1 We thank reviewers for detailed comments and suggestions. We will address all comments in the revision. In this work,
2 we consider a novel unsupervised stochastic contextual bandits problem. In follow-on work, we will study relevant
3 open questions like, lower bounds, private information, and real-valued feedback, pointed out by reviewers.

4 **Novelty.** The novelty of our work, as such, is a combination of novel modeling principles to account for unsupervised
5 contextual sequential selection, as well as subsequent method and analysis, and experimentation. The earlier work
6 (Verma et al. AIStats'19) considered the problem of learning an optimal action but ignored the contextual information.
7 In this work, we incorporated the contextual information, which is readily available in many applications. Exploiting
8 the *real-valued* contextual information (features) for improving the arm selection strategy is non-trivial due to the
9 unsupervised nature of the problem where the standard analysis of contextual bandits does not apply. We made necessary
10 modeling assumptions leveraging GLM models and extended the existing definitions to address the learnability issues
11 in the new setup. However, the problem still requires new ideas and analysis methods to derive an efficient algorithm
12 and poses new technical challenges for analysis.

13 **Response to common comments of Reviewers $2$ and $4$:**
14 `The idea might look incremental.  What are the main challenges solved by this work?  :`
15 We respectfully disagree. The new ideas and technical challenges addressed in this work are as follows:

16 **New ideas:** USS-UCB (Verma et al. AIStats'19) uses two-sided test derived from Eq. (6) and Eq. (7) to identify an
17 optimal arm. In contrast, our work identifies that a one-sided test is enough to learn an optimal arm with contextual
18 information. We exploit this idea to come up with simpler algorithms. The difference in tests used in earlier and our
19 work becomes apparent by comparing arm selection strategy in USS-PD (line 9) with that of USS-UCB (lines 7-9)
20 where two sets $\hat{B}_t^h$ and $\hat{B}_t^l$ and their interaction needs to be computed. However, this simplification in USS-UCB throws
21 some challenges in the analysis that do not arise with two-sided test, but we carefully handle it (see lines 488-492).

22 **Technical Challenges**
23   1. GLM bandits are well studied but require reward or loss information. In the USS setup, loss of selected arm can
24      not be observed; hence finding the optimal arm is challenging. We have shown that if problem instance satisfies
25      contextual weak dominance (CWD) property, then the pairwise disagreement between arms can be used to estimate
26      context-dependent disagreement probability, and that can be used to find an optimal arm for a given context.
27   2. Regret analysis of GLM bandits hinges on bounding the instantaneous regret in each round, which is tied to the
28      estimation error of the GLM parameters. Due to the unsupervised setting and cascade structure, this way of
29      bounding regret does not work in our setup. Our analysis goes by bounding the number of pulls of the sub-optimal
30      arms. However, unlike standard bandits, we have to distinguish whether the sub-optimal arm pulled by USS-PD is
31      on the 'left' or 'right' of the optimal arm in the cascade. It requires our analysis to carefully handle both the cases
32      (see Lemma 5 and 6). Since USS-PD uses a similar MLE estimator for parameter estimation as in GLM bandits,
33      we only adapt their asymptotic normality results, the other steps of bounds are new in our work.
34   3. Though it is not reported in our work, we did try several other models and analysis approaches to solve the USS
35      problem with contextual information. However, due to the weak feedback structure of the problem, the other
36      methods are not amenable for analysis. Our final presentation is a model and analysis that is clean and complete.
37      For example, in the appendix, we point out that analysis based on Optimism in the Face of Uncertainty for Linear
38      bandits (OFUL) method with a regularizer did not go through without making more assumptions.

39 **Response to Reviewer $1$:**
40 `Lemma 1 should be equality, and does not require` $j > i$: Thanks for catching the typo. We will fix it.

41 **Response to Reviewer $2$:**
42 `Experiments lacks comparison with baselines and other SOTA methods:`
43 To the best of our knowledge, we are first to consider the contextual USS problem, so there is no state-of-the-art (SOTA)
44 method. In our experiments, we have considered baseline policies that select either a fixed arm or a uniformly random
45 arm in each round (see Figure 1c). We also compare USS-PD's performance with the policy that can observe the true
46 loss (see Figure 1b).

47 **Response to Reviewer $3$:**
48 `In equation (2),` $\lambda_{I_t}$ `is missing:` Thank you for catching the mistake. We will correct it.

49 **Response to Reviewer $4$:**
50 `In [1], the exactly same weak dominance property is introduced and Theorem 1 is same:`
51 We **disagree** on both. In this paper, the WD property is context-dependent, whereas the contextual information is
52 ignored in [1]. This new definition has nuances. First, Eq. 4 points to the fact that examples can be partitioned
53 based on strength of CWD, a situation that does not arise in [1]. Additionally, we allow for instances to violate CWD
54 property, and present algorithms that are agnostic to the presence of these instances. As for Theorem 1, although it
55 bears similarities with the previous work but it requires adaptation to the contextual setting.