
Sum-max Submodular Bandits

Anonymous Authors¹

Abstract

Many online decision-making problems correspond to maximizing a sequence of submodular functions. In this work, we introduce sum-max functions, a subclass of monotone submodular functions capturing several interesting problems, including best-of- K -bandits, combinatorial bandits, and the bandit versions on M -medians and hitting sets. We show that all functions in this class satisfy a key property that we call pseudo-concavity. This allows us to prove $(1 - \frac{1}{e})$ -regret bounds for bandit feedback in the nonstochastic setting of the order of \sqrt{MKT} (ignoring log factors), where T is the time horizon and M is a cardinality constraint. This bound, attained by a simple and efficient algorithm, significantly improves on the $\tilde{O}(T^{2/3})$ regret bound for online monotone submodular maximization with bandit feedback. We also extend our results to a bandit version of the facility location problem.

1. INTRODUCTION

In many concrete settings of sequential decision-making, decisions are subsets of a finite set $[K]$ (possibly with cardinality constraints) and utilities, or rewards, are non-linear set functions over $[K]$. Although we may know that utility functions have some specific structure, e.g., they are submodular, the feedback may not reveal anything beyond the utility of the current decision. For example, consider an advertising campaign over $[K]$ digital channels (e.g., web, apps, and social media). Due to budget constraints, the campaign can show ads only on a subset of M channels for every user. If a user ends up buying the advertised product, we observe that

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the Workshop on Foundations of Reinforcement Learning and Control at the International Conference on Machine Learning (ICML). Do not distribute.

a sale occurred, but we may not know which of the M channels triggered the purchase. The advertiser's goal is to choose the subset of channels for each new user in order to maximize the number of sales.

The same problem was studied (with a different motivation) by Simchowitz et al. (2016) under stochastic assumptions on the generation of the Bernoulli random variables each indicating whether displaying an ad on a certain channel triggers a purchase for the current user. In this work, we study the nonstochastic variant of this problem, where the binary variables associated with the channels are chosen, for each user, by an oblivious adversary. Our main result is an efficient algorithm minimizing regret in a much larger class of problems containing the multichannel advertising problem as a special case. In particular, our regret analysis applies to any sequential decision-making problem where reward functions belong to a subclass of all monotone submodular functions called *sum-max*.

A sum-max function is defined by a nonnegative matrix with K columns and an arbitrary number of rows. The value of the function evaluated at a subset $\mathcal{S} \subset [K]$ of columns is the sum over the rows of the maximum row element over the subset \mathcal{S} of columns. In the multichannel campaign example, the matrix is binary with a single row. The j -th entry indicates whether the current user would buy the product if advertised on channel j . If the matrix is square and symmetric, then we recover the non-metric facility location problem as a special case.

As we said earlier, our analysis of regret for sum-max functions assumes bandit feedback: at each time t we only observe the reward $r_t(\mathcal{A}_t)$ associated with our decision \mathcal{A}_t , where r_t is the sum-max function chosen by the adversary at time t . Hence, the reward $r_t(\mathcal{S})$ that we would have obtained by choosing any $\mathcal{S} \neq \mathcal{A}_t$ remains unknown. We also consider cardinality constraints, in the form of a parameter M requiring that the decision \mathcal{A}_t at each time t satisfy $|\mathcal{A}_t| \leq M$. Note that when $M = 1$ we recover the adversarial K -armed bandit problem.

Our main result is an efficient algorithm, MSE3, achieving a $\tilde{O}(\sqrt{MKT})$ bound on the γ_M -regret for $\gamma_M = 1 -$

(1 - 1/M)^M. For comparison, for the class of all monotone submodular functions, Niazadeh et al. (2021) obtain a (1 - 1/e)-regret bound of O((ln K)^{1/3} M (KT)^{2/3}). As gamma_M > 1 - 1/e for all M > 1, this bound is worse than ours in both approximation factor and regret.

When M = 1, algorithm MSE3 reduces to the standard EXP3 algorithm for K-armed bandits and our result specializes to the standard O(sqrt(K ln K T)) regret bound of EXP3. This implies that the sqrt(KT) dependence in the regret bound is not improvable, even disregarding efficiency. Moreover, we show that improving on the approximation factor gamma_M with an efficient algorithm would give an efficient randomized algorithm for solving set cover on [K] with an approximation ratio of (1 - epsilon) ln K, which is NP-hard for any epsilon > 0 (Dinur and Steurer, 2014).

In many real world problems, including an element i in the decision A_t at round t invokes a cost (i.e., a negative reward) c_{t,i} >= 0. When this is the case we would like to maximize the cumulative reward:

$$\sum_{t \in [T]} \left(r_t(\mathcal{A}_t) - \sum_{i \in \mathcal{A}_t} c_{t,i} \right).$$

We show that MSE3 can handle this generalized problem if it receives, at the end of each round t, the values of c_{t,i} for all i in A_t. We note that the bandit MSE3 without costs is a special case of MSE3 with costs.

The inclusion of costs creates a tension between including arms in A_t to increase the reward and, simultaneously, avoid including too many arms to control the costs. We address this trade-off in Section 5 by introducing and analyzing a variant of MSE3 for regret minimization with costs and bandit feedback where the rewards are sum-max functions without cardinality constraints. We call this setting the bandit facility location problem because it is a bandit version of the online facility location problem studied by Pasteris et al. (2021).

For M > 1 and arbitrary costs, MSE3 selects A_t by performing M independent draws a_{t,1}, ..., a_{t,M} from a distribution p_t = (p_{t,1}, ..., p_{t,K}) in Delta_K. Then, a reward estimate for each i in [K] is computed using

$$g_{t,i} = \frac{r_t(\mathcal{A}_t) - c_{t,i}}{p_{t,i}} \sum_{j \in [M]} \mathbb{I}[a_{t,j} = i], \quad (1)$$

where, for any statement S, the Iverson bracket notation [[S]] is defined as [[S]] = 1 if S is true and [[S]] = 0 otherwise. Note that for M = 1 and c_{t,i} = 0, the above reduces to the standard reward estimate of EXP3.

We now give an overview of how MSE3 works when we have no costs (i.e., c_{t,i} = 0). For all set functions r,

we construct a function Phi^r : R_+^K -> R such that for all q in Delta_K we have that Phi^r(q) is the expected value of r(A) when A is constructed by drawing M arms i.i.d. with replacement from q. Specifically, we first show that there exists a function tilde{r} : 2^{[K]} -> R such that for all Q subseteq [K] we have r(Q) = sum_{S subseteq [K]} [[Q subseteq S]] tilde{r}(S). For all q in R_+^K we then define:

$$\Phi^r(q) = \sum_{S \subseteq [K]} \tilde{r}(S) \left(\sum_{i \in [K]} \mathbb{I}[i \in S] q_i \right)^M.$$

We learn via online exponentiated gradient ascent using the unbiased estimates (1) of the gradient of Phi^{r_t}. Clearly, for exponentiated gradient ascent to work we must have that, for all rounds t, our objective function Phi^{r_t} is concave over the simplex. We show that a sufficient condition for this to hold is that the function r_t is pseudo-concave, see Section 2 for a formal definition.

Next, we bound the regret with respect to any vector p^* in Delta_K. Namely, we bound the expected reward of our algorithm relative to sum_{t in [T]} Phi^{r_t}(p^*). By taking p^* such that p_i^* = [[i in S]]/|S| for some set S we show that, because r_t is submodular, we have Phi^{r_t}(S) >= (1 - alpha^M) r_t(S) where alpha = (|S| - 1)/|S|. By bounding the variance of the gradient estimate we show that the regret term is O(sqrt(MKT)).

We have provided an overview of how, when we have no costs, MSE3 works and why we require r_t to be pseudo-concave and submodular. We now describe how costs are incorporated. This is done by using, instead of Phi^{r_t}, the function Psi_t defined by:

$$\Psi_t(q) = \Phi^{r_t}(q) - M \sum_{i \in [K]} q_i c_{t,i},$$

so that Psi_t(p_t) lower bounds the expected profit on trial t. Since Psi_t differs from Phi^{r_t} by a linear function it is straightforward to extend the above methodology to this new objective function.

2. SUM-MAX FUNCTIONS

We now introduce sum-max functions and define the key property of this class that allows us to learn it with bandit feedback.

Definition 2.1. A set function r : 2^{[K]} -> R is *sum-max* if and only if there exists some N in N and some matrix V in R^{N x K} such that for all S subseteq [K] with S not= empty we have:

$$r(S) = \sum_{k \in [N]} \max_{i \in S} V_{k,i} \quad \text{and} \quad r(\emptyset) \leq \sum_{k \in [N]} \min_{i \in [K]} V_{k,i}$$

For example, consider a marketplace with N buyers and K sellers. The value $V_{k,i}$ is the combined utility of buyer k going to seller i . The value $r(\mathcal{S})$ is the social welfare when a subset \mathcal{S} of sellers participate in the marketplace, and buyers match up with sellers to optimize their combined utilities. When there is only one buyer ($N = 1$), \mathbf{V} is a vector (V_1, \dots, V_K) and we view each $i \in [K]$ as an arm with reward V_i . Then $r(\mathcal{S}) = \max_{i \in \mathcal{S}} V_i$, the maximum reward of an arm in the chosen set \mathcal{S} .

As sum-max functions are sums of monotone submodular functions, they are monotone submodular. We now list a number of sequential decision-making problems that can be expressed as regret minimization of specific sum-max functions under bandit feedback.

The multichannel campaign problem. This is our nonstochastic variant of the best-of- k bandit problem of Simchowitz et al. (2016). To view it as an instance of sum-max optimization, set $N = 1$ and let $V_i \in \{0, 1\}$ indicate whether a user makes a purchase when the ad is displayed on channel i . Then (V_1, \dots, V_K) can be viewed as the incidence vector of a subset $\mathcal{D} \subseteq [K]$ of channels, and the reward is defined by $r(\mathcal{S}) = \llbracket \mathcal{S} \cap \mathcal{D} \neq \emptyset \rrbracket$. The feedback is bandit because we do not know what channel triggered the sale for that user.

Bandit hitting sets. This is a generalization of the previous example where $N \geq 1$ and \mathbf{V} is a boolean matrix. Each row of \mathbf{V} denotes a subset \mathcal{C}_k of $[K]$ and $V_{k,i}$ indicates whether $i \in \mathcal{C}_k$. The value $r(\mathcal{S})$ then counts how many sets \mathcal{C}_k have a non-empty intersection with \mathcal{S} . Bandit setting occurs when the sets remain unknown and each time we only observe the number of intersected sets.

Combinatorial bandits. Another important special case is when we receive the sum of the rewards r_i of the arms $i \in \mathcal{S}$. In this case $N = K$ and $V_{k,i} = \llbracket k = i \rrbracket r_i$. The problem is then equivalent to a combinatorial bandit (with full bandit feedback) over the class of M -sized subsets (Cesa-Bianchi and Lugosi, 2012).

Bandit k -medians. Given $\mathbf{x}_1, \dots, \mathbf{x}_N$ points in a metric space (\mathcal{X}, d) , consider the version of the k -medians problem (for $k = M$) where the M centroids have to be chosen in the given set of points. The value of the objective function at a candidate solution $\mathcal{S} \subset [K]$ with $|\mathcal{S}| \leq M$ can be written as

$$r(\mathcal{S}) = - \sum_{k \in [N]} \min_{i \in \mathcal{S}} d(\mathbf{x}_k, \mathbf{x}_i).$$

Clearly, this is a sum-max function for \mathbf{V} with elements $V_{k,i} := -d(\mathbf{x}_k, \mathbf{x}_i)$. The feedback is bandit when we do not know the metric, but we can observe the value

of the objective function.

Next, we introduce an important property of sum-max functions.

Definition 2.2. Suppose we have a set function $r : 2^{[K]} \rightarrow \mathbb{R}$. For any $\mathcal{S} \subseteq [K]$ define the matrix $\mathbf{U}^{r,\mathcal{S}} \in \mathbb{R}^{K \times K}$ such that $U_{i,j}^{r,\mathcal{S}} = r(\mathcal{S} \cup \{i, j\})$ for all $i, j \in [K]$. We call the function r *pseudo-concave* if and only if $\mathbf{x}^\top \mathbf{U}^{r,\mathcal{S}} \mathbf{x} \leq 0$ for all $\mathcal{S} \subseteq [K]$ and all $\mathbf{x} \in \mathbb{R}^K$ with $\mathbf{x} \cdot \mathbf{1} = 0$.

In Appendix C, we show that there are monotone submodular functions that are not pseudo-concave. As a consequence, sum-max functions are indeed a proper subset of the class of monotone submodular functions. The following theorem confirms that all sum-max functions are pseudo-concave:

Theorem 2.3. *Any sum-max set function is pseudo-concave.*

Proof. Suppose we have some sum-max function $r : 2^{[K]} \rightarrow [0, 1]$. Let \mathbf{V} be as in Definition 2.1. Without loss of generality, we will assume that all components of \mathbf{V} are non-negative and $r(\emptyset) = 0$ (since any sum-max function can be transformed into this form by the addition of a constant).

Define, for any $\mathcal{Q} \subseteq [K]$, the set function $r^\mathcal{Q} : 2^{[K]} \rightarrow [0, 1]$ such that for all $\mathcal{S} \subseteq [K]$ we have

$$r^\mathcal{Q}(\mathcal{S}) := \llbracket \mathcal{S} \cap \mathcal{Q} \neq \emptyset \rrbracket.$$

We shall now show that for all such \mathcal{Q} we have that $r^\mathcal{Q}$ is pseudo-concave. Choose any $\mathbf{x} \in \mathbb{R}^K$ with $\mathbf{x} \cdot \mathbf{1} = 0$ and any $\mathcal{S} \subseteq [K]$. We have two cases:

1. If $\mathcal{S} \cap \mathcal{Q} \neq \emptyset$, for all $i, j \in [K]$ we have $r^\mathcal{Q}(\mathcal{S} \cup \{i, j\}) = 1$, this implies $\mathbf{U}^{r^\mathcal{Q},\mathcal{S}} = \mathbf{1}\mathbf{1}^\top$ and hence $\mathbf{x}^\top \mathbf{U}^{r^\mathcal{Q},\mathcal{S}} \mathbf{x} = 0$.
2. If $\mathcal{S} \cap \mathcal{Q} = \emptyset$ then for all $i, j \in [K]$ we have

$$r^\mathcal{Q}(\mathcal{S} \cup \{i, j\}) = \llbracket (i \in \mathcal{Q}) \vee (j \in \mathcal{Q}) \rrbracket.$$

Let $\mathbf{z} \in \{0, 1\}^K$ be such that for all $k \in [K]$ we have $z_k := \llbracket k \in \mathcal{Q} \rrbracket$. Then for all $i, j \in [K]$ we have

$$\llbracket (i \in \mathcal{Q}) \vee (j \in \mathcal{Q}) \rrbracket = 1 - z_i z_j,$$

so that, by above, we have $\mathbf{U}^{r^\mathcal{Q},\mathcal{S}} = \mathbf{1}\mathbf{1}^\top - \mathbf{z}\mathbf{z}^\top$. This implies that: $\mathbf{x}^\top \mathbf{U}^{r^\mathcal{Q},\mathcal{S}} \mathbf{x} = -(\mathbf{x} \cdot \mathbf{z})^2 \leq 0$.

And therefore, $r^\mathcal{Q}$ is pseudo-concave.

Now suppose we have a vector $\mathbf{v} \in \mathbb{R}_+^K$ and define the set function $r^\mathbf{v} : 2^{[K]} \rightarrow \mathbb{R}_+$ such that for all $\mathcal{S} \subseteq [K]$ we have

$$r^\mathbf{v}(\mathcal{S}) := \max_{i \in \mathcal{S}} v_i,$$

where the maximum of the empty set is defined as equal to zero. We can order the set $[K]$ into the sequence $\langle j_i \mid i \in [K] \rangle = [K]$ where $v_{j_{i+1}} \leq v_{j_i}$ for all $i \in [K-1]$. For all $i \in [K]$ we can define $\mathcal{Q}_i := \{j_k \mid k \leq i\}$. Now note then that for all $\mathcal{S} \subseteq [K]$ the set function $r^{\mathbf{v}}(\mathcal{S})$ can be expressed as

$$\begin{aligned} & \sum_{i \in [K-1]} (v_{j_i} - v_{j_{i+1}}) \mathbb{1}[\mathcal{S} \cap \mathcal{Q}_i \neq \emptyset] + v_{j_K} \mathbb{1}[\mathcal{S} \cap \mathcal{Q}_K \neq \emptyset] \\ &= \sum_{i \in [K-1]} (v_{j_i} - v_{j_{i+1}}) r^{\mathcal{Q}_i}(\mathcal{S}) + v_{j_K} r^{\mathcal{Q}_K}(\mathcal{S}), \end{aligned}$$

so, by above, $r^{\mathbf{v}}$ is a positive sum of pseudo-concave functions and is hence itself pseudo-concave. Note also that $r^{\mathbf{v}}$ is clearly submodular. Noting that r is a positive sum of functions of the form $r^{\mathbf{v}}$ we have now shown that it is both pseudo-concave and submodular as required. \square

3. ADDITIONAL RELATED WORK

The work closest to ours is Pasteris et al. (2021), where they study online facility location with full information feedback. Our work improves on theirs in many respects: First, we solve the problem with bandit feedback, which requires designing an entirely different algorithm based on our discovery of an unbiased estimator for the gradient of our expected reward (we find it remarkable that such an estimator exists). As a consequence, our algorithm is also applicable to the full-information setting, where we obtain a per-trial running time of $\mathcal{O}(MK)$ when given an oracle for the reward function. When considering general sum-max functions, the methodology of Pasteris et al. (2021) would instead require a per-trial running time exponential in K^1 . Second, our algorithm can efficiently learn classes that are even more general than sum-max functions. Third, we obtain tighter approximation ratios and show optimality for the multichannel campaign problem (and thus optimality in general).

Sum-max functions are a special case of linear submodular functions (Yue and Guestrin, 2011), which are of the form $r(\mathcal{S}) = \sum_{i \in [N]} w_i F_i(\mathcal{S})$ for F_1, \dots, F_N monotone submodular functions and w_1, \dots, w_N non-negative coefficients. However, linear submodular functions have been only studied in stochastic settings, assuming preliminary knowledge of F_1, \dots, F_N , and using a feedback model more informative than our bandit feedback.

Click-models (Lattimore and Szepesvári, 2020; Lattimore et al., 2018; Kveton et al., 2015) provide a different

¹The work of Pasteris et al. (2021) only considered single-user cases, but it is straightforward to extend their methodology to general sum-max functions.

Algorithm 1 MSE3

Set $\eta := \ln(K)/R$ and $p_{1,i} := 1/K$ for $i \in [K]$

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

1. For all $j \in [M]$ draw $a_{t,j} \in [K]$ from distribution \mathbf{p}_t
2. Define $\mathcal{A}_t := \{a_{t,j} \mid j \in [M]\}$
3. Receive $r_t(\mathcal{A}_t)$ and $\{c_{t,i} \mid i \in \mathcal{A}_t\}$
4. For all $i \in [K]$ set

$$g_{t,i} := \frac{r_t(\mathcal{A}_t) - c_{t,i}}{p_{t,i}} \sum_{j \in [M]} \mathbb{1}[a_{t,j} = i]$$

5. For all $i \in [K]$ define $\tilde{p}_{t,i} := p_{t,i} \exp(\eta g_{t,i})$
 6. Define $\mathbf{p}_{t+1} := \tilde{\mathbf{p}}_t / \|\tilde{\mathbf{p}}_t\|_1$
-

stochastic formalization of the best-of- k bandit problem. Here the user is presented with an ordered list of items, and the learner receives a positive reward if the user clicks on one of the presented items. The difference with our multichannel campaign problem is that the items are ordered, and the likelihood of clicking an item is also affected by the position of the item within the list.

4. MAIN RESULT

Our learning problem is formally defined as follows. The values $M, K \in \mathbb{N}$ and $C \in \mathbb{R}_+$ are all preliminarily known to the learner. Hidden from the learner, the adversary selects a sequence of set functions $\langle r_t \mid t \in [T] \rangle$, each with domain $2^{[K]}$ and a sequence of vectors $\langle \mathbf{c}_t \mid t \in [T] \rangle$ each in $[0, C]^K$. On each trial $t \in [T]$:

1. The learner chooses some $\mathcal{A}_t \subseteq [K]$ with $|\mathcal{A}_t| \leq M$.
2. The value $r_t(\mathcal{A}_t)$ is revealed.
3. For all $i \in \mathcal{A}_t$ the value $c_{t,i}$ is also revealed.

The learner maintains a probability vector $\mathbf{p}_t \in \Delta_K$, and behaves as described in Algorithm 1.

To aid our theorem statement we add the following definitions. For all $t \in [T]$ and $\mathcal{Q} \subseteq [K]$ we define $\hat{r}_t(\mathcal{Q}) := r_t(\mathcal{Q}) - r_t(\emptyset)$, which is the difference between the learner's profit on trial t and that which it would have obtained by selecting the empty set, $\psi_t := \hat{r}_t(\mathcal{A}_t) - \gamma_t(\mathcal{A}_t)$, and $\gamma_t(\mathcal{Q}) := \sum_{i \in \mathcal{Q}} c_{t,i}$. We note that by considering \hat{r}_t instead of r_t our bounds do not change when r_t is shifted by an additive constant (which can be different for different trials t) as long as the range of r_t falls within the bounds described as follows.

We assume that the learner knows upper and lower bounds on the range r_t for all trials t . Hence, without loss of generality, assume that $r_t(\mathcal{Q}) \in [-1, 0]$ for all

$t \in [T]$ and $\mathcal{Q} \setminus [K]$ (otherwise scale and shift r_t and C). Let

$$R := (1 + C)\sqrt{2\ln(K)M(K + M - 1)T}.$$

Our results hold for a relaxed notion of submodularity, which we call *pseudo-submodularity*.

Definition 4.1. A set function $r : 2^{[K]} \rightarrow \mathbb{R}$ is *pseudo-submodular* if and only if for every set $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$ there exists some $i \in \mathcal{S}$ such that for all $\mathcal{Q} \subseteq \mathcal{S} \setminus \{i\}$ we have $r(\mathcal{Q} \cup \{i\}) - r(\mathcal{Q}) \geq r(\mathcal{S}) - r(\mathcal{S} \setminus \{i\})$.

Note that all pseudo-submodular set functions are also submodular. We now present our main result.

Theorem 4.2. *Given r_t is pseudo-concave and pseudo-submodular for all $t \in [T]$, then for any set $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$ we have*

$$\sum_{t \in [T]} \mathbb{E}[\psi_t] \geq (1 - \alpha^M) \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) - \frac{M}{|\mathcal{S}|} \sum_{t \in [T]} \gamma_t(\mathcal{S}) - R,$$

where

$$\alpha := 1 - \frac{1}{|\mathcal{S}|}.$$

Proof. See Section 6 □

We note that both the standard facility location and k -medians problems are often phrased as the minimization of a loss rather than a maximization of a profit. Our results easily capture this by considering the reward as a negative loss.

We now show that the approximation ratio $1 - \alpha^M$ is not improvable in general in the class of sum-max functions. In particular, we show that obtaining an efficient online learning algorithm for the multichannel advertising problem with a sublinear γ -regret with $\gamma < 1 - \alpha^M$ would give an efficient randomized algorithm for solving set cover on $[K]$ with an approximation better than $\ln K$. As shown in (Dinur and Steurer, 2014), obtaining an approximation of $(1 - \varepsilon) \ln K$ for set cover is NP-hard for any $\varepsilon > 0$.

Recall that an instance of the multichannel campaign problem over K ads is defined by a sequence $\langle r_t \mid t \in [T] \rangle$ of set functions over $[K]$ such that for all $t \in [T]$ there exists some $\mathcal{D}_t \subseteq [K]$ with $r_t(\mathcal{Q}) = \mathbb{1}[\mathcal{Q} \cap \mathcal{D}_t \neq \emptyset]$ for all $\mathcal{Q} \subseteq [K]$.

Theorem 4.3. *Suppose that there exists some $d \in \mathbb{N}$, $s \in (0, 1)$, $\gamma > 1$, and a randomized polynomial time algorithm for the learner such that for all $K, M \in \mathbb{N}$ and for any instance of the multichannel advertising problem, it holds that $|\mathcal{A}_t| \leq M$ for all $t = 1, \dots, T$*

Algorithm 2 FLE3

Run MSE3 with $L = 2K$ arms and $M := \frac{K}{2} \ln(T/K^2)$.

On each trial $t \in [T]$:

1. Let \mathcal{A}'_t be the output of MSE3
 2. Output $\mathcal{A}_t := \mathcal{A}'_t \cap [K]$
 3. Receive $r_t(\mathcal{A}_t)$ and $\{c_{t,i} \mid i \in [K]\}$
 4. For all $i \in [L] \setminus [K]$ set $c_{t,i} := 0$
 5. Feed $r_t(\mathcal{A}_t)$ and $\{c_{t,i} \mid i \in [L]\}$ back to MSE3
-

and, for any subset $\mathcal{S} \subseteq [K]$,

$$\mathbb{E} \left[\sum_{t \in [T]} r_t(\mathcal{A}_t) \right] \geq (1 - \alpha^{\gamma M}) \sum_{t \in [T]} r_t(\mathcal{S}) - R',$$

where $R' \in \mathcal{O}(K^d T^s)$ and $\alpha := 1 - \frac{1}{|\mathcal{S}|}$. Then, for all $\varepsilon \in (0, 1 - 1/\gamma)$ and $B > 4^{1/((1-\varepsilon)\gamma-1)}$, there exists a randomized polynomial-time algorithm for the set cover problem on $[B]$ that, with probability at least $\frac{1}{2}$, achieves approximation ratio at least $(1 - \varepsilon) \ln(B)$.

The proof can be found in Appendix B.

5. BANDIT FACILITY LOCATION

We can view this setting as a generalization of the marketplace example where sellers pay a known cost to enter the market. At each round, the platform admits a subset \mathcal{A}_t of sellers and only observes the resulting social welfare (bandit feedback).

In this application, there are no restrictions on the set of arms \mathcal{A}_t that we choose. We seek to maximize $r(\mathcal{A}_t) - \gamma(\mathcal{A}_t)$ where r is the sum-max reward function and γ is the linear and positive cost function.

For the facility location problem we must choose M , noting that although a high value of M increases the approximation ratio on the reward, it also increases that on the costs. To decrease the potentially large approximation ratio on the costs, we borrow from (Pasteris et al., 2021) the idea of *dummy arms* and the tuning of M . This leads to our algorithm FLE3 described in Algorithm 2. The bound on the total profit of FLE3 is given in the following theorem.

Theorem 5.1. *Given that $C = 1$ and $r_t : 2^K \rightarrow [-1, 0]$ is pseudo-concave and pseudo-submodular for all $t \in [T]$, we have that the algorithm FLE3 obtains the following bound for all $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$:*

$$\sum_{t \in [T]} \mathbb{E}[\psi_t] \geq \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) - \frac{1}{2} \ln \left(\frac{T}{K^2} \right) \sum_{t \in [T]} \gamma_t(\mathcal{S}) - R'',$$

where $R'' \in \tilde{\mathcal{O}}(K\sqrt{T})$.

275 *Proof.* For all $t \in [T]$ define the set function $r'_t : 2^L \rightarrow$
 276 $[0, 1]$ such that for all $\mathcal{Q} \subseteq [L]$, $r'_t(\mathcal{Q}) := r_t(\mathcal{Q} \cap [K])$
 277 and, as consequence, $\hat{r}'_t(\mathcal{Q}) := r'_t(\mathcal{Q}) - r'_t(\emptyset)$. Now
 278 taking into consideration any possible comparator set
 279 $\mathcal{S} \subseteq [K]$, we define

$$280 \quad \mathcal{S}' := \mathcal{S} \cup \{K + i \mid i \in [K - |\mathcal{S}|\}],$$

282 noting that $|\mathcal{S}'| = K$. Note that r'_t is sum-max and
 283 hence, by Theorem 2.3, is pseudo-concave and submod-
 284 ular for all $t \in [T]$. This allows us to apply Theorem
 285 4.2, that gives us:

$$\begin{aligned}
 287 \quad & \sum_{t \in [T]} \mathbb{E}[\psi_t] \geq (1 - \alpha^M) \sum_{t \in [T]} \hat{r}'_t(\mathcal{S}') - \frac{M}{|\mathcal{S}'|} \sum_{t \in [T]} \gamma_t(\mathcal{S}') - R \\
 288 \quad & = (1 - \alpha^M) \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) - \frac{M}{|\mathcal{S}'|} \sum_{t \in [T]} \gamma_t(\mathcal{S}) - R \quad (2) \\
 290 \quad & = (1 - \alpha^M) \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) - \frac{1}{2} \ln \frac{T}{K^2} \sum_{t \in [T]} \gamma_t(\mathcal{S}) - R, \\
 292 \quad & \quad \quad \quad (3)
 \end{aligned}$$

297 where equation (2) comes from the contribution of the
 298 dummy arms and equation (3) from the definition of
 299 M . Given that

$$300 \quad \alpha := \frac{|\mathcal{S}'| - 1}{|\mathcal{S}'|} = \frac{K - 1}{K} \leq \exp(-1/K),$$

303 we can therefore see that

$$\begin{aligned}
 305 \quad & \alpha^M \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) \leq \alpha^M T = \exp(-M/K) T \\
 306 \quad & = \frac{1}{\sqrt{T/K^2}} T = \sqrt{TK^2}, \quad (4)
 \end{aligned}$$

310 where we used the definition of M given in Algorithm 2.
 311 Putting together (3) and (4) gives us the result, where
 312 $R'' = R + \sqrt{TK^2}$. \square

314 6. ANALYSIS

316 We now give an overview of the proof of Theorem 4.2.

318 We first consider the case that we have no costs (i.e.
 319 $\mathbf{c}_t = \mathbf{0}$). MSE3 works by maintaining a probability
 320 distribution over the set of arms. Specifically, $\mathbf{p}_t \in \Delta_K$
 321 is the vector whose components are the probabilities of
 322 drawing the actions on trial t . On trial t the algorithm
 323 constructs the set \mathcal{A}_t by drawing a sequence $\langle a_{t,j} \mid j \in$
 324 $[M] \rangle$ of arms i.i.d. with replacement from \mathbf{p}_t and then
 325 setting $\mathcal{A}_t := \{a_{t,j} \mid j \in [M]\}$.

326 This stochastic draw of a sequence and set from a
 327 probability vector will be represented by the following
 328 notation.

329

Definition 6.1. For all $\mathbf{q} \in \Delta_K$ let $\langle b_j(\mathbf{q}) \mid j \in [M] \rangle$
 be a sequence of stochastic quantities drawn i.i.d. at
 random from (the probability distribution characterised
 by) \mathbf{q} . In addition, let $\mathcal{B}(\mathbf{q}) := \{b_j(\mathbf{q}) \mid j \in [M]\}$.

Note that our expected reward on trial t is $\mathbb{E}[r_t(\mathcal{B}(\mathbf{p}_t))]$
 and hence, for all set functions r we shall construct
 a differentiable function $\Phi^r : \mathbb{R}^K \rightarrow \mathbb{R}$ such that for
 all $\mathbf{q} \in \Delta_K$ we have $\Phi^r(\mathbf{q}) = \mathbb{E}[r(\mathcal{B}(\mathbf{q}))]$. This con-
 struction is based on the following notion of a *subset*
decomposition.

Definition 6.2. Given a function $r : 2^{[K]} \rightarrow \mathbb{R}$, we
 call a function $\tilde{r} : 2^{[K]} \rightarrow \mathbb{R}$ a *subset decomposition* of
 r if and only if for all $\mathcal{Q} \subseteq [K]$ we have

$$r(\mathcal{Q}) = \sum_{\mathcal{S} \subseteq [K]} \mathbb{1}[\mathcal{Q} \subseteq \mathcal{S}] \tilde{r}(\mathcal{S}).$$

The following lemma confirms that every set function
 has a unique subset decomposition.

Lemma 6.3. *Given a function $r : 2^{[K]} \rightarrow \mathbb{R}$ there*
exists a unique subset decomposition \tilde{r} of r .

Proof. See Appendix A.1. \square

Now we can define our function Φ^r .

Definition 6.4. For all $r : 2^K \rightarrow \mathbb{R}$ and all $\mathbf{q} \in \mathbb{R}^K$
 define

$$\Phi^r(\mathbf{q}) := \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \left(\sum_{i \in [K]} \mathbb{1}[i \in \mathcal{S}] q_i \right)^M,$$

where, by Lemma 6.3, \tilde{r} is the unique subset decompo-
 sition of r .

The following lemma confirms that our function Φ^r
 indeed satisfies our condition.

Lemma 6.5. *For all $r : 2^K \rightarrow \mathbb{R}$ and all $\mathbf{q} \in \Delta_K$ we*
have $\Phi^r(\mathbf{q}) = \mathbb{E}[r(\mathcal{B}(\mathbf{q}))]$.

Proof. See Appendix A.2. \square

Drawing inspiration from (Auer et al., 2001) we will
 learn via online exponentiated gradient ascent with the
 functions Φ^{r_t} using unbiased gradient estimates. Of
 course, this means that we must be able to construct
 unbiased gradient estimates. Remarkably, we now show
 that we can use our sequence $\langle a_{t,j} \mid j \in [M] \rangle$ and the
 observed reward $r_t(\mathcal{A}_t)$ to construct an unbiased gradi-
 ent estimate \mathbf{g}_t defined in Algorithm 1 of the function
 Φ^{r_t} at \mathbf{p}_t .

Lemma 6.6. For all $r : 2^K \rightarrow \mathbb{R}$, all $\mathbf{q} \in \Delta_K$ and all $i \in [K]$ we have

$$\partial_i \Phi^r(\mathbf{q}) = \mathbb{E} \left[\frac{r(\mathcal{B}(\mathbf{q}))}{q_i} \sum_{j \in [M]} \mathbb{1}[b_j(\mathbf{q}) = i] \right].$$

Proof. See Appendix A.3. \square

For exponentiated gradient ascent to work, we must have that, for all trials t , our objective function Φ^{r_t} is concave over the simplex. We now show that a sufficient condition for this to hold is that the function r_t is pseudo-concave.

Lemma 6.7. For all pseudo-concave set functions $r : 2^K \rightarrow \mathbb{R}$ we have that Φ^r is concave over the simplex Δ_K .

Proof. See Appendix A.4. \square

Now that we have all the underpinnings for exponentiated gradient ascent to function properly, we can establish a bound on the regret relative to any vector $\mathbf{p}^* \in \Delta_K$ via the following classic result.

Lemma 6.8. For any vector $\mathbf{p}^* \in \Delta_K$ we have

$$\begin{aligned} \sum_{t \in [T]} (\mathbf{p}^* - \mathbf{p}_t) \cdot \mathbf{g}_t &\leq \frac{1}{\eta} \sum_{i \in [K]} p_i^* \ln(K p_i^*) \\ &\quad + \eta \sum_{t \in [T]} \sum_{i \in [K]} p_{t,i} g_{t,i}^2. \end{aligned}$$

Proof. A classic result from the analysis of HEDGE. \square

This lemma gives a bound on the regret since, because we have shown that \mathbf{g}_t is an unbiased estimate of the gradient and the objective function is concave over the simplex, the term $(\mathbf{p}^* - \mathbf{p}_t) \cdot \mathbf{g}_t$ is bounded below by $\Phi^{r_t}(\mathbf{p}^*) - \Phi^{r_t}(\mathbf{p}_t)$. Note that we have shown above that $\Phi^{r_t}(\mathbf{p}_t)$ is equal to $\mathbb{E}[r_t(\mathcal{A}_t)]$.

We will later discuss the bounding of the regret itself, but first we shall show how to choose \mathbf{p}^* such that we can bound $\Phi^{r_t}(\mathbf{p}^*)$ relative to $r_t(\mathcal{S})$ for some set $\mathcal{S} \subseteq [K]$. Specifically, we will choose \mathbf{p}^* equal to $\mathbf{p}^{\mathcal{S}}$ in the following definition.

Definition 6.9. For all $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$ define $\mathbf{p}^{\mathcal{S}} \in \Delta_K$ such that for all $i \in [K]$ we have

$$p_i^{\mathcal{S}} := \frac{\mathbb{1}[i \in \mathcal{S}]}{|\mathcal{S}|}.$$

We use the following lemma will to bound $\Phi^{r_t}(\mathbf{p}^{\mathcal{S}})$, and it explains why we require r_t to be pseudo-submodular.

Lemma 6.10. Let $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$, $r : 2^{[K]} \rightarrow \mathbb{R}$ be a pseudo-submodular function, and $\mathcal{Z} \subseteq [K]$ be a set formed by drawing M elements uniformly at random (with replacement) from \mathcal{S} . Then we have

$$\mathbb{E}[r(\mathcal{Z}) - r(\emptyset)] \geq \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|} \right)^M \right) (r(\mathcal{S}) - r(\emptyset)).$$

Proof. See Appendix A.5. \square

With this lemma in hand, we can now bound $\Phi^{r_t}(\mathbf{p}^{\mathcal{S}})$.

Lemma 6.11. Given any $\mathcal{S} \subseteq [K]$ and any pseudo-submodular set function $r : 2^{[K]} \rightarrow \mathbb{R}$ we have

$$\Phi^r(\mathbf{p}^{\mathcal{S}}) \geq r(\emptyset) + \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|} \right)^M \right) (r(\mathcal{S}) - r(\emptyset)).$$

Proof. See Appendix A.6. \square

Before we bound the regret term, we show how to incorporate the costs, so that \mathbf{c}_t can be non-zero. This is done by choosing, instead of Φ^{r_t} , the objective function Ψ_t defined as follows.

Definition 6.12. For all trials $t \in [T]$ define $\Psi_t : \mathbb{R}^K \rightarrow \mathbb{R}$ such that for all $\mathbf{q} \in \mathbb{R}^K$ we have

$$\Psi_t(\mathbf{q}) := \Phi^{r_t}(\mathbf{q}) - M \mathbf{q} \cdot \mathbf{c}_t.$$

Note that by Lemma 6.5 we have that $\Psi_t(\mathbf{p}_t)$ is a lower bound on the expected profit and by Lemma 6.7 Ψ_t is concave over the simplex. It can hence serve as a surrogate concave objective function.

Lemma 6.6 leads to the following lemma, which confirms that \mathbf{g}_t is an unbiased gradient estimate of Ψ_t at \mathbf{p}_t .

Lemma 6.13. For all trials $t \in [T]$ and Ψ_t as defined in Definition 6.12, we have

$$\nabla \Psi_t(\mathbf{p}_t) = \mathbb{E}[\mathbf{g}_t \mid \mathbf{p}_t],$$

Proof. See Appendix A.7. \square

Now we have shown that our results carry over to the case of non-zero costs, we can finally bound the regret via Lemma 6.8 and the following lemma.

Lemma 6.14. For all trials $t \in [T]$ we have

$$\mathbb{E} \left[\sum_{i \in [K]} p_{t,i} g_{t,i}^2 \right] \leq (1 + C)^2 M (K + M - 1).$$

Proof. See Appendix A.8. \square

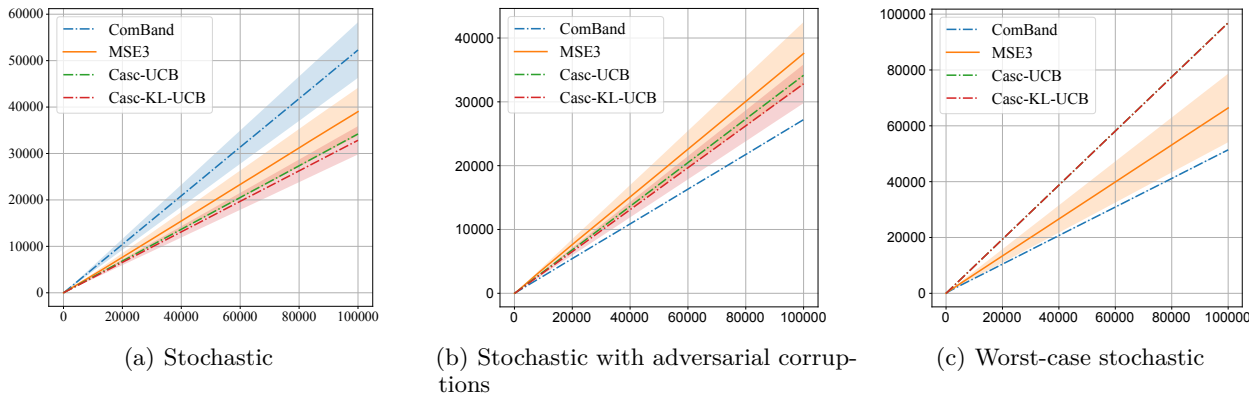


Figure 1. Cumulative reward over time in the three environment settings is described. The results also display the 95% confidence intervals over 35 runs with an Intel Xeon Gold 6312U, calculated using the standard error multiplied by the z-score of 1.96.

This completes the analysis. Although discussed here, Appendix A.9 formally shows how to piece the lemmas together in order to prove Theorem 4.2.

7. EXPERIMENTS

We experimentally evaluated the performance of our method by comparing it against two baselines: CASCADEBANDIT from (Kveton et al., 2015) (in both the UCB and KL settings) and COMBAND from (Cesa-Bianchi and Lugosi, 2012) for M -sized subsets, whose efficient implementation is described in Appendix D. We conducted our experiments in various synthetic settings. In each of these environments, a hidden vector $\theta \in \mathbb{R}^K$ is maintained. For each $k \in [K]$, the entry θ_k represents the probability of obtaining a unit reward. These values can be viewed as *attraction probabilities*: the probability that a user clicks on the specific item. After presenting a subset of M elements, the learner gets a unit reward if any of the selected items returns a 1, and 0 otherwise. It is worth emphasizing that our model does not necessitate binary rewards; it offers the flexibility to accommodate any sum-max reward function (as discussed in Section 2). The use of a binary reward model is specifically required for comparisons with click models such as CASCADEBANDIT.

Environments for the experiments. We experimentally evaluated our method in three different synthetic environments. We conducted experiments across a wide range of values for K , M , T , and for the probabilities associated with both optimal and suboptimal arms. In Figure 1, we display the cumulative reward over time obtained with $T = 10^5$, $K = 20$, $M = 3$ when the environments are set as follows:

- Stochastic** (Figure 1(a)): we randomly select M *good actions* to which we assign a reward probability of 0.3. The reward probabilities of the remaining $k - M$ arms are set to 0.1.
- Stochastic with adversarial corruptions** (Figure 1(b)): the rewards are generated as in the stochastic setting. However, in the first \sqrt{T} rounds all good actions have a deterministic reward of 0.
- Worst-case stochastic** (Figure 1(c)): this setting is inspired by the lower bound of Cohen et al. (2017). Here the set $\mathcal{M} \subset [K]$ of M *good actions* is drawn uniformly at random. Then, for each $k \in [K]$, the probabilities are assigned as follows:

$$\theta_k = \begin{cases} X_k + \epsilon & \text{if } k \in \mathcal{M} \\ X_k & \text{otherwise,} \end{cases} \quad \text{where } X_k \sim N\left(\frac{1}{2}, \sigma^2\right),$$

$$\sigma^2 = \frac{1}{192 + 96 \log T} \quad \text{and} \quad \epsilon = \sigma \sqrt{\frac{KM}{8T}}.$$

In Appendix E we present also results obtained varying the subset size M .

Results As expected, our most compelling results were achieved in the adversarial setting, where our approach demonstrated its superiority. In the two stochastic settings, we observed results that were on par with the established baseline methods, affirming the competitiveness of our proposed approach. These findings collectively underscore the effectiveness of our method, particularly in the challenging adversarial context, while also highlighting its versatility in stochastic scenarios. We emphasize that our method is the most efficient one, as each prediction only requires sampling M times from a probability distribution over the K available actions.

References

- P. Auer, Y. Freund, and R. E. Schapire. The non-stochastic multi-armed bandit problem. 2001. URL <https://api.semanticscholar.org/CorpusID:7732525>.
- N. Cesa-Bianchi and G. Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5): 1404–1422, 2012.
- A. Cohen, T. Hazan, and T. Koren. Tight bounds for bandit combinatorial optimization. In *Conference on Learning Theory*, pages 629–642. PMLR, 2017.
- I. Dinur and D. Steurer. Analytical approach to parallel repetition. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 624–633, 2014.
- D. Harvey and J. Van Der Hoeven. Integer multiplication in time $o(n \log n)$. *Annals of Mathematics*, 193(2):563–617, 2021.
- B. Kveton, C. Szepesvari, Z. Wen, and A. Ashkan. Cascading bandits: Learning to rank in the cascade model. In *International conference on machine learning*, pages 767–776. PMLR, 2015.
- T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- T. Lattimore, B. Kveton, S. Li, and C. Szepesvari. Toprank: A practical algorithm for online stochastic ranking. *Advances in Neural Information Processing Systems*, 31, 2018.
- R. Niazadeh, N. Golrezaei, J. R. Wang, F. Susan, and A. Badanidiyuru. Online learning via offline greedy algorithms: Applications in market design and optimization. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pages 737–738, 2021.
- S. Pasteris, T. He, F. Vitale, S. Wang, and M. Herbster. Online learning of facility locations. In *Algorithmic Learning Theory*, pages 1002–1050. PMLR, 2021.
- M. Simchowitz, K. Jamieson, and B. Recht. Best-of-k bandits. In *Conference on Learning Theory*, pages 1440–1489. PMLR, 2016.
- Y. Yue and C. Guestrin. Linear submodular bandits and their application to diversified retrieval. *Advances in Neural Information Processing Systems*, 24, 2011.

A. ANALYSIS PROOFS (PROOF OF THEOREM 4.2)

A.1. Lemma 6.3

Lemma A.1. *Given a function $r : 2^{[K]} \rightarrow \mathbb{R}$ there exists a unique subset decomposition \tilde{r} of r .*

Proof. For all $k \in [K] \cup \{0\}$ define $\mathcal{V}_k := \{\mathcal{S} \in 2^{[K]} \mid k \leq |\mathcal{S}|\}$. We take the inductive hypothesis that for all $k \in [K] \cup \{0\}$ there exists a unique function $\tilde{r}_k : \mathcal{V}_k \rightarrow \mathbb{R}$ such that for all $\mathcal{Q} \in \mathcal{V}_k$ we have

$$r(\mathcal{Q}) = \sum_{\mathcal{S} \in \mathcal{V}_k} \mathbb{I}[\mathcal{Q} \subseteq \mathcal{S}] \tilde{r}_k(\mathcal{S}).$$

We will prove the inductive hypothesis via reverse induction on k (i.e. from $k = K$ to $k = 0$).

The inductive hypothesis holds for $k = K$ since the only element of \mathcal{V}_K is $[K]$ so we must have $\tilde{r}_K([K]) := r([K])$

Now suppose, for some $i \in [K]$, the inductive hypothesis holds when $k = i$. Now consider the case that $k = i - 1$. Note that for all $\mathcal{Q} \in \mathcal{V}_i$ and $\mathcal{S} \in \mathcal{V}_{i-1} \setminus \mathcal{V}_i$ we must have that $\mathcal{Q} \not\subseteq \mathcal{S}$ and hence we must have that:

$$r(\mathcal{Q}) = \sum_{\mathcal{S} \in \mathcal{V}_i} \mathbb{I}[\mathcal{Q} \subseteq \mathcal{S}] \tilde{r}_{i-1}(\mathcal{S}),$$

so, by the inductive hypothesis, the restriction of \tilde{r}_{i-1} to \mathcal{V}_i is equal to \tilde{r}_i . Now choose some arbitrary $\mathcal{Q} \in \mathcal{V}_{i-1} \setminus \mathcal{V}_i$ and define:

$$v(\mathcal{Q}) := \sum_{\mathcal{S} \in \mathcal{V}_i} \mathbb{I}[\mathcal{Q} \subseteq \mathcal{S}] \tilde{r}_{i-1}(\mathcal{S})$$

which, by above, is uniquely defined. Note that for all $\mathcal{S} \in \mathcal{V}_{i-1} \setminus \mathcal{V}_i$ we have that $\mathcal{Q} \subseteq \mathcal{S}$ if and only if $\mathcal{S} = \mathcal{Q}$ and hence we must have that:

$$r(\mathcal{Q}) = \sum_{\mathcal{S} \in \mathcal{V}_i} \mathbb{I}[\mathcal{Q} \subseteq \mathcal{S}] \tilde{r}_{i-1}(\mathcal{S}) + \tilde{r}_{i-1}(\mathcal{Q}) = v(\mathcal{Q}) + \tilde{r}_{i-1}(\mathcal{Q}),$$

so that $\tilde{r}_{i-1}(\mathcal{Q}) = r(\mathcal{Q}) - v(\mathcal{Q})$ which is unique.

We have hence shown that the inductive hypothesis holds for $k = i - 1$ and hence holds always. Noting that $\mathcal{V}_0 = 2^{[K]}$ we then get the result by necessarily setting $\tilde{r} = \tilde{r}_0$. \square

A.2. Lemma 6.5

Lemma A.2. *For all $r : 2^K \rightarrow \mathbb{R}$ and all $\mathbf{q} \in \Delta_K$ we have $\Phi^r(\mathbf{q}) = \mathbb{E}[r(\mathcal{B}(\mathbf{q}))]$.*

Proof. Let \tilde{r} be a subset decomposition of r . We have

$$\begin{aligned} \mathbb{E}[r(\mathcal{B}(\mathbf{q}))] &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{P}[\mathcal{B}(\mathbf{q}) \subseteq \mathcal{S}] \\ &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \prod_{j \in [M]} \mathbb{P}[b_j(\mathbf{q}) \in \mathcal{S}] \\ &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \prod_{j \in [M]} \sum_{i \in [K]} \mathbb{I}[i \in \mathcal{S}] q_i \\ &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \left(\sum_{i \in [K]} \mathbb{I}[i \in \mathcal{S}] q_i \right)^M \\ &= \Phi^r(\mathbf{q}) \end{aligned}$$

as required. \square

A.3. Lemma 6.6

Lemma A.3. For all $r : 2^K \rightarrow \mathbb{R}$, all $\mathbf{q} \in \Delta_K$ and all $i \in [K]$ we have

$$\partial_i \Phi^r(\mathbf{q}) = \mathbb{E} \left[\frac{r(\mathcal{B}(\mathbf{q}))}{q_i} \sum_{j \in [M]} \mathbb{1}[b_j(\mathbf{q}) = i] \right].$$

Proof. Let \tilde{r} be a subset decomposition of r . For all $\mathbf{q}' \in \mathbb{R}^K$ and $\mathcal{S} \subseteq [K]$ define

$$\Lambda^{\mathcal{S}}(\mathbf{q}') := \left(\sum_{k \in [K]} \mathbb{1}[k \in \mathcal{S}] q'_k \right)^M.$$

Fix some $j \in [M]$. Note that

$$\begin{aligned} \partial_i \Lambda^{\mathcal{S}}(\mathbf{q}) &= M \mathbb{1}[i \in \mathcal{S}] \left(\sum_{k \in [K]} \mathbb{1}[k \in \mathcal{S}] q_k \right)^{M-1} \\ &= M \mathbb{1}[i \in \mathcal{S}] \prod_{j' \in [M] \setminus \{j\}} \sum_{k \in [K]} \mathbb{1}[k \in \mathcal{S}] q_k \\ &= M \mathbb{1}[i \in \mathcal{S}] \prod_{j' \in [M] \setminus \{j\}} \mathbb{P}[b_{j'}(\mathbf{q}) \in \mathcal{S}] \\ &= \frac{M}{q_i} \mathbb{P}[b_j(\mathbf{q}) = i] \mathbb{1}[i \in \mathcal{S}] \prod_{j' \in [M] \setminus \{j\}} \mathbb{P}[b_{j'}(\mathbf{q}) \in \mathcal{S}] \\ &= \frac{M}{q_i} \mathbb{P}[(b_j(\mathbf{q}) = i) \wedge (i \in \mathcal{S})] \prod_{j' \in [M] \setminus \{j\}} \mathbb{P}[b_{j'}(\mathbf{q}) \in \mathcal{S}] \\ &= \frac{M}{q_i} \mathbb{P}[(b_j(\mathbf{q}) = i) \wedge (b_j(\mathbf{q}) \in \mathcal{S})] \prod_{j' \in [M] \setminus \{j\}} \mathbb{P}[b_{j'}(\mathbf{q}) \in \mathcal{S}] \\ &= \frac{M}{q_i} \mathbb{P}[(b_j(\mathbf{q}) = i) \wedge (\forall j' \in [M], b_{j'}(\mathbf{q}) \in \mathcal{S})] \\ &= \frac{M}{q_i} \mathbb{P}[(b_j(\mathbf{q}) = i) \wedge (\mathcal{B}(\mathbf{q}) \subseteq \mathcal{S})] \\ &= \frac{M}{q_i} \mathbb{E}[\mathbb{1}[b_j(\mathbf{q}) = i] \mathbb{1}[\mathcal{B}(\mathbf{q}) \subseteq \mathcal{S}]], \end{aligned}$$

so since:

$$\Phi^r(\mathbf{q}) = \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \Lambda^{\mathcal{S}}(\mathbf{q}),$$

we have

$$\begin{aligned} \partial_i \Phi^r(\mathbf{q}) &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \partial_i \Lambda^{\mathcal{S}}(\mathbf{q}) \\ &= \frac{M}{q_i} \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{E}[\mathbb{1}[b_j(\mathbf{q}) = i] \mathbb{1}[\mathcal{B}(\mathbf{q}) \subseteq \mathcal{S}]] \\ &= \frac{M}{q_i} \mathbb{E} \left[\mathbb{1}[b_j(\mathbf{q}) = i] \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{1}[\mathcal{B}(\mathbf{q}) \subseteq \mathcal{S}] \right] \\ &= \frac{M}{q_i} \mathbb{E}[\mathbb{1}[b_j(\mathbf{q}) = i] r(\mathcal{B}(\mathbf{q}))]. \end{aligned}$$

605 Summing over all $j \in [M]$ and dividing by M then gives us

$$606 \quad \partial_i \Phi^r(\mathbf{q}) = \mathbb{E} \left[\frac{r(\mathcal{B}(\mathbf{q}))}{q_i} \sum_{j \in [M]} \mathbb{1}[b_j(\mathbf{q}) = i] \right]$$

607 as required. □

612 A.4. Lemma 6.7

613 **Lemma A.4.** For all pseudo-concave set functions $r : 2^K \rightarrow \mathbb{R}$ we have that Φ^r is concave over the simplex Δ_K .

614 *Proof.* Choose any $\mathbf{q} \in \Delta_K$. Define $\langle b'_j(\mathbf{q}) \mid j \in [M-2] \rangle$ to be a sequence of stochastic quantities drawn i.i.d. at
615 random from (the probability distribution characterised by) \mathbf{q} . In addition, let:

$$616 \quad \mathcal{B}'(\mathbf{q}) := \{b'_j(\mathbf{q}) \mid j \in [M-2]\}.$$

617 Direct from the definition of Φ^r we have, for all $i, i' \in [K]$, that

$$\begin{aligned} 618 \quad \partial_i \partial_{i'} \Phi^r &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{1}[i \in \mathcal{S}] \mathbb{1}[i' \in \mathcal{S}] \left(\sum_{k \in [K]} \mathbb{1}[k \in \mathcal{S}] q_k \right)^{M-2} \\ 619 &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{1}[i \in \mathcal{S}] \mathbb{1}[i' \in \mathcal{S}] \prod_{j \in [M-2]} \sum_{k \in [K]} \mathbb{1}[k \in \mathcal{S}] q_k \\ 620 &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{1}[i \in \mathcal{S}] \mathbb{1}[i' \in \mathcal{S}] \prod_{j \in [M-2]} \mathbb{P}[b'_j(\mathbf{q}) \in \mathcal{S}] \\ 621 &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{1}[i \in \mathcal{S}] \mathbb{1}[i' \in \mathcal{S}] \mathbb{P}[\mathcal{B}'(\mathbf{q}) \subseteq \mathcal{S}] \\ 622 &= \sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{P}[\mathcal{B}'(\mathbf{q}) \cup \{i, i'\} \subseteq \mathcal{S}] \\ 623 &= \mathbb{E} \left[\sum_{\mathcal{S} \subseteq [K]} \tilde{r}(\mathcal{S}) \mathbb{1}[\mathcal{B}'(\mathbf{q}) \cup \{i, i'\} \subseteq \mathcal{S}] \right] \\ 624 &= \mathbb{E}[r(\mathcal{B}'(\mathbf{q}) \cup \{i, i'\})] \\ 625 &= \sum_{\mathcal{S} \subseteq [K]} \mathbb{P}[\mathcal{B}'(\mathbf{q}) = \mathcal{S}] r(\mathcal{S} \cup \{i, i'\}) \\ 626 &= \sum_{\mathcal{S} \subseteq [K]} \mathbb{P}[\mathcal{B}'(\mathbf{q}) = \mathcal{S}] U_{i, i'}^{r, \mathcal{S}}. \end{aligned}$$

627 So for all $\mathbf{x} \in \mathbb{R}^K$ with $\mathbf{x} \cdot \mathbf{1} = 0$ we have

$$\begin{aligned} 628 \quad \mathbf{x}^\top (\nabla^2 \Phi^r(\mathbf{q})) \mathbf{x} &= \sum_{i, i' \in [K]} x_i (\partial_i \partial_{i'} \Phi^r) x_{i'} \\ 629 &= \sum_{i, i' \in [K]} x_i x_{i'} \sum_{\mathcal{S} \subseteq [K]} \mathbb{P}[\mathcal{B}'(\mathbf{q}) = \mathcal{S}] U_{i, i'}^{r, \mathcal{S}} \\ 630 &= \sum_{\mathcal{S} \subseteq [K]} \mathbb{P}[\mathcal{B}'(\mathbf{q}) = \mathcal{S}] \sum_{i, i' \in [K]} x_i U_{i, i'}^{r, \mathcal{S}} x_{i'} \\ 631 &= \sum_{\mathcal{S} \subseteq [K]} \mathbb{P}[\mathcal{B}'(\mathbf{q}) = \mathcal{S}] (\mathbf{x}^\top U^{r, \mathcal{S}} \mathbf{x}) \\ 632 &\leq 0, \end{aligned}$$

633 which means that Φ^r is concave on Δ_K as required. □

A.5. Lemma 6.10

Lemma A.5. *Let $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$, $r : 2^{[K]} \rightarrow \mathbb{R}$ be a pseudo-submodular function, and $\mathcal{Z} \subseteq [K]$ be a set formed by drawing M elements uniformly at random (with replacement) from \mathcal{S} . Then we have*

$$\mathbb{E}[r(\mathcal{Z}) - r(\emptyset)] \geq \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) (r(\mathcal{S}) - r(\emptyset)).$$

Proof. Without loss of generality assume that $r(\emptyset) = 0$.

We prove by induction on m that the lemma holds whenever $M \leq m$. In the case that $m = 0$ we have $\mathbb{E}[r(\mathcal{Z})] = r(\emptyset) = 0$ and $M = 0$ so the result holds. Now assume that it holds for all $M \leq m$ and consider the case that $M = m + 1$.

Since r is pseudo-submodular choose $i \in \mathcal{S}$ such that

$$r(\mathcal{Q} \cup \{i\}) - r(\mathcal{Q}) \geq r(\mathcal{S}) - r(\mathcal{S} \setminus \{i\}) \quad (5)$$

for all $\mathcal{Q} \subseteq \mathcal{S} \setminus \{i\}$. Define $\sigma := |\mathcal{S}|$ and

$$\phi := r(\mathcal{S}) - r(\mathcal{S} \setminus \{i\}). \quad (6)$$

Let $\langle z_s \mid s \in [M] \rangle$ be a sequence of M elements drawn uniformly at random from \mathcal{S} such that $\mathcal{Z} = \{z_s \mid s \in [M]\}$. Define

$$\mu := \sum_{s \in [M]} \mathbb{1}[z_s \neq i].$$

For all $j \in [M] \cup \{0\}$ let \mathcal{Z}_j be a set formed by sampling j actions independently and uniformly at random from $\mathcal{S} \setminus \{i\}$.

Note that by the inductive hypothesis, we have

$$\begin{aligned} \mathbb{E}[\mathbb{1}[i \notin \mathcal{Z}]r(\mathcal{Z})] &= \mathbb{P}[i \notin \mathcal{Z}] \mathbb{E}[r(\mathcal{Z}) \mid i \notin \mathcal{Z}] \\ &= \mathbb{P}[i \notin \mathcal{Z}] \left(1 - \left(\frac{|\mathcal{S} \setminus \{i\}| - 1}{|\mathcal{S} \setminus \{i\}|}\right)^M\right) r(\mathcal{S} \setminus \{i\}) \\ &= \mathbb{P}[i \notin \mathcal{Z}] \left(1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^M\right) r(\mathcal{S} \setminus \{i\}) \\ &= \mathbb{P}[\mu = M] \left(1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^M\right) r(\mathcal{S} \setminus \{i\}). \end{aligned} \quad (7)$$

Note also that

$$\mathbb{E}[\mathbb{1}[i \in \mathcal{Z}]r(\mathcal{Z})] = \sum_{j \in [m] \cup \{0\}} \mathbb{P}[\mu = j] \mathbb{E}[r(\mathcal{Z}_j \cup \{i\})]. \quad (8)$$

By equations (5) and (6) and the inductive hypothesis we have, for all $j \in [m] \cup \{0\}$, that

$$\begin{aligned} \mathbb{E}[r(\mathcal{Z}_j \cup \{i\})] &\geq \mathbb{E}[\phi + r(\mathcal{Z}_j)] \\ &= \phi + \mathbb{E}[r(\mathcal{Z}_j)] \\ &\geq \phi + \left(1 - \left(\frac{|\mathcal{S} \setminus \{i\}| - 1}{|\mathcal{S} \setminus \{i\}|}\right)^j\right) r(\mathcal{S} \setminus \{i\}) \\ &= \phi + \left(1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^j\right) r(\mathcal{S} \setminus \{i\}). \end{aligned} \quad (9)$$

We also have that

$$\sum_{j \in [m] \cup \{0\}} \mathbb{P}[\mu = j] = \mathbb{P}[i \in \mathcal{Z}]. \quad (10)$$

Substituting equations (9) and (10) into Equation (8) gives us

$$\mathbb{E}[\mathbb{1}[i \in \mathcal{Z}]r(\mathcal{Z})] = \mathbb{P}[i \in \mathcal{Z}]\phi + \sum_{j \in [m] \cup \{0\}} \mathbb{P}[\mu = j] \left(1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^j\right) r(\mathcal{S} \setminus \{i\}).$$

Adding this equation to Equation (7) gives us

$$\mathbb{E}[r(\mathcal{Z})] = \mathbb{P}[i \in \mathcal{Z}]\phi + \sum_{j \in [M] \cup \{0\}} \mathbb{P}[\mu = j] \left(1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^j\right) r(\mathcal{S} \setminus \{i\}). \quad (11)$$

Take any $k \in \mathcal{S} \setminus \{i\}$. Note that

$$1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^j = 1 - (1 - 1/(\sigma - 1))^j = 1 - \mathbb{P}[k \notin \mathcal{Z}_j] = \mathbb{P}[k \in \mathcal{Z}_j],$$

so that

$$\begin{aligned} \sum_{j \in [M] \cup \{0\}} \mathbb{P}[\mu = j] \left(1 - \left(\frac{\sigma - 2}{\sigma - 1}\right)^j\right) &= \sum_{j \in [M] \cup \{0\}} \mathbb{P}[\mu = j] \mathbb{P}[k \in \mathcal{Z}_j] \\ &= \sum_{j \in [M] \cup \{0\}} \mathbb{P}[\mu = j] \mathbb{P}[k \in \mathcal{Z} \setminus \{i\} \mid \mu = j] \\ &= \sum_{j \in [M] \cup \{0\}} \mathbb{P}[\mu = j] \mathbb{P}[k \in \mathcal{Z} \mid \mu = j] \\ &= \mathbb{P}[k \in \mathcal{Z}]. \end{aligned}$$

Substituting into Equation (11) gives us:

$$\begin{aligned} \mathbb{E}[r(\mathcal{Z})] &\geq \mathbb{P}[i \in \mathcal{Z}]\phi + \mathbb{P}[k \in \mathcal{Z}]r(\mathcal{S} \setminus \{i\}) \\ &= \mathbb{P}[i \in \mathcal{Z}](\phi + r(\mathcal{S} \setminus \{i\})) \\ &= \mathbb{P}[i \in \mathcal{Z}]r(\mathcal{S}) \\ &= (1 - \mathbb{P}[i \notin \mathcal{Z}])r(\mathcal{S}) \\ &= (1 - (1 - 1/\sigma)^M)r(\mathcal{S}) \\ &= \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) r(\mathcal{S}). \end{aligned}$$

So the inductive hypothesis holds for all $M \in [m + 1]$ and hence holds always. \square

A.6. Lemma 6.11

Lemma A.6. *Given any $\mathcal{S} \subseteq [K]$ and any pseudo-submodular set function $r : 2^{[K]} \rightarrow \mathbb{R}$ we have*

$$\Phi^r(\mathbf{p}^{\mathcal{S}}) \geq r(\emptyset) + \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) (r(\mathcal{S}) - r(\emptyset)).$$

Proof. Let \mathcal{Z} be a set formed by drawing M elements i.i.d. with replacement from \mathcal{S} . Let z be an element drawn i.i.d. from \mathcal{S} . Let \tilde{r} be a subset-decomposition of r . Note that for all $i \in [K]$ we have

$$p_i^{\mathcal{S}} = \mathbb{P}[z = i].$$

770 Hence, we have

$$\begin{aligned}
 771 \quad \Phi^r(\mathbf{p}^S) &= \sum_{\mathcal{Q} \subseteq [K]} \tilde{r}(\mathcal{Q}) \left(\sum_{i \in [K]} \mathbb{1}[i \in \mathcal{Q}] p_i^S \right)^M \\
 772 \quad &= \sum_{\mathcal{Q} \subseteq [K]} \tilde{r}(\mathcal{Q}) \left(\sum_{i \in [K]} \mathbb{1}[i \in \mathcal{Q}] \mathbb{P}[z = i] \right)^M \\
 773 \quad &= \sum_{\mathcal{Q} \subseteq [K]} \tilde{r}(\mathcal{Q}) \mathbb{P}[z \in \mathcal{Q}]^M \\
 774 \quad &= \sum_{\mathcal{Q} \subseteq [K]} \tilde{r}(\mathcal{Q}) \mathbb{P}[\mathcal{Z} \subseteq \mathcal{Q}] \\
 775 \quad &= \sum_{\mathcal{Q} \subseteq [K]} \tilde{r}(\mathcal{Q}) \mathbb{E}[\mathbb{1}[\mathcal{Z} \subseteq \mathcal{Q}]] \\
 776 \quad &= \mathbb{E} \left[\sum_{\mathcal{Q} \subseteq [K]} \tilde{r}(\mathcal{Q}) \mathbb{1}[\mathcal{Z} \subseteq \mathcal{Q}] \right] \\
 777 \quad &= \mathbb{E}[r(\mathcal{Z})].
 \end{aligned}$$

778 So

$$779 \quad \Phi^r(\mathbf{p}^S) - r(\emptyset) = \mathbb{E}[r(\mathcal{Z}) - r(\emptyset)],$$

780 Lemma 6.10 then gives us the result. □

781 A.7. Lemma 6.13

782 **Lemma A.7.** For all trials $t \in [T]$ and Ψ_t as defined in Definition 6.12, we have

$$783 \quad \nabla \Psi_t(\mathbf{p}_t) = \mathbb{E}[\mathbf{g}_t | \mathbf{p}_t],$$

784 *Proof.* Take any $i \in [K]$. For any $j \in [M]$ we have

$$\begin{aligned}
 785 \quad c_{t,i} &= p_{t,i} c_{t,i} / p_{t,i} \\
 786 \quad &= \mathbb{P}[a_{t,j} = i | \mathbf{p}_t] c_{t,i} / p_{t,i} \\
 787 \quad &= \mathbb{E}[\mathbb{1}[a_{t,j} = i] c_{t,i} / p_{t,i} | \mathbf{p}_t].
 \end{aligned}$$

788 So:

$$\begin{aligned}
 789 \quad M c_{t,i} &= \sum_{j \in [M]} \mathbb{E}[\mathbb{1}[a_{t,j} = i] c_{t,i} / p_{t,i} | \mathbf{p}_t] \\
 790 \quad &= \mathbb{E} \left[\frac{c_{t,j}}{p_{t,i}} \sum_{j \in [M]} \mathbb{1}[a_{t,j} = i] \middle| \mathbf{p}_t \right].
 \end{aligned}$$

791 Hence, by Lemma 6.6, we have

$$\begin{aligned}
 792 \quad \partial_i \Psi_t(\mathbf{p}_t) &= \partial_i \Phi^{r_t}(\mathbf{p}_t) - M c_{t,i} \\
 793 \quad &= \mathbb{E} \left[\frac{r_t(\mathcal{B}(\mathbf{p}_t))}{p_{t,i}} \sum_{j \in [M]} \mathbb{1}[b_j(\mathbf{p}_t) = i] \right] - M c_{t,i} \\
 794 \quad &= \mathbb{E} \left[\frac{r_t(\mathcal{A}_t)}{p_{t,i}} \sum_{j \in [M]} \mathbb{1}[a_{t,j} = i] \middle| \mathbf{p}_t \right] - M c_{t,i} \\
 795 \quad &= \mathbb{E}[\mathbf{g}_{t,i} | \mathbf{p}_t]
 \end{aligned}$$

796 as required. □

A.8. Lemma 6.14

Lemma A.8. For all trials $t \in [T]$ we have

$$\mathbb{E} \left[\sum_{i \in [K]} p_{t,i} g_{t,i}^2 \right] \leq (1+C)^2 M(K+M-1).$$

Proof. Given $i \in [K]$ we have that

$$\begin{aligned} \frac{\mathbb{E}[g_{t,i}^2]}{(1+C)^2} &= \frac{1}{(1+C)^2} \mathbb{E} \left[(r_t(\mathcal{A}_t) - c_{t,i})^2 \sum_{j,j' \in [M]} \frac{\mathbb{1}[a_{t,j} = i] \mathbb{1}[a_{t,j'} = i]}{p_{t,i}^2} \right] \\ &\leq \sum_{j,j' \in [M]} \mathbb{E} \left[\frac{\mathbb{1}[a_{t,j} = i] \mathbb{1}[a_{t,j'} = i]}{p_{t,i}^2} \right] \\ &= \sum_{j \in [M]} \mathbb{E} \left[\frac{\mathbb{1}[a_{t,j} = i]}{p_{t,i}^2} \right] + \sum_{j,j' \in [M]} \mathbb{1}[j \neq j'] \mathbb{E} \left[\frac{\mathbb{1}[a_{t,j} = i] \mathbb{1}[a_{t,j'} = i]}{p_{t,i}^2} \right] \\ &= \sum_{j \in [M]} \frac{\mathbb{P}[a_{t,j} = i]}{p_{t,i}^2} + \sum_{j,j' \in [M]} \mathbb{1}[j \neq j'] \frac{\mathbb{P}[a_{t,j} = i] \mathbb{P}[a_{t,j'} = i]}{p_{t,i}^2} \\ &= \sum_{j \in [M]} \frac{1}{p_{t,i}} + \sum_{j,j' \in [M]} \mathbb{1}[j \neq j'] \\ &= \frac{M}{p_{t,i}} + M(M-1), \end{aligned}$$

and hence

$$\mathbb{E} \left[\sum_{i \in [K]} p_{t,i} g_{t,i}^2 \right] = \sum_{i \in [K]} p_{t,i} \mathbb{E}[g_{t,i}^2] \leq (1+C)^2 M(K+M-1)$$

as required. \square

A.9. Theorem 4.2

Theorem 4.2. Given r_t is pseudo-concave and pseudo-submodular for all $t \in [T]$, then for any set $\mathcal{S} \subseteq [K]$ with $\mathcal{S} \neq \emptyset$ we have

$$\sum_{t \in [T]} \mathbb{E}[\psi_t] \geq (1 - \alpha^M) \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) - \frac{M}{|\mathcal{S}|} \sum_{t \in [T]} \gamma_t(\mathcal{S}) - R,$$

where

$$\alpha := 1 - \frac{1}{|\mathcal{S}|}.$$

Proof. Consider some trial $t \in [T]$. By Lemma 6.7 and the definition of Ψ_t we have that Ψ_t is concave over Δ_K . Hence, by Lemma 6.13, we have

$$\begin{aligned} \mathbb{E}[(\mathbf{p}^{\mathcal{S}} - \mathbf{p}_t) \cdot \mathbf{g}_t | \mathbf{p}_t] &= (\mathbf{p}^{\mathcal{S}} - \mathbf{p}_t) \cdot \mathbb{E}[\mathbf{g}_t | \mathbf{p}_t] \\ &= (\mathbf{p}^{\mathcal{S}} - \mathbf{p}_t) \cdot \nabla \Psi_t(\mathbf{p}_t) \\ &\geq \Psi_t(\mathbf{p}^{\mathcal{S}}) - \Psi_t(\mathbf{p}_t). \end{aligned} \tag{12}$$

880 Lemma 6.11 gives us:

$$\begin{aligned}
 881 & \\
 882 & \Psi_t(\mathbf{p}^S) = \Phi^{r_t}(\mathbf{p}^S) - M \sum_{i \in [K]} p_{t,i}^S c_{t,i} \\
 883 & \\
 884 & \geq r(\emptyset) + \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) (r(\mathcal{S}) - r(\emptyset)) - \frac{M}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} c_{t,i} \\
 885 & \\
 886 & \end{aligned} \tag{13}$$

887 and Lemma 6.5 gives us:

$$\begin{aligned}
 889 & \Psi_t(\mathbf{p}_t) = \Phi^{r_t}(\mathbf{p}_t) - M \sum_{i \in [K]} p_{t,i} c_{t,i} \\
 890 & \\
 891 & = \mathbb{E}[r_t(\mathcal{B}(\mathbf{p}_t))] - M \sum_{i \in [K]} p_{t,i} c_{t,i} \\
 892 & \\
 893 & = \mathbb{E}[r_t(\mathcal{A}_t) | \mathbf{p}_t] - M \sum_{i \in [K]} p_{t,i} c_{t,i} \\
 894 & \\
 895 & = \mathbb{E}[r_t(\mathcal{A}_t) | \mathbf{p}_t] - \sum_{j \in [M]} \sum_{i \in [K]} \mathbb{P}[a_{t,j} = i | \mathbf{p}_t] c_{t,i} \\
 896 & \\
 897 & = \mathbb{E}[r_t(\mathcal{A}_t) | \mathbf{p}_t] - \sum_{j \in [M]} \mathbb{E}[c_{t,a_{t,j}} | \mathbf{p}_t] \\
 898 & \\
 899 & = \mathbb{E}[r_t(\mathcal{A}_t) | \mathbf{p}_t] - \mathbb{E}\left[\sum_{j \in [M]} c_{t,a_{t,j}} \mid \mathbf{p}_t\right] \\
 900 & \\
 901 & \leq \mathbb{E}[r_t(\mathcal{A}_t) | \mathbf{p}_t] - \mathbb{E}\left[\sum_{i \in \mathcal{A}_t} c_{t,a_{t,j}} \mid \mathbf{p}_t\right] \\
 902 & \\
 903 & = \mathbb{E}[\psi_t | \mathbf{p}_t] + r_t(\emptyset). \\
 904 & \\
 905 & \\
 906 & \\
 907 & \\
 908 & \end{aligned} \tag{14}$$

909 Substituting equations (13) and (14) into Equation (12) gives us:

$$910 \mathbb{E}[(\mathbf{p}^S - \mathbf{p}_t) \cdot \mathbf{g}_t | \mathbf{p}_t] \geq -\mathbb{E}[\psi_t | \mathbf{p}_t] + \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) \hat{r}_t(\mathcal{S}) - \frac{M}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} c_{t,i}$$

911 and hence:

$$912 \mathbb{E}[(\mathbf{p}^S - \mathbf{p}_t) \cdot \mathbf{g}_t] \geq -\mathbb{E}[\psi_t] + \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) \hat{r}_t(\mathcal{S}) - \frac{M}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} c_{t,i}. \tag{15}$$

913 Lemma 6.14 gives us:

$$914 \mathbb{E}\left[\sum_{i \in [K]} p_{t,i} g_{t,i}^2\right] \leq \frac{R^2}{T}, \tag{16}$$

915 Lemma 6.8 gives us:

$$916 \sum_{t \in [T]} \mathbb{E}[(\mathbf{p}^S - \mathbf{p}_t) \cdot \mathbf{g}_t] \leq \frac{\ln(K)}{\eta} + \eta \sum_{t \in [T]} \mathbb{E}\left[\sum_{i \in [K]} p_{t,i} g_{t,i}^2\right]. \tag{17}$$

917 Substituting equations (15) and (16) into Equation (17) gives us:

$$918 - \sum_{t \in [T]} \mathbb{E}[\psi_t] + \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|}\right)^M\right) \sum_{t \in [T]} \hat{r}_t(\mathcal{S}) - \frac{M}{|\mathcal{S}|} \sum_{t \in [T]} \sum_{i \in \mathcal{S}} c_{t,i} \leq \frac{\ln(K)}{\eta} + \eta R^2.$$

919 Since $\eta = \ln(K)/R$ this implies the result. \square

B. PROOF OF THEOREM 4.3

Theorem 4.3. *Suppose that there exists some $d \in \mathbb{N}$, $s \in (0, 1)$, $\gamma > 1$, and a randomized polynomial time algorithm for the learner such that for all $K, M \in \mathbb{N}$ and for any instance of the multichannel advertising problem, it holds that $|\mathcal{A}_t| \leq M$ for all $t = 1, \dots, T$ and, for any subset $\mathcal{S} \subseteq [K]$,*

$$\mathbb{E} \left[\sum_{t \in [T]} r_t(\mathcal{A}_t) \right] \geq (1 - \alpha^{\gamma M}) \sum_{t \in [T]} r_t(\mathcal{S}) - R',$$

where $R' \in \mathcal{O}(K^d T^s)$ and $\alpha := 1 - \frac{1}{|\mathcal{S}|}$. Then, for all $\varepsilon \in (0, 1 - 1/\gamma)$ and $B > 4^{1/((1-\varepsilon)\gamma-1)}$, there exists a randomized polynomial-time algorithm for the set cover problem on $[B]$ that, with probability at least $\frac{1}{2}$, achieves approximation ratio at least $(1 - \varepsilon) \ln(B)$.

Proof. Suppose we have such an algorithm. Let $c > 0$ and $\gamma > 1$ be such that

$$\mathbb{E} \left[\sum_{t \in [T]} r_t(\mathcal{A}_t) \right] \geq \left(1 - \left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|} \right)^{\gamma M} \right) \sum_{t \in [T]} r_t(\mathcal{S}) - cK^d T^s. \quad (18)$$

Choose any $\rho \in (1/\gamma, 1)$ and then consider any $B \in \mathbb{N}$ such that $B > 4^{1/(\rho\gamma-1)}$. Consider also any collection of sets $\{\mathcal{C}_k \mid k \in [K]\} \subseteq 2^{[B]}$ such that

$$\bigcup_{k \in [K]} \mathcal{C}_k = [B].$$

Let \mathcal{S} be a subset of $[K]$ of minimum cardinality such that

$$\bigcup_{k \in \mathcal{S}} \mathcal{C}_k = [B].$$

Now choose

$$T := \left\lceil (4cK^d B)^{1/(1-s)} \right\rceil.$$

and choose any $M \in \mathbb{N}$ such that $M \geq \rho \ln(B) |\mathcal{S}|$. For all $t \in [T]$ draw \mathcal{D}_t randomly as follows. First draw β_t uniformly at random from $[B]$ and then define

$$\mathcal{D}_t := \{k \in [K] \mid \beta_t \in \mathcal{C}_k\}.$$

It is a classic result that

$$\left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|} \right)^{|\mathcal{S}|} \leq e^{-1}.$$

so by the conditions on B and M we have

$$\left(\frac{|\mathcal{S}| - 1}{|\mathcal{S}|} \right)^{\gamma M} \leq \exp(-\gamma M / |\mathcal{S}|) = B^{-\rho\gamma} = \frac{B^{1-\rho\gamma}}{B} < \frac{1}{4B}. \quad (19)$$

By definition of \mathcal{S} we have, for all $t \in [T]$, that there exists some $k \in \mathcal{S}$ such that $\beta_t \in \mathcal{C}_k$ so that $\mathcal{D}_t \cap \mathcal{S} \neq \emptyset$. This implies

$$\sum_{t \in [T]} r_t(\mathcal{S}) = T,$$

and hence, by (18) and (19), we have

$$\mathbb{E} \left[\sum_{t \in [T]} (1 - r_t(\mathcal{A}_t)) \right] \leq T - T + \frac{T}{4B} + cK^d T^s \leq \frac{T}{4B} + \frac{TcK^d}{T^{1-s}} \leq \frac{T}{2B}. \quad (20)$$

Fix t and a realization of \mathcal{A}_t . If we have

$$\bigcup_{k \in \mathcal{A}_t} \mathcal{C}_k \neq [B],$$

then we must also have that

$$\begin{aligned} \mathbb{E}[1 - r_t(\mathcal{A}_t) \mid \mathcal{A}_t] &= \mathbb{P}[\mathcal{A}_t \cap \mathcal{D}_t = \emptyset \mid \mathcal{A}_t] \\ &= \mathbb{P}[\forall k \in \mathcal{A}_t, \beta_t \notin \mathcal{C}_k \mid \mathcal{A}_t] \\ &= \mathbb{P}\left[\beta_t \notin \bigcup_{k \in \mathcal{A}_t} \mathcal{C}_k \mid \mathcal{A}_t\right] \geq \frac{1}{B}. \end{aligned}$$

Hence, by taking the randomness of $\mathcal{A}_1, \dots, \mathcal{A}_T$ into account,

$$\begin{aligned} &\mathbb{P}\left[\sum_{t \in [T]} \mathbb{I}\left[\bigcup_{k \in \mathcal{A}_t} \mathcal{C}_k \neq [B]\right] = T\right] \\ &\leq \mathbb{P}\left[\mathbb{E}\left[\sum_{t \in [T]} (1 - r_t(\mathcal{A}_t)) \mid \mathcal{A}_1, \dots, \mathcal{A}_T\right] \geq \frac{T}{n}\right] \leq \frac{1}{2} \end{aligned}$$

by (20). Since T is polynomial in KB and $|\mathcal{A}_t| \leq M$, we have a randomized polynomial-time algorithm that, with probability at least $\frac{1}{2}$, solves the set cover problem on $[B]$ with approximation ratio $(1 - \varepsilon) \ln(B)$ for $\varepsilon = 1 - \rho \in (0, 1 - 1/\gamma)$. \square

C. SUBMODULAR MONOTONE NON-PSEUDOCONCAVE FUNCTIONS

We provide a function counterexample to show that there are submodular monotone functions which are not pseudoconcave.

Let $K = 8$, $\mathcal{P} = 2^{[K]}$, $\mathcal{S} = \{K\}$, and $\alpha > 0$. We define $U^{r, \mathcal{S}}$ as follows:

$$U^{r, \mathcal{S}} := \begin{pmatrix} 1 & 2 & 2 & 2 & 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 \\ 2 & 1 & 2 & 2 & 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 \\ 2 & 2 & 1 & 2 & 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 \\ 2 & 2 & 2 & 1 & 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 \\ 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 & 2 & 2 & 1 \\ 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 + \alpha & 2 & 1 & 2 & 1 \\ 1 + \alpha & 1 + \alpha & 1 + \alpha & 1 + \alpha & 2 & 2 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 \end{pmatrix}.$$

Now, let $\mathbf{x} = (1, 1, 1, 1, -1, -1, -1, -1)^\top$. Note that we have $\langle \mathbf{x}, \mathbf{1} \rangle = 0$ as required by the pseudoconcavity definition, and $\mathbf{x}^\top U^{r, \mathcal{S}} \mathbf{x} = 17 - 24\alpha$, which is positive for $\alpha \in (0, \frac{17}{24})$, implying therefore the non-pseudoconcavity of r for such values of α .

We now show how to define r starting from $U^{r, \mathcal{S}}$ in such a way that it is both monotone and submodular while being therefore also non-pseudoconcave.

We have $|\mathcal{P}| = 2^K = 256$ possible subsets as the arguments of r , 29 of which are already defined by the above matrix $U^{r, \mathcal{S}}$:

- 1 subset ($\{K\}$) with cardinality 1,
- 7 subsets ($\{j, K\}_{j \in [K-1]}$) with cardinality 2,

1045 • 21 subsets $(\{i, j, K\}_{1 \leq j < i \leq K-1})$ with cardinality 3.

1046
 1047 For any $i \in [K]$, let δ_i and Δ_i be equal respectively to the minimum and the maximum difference (gain) over
 1048 all values of r for subsets with cardinality i and all the ones for subsets with cardinality $i - 1$. As previously
 1049 anticipated, we construct function r starting from the above matrix $U^{r,S}$ in such a way that for all $i \in [K - 1]$,
 1050 we have

$$1051 \quad \delta_i \geq \Delta_{i+1},$$

1052 which is a sufficient condition for submodularity because, for all $i \in [K]$, *each* subset $S_i \in \mathcal{P}$ with cardinality i
 1053 can be generated by adding one of its element *only* from a subset $S_{i-1} \subset S_i$ with cardinality $i - 1$.

1054 We set $\alpha = \frac{2}{3} < \frac{17}{24}$, which guarantees the non-pseudoconcavity of r . To ensure monotonicity and submodularity,
 1055 we define

- 1056 • $r(S_1) := 0$ for all subsets $S_1 \in \mathcal{P}$ with $|S_1| = 1$ (consistently with $U_{K,K}^{r,S}$);
 - 1057 • $r(S_2) := 1$ for all subsets $S_2 \in \mathcal{P}$ with $|S_2| = 2$ (consistently with $U_{K,j}^{r,S}, U_{j,j}^{r,S}, U_{j,1}^{r,S}$ for all $j \in [K - 1]$);
 - 1058 • $r(S_3) := 1 + \frac{2}{3} = \frac{5}{3}$ for all subsets $S_3 \in \mathcal{P}$ with $|S_3| = 3$ that are not already defined by $U^{r,S}$;
 - 1059 • $r(S_4) := r(S_3) + \frac{1}{2} = \frac{5}{3} + \frac{1}{2} = 2 + \frac{1}{6} > \max_{i,j} U_{i,j}^{r,S} = 2$ for all subsets $S_4 \in \mathcal{P}$ with $|S_4| = 4$;
 - 1060 • $r(S_5) := r(S_4) + \frac{1}{6} = 2 + \frac{2}{6}$,
 - 1061 • $r(S_6) := r(S_5) + \frac{1}{6} = 2 + \frac{3}{6}$,
 - 1062 • $r(S_7) := r(S_6) + \frac{1}{6} = 2 + \frac{4}{6}$,
 - 1063 • $r(S_8) := r(S_7) + \frac{1}{6} = 2 + \frac{5}{6}$
- 1064 for all subsets $S_5, S_6, S_7, S_8 \in \mathcal{P}$ such that $|S_5| = 5, |S_6| = 6, |S_7| = 7, |S_8| = 8$.

1065 Finally, we also set $r(\emptyset) = -1$. Note that, to ensure that submodularity is not violated, for each subset S_3^U with
 1066 $|S_3^U| = 3$ defined by $U^{r,S}$, we have that the difference $r(S_3^U) - r(S_2)$ for any subset $S_2 \subset S_3^U$ with $|S_2| = 2$ is
 1067 either equal to $\alpha = \frac{2}{3}$ or 1, that is not smaller than the maximum difference $r(S_4) - r(S_3)$ over all $S_3, S_4 \in \mathcal{P}$
 1068 with $|S_3| = 3$ and $|S_4| = 4$, which in turn is equal to $\frac{1}{2} < \frac{2}{3}$. Furthermore, $r(S_4) = 2 + \frac{1}{6}$ is never smaller than
 1069 any values of $r(S_3^U)$ for all subsets $S_3^U \in \mathcal{P}$ with $|S_3| = 3$ that are already defined by $U^{r,S}$, because we have
 1070 $r(S_3^U) \leq 2$, thereby preserving monotonicity for all subsets in \mathcal{P} with cardinality smaller or equal to 4.

1071 Now, we recall that for any $i \in [K]$, δ_i and Δ_i are defined to be respectively equal to the minimum and the
 1072 maximum difference (gain) over all values of r for subsets with cardinality i and all the ones for subsets with
 1073 cardinality $i - 1$. Since we have

- 1074 • $\delta_1, \Delta_1, \delta_2, \Delta_2 = 1$ (which immediately implies $\Delta_2 \leq \delta_1$),
- 1075 • $\delta_3 = \frac{2}{3}; \quad \Delta_3 = 1 \leq \delta_2$,
- 1076 • $\delta_4 = \frac{1}{6}; \quad \Delta_4 = \frac{1}{2} \leq \delta_3$,
- 1077 • $\delta_5, \Delta_5, \delta_6, \Delta_6, \delta_7, \Delta_7, \delta_8, \Delta_8 = \frac{1}{6} \leq \delta_4$,

1078 then $\delta_i \geq \Delta_{i+1}$ for all $i \in [K - 1]$ which guarantees the submodularity of r . Finally, it is immediate to verify that
 1079 r is monotone also for all subsets in \mathcal{P} with cardinality larger than 4. Hence, we conclude that r is monotone
 1080 submodular and non-pseudoconcave.

1081 □

D. EFFICIENT IMPLEMENTATION OF COMBAND

To implement the algorithm the COMBAND presented in (Cesa-Bianchi and Lugosi, 2012), it is necessary to devise an efficient method for sampling from a set whose size can be exponential in K . In fact, at each trial, given a set \mathcal{S} of positive real numbers, we need to select any of the subsets S with a given size m from \mathcal{S} with a probability proportional to the product of the elements contained in S itself.

To be consistent with the notation used in (Cesa-Bianchi and Lugosi, 2012), henceforth we use the symbol d in place of K .

Given a set $\mathcal{S} = \{q_1, q_2, \dots, q_d\}$ of real positive numbers, we now show how to select a m -sized subset of indices $\{u_1, \dots, u_m\}$ from $[d]$ with a probability proportional to $\prod_{i=1}^m q_{u_i}$ by using dynamic programming. The running time of this sampling method is always linear² in $m \cdot d$.

For each sampling operation, consider the sequence of element indices u_1, u_2, \dots, u_m ordered according to the elements in $[d]$, i.e., $u_i < u_{i+1}$ for all $i \in [m-1]$.

The main idea of this method is to sample first u_m , and then u_{m-1}, \dots, u_1 (i.e., in reverse order) having derived in a preliminary phase via dynamic programming all the probabilities that $u_m = j$ for all $m \leq j \leq d$, and the conditional probabilities that $u_{m'} = j$ given that $u_{m'+1} = j'$, for all $m' \in [m-1]$ and $m' \leq j < j' \leq d - m + m'$.

We denote the conditional probability that $u_{m'} = j$ given that $u_{m'+1} = j'$, where $m' \in [m-1]$ and $m' \leq j < j' \leq d - m + m'$ by

$$P_{m',j|j'} := \mathbb{P}(u_{m'} = j | u_{m'+1} = j'),$$

and, for the selection of u_m , we define for all $j \in [d]$

$$P_{m,j} := \mathbb{P}(u_m = j),$$

because there is no element $u_{j'} > u_m$ (with $j' > m$) in the sequence of selected indices from $[d]$. We clearly have $\sum_{j=m'}^{j'-1} P_{m',j|j'} = 1$ and $\sum_{j=m}^d P_{m,j} = 1$.

For each $m' \in [m]$ and $m' \leq j \leq d - m + m'$ let $z_{m',j}$, be the the sum of the products of numbers of \mathcal{S} with indices $u_1, u_2, \dots, u_{m'}$ contained in each m' -sized subset of $[j]$ such that $u_{m'} = j$. We define $Z_{m',k} := \sum_{i=m'}^k z_{m',i}$ for any integer k such that $m' \leq k \leq d - m + m'$. Thus, for all $m' \in [m-1]$ and $m' \leq j < j' \leq d - m + m'$ we have

$$P_{m',j|j'} = \frac{z_{m',j}}{Z_{m',j'-1}}.$$

Analogously, for the selection of u_m , for all $m \leq j \leq d$ we can write

$$P_{m,j} = \frac{z_{m,j}}{Z_{m,d}}.$$

Hence, once we obtain $z_{m',j}$ and $Z_{m',j'-1}$ for all $m' \in [m-1]$ and $m' \leq j < j' \leq d - m + m'$, $z_{m,j}$ for all $m \leq j \leq d$, and $Z_{m,d}$, we can immediately compute the desired probabilities to sample u_m, u_{m-1}, \dots, u_1 in this (reverse) order.

We now show how to calculate these values. To this goal, since $Z_{m',k} := \sum_{i=m'}^k z_{m',i}$, we only need to show how to compute the values appearing at the numerator in the above probability formulas.

The possibility to efficiently the above probabilities is given by the following observation:

²We assume that multiplying two numbers requires a constant time. Removing this assumption, since it is known that it is possible to multiply two numbers represented by at most m bits in time equal to $\tilde{\mathcal{O}}(m)$ when $m \gg 1$ (Harvey and Van Der Hoeven, 2021), the total sampling time would be $\tilde{\mathcal{O}}(m^2 d)$ instead of $\mathcal{O}(md)$.

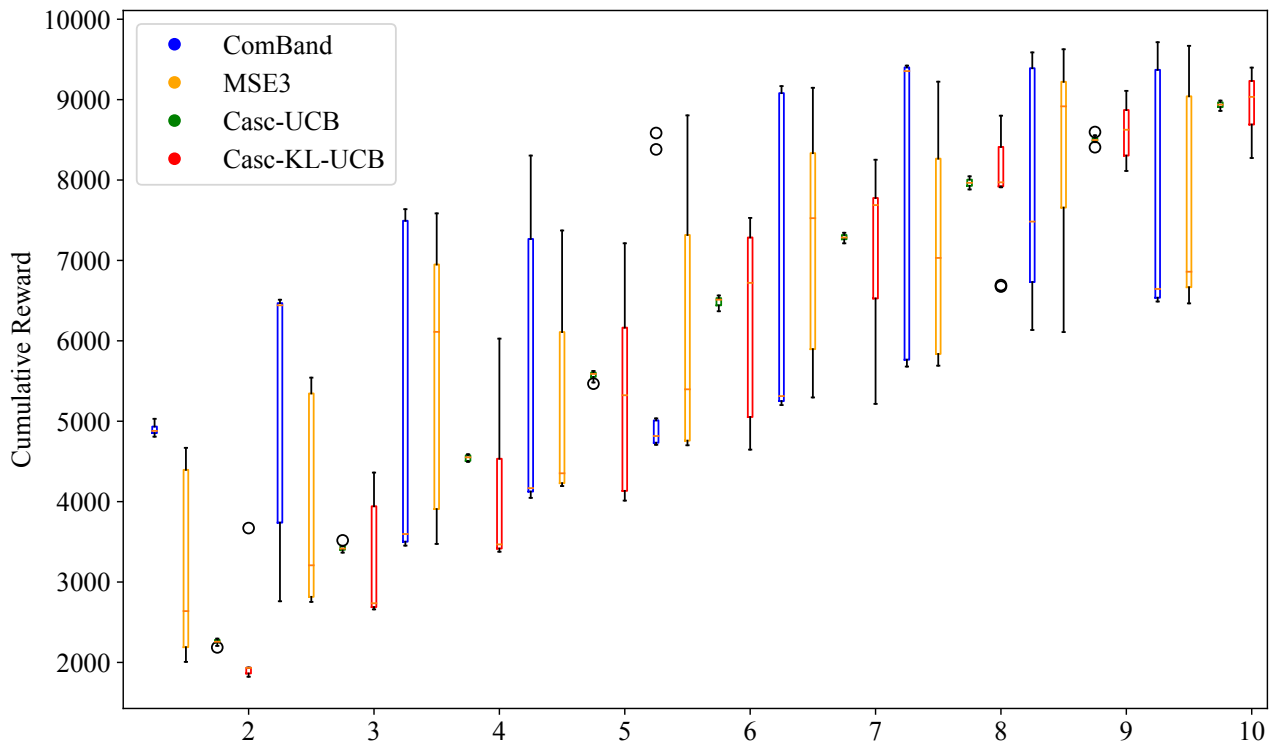


Figure 2. Stochastic environment, cumulative reward with respect to the M parameter

$$z_{m',j} = Z_{m'-1,j-1} \cdot q_j.$$

Note that $Z_{m'-1,j-1}$ can be in turn defined in terms of $z_{m'-1,m'-1}, z_{m'-1,m'}, z_{m'-1,m'+1}, \dots, z_{m'-1,j-2}, z_{m'-1,j-1}$. This recurrence relation allows us to compute all these values once we know $z_{1,1}, z_{1,2}, \dots, z_{1,d}$. Since we clearly have $z_{1,j} = q_j$ for all $j \in [d]$, we can therefore compute all these values and the above probabilities to efficiently accomplish this sampling operation by finding the indices u_m, u_{m-1}, \dots, u_1 in this order. It is immediate to verify that both the number of sum and multiplication operations are equal to $\Theta(md)$.

E. ADDITIONAL EXPERIMENTAL RESULTS

1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249

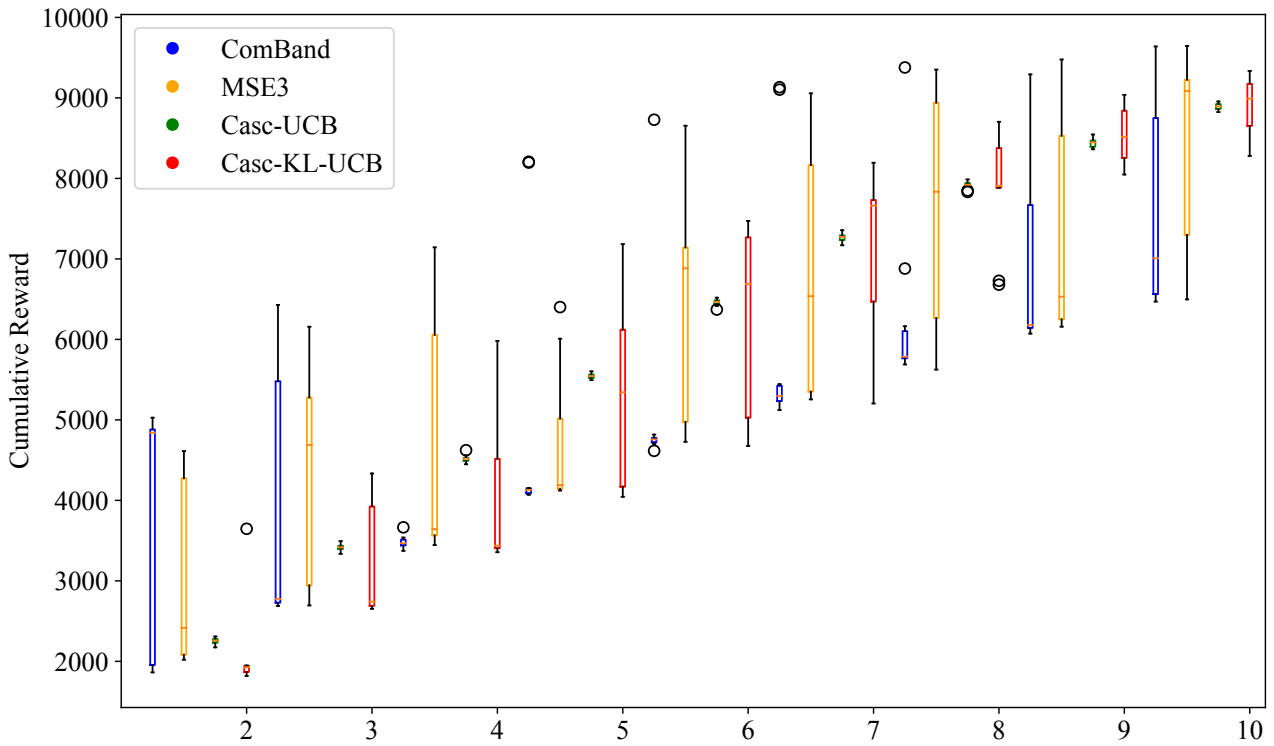


Figure 3. Stochastic with adversarial corruptions environment, cumulative reward with respect to the M parameter

1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264

1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319

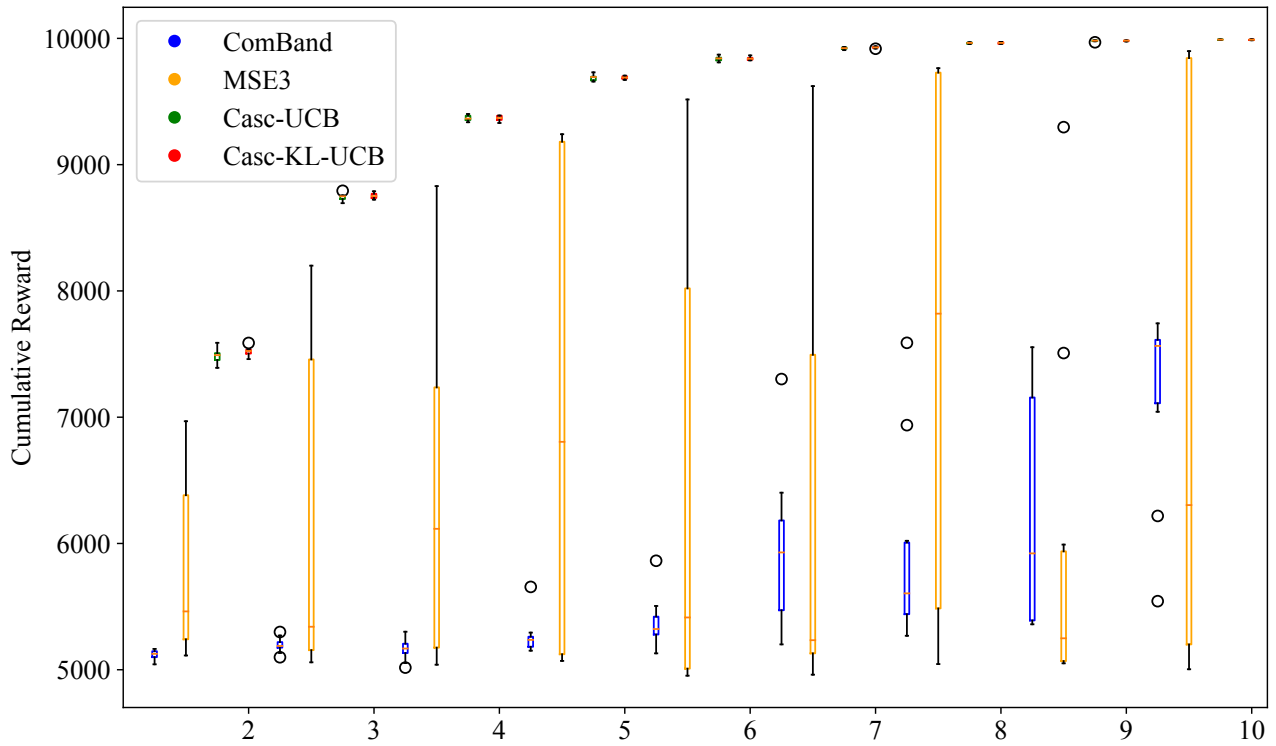


Figure 4. Worst case stochastic environment, cumulative reward with respect to the M parameter