

1 We thank the reviewers for the detailed comments. We will address all the minor issues and do not discuss them  
2 individually here. We will also add high-level pictures and proof sketches in the revision.

3 **Novelty and Originality:** We would like to first describe briefly the key differences between our work and the existing  
4 methods.

5 (1) [17, 18] are the only prior work on OCSM where up to  $K = T^{3/2}$  gradients are required per iteration. Note that  $K$  is  
6 not a constant, but a function of the total number of rounds  $T$ . Thus reducing  $K$  to 1 is an important step. To do so, we  
7 proposed a series of novel methods including the blocking procedure and the permutation methods (L114-131). Besides  
8 these, as noted in L137-144, a novel error analysis was performed although we relied on the same averaging technique  
9 proposed as in [37, 38].

10 (2) The extension from online setting to bandit is far from trivial given the one-point estimator. Indeed, the bandit  
11 information setting is far more challenging than the online (full information) setting and several novel steps are  
12 required to design a low-regret algorithm: First, in the bandit algorithm, the point at which we play and the point at  
13 which we get the gradient estimation are different (L206-209). To circumvent this issue, we proposed the biphasic  
14 (exploration/exploitation) method (L210-218). Second, the point for estimation may fall out of the constraint set  
15 (L178-179). In [23], this issue was resolved by assuming that  $rB^d \subset \mathcal{K} \subset RB^d$ , which does not hold for many  
16 DR-continuous submodular function, whose domain is defined to be a subset of the non-negative orthant. Therefore we  
17 introduced the definition of  $\delta$ -interior, and explained how it can help us address the bandit problem (L180-188). We  
18 also proposed a method to construct proper  $\delta$ -interior (L195-197) and prove the result by a geometric analysis (Lemma  
19 1). We established the regret bound based on Lemma 1 (Theorem 2), but also provided the result for general constraint  
20 (Theorem 4, Lines 219-221).

21 (3) By providing the hardness result (Lemma 2), we showed that it is difficult to extend our method directly from  
22 continuous settings to discrete case. Then we considered the RBSM model, and established a sublinear regret bound.

23 **Response to Reviewer #1:**

24 **Q1:** Using other gradient estimation techniques may lead to better regret bounds. **A1:** We agree with the reviewer and  
25 are grateful for suggesting an interesting future research direction.

26 **Q2:** Provide results for the case that the horizon  $T$  is not available offline. **A2:** Using the doubling trick (Auer et al.,  
27 1995) will easily extend our methods to the cases where  $T$  is unknown. The exploration and exploitation phases are  
28 similar to those in Alg. 2, 3.

29 **Q3:** Provide numerical experiments and compare its running time with previous algorithms for this problem. **A3:** We  
30 thank the reviewer for the suggestion. However, since our algorithms are the first with a sublinear regret bound for both  
31 one-shot online learning and continuous bandit problems, there is no previous work to compare with.

32 **Q4:** Extension to continuous submodular functions. **A4:** Our results only apply to DR-submodular functions. We tried  
33 to make it clear in the abstract and throughout the paper. We are sorry for the confusion.

34 **Response to Reviewer #2:**

35 **Q1:** Some detailed comments. **A1:** The definition of radius uses the assumption that the constraint contains 0.  
36 Assumption 5 is an analogue of the assumption in [23], which assumes that  $rB^d \subset \mathcal{K} \subset RB^d$ . This assumption is  
37 listed in the statement of Theorem 2 (4 to 6 on the first line). We assume that  $r$  is given.

38 **Q2:** Applications of RBSM. **A2:** In theory, RBSM can be regarded as a relaxation of BSM, which helps us to better  
39 understand the nature of BSM. In practice, the responsive model (not only for submodular maximization or bandit) has  
40 potentially many applications when a decision cannot be committed, while we can still get the potential outcome of the  
41 decision as feedback. For example, suppose that we have an replenishable inventory of  $n$  items. Each time there is  
42 a customer coming with a utility function unknown to us. We allocate a collection of items to her, and the goal is to  
43 maximize the total utility (reward) of all the customers. We may use a partition matroid to model diversity (in terms of  
44 category, time, etc). In the RBSM model, we cannot allocate the collection of items which violates the constraint to the  
45 customer, but we can use it as a questionnaire, and the customer will tell us the potential utility if she received those  
46 items. The feedback will help us to make better decisions in the future. Similar examples include portfolio selection:  
47 when the investment choice is too risky, i.e., violates the recommended constraint set, we may stop trading and thus get  
48 no reward on that trading period, but at the same time observe the potential reward if we invested in that way. We will  
49 add more examples in the revision.

50 **Response to Reviewer #3:**

51 **Q1:** The resulting regret bounds are worse than the previous one, and shows a limitation of the techniques. **A1:** The  
52 previous  $O(\sqrt{T})$  bound is achieved by using  $\sqrt{T}$  exact gradients or  $T^{3/2}$  stochastic gradients per round, while our  
53 method only needs one single gradient. Our result opens the possibility of achieving sub-linear regret with only one  
54 gradient evaluation. We agree that it is an interesting future work to achieve the same  $O(\sqrt{T})$  regret bound.