

1 We would like to thank the reviewers for their insightful reviews. We are glad that the feedback was generally positive,  
2 as the reviewers were happy with the motivation and practicality of our work.

3 The primary weakness that several reviewers brought up was that the methods and analysis were straightforward. We  
4 believe that our novelty is in proposing Thompson sampling latent bandit algorithms using offline-learned graphical  
5 models (PGMs) as side information, with full regret analysis. Reviewers are correct in noting that our analysis relies  
6 on insights made earlier by Russo and Van Roy. However, we are the first to apply it to bandit problems with latent  
7 variables, which is very common in many applications e.g., personalized recommender systems. Our near-optimal  
8 analysis is very general and we are working on applying it to problems with more complicated and problem-dependent  
9 PGMs (e.g., PGMs with latent state transition dynamics or with factored, continuous, latent state structures).

10 **Reviewer #1** “*algorithm ... is not designed keeping short horizons in mind*”: Our algorithms quickly personalize by  
11 assuming users can be clustered into a finite set of latent states. When the set is small, identifying which state best  
12 describes the user’s preferences can happen much quicker than learning the user’s preferences from scratch.

13 **Reviewer #2** “*suffers an exploration-exploitation tradeoff*”: You are correct in noting that our algorithm depends on  
14 exploration of the offline data to learn good models. If the models are good though, our method achieves much greater  
15 sample-efficiency than baselines that ignore the offline data.

16 “*unified analyses cannot cover instance-dependent bounds*”: We derive Bayes regret bounds, which contain an expectation  
17 over possible instances. We agree that instance-dependent regret bounds are more informative, but more difficult to  
18 derive and interesting future work. Our view our work as a first step to achieving this.

19 **Reviewer #3** Thank you for your detailed corrections! We will update the paper with your clarifications.

20 “*the available epsilon-bounds are wildly pessimistic*”: You are correct in noting that our regret bounds require that the  
21 offline-learned model has low prediction error  $\varepsilon$ . Choosing  $\varepsilon$  using model-learning guarantees will likely lead to an  
22 overly conservative mmUCB algorithm. However, we view mmUCB not as a practical algorithm, but one that can be used to  
23 analyze mmTS, which is practical and general. The mmTS algorithm additionally computes posterior model parameters  
24 using online interactions. This means the  $\varepsilon$  term should not affect its long-term performance, and the regret bound we  
25 derived is likely too conservative. Updating our regret bounds to reflect this is a future line of work.

26 “*no attempt to analyze the problem in computational terms*”: Given a latent variable model, the optimal policy would plan  
27 out the entire sequence of actions using the model, akin to solving a special case of a POMDP. It is unclear how to do  
28 so feasibly in our setting, even on short horizons. Gittins index will compute an optimal strategy for Bayesian bandits,  
29 but we are also unsure how to generalize the method to complex latent models. We instead compare our algorithms to a  
30 “post hoc clairvoyant” algorithm that performs at least as well

31 *Relation to value-of-information problem*: A key distinction is that prior value-of-information work seems more  
32 concerned with discovering the optimal solution, similar to best-arm-identification, whereas our setting deals with  
33 regret minimization. We will cite such work by Krause et. al and discuss its differences in our work.

34 **Reviewer #4** “*regret bounds for the misspecified model case ... are linear in n*”: It is true that the performance of  
35 our algorithms depend on the quality of the offline-learned model. Past work in offline model-learning e.g. spectral  
36 methods, can give guarantees that are  $\varepsilon = O(1/\sqrt{n_{\text{offline}}})$  where  $n_{\text{offline}}$  is the size of the offline dataset. On short  
37 horizons  $n \leq n_{\text{offline}}$ , the contribution of  $\varepsilon$  to overall regret is small. It is also important to note that our proposed mmTS  
38 algorithm computes posterior model parameters given online interactions. This means the  $\varepsilon$  term due to the prior is too  
39 conservative. Having a regret bound that reflects this property remains as future work.

40 *Questions about latent states*: We consider the latent bandits setting in Maillard and Mannor where the underlying  
41 latent state is fixed, and comes from a finite set of known size. In practice, the number of latent states could be be tuned  
42 during offline learning via cross validation. Our algorithms can also work with non-parametric latent models, though  
43 parametric ones are more well-studied in literature and have recovery guarantees.

44 *Questions about MovieLens experiment*: For each episode/user, actions are movies to recommend, and context is the  
45 concatenation of feature vectors for each movie; the context was learned from the training set and known beforehand.  
46 The mean reward for recommending movie  $i$  to user  $j$  is the dot product between movie  $i$  and user  $j$ ’s feature vectors,  
47 both of which derived from the test set and not known to the learning method. Our hypothesis was that using a Gaussian  
48 mixture-model (GMM) over user features fitted on the train set would improve performance by allowing for quickly  
49 associating each user’s hidden preferences with a cluster in the GMM. We will clarify the ambiguities in our work.

50 “*Exp4 performance seems a bit surprising*”: Exp4 can solve the latent bandits setting if we interpret the conditional  
51 model for each latent state as an expert; however, Exp4 is designed for adversarial rather than stochastic settings, which  
52 is why our algorithms greatly outperformed Exp4 using the same GMM model.