequality's assumption.
 471 Also note that the deduction above makes no assumption about the value of $t \in N$, thus the desired
 472 lemma is proved. \square

473 C.2 Discretization error

474 **Lemma 6.** If f^α satisfies $\|f - f^\alpha\|_\infty \leq \alpha$, and $m_f > 0, M_f < 1$, we have

$$\left|\frac{1}{2} \|f^\alpha - f^*\|_{E,t} - \frac{1}{2} \|f - f^*\|_{E,t} + L_{2,t}(f) - L_{2,t}(f^\alpha)\right| \leq \alpha \eta_t \quad (25)$$

475 where $\eta_t := (4M_f + \frac{1}{\min\{m_f, 1 - M_f\}})t$.

476 *Proof.* It is sufficient to consider a single sample's discretization error and then sum them up over
 477 time t . Thus we omit the subscript t in the following derivation.

478 We have

$$\begin{aligned} |(f^\alpha - f^*)^2 - (f - f^*)^2| &= |(f^\alpha - f)(f^\alpha + f) + 2f^*(f - f^\alpha)| \\ &\leq |(f^\alpha - f)(f^\alpha + f)| + |2f^*(f - f^\alpha)| \\ &\leq 2M_f \alpha + 2M_f \alpha \\ &= 4M_f \alpha \end{aligned} \quad (26)$$

479 and

$$\begin{aligned} |l(f) - l(f^\alpha)| &= |C(\log f^\alpha - \log f) + (1 - C)(\log(1 - f^\alpha) - \log(1 - f))| \\ &\leq C|\log f^\alpha - \log f| + |(1 - C)(\log(1 - f^\alpha) - \log(1 - f))| \\ &\leq \frac{C}{m_f} |f^\alpha - f| + \frac{1 - C}{1 - M_f} |f^\alpha - f| \\ &\leq \frac{1}{\min\{m_f, 1 - M_f\}} \alpha \end{aligned} \quad (27)$$

480 Summing over t and using triangular inequality, we get the desired result. \square

481 **C.3 Proof of Proposition 1**

Proof. Let \mathcal{F}^α be an α -covering of \mathcal{F} with respect to the sup norm such that for any $f \in \mathcal{F}$, there exists an $f^\alpha \in \mathcal{F}^\alpha$ such that $\|f - f^\alpha\|_\infty \leq \alpha$. By a union bound to Lemma 5, with probability at least $1 - \delta$,

$$L_{2,t}(f^\alpha) - L_{2,t}(f^*) \geq \frac{1}{2}\|f^\alpha - f^*\|_{E,t} - \frac{1}{\lambda_0} \log(|\mathcal{F}^\alpha|/\delta) \quad \forall t \in N, f^\alpha \in \mathcal{F}^\alpha.$$

482 Thus, with probability at least $1 - \delta$ for all $t \in N$ and $f \in \mathcal{F}$,

$$\begin{aligned} L_{2,t}(f) - L_{2,t}(f^*) &\geq \frac{1}{2}\|f - f^*\|_{E,t} - \frac{1}{\lambda_0} \log(|\mathcal{F}^\alpha|/\delta) \\ &\quad + \underbrace{\min_{f^\alpha \in \mathcal{F}^\alpha} \left\{ \frac{1}{2}\|f^\alpha - f^*\|_{E,t} - \frac{1}{2}\|f - f^*\|_{E,t} + L_{2,t}(f) - L_{2,t}(f^\alpha) \right\}}_{\text{Discretization error}} \end{aligned} \quad (28)$$

483 Since $\hat{f}_t^{\text{LOGLOSS}} \in \arg \min_{f \in \mathcal{F}} L_{2,t}(f)$ and $f^* \in \mathcal{F}$, we have

$$L_{2,t}(\hat{f}_t^{\text{LOGLOSS}}) - L_{2,t}(f^*) \leq 0. \quad (29)$$

484 Using the two-side bound of the discretization error term established by Lemma 6, we find with
485 probability at least $1 - \delta$,

$$\frac{1}{2}\|\hat{f}_t^{\text{LOGLOSS}} - f^*\|_{E,t} \leq \frac{1}{\lambda_0} \log(|\mathcal{F}^\alpha|/\delta) + \alpha\eta_t \quad (30)$$

Taking the infimum over the size of α -covers implies

$$\|f^* - \hat{f}_t^{\text{LOGLOSS}}\|_{E,t} \leq \sqrt{\beta_t^*(\mathcal{F}, \delta, \alpha)}$$

486 where $\beta_t^*(\mathcal{F}, \delta, \alpha) := \frac{2}{\lambda_0} \log(N(\mathcal{F}, \alpha, \|\cdot\|_\infty)/\delta) + 2\alpha\eta_t$. □

487 **Appendix D The sub-Gaussian property of $C_t - f_t^*$**

488 In this section, we prove $C_t - f_t^*$ is a sub-Gaussian random variable.

Lemma 7. Let $\mathcal{H}_t^x = \mathcal{H}_t \cup \{\mathcal{A}_t, \mathbf{x}_t\}$ be a new filtration constructed on top of \mathcal{H}_t , where $\mathcal{H}_t := \{(\mathcal{A}_i, \mathbf{x}_i, C_i)\}_{i=1}^{t-1}$ and C_t is a binary random variable such that $\mathbb{E}[C_t | \mathbf{x}_t] = f_t^*$. The random variable $\zeta_t := C_t - f_t^*$ is σ -sub-Gaussian conditioned on \mathcal{H}_t^x . That is to say, we have

$$\mathbb{E}[\zeta_t | \mathcal{H}_t^x] = 0$$

and there exists a $\sigma > 0$ such that

$$\mathbb{E}[e^{\lambda\zeta_t} | \mathcal{H}_t^x] \leq e^{\frac{\lambda^2\sigma^2}{2}}$$

489 for any $\lambda \in \mathcal{R}$.

Proof. Since $\mathbb{E}[C_t | \mathbf{x}_t] = f_t^*$ and $f_t^* \in \mathcal{H}_t^x$, we have $\mathbb{E}[\zeta_t | \mathcal{H}_t^x] = \mathbb{E}[C_t - f_t^* | \mathcal{H}_t^x] = 0$. Since ζ_t is bounded in $(-1, 1)$, by Hoeffding's lemma,

$$\mathbb{E}[e^{\lambda\zeta_t} | \mathcal{H}_t^x] \leq e^{\frac{\lambda^2}{2}}.$$

490 Thus conditioned on \mathcal{H}_t^x , ζ_t is σ -sub-Gaussian, where one can choose $\sigma = 1$. □

491 **Appendix E Alpha-covering number of function space**

492 We provide an α -covering number for the function space used in E-C Bandit.

Lemma 8. Let $\boldsymbol{\theta} := [\boldsymbol{\theta}_C^\top, \boldsymbol{\theta}_E^\top]^\top \in \mathcal{B}_s$, $\mathbf{x} := [\mathbf{x}_C^\top, \mathbf{x}_E^\top]^\top \in \mathcal{B}_x$, where $\mathcal{B}_s := \{\boldsymbol{\theta} \in \mathcal{R}^d : \|\boldsymbol{\theta}\|_2 \leq s\}$ and $\mathcal{B}_x := \{\mathbf{x} \in \mathcal{R}^d : \|\mathbf{x}\|_2 \leq x\}$. Let $\mathcal{F} = \{f : \mathcal{B}_x \rightarrow \mathcal{R} | f = \rho(\mathbf{x}_C^\top \boldsymbol{\theta}_C) \rho(\mathbf{x}_E^\top \boldsymbol{\theta}_E), [\boldsymbol{\theta}_C^\top, \boldsymbol{\theta}_E^\top]^\top \in \mathcal{B}_s\}$. Let $\gamma = 2M_\rho k_\rho x$. We have that

$$N(\mathcal{F}, \alpha, \|\cdot\|_\infty) = (\gamma/\alpha)^d$$

493 holds for any $\alpha > 0$.

494 *Proof.* For any $\theta_1, \theta_2 \in \mathcal{B}_s$, we have

$$\begin{aligned}
& \|f(\theta_1) - f(\theta_2)\| \\
&= \left\| \rho(\mathbf{x}_C^\top \theta_{C,1}) \rho(\mathbf{x}_E^\top \theta_{E,1}) - \rho(\mathbf{x}_C^\top \theta_{C,1}) \rho(\mathbf{x}_E^\top \theta_{E,2}) + \rho(\mathbf{x}_C^\top \theta_{C,1}) \rho(\mathbf{x}_E^\top \theta_{E,2}) - \rho(\mathbf{x}_C^\top \theta_{C,2}) \rho(\mathbf{x}_E^\top \theta_{E,2}) \right\| \\
&\leq \left\| \rho(\mathbf{x}_C^\top \theta_{C,1}) \rho(\mathbf{x}_E^\top \theta_{E,1}) - \rho(\mathbf{x}_C^\top \theta_{C,1}) \rho(\mathbf{x}_E^\top \theta_{E,2}) \right\| + \left\| \rho(\mathbf{x}_C^\top \theta_{C,1}) \rho(\mathbf{x}_E^\top \theta_{E,2}) - \rho(\mathbf{x}_C^\top \theta_{C,2}) \rho(\mathbf{x}_E^\top \theta_{E,2}) \right\| \\
&\leq M_\rho (\left\| \rho(\mathbf{x}_E^\top \theta_{E,1}) - \rho(\mathbf{x}_E^\top \theta_{E,2}) \right\| + \left\| \rho(\mathbf{x}_C^\top \theta_{C,1}) - \rho(\mathbf{x}_C^\top \theta_{C,2}) \right\|) \\
&\leq M_\rho k_\rho \left\| \mathbf{x}_E^\top (\theta_{E,1} - \theta_{E,2}) \right\| + \left\| \mathbf{x}_C^\top (\theta_{C,1} - \theta_{C,2}) \right\| \\
&\leq 2M_\rho k_\rho x \|\theta_1 - \theta_2\| \\
&= \gamma \|\theta_1 - \theta_2\|
\end{aligned} \tag{31}$$

495 where all norms are infinity norm. Thus a α -covering of \mathcal{F} can be achieved by a (α/γ) -covering of
496 \mathcal{B}_s . Evenly divide \mathcal{B}_s in each dimension, we get the desired result that $N(\mathcal{F}, \alpha, \|\cdot\|_\infty) = (\gamma/\alpha)^d$,
497 where we omit the ceiling function for simplicity. \square

498 Appendix F Linear regret of a mis-specified model

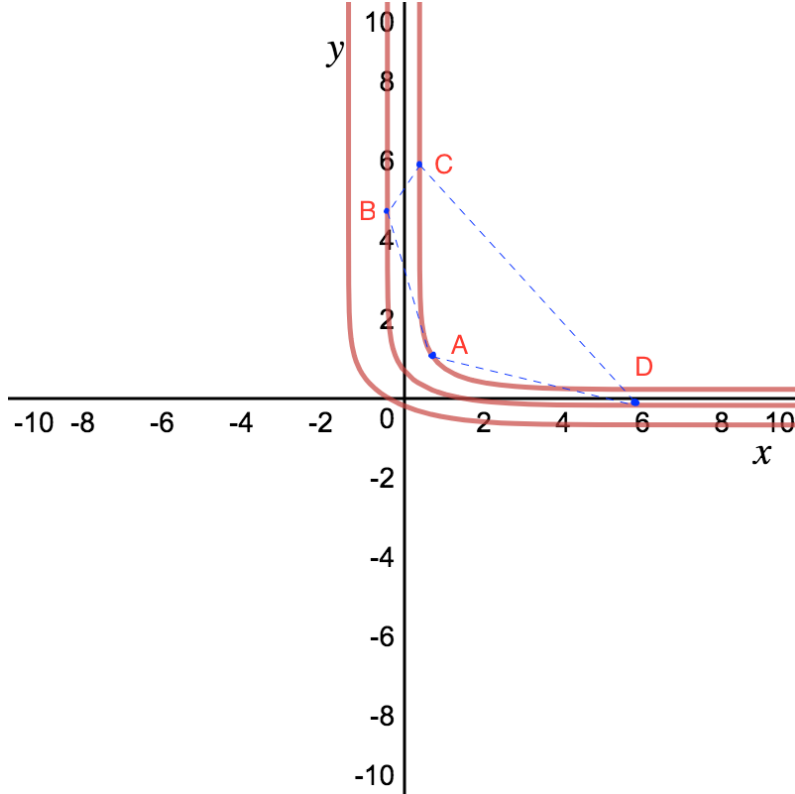


Figure 5: A contour plot of $f(x, y) = \rho(-x)\rho(-2y)$. Four arms are labeled by A, B, C, D such that the four points forms a convex quadrilateral and both pairs of opposite vertexes have the same function value(reward). One can always find such four points due to the non-convexity of f .

499 In this section, we construct an example when a Logistic Bandit suffers a linear regret under the
500 examination's hypothesis. It is enough to consider a simple case with $d_C = d_E = 1$ and $|\mathcal{A}| = 4$.
501 Figure 5 plots a simple but representative case when $\theta = [1, 2]^T$. One can think of the x -axis as the
502 context feature of result examination and y -axis as the context feature of relevance judgment. Denote
503 the reward at arm a by R_a , and the corresponding coordinate by (x_a, y_a) , where $a \in \{A, B, C, D\}$.
504 Then the figure indicates that

$$R_A = R_C > R_B = R_D. \tag{32}$$

505 Note that for a Logistic Bandit, which does not consider the decomposition of implicit feedback thus
 506 combines x and y in a linear predictor, will not be able to learn the correct order between those arms,
 507 i.e., Equation 32. This can be easily verified through a necessary condition. In fact, assume there is
 508 $[a, b]^T$ such that

$$\begin{aligned} ax_A + by_A &= ax_C + by_C \\ ax_B + by_B &= ax_D + by_D \end{aligned} \tag{33}$$

509 which is a necessary condition for a Logistic Bandit to learn Equation 32 exactly. As long as
 510 $x_B < x_A < x_C < x_D$ and $y_A < y_B < y_D < y_C$ (the case we plot in the figure), the determinant
 511 of coefficient matrix is less than 0 so that the only solution is $a = b = 0$. In this case, Logistic
 512 Bandit fails to differentiate A, C from B, D , and would suffer a regret of $R_A - R_B$ with probability
 513 $\frac{1}{2}$ when ties break evenly. Without loss of generality, we assume $a \neq 0$. If $a > 0, b \geq 0$,
 514 then $ax_D + by_D > ax_A + by_A$, thus a regret of $R_A - R_D$ is inevitable when the available arm set
 515 is $\{A, D\}$, since logistic function is monotonically increasing with respect to its linear predictor.
 516 So is case when $a < 0, b \leq 0$ with the inequality $ax_B + by_B > ax_C + by_C$. The cases when
 517 $a > 0, b < 0$ with the inequality $ax_D + by_D > ax_C + by_C$, and when $a < 0, b > 0$ with the
 518 inequality $ax_B + by_B > ax_A + by_A$ follow the same logic. Thus, we have constructed a case with
 519 an adversary environment when a linear regret is inevitable for a Logistic Bandit lacking in the
 520 decomposition of implicit feedback.

521 Appendix G Features used for the modeling of click and examination

522 We list the features we use for the modeling of examination and click in E-C Bandit in Table 1.
 523 Table 2 gives detailed description of these features. The main criteria of separating features between
 examination and relevance is that whether this feature is question-related or not.

Table 1: Feature usage in E-C Bandit

Feature	Used for examination?	Used for click(relevance)?
Playing speed	✓	×
Device	✓	×
Position	✓	×
Content Related	×	✓

Table 2: Feature description

Features	Description
Playing speed	Students are able to play the video at five different speeds. They are 1, 1.25, 1.5, 1.75 and 2 times the default speed, respectively. The faster the playing speed, the less probable the examination.
Device	A student may watch a video using different devices, such as a PC, iPad or smart cellphones. Different device should have different probability of examination.
Position	The question may be displayed at the beginning of the video, or somewhere middle of the video. Different CTRs for the same question at different positions are observed in real data, which we assume is caused by examination.
Content Related	A question's semantic meaning and its similarity with the subtitle of a video. Distributed embedding or tf-idf score may be used to gain a numeric representation of these features.

524