

# Max-Sum Diversification, Monotone Submodular Functions and Dynamic Updates

[Extended Abstract] \*

Allan Borodin  
Department of Computer  
Science  
University of Toronto  
Toronto, Ontario, Canada  
M5S 3G4  
bor@cs.toronto.edu

Hyun Chul Lee  
Linkedin Corporation  
2025 Stierlin Court,  
Mountain View, CA, 94043  
culee@linkedin.com

Yuli Ye  
Department of Computer  
Science  
University of Toronto  
Toronto, Ontario, Canada  
M5S 3G4  
y3ye@cs.toronto.edu

## ABSTRACT

Result diversification has many important applications in databases, operations research, information retrieval, and finance. In this paper, we study and extend a particular version of result diversification, known as max-sum diversification. More specifically, we consider the setting where we are given a set of elements in a metric space and a set valuation function  $f$  defined on every subset. For any given subset  $S$ , the overall objective is a linear combination of  $f(S)$  and the sum of the distances induced by  $S$ . The goal is to find a subset  $S$  satisfying some constraints that maximizes the overall objective.

This problem is first studied by Gollapudi and Sharma in [17] for modular set functions and for sets satisfying a cardinality constraint (uniform matroids). In their paper, they give a 2-approximation algorithm by reducing to an earlier result in [20]. The first part of this paper considers an extension of the modular case to the monotone submodular case, for which the algorithm in [17] no longer applies. Interestingly, we are able to maintain the same 2-approximation using a natural, but different greedy algorithm. We then further extend the problem by considering any matroid constraint and show that a natural single swap local search algorithm provides a 2-approximation in this more general setting. This extends the Nemhauser, Wolsey and Fisher approximation result [29] for the problem of submodular function maximization subject to a matroid constraint (without the distance function component).

The second part of the paper focuses on dynamic updates for the modular case. Suppose we have a good initial approx-

imate solution and then there is a single weight-perturbation either on the valuation of an element or on the distance between two elements. Given that users expect some stability in the results they see, we ask how easy is it to maintain a good approximation without significantly changing the initial set. We measure this by the number of updates, where each update is a swap of a single element in the current solution with a single element outside the current solution. We show that we can maintain an approximation ratio of 3 by just a single update if the perturbation is not too large.

## Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Selection process*

## General Terms

Algorithm, Design, Performance, Theory

## Keywords

Diversification, information retrieval, ranking, submodular functions, matroids, greedy algorithm, local search, approximation algorithm, dynamic update

## 1. INTRODUCTION

The objective in many optimization problems is to find the “best” subset amongst a set of given items. While the definition of “best” is often vague, one common approach is to quantify the desired property for each element in the set and then select a subset of elements accordingly. Although this is a viable approach for many problems, for some applications, this does not yield good results. For example, in portfolio management, allocating equities only according to their expected returns might lead to a large potential risk as the portfolio is not diversified. A similar situation occurs in databases, for example, query result handling. When knowledge of the user’s intent is not fully available, it is actually better for a database system to diversify its displayed query results to improve user satisfaction. In many such scenarios, *diversity* is a necessary criterion.

We focus on a particular form of result diversification: max-sum diversification. We design algorithms for computing a “quality” subset, while also taking into account diver-

\*This research is supported by MITACS Accelerate program, Thoor Inc., Natural Sciences and Engineering Research Council of Canada and University of Toronto, Department of Computer Science.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$10.00.

sity which is defined as the sum of pairwise distances between elements of the set being returned. We also show how to gradually change such a set in a dynamically changing environment.

We consider the case where quality is measured by a monotone submodular function  $f(S)$  of the returned set  $S$ . In this way, we are extending beyond the linear (i.e., modular) case considered in [17]. Submodular functions have been extensively considered since they model many natural phenomena. For example, in terms of keyword based search in database systems, it is well understood that users begin to gradually (or sometimes abruptly) lose interest the more results they have to consider [37, 38]. But on the other hand, as long as a user continues to gain some benefit, additional query results can improve the overall quality but at a decreasing rate. As in [17], we consider the case of maximizing a linear combination of the quality  $f(S)$  and the (distance based) diversity subject to a cardinality constraint (i.e.,  $|S| \leq p$  for some given  $p$ ). We present a greedy algorithm that is somewhat unusual in that it does not try to optimize the objective in each iteration but rather optimizes a closely related potential function. We show that our greedy approach matches the 2-approximation [17] obtained for the modular case.

Our next result continues with the submodular case but now we go beyond a cardinality constraint (i.e., the uniform matroid) on  $S$  and allow the constraint to be that  $S$  is independent in a given matroid. This allows a substantial increase in generality. For example, while diversity might represent the distance between retrieved database tuples under a given criterion (for instance, a kernel based diversity measure called *answer tree kernel* is used in [43]), we could use a partition matroid to insure that (for example) the retrieved database tuples come from a variety of different sources. That is, we may wish to have  $n_i$  tuples from a specific database field  $i$ . This is, of course, another form of diversity but one orthogonal to diversity based on the given criterion. Similarly in the stock portfolio example, we might wish to have a balance of stocks in terms of say risk and profit margins (using some statistical measure of distances) while using a partition matroid to insure that different sectors of the economy are well represented. Another important class of matroids (relevant to our application) is that of transversal matroids. Suppose we have a collection  $\{C_1, C_2, \dots, C_m\}$  of (possibly) *overlapping* sets (i.e., the collection is not a partition) of database tuples (or stocks). Our goal might be to derive a set  $S$  such that the database tuples in  $S$  form a set of representatives for the collection; that is, every database tuple in  $S$  represents (and is in) a unique set  $C_i$  in the collection. The set  $S$  is then an independent set in the transversal matroid induced by the collection. We also note [35] that the intersection of any matroid with a uniform matroid is still a matroid so that in the above examples, we could further impose the constraint that the set  $S$  has at most  $p$  elements.

Our final theoretical result concerns dynamic updates. Here we restrict attention to a modular set function  $f(S)$ ; that is, we now have weights on the elements and  $f(S) = \sum_{u \in S} w(u)$  where  $w(u)$  is the weight of element  $u$ . This allows us to consider changes to the weight of a single element as well as changes to the distance function. We also mention some preliminary experiments on synthetic data that in the modular set function case suggest that our greedy and

local search algorithms may perform significantly better in practice than the proven worst case bounds.

The rest of the paper is organized as follows. In Section 2, we discuss related work in result diversification. We formulate the problem into a combinatorial optimization problem and show its connection to the dispersion problem in location theory in Section 3. In Section 4, we discuss max-sum diversification with monotone submodular set functions and give a simple greedy algorithm that achieves a 2-approximation when the set is of bounded cardinality. We extend the problem to the matroid case in Section 5 and discuss dynamic updates in Section 6. Section 7 carries out two preliminary experiments and Section 8 concludes the paper.

## 2. RELATED WORK

With the proliferation of today’s social media, database and web content, ranking becomes an important problem as it decides what gets selected and what does not, and what to be displayed first and what to be displayed last. Many early ranking algorithms, for example in web search, are based on the notion of “relevance”, i.e., the closeness of the object to the search query. However, there has been a rising interest to incorporate some notion of “diversity” into measures of quality.

One early work in this direction is the notion of “Maximal Marginal Relevance” (MMR) introduced by Carbonell and Goldstein in [6]. More specifically, MMR is defined as follows:

$$\text{MMR} = \max_{D_i \in R \setminus S} [\lambda \cdot \text{sim}_1(D_i, Q) - (1 - \lambda) \max_{D_j \in S} \text{sim}_2(D_i, D_j)],$$

where  $Q$  is a query;  $R$  is the ranked list of documents retrieved;  $S$  is the subset of documents in  $R$  already selected;  $\text{sim}_1$  is the similarity measure between a document and a query, and  $\text{sim}_2$  is the similarity measure between two documents. The parameter  $\lambda$  controls the trade-off between novelty (a notion of diversity) and relevance. The MMR algorithm iteratively selects the next document with respect to the MMR objective function until a given cardinality condition is met. The MMR heuristic has been widely used, but to the best of our knowledge, it has not been theoretically justified. Our paper provides some theoretical evidence why MMR is a legitimate approach for diversification. The greedy algorithm we propose in this paper can be viewed as a natural extension of MMR.

There is extensive research on how to diversify returned ranking results to satisfy multiple users. Namely, the result diversity issue occurs when many facets of queries are discovered and a set of multiple users expect to find their desired facets in the first page of the results. Thus, the challenge is to find the best strategy for ordering the results such that many users would find their relevant pages in the top few slots.

Rafei et al. [32] modeled this as a continuous optimization problem. They introduce a weight vector  $W$  for the search results, where the total weight sums to one. They define the portfolio variance to be  $W^T C W$ , where  $C$  is the covariance matrix of the result set. The goal then is to minimize the portfolio variance while the expected relevance is fixed at a certain level. They report that their proposed algorithm can improve upon Google in terms of the diversity on random queries, retrieving 14% to 38% more aspects of queries in

top five, while maintaining a precision very close to Google.

Bansal et al. [2] considered the setting in which various types of users exist and each is interested in a subset of the search results. They use a performance measure based on *discounted cumulative gain*, which defines the usefulness (gain) of a document as its position in the resulting list. Based on this measure, they suggest a general approach to develop approximation algorithms for ranking search results that captures different aspects of users’ intents. They also take into account that the relevance of one document cannot be treated independent of the relevance of other documents in a collection returned by a search engine. They consider both the scenario where users are interested in only a single search result (e.g., navigational queries) and the scenario where users have different requirements on the number of search results, and develop good approximation solutions for them.

The database community has recently studied the query diversification problem, which is mainly for keyword search in databases [27, 40, 12, 38, 43, 37, 10]. Given a very large database, an exploratory query can easily lead to a vast answer set. Typically, an answer’s relevance to the user query is based on *top-k* or *tf-idf*. As a way of increasing user satisfaction, different query diversification techniques have been proposed including some system based ones taking into account query parameters, evaluation algorithms, and dataset properties. For many of these, a max-sum type objective function is usually used.

Other than those discussed above, there are many recent papers studying result diversification in different settings, via different approaches and through different perspectives, for example [42, 9, 44, 41, 30, 1, 3, 34, 11, 36]. The reader is referred to [1, 13] for a good summary of the field. Most relevant to our work is the paper by Gollapudi and Sharma [17], where they develop an axiomatic approach to characterize and design diversification systems. Furthermore, they consider three different diversification objectives and using earlier results in facility dispersion, they are able to give algorithms with good approximation guarantees. This paper is a continuation of research along this line.

Recently, Minack et al. [28] have studied the problem of incremental diversification for very large data sets. Instead of viewing the input of the problem as a set, they consider the input as a stream, and use a simple online algorithm to process each element in an incremental fashion, maintaining a near-optimal diverse set at any point in the stream. Although their results are largely experimental, this approach significantly reduces CPU and memory consumption, and hence is applicable to large data sets. Our dynamic update algorithm deals with a problem of a similar nature, but instead of relying on experimental results, we prove theoretical guarantees. To the best of our knowledge, our work is the first of its kind to obtain a near-optimality condition for result diversification in a dynamically changing environment.

### 3. PROBLEM FORMULATION

Although the notion of “diversity” naturally arises in the context of databases, social media and web search, the underlying mathematical object is not new. As presented in [17], there is a rich and long line of research in location theory dealing with a similar concept; in particular, one objective is the placement of facilities on a network to maximize some function of the distances between facilities. The situation

arises when proximity of facilities is undesirable, for example, the distribution of business franchises in a city. Such location problems are often referred to as *dispersion* problems; for more motivation and early work, see [15, 16, 23].

Analytical models for the dispersion problem assume that the given network is represented by a set  $V = \{v_1, v_2, \dots, v_n\}$  of  $n$  vertices with metric distance between every pair of vertices. The objective is to locate  $p$  facilities ( $p \leq n$ ) among the  $n$  vertices, with at most one facility per vertex, such that some function of distances between facilities is maximized. Different objective functions are considered for the dispersion problems in the literature including: the max-sum criterion (maximize the total distances between all pairs of facilities) in [39, 15, 33], the max-min criterion (maximize the minimum distance between a pair of facilities) in [23, 15, 33], the max-mst (maximize the minimum spanning tree among all facilities) and many other related criteria in [18, 8]. The general problem (even in the metric case) for most of these criteria is NP-hard, and approximation algorithms have been developed and studied; see [8] for a summary of known results. Most relevant to this paper is the max-sum dispersion problem. The problem is known to be NP-hard [19], but it is not known whether or not it admits a PTAS. In [33], Ravi, Rosenkrantz and Tayi give a greedy algorithm and show it has an approximation ratio within a factor of 4. This is later improved by Hassin, Rubinstein and Tamir [20], who show a different algorithm with an approximation ratio of 2. This is the best known ratio today.

The dispersion problem is related to the diversification problem as both are trying to select a subset of elements which are element-wise far apart. The difference is that the diversification problem also considers vertex weight, so it is a bi-criteria optimization problem.

#### PROBLEM 1. Max-Sum Diversification

Let  $U$  be the underlying ground set, and let  $d(\cdot, \cdot)$  be a metric distance function on  $U$ . For any subset of  $U$ , let  $f(\cdot)$  be a non-negative set function measuring the value of a subset. Given a fixed integer  $p$ , the goal of the problem is to find a subset  $S \subseteq U$  that:

$$\begin{aligned} &\text{maximizes} && f(S) + \lambda \sum_{\{u,v\}:u,v \in S} d(u,v) \\ &\text{subject to} && |S| = p, \end{aligned}$$

where  $\lambda$  is a parameter specifying a desired trade-off between the two objectives.

The max-sum diversification problem is first proposed and studied in the context of result diversification in [17]<sup>1</sup>, where the function  $f(\cdot)$  is modular. In their paper, the value of  $f(S)$  measures the relevance of a given subset to a search query, and the value  $\sum_{\{u,v\}:u,v \in S} d(u,v)$  gives a diversity measure on  $S$ . The parameter  $\lambda$  specifies a desired trade-off between diversity and relevance. They reduce the problem to the max-sum dispersion problem, and using an algorithm in [20], they obtain an approximation ratio of 2.

In this paper, we first study the problem with more general valuation functions: normalized, monotone submodular set functions. For notational convenience, for any two sets  $S, T$  and an element  $e$ , we write  $S \cup \{e\}$  as  $S + e$ ,  $S \setminus \{e\}$  as  $S - e$ ,  $S \cup T$  as  $S + T$ , and  $S \setminus T$  as  $S - T$ . A set function

<sup>1</sup>In fact, they have a slightly different but equivalent formulation.

$f$  is *normalized* if  $f(\emptyset) = 0$ . The function is *monotone* if for any  $S, T \subseteq U$  and  $S \subseteq T$ ,

$$f(S) \leq f(T).$$

It is *submodular* if for any  $S, T \subseteq U$ ,  $S \subseteq T$  with  $u \in U$ ,

$$f(T + u) - f(T) \leq f(S + u) - f(S).$$

In the remainder of paper, all functions considered are normalized.

We proceed to our first contribution, a greedy algorithm (different than the one in [17]) that obtains a 2-approximation for monotone submodular set functions.

## 4. SUBMODULAR FUNCTIONS

Submodular set functions can be characterized by the property of a decreasing marginal gain as the size of the set increases. As such, submodular functions are well-studied objects in economics, game theory and combinatorial optimization. More recently, submodular functions have attracted attention in many practical fields of computer science. For example, Kempe et al. [21] study the problem of selecting a set of most influential nodes to maximize the total information spread in a social network. They have shown that under two basic diffusion models, the amount of influence of a set is submodular, hence the problem admits a good approximation algorithm. In natural language processing, Lin and Bilmes [26, 24, 25] have studied a class of submodular functions for document summarization. These functions each combine two terms, one which encourages the summary to be representative of the corpus, and the other which positively rewards diversity. Their experimental results show that a greedy algorithm with the objective of maximizing these submodular functions outperforms the existing state-of-art results in both generic and query-focused document summarization.

Both of the above mentioned results are based on the fundamental work of Nemhauser, Wolsey and Fisher [29], which has shown an  $\frac{e}{e-1}$ -approximation for maximizing monotone submodular set functions over a uniform matroid; and this bound is known to be tight even for a general matroid [5]. Our max-sum diversification problem with monotone submodular set functions can be viewed as an extension of that problem: the objective function now not only contains a submodular part, but also has a supermodular part: the sum of distances.

Since the max-sum diversification problem with modular set functions studied in [17] admits a 2-approximation algorithm, it is natural to ask what approximation ratio is obtainable for the same problem with monotone submodular set functions. Note that the algorithm in [17] does not apply to the submodular case. In what follows we assume (as is standard when considering submodular function) access to an oracle for finding an element  $u \in U - S$  that maximizes  $f(S + u) - f(S)$ . When  $f$  is modular, this simply means accessing the element  $u \in U - S$  having maximum weight.

**THEOREM 1.** *There is a simple linear time greedy algorithm that achieves a 2-approximation for the max-sum diversification problem with monotone submodular set functions satisfying a cardinality constraint.*

Before giving the proof of Theorem 1, we first introduce our notation. We extend the notion of distance function

to sets. For disjoint subsets  $S, T \subseteq U$ , we let  $d(S) = \sum_{\{u,v\}:u,v \in S} d(u,v)$ , and  $d(S, T) = \sum_{\{u,v\}:u \in S, v \in T} d(u,v)$ .

Now we define various types of marginal gain. For any given subset  $S \subseteq U$  and an element  $u \in U - S$ : let  $\phi(S)$  be the value of the objective function,  $d_u(S) = \sum_{v \in S} d(u,v)$  be the marginal gain on the distance,  $f_u(S) = f(S + u) - f(S)$  be the marginal gain on the weight, and  $\phi_u(S) = f_u(S) + \lambda d_u(S)$  be the total marginal gain on the objective function. Let  $f'_u(S) = \frac{1}{2}f_u(S)$ , and  $\phi'_u(S) = f'_u(S) + \lambda d_u(S)$ . We consider the following simple greedy algorithm:

GREEDY ALGORITHM

```

S = ∅
while |S| < p
    find u ∈ U - S maximizing φ'_u(S)
    S = S + u
end while
return S

```

Note that the above greedy algorithm is “non-oblivious” (in the sense of [22]) as it is not selecting the next element with respect to the objective function  $\phi(\cdot)$ . This might be of an independent interest. We utilize the following lemma in [33].

**LEMMA 1.** *Given a metric distance function  $d(\cdot, \cdot)$ , and two disjoint sets  $X$  and  $Y$ , we have the following inequality:*

$$(|X| - 1)d(X, Y) \geq |Y|d(X).$$

Now we are ready to prove Theorem 1.

**PROOF.** Let  $O$  be the optimal solution, and  $G$ , the greedy solution at the end of the algorithm. Let  $G_i$  be the greedy solution at the end of step  $i$ ,  $i < p$ ; and let  $A = O \cap G_i$ ,  $B = G_i - A$  and  $C = O - A$ . By lemma 1, we have the following three inequalities:

$$(|C| - 1)d(B, C) \geq |B|d(C) \tag{1}$$

$$(|C| - 1)d(A, C) \geq |A|d(C) \tag{2}$$

$$(|A| - 1)d(A, C) \geq |C|d(A) \tag{3}$$

Furthermore, we have

$$d(A, C) + d(A) + d(C) = d(O) \tag{4}$$

Note that the algorithm clearly achieves the optimal solution if  $p = 1$ . If  $|C| = 1$ , then  $i = p - 1$  and  $G_i \subset O$ . Let  $v$  be the element in  $C$ , and let  $u$  be the element taken by the greedy algorithm in the next step, then  $\phi'_u(G_i) \geq \phi'_v(G_i)$  for all  $v \in U - S$ . Therefore,

$$\frac{1}{2}f_u(G_i) + \lambda d_u(G_i) \geq \frac{1}{2}f_v(G_i) + \lambda d_v(G_i),$$

which implies

$$\begin{aligned} \phi_u(G_i) &= f_u(G_i) + \lambda d_u(G_i) \\ &\geq \frac{1}{2}f_u(G_i) + \lambda d_u(G_i) \\ &\geq \frac{1}{2}f_v(G_i) + \lambda d_v(G_i) \\ &\geq \frac{1}{2}\phi_v(G_i); \end{aligned}$$

and hence  $\phi(G) \geq \frac{1}{2}\phi(O)$ .

Now we can assume that  $p > 1$  and  $|C| > 1$ . We apply the following non-negative multipliers to equations (1), (2),

(3), (4) and add them:  $(1) * \frac{1}{|C|-1} + (2) * \frac{|C|-|B|}{p(|C|-1)} + (3) * \frac{i}{p(p-1)} + (4) * \frac{i|C|}{p(p-1)}$ ; we then have

$$d(A, C) + d(B, C) - \frac{i|C|(p - |C|)}{p(p-1)(|C|-1)}d(C) \geq \frac{i|C|}{p(p-1)}d(O).$$

Since  $p > |C|$ ,

$$d(C, G_i) \geq \frac{i|C|}{p(p-1)}d(O).$$

By submodularity and monotonicity of  $f'(\cdot)$ , we have

$$\sum_{v \in C} f'_v(G_i) \geq f'(C \cup G_i) - f'(G_i) \geq f'(O) - f'(G).$$

Therefore,

$$\begin{aligned} \sum_{v \in C} \phi'_v(G_i) &= \sum_{v \in C} [f'_v(G_i) + \lambda d(\{v\}, G_i)] \\ &= \sum_{v \in C} f'_v(G_i) + \lambda d(C, G_i) \\ &\geq [f'(O) - f'(G)] + \frac{\lambda i |C|}{p(p-1)}d(O). \end{aligned}$$

Let  $u_{i+1}$  be the element taken at step  $(i+1)$ , then we have

$$\phi'_{u_{i+1}}(G_i) \geq \frac{1}{p}[f'(O) - f'(G)] + \frac{\lambda i}{p(p-1)}d(O).$$

Summing over all  $i$  from 0 to  $p-1$ , we have

$$\phi'(G) = \sum_{i=0}^{p-1} \phi'_{u_{i+1}}(G_i) \geq [f'(O) - f'(G)] + \frac{\lambda}{2}d(O).$$

Hence,

$$f'(G) + \lambda d(G) \geq f'(O) - f'(G) + \frac{\lambda}{2}d(O),$$

and

$$\phi(G) = f(G) + \lambda d(G) \geq \frac{1}{2}[f(O) + \lambda d(O)] = \frac{1}{2}\phi(O).$$

This completes the proof.  $\square$

The greedy algorithm runs in linear time when  $p$  is a fixed constant. Note that the approximation ratio of two is tight for this particular greedy algorithm. To see this, considering the special case where  $f(S) = 0$  for any subset  $S$  of  $U$ . Let  $A, B$  be a bipartition of  $U$ , each having size  $p$ . The distance between any two elements in  $A$  is one, and the distance between any two elements in  $B$  is two. The distance between an element in  $A$  and an element in  $B$  is one. It is not hard to see that this distance function is a metric and it is possible for the greedy algorithm to return  $A$  as a solution while the optimal solution is  $B$ .

## 5. MATROIDS AND LOCAL SEARCH

Theorem 1 provides a 2-approximation for max-sum diversification when the set function is submodular and the set constraint is a cardinality constraint, i.e., a uniform matroid. It is natural to ask if the same approximation guarantee can be obtained for an arbitrary matroid. In this section, we show that the max-sum diversification problem with monotone submodular function admits a 2-approximation subject to a general matroid constraint.

Matroids are well studied objects in combinatorial optimization. A matroid  $\mathcal{M}$  is a pair  $\langle U, \mathcal{F} \rangle$ , where  $U$  is a

set of ground elements and  $\mathcal{F}$  is a collection of subsets of  $U$ , called *independent sets*, with the following properties :

- **Hereditary:** The empty set is independent and if  $S \in \mathcal{F}$  and  $S' \subset S$ , then  $S' \in \mathcal{F}$ .
- **Augmentation:** If  $A, B \in \mathcal{F}$  and  $|A| > |B|$ , then  $\exists e \in A - B$  such that  $B \cup \{e\} \in \mathcal{F}$ .

The maximal independent sets of a matroid are called *bases* of  $\mathcal{M}$ . Note that all bases have the same number of elements, and this number is called the *rank* of  $\mathcal{M}$ . The definition of a matroid captures the key notion of independence from linear algebra and extends that notion so as to apply to many combinatorial objects. We have already mentioned two classes of matroids relevant to our results, namely partition matroids and transversal matroids. In a partition matroid, the universe  $U$  is partitioned into sets  $C_1, \dots, C_m$  and the independent sets  $S$  satisfy  $S = \cup_{1 \leq i \leq m} S_i$  with  $|S_i| \leq k_i$  for some given bounds  $k_i$  on each part of the partition. A uniform matroid is a special case of a partition matroid with  $m = 1$ . In a transversal matroid, the universe  $U$  is a collection of (possibly) non-intersecting sets  $\mathcal{C} = C_1, \dots, C_m$  and a set  $S$  is independent if there is an injective function  $\phi$  from  $S$  into  $\mathcal{C}$  with say  $\phi(s_i) = S_i$  and  $\phi(s) \in C_i$ . That is,  $S$  forms a set of representatives for each set  $C_i$ . (Note that a given  $s_i$  could occur in other sets  $C_j$ .)

### PROBLEM 2. Max-Sum Diversification for Matroids

Let  $U$  be the underlying ground set, and  $\mathcal{F}$  be the set of independent subsets of  $U$  such that  $\mathcal{M} = \langle U, \mathcal{F} \rangle$  is a matroid. Let  $d(\cdot, \cdot)$  be a (non-negative) metric distance function measuring the distance on every pair of elements. For any subset of  $U$ , let  $f(\cdot)$  be a non-negative monotone submodular set function measuring the weight of the subset. The goal of the problem is to find a subset  $S \in \mathcal{F}$  that:

$$\text{maximizes } f(S) + \lambda \sum_{\{u,v\}: u,v \in S} d(u,v)$$

where  $\lambda$  is a parameter specifying a desired trade-off between the two objectives. As before, we let  $\phi(S)$  be the value of the objective function. Note that since the function  $\phi(\cdot)$  is monotone,  $S$  is essentially a basis of the matroid  $\mathcal{M}$ . The greedy algorithm in Section 4 still applies, but it fails to achieve any constant approximation ratio. This is in contrast to the greedy algorithm of Nemhauser, Wolsey and Fisher, which achieves 2-approximation for general matroids.

Note that the problem is trivial if the rank of the matroid is less than two. Therefore, without loss of generality, we assume the rank is greater or equal to two. Let

$$\{x, y\} = \arg \max_{\{x,y\} \in \mathcal{F}} [f(\{x, y\}) + \lambda d(x, y)].$$

We now consider the following oblivious local search algorithm:

#### LOCAL SEARCH ALGORITHM

let  $S$  be a basis of  $\mathcal{M}$  containing both  $x$  and  $y$   
while there is an  $u \in U - S$  and  $v \in S$  such that  $S + u - v \in \mathcal{F}$   
and  $\phi(S + u - v) > \phi(S)$   
 $S = S + u - v$

end while

return  $S$

**THEOREM 2.** *The local search algorithm achieves an approximation ratio of 2 for max-sum diversification with a matroid constraint.*

Note that if the rank of the matroid is two, then the algorithm is clearly optimal. From now on, we assume the rank of the matroid is greater than two. Before we prove the theorem, we first give several lemmas. All the lemmas assume the problem and the underlying matroid without explicitly mentioning it. Let  $O$  be the optimal solution, and  $S$ , the solution at the end of the local search algorithm. Let  $A = O \cap S$ ,  $B = S - A$  and  $C = O - A$ .

**LEMMA 2.** *For any two sets  $X, Y \in \mathcal{F}$  with  $|X| = |Y|$ , there is a bijective mapping  $g : X \rightarrow Y$  such that  $X - x + g(x) \in \mathcal{F}$  for any  $x \in X$ .*

This is a known property of a matroid and its proof can be found in [4]. Since both  $S$  and  $O$  are bases of the matroid, they have the same cardinality. Therefore,  $B$  and  $C$  have the same cardinality. By Lemma 2, there is a bijective mapping  $g : B \rightarrow C$  such that  $S - b + g(b) \in \mathcal{F}$  for any  $b \in B$ . Let  $B = \{b_1, b_2, \dots, b_t\}$ , and let  $c_i = g(b_i)$  for all  $i$ . Without loss of generality, we assume  $t \geq 2$ , for otherwise, the algorithm is optimal by the local optimality condition.

**LEMMA 3.**  $f(S) + \sum_{i=1}^t f(S - b_i + c_i) \geq f(S - \sum_{i=1}^t b_i) + \sum_{i=1}^t f(S + c_i)$ .

**PROOF.** Since  $f$  is submodular,

$$f(S) - f(S - b_1) \geq f(S + c_1) - f(S + c_1 - b_1)$$

$$f(S - b_1) - f(S - b_1 - b_2) \geq f(S + c_2) - f(S + c_2 - b_2)$$

⋮

$$f(S - \sum_{i=1}^{t-1} b_i) - f(S - \sum_{i=1}^t b_i) \geq f(S + c_t) - f(S + c_t - b_t).$$

Summing up these inequalities, we have

$$f(S) - f(S - \sum_{i=1}^t b_i) \geq \sum_{i=1}^t f(S + c_i) - \sum_{i=1}^t f(S - b_i + c_i),$$

and the lemma follows.  $\square$

**LEMMA 4.**  $\sum_{i=1}^t f(S + c_i) \geq (t-1)f(S) + f(S + \sum_{i=1}^t c_i)$ .

**PROOF.** Since  $f$  is submodular,

$$f(S + c_t) - f(S) = f(S + c_t) - f(S)$$

$$f(S + c_{t-1}) - f(S) \geq f(S + c_t + c_{t-1}) - f(S + c_t)$$

$$f(S + c_{t-2}) - f(S) \geq f(S + c_t + c_{t-1} + c_{t-2}) - f(S + c_t + c_{t-1})$$

⋮

$$f(S + c_1) - f(S) \geq f(S + \sum_{i=1}^t c_i) - f(S + \sum_{i=2}^t c_i)$$

Summing up these inequalities, we have

$$\sum_{i=1}^t f(S + c_i) - tf(S) \geq f(S + \sum_{i=1}^t c_i) - f(S),$$

and the lemma follows.  $\square$

**LEMMA 5.**  $\sum_{i=1}^t f(S - b_i + c_i) \geq (t-2)f(S) + f(O)$ .

**PROOF.** Combining Lemma 3 and Lemma 4, we have

$$\begin{aligned} & f(S) + \sum_{i=1}^t f(S - b_i + c_i) \\ & \geq f(S - \sum_{i=1}^t b_i) + \sum_{i=1}^t f(S + c_i) \\ & \geq (t-1)f(S) + f(S + \sum_{i=1}^t c_i) \\ & = (t-1)f(S) + f(S + C) \\ & \geq (t-1)f(S) + f(O). \end{aligned}$$

Therefore, the lemma follows.  $\square$

**LEMMA 6.** *If  $t > 2$ ,  $d(B, C) - \sum_{i=1}^t d(b_i, c_i) \geq d(C)$ .*

**PROOF.** For any  $b_i, c_j, c_k$ , we have

$$d(b_i, c_j) + d(b_i, c_k) \geq d(c_j, c_k).$$

Summing up these inequalities over all  $i, j, k$  with  $i \neq j$ ,  $i \neq k$ ,  $j \neq k$ , we have each  $d(b_i, c_j)$  with  $i \neq j$  is counted  $(t-2)$  times; and each  $d(c_i, c_j)$  with  $i \neq j$  is counted  $(t-2)$  times. Therefore

$$(t-2)[d(B, C) - \sum_{i=1}^t d(b_i, c_i)] \geq (t-2)d(C),$$

and the lemma follows.  $\square$

**LEMMA 7.**  $\sum_{i=1}^t d(S - b_i + c_i) \geq (t-2)d(S) + d(O)$ .

**PROOF.**

$$\begin{aligned} & \sum_{i=1}^t d(S - b_i + c_i) \\ & = \sum_{i=1}^t [d(S) + d(c_i, S - b_i) - d(b_i, S - b_i)] \\ & = td(S) + \sum_{i=1}^t d(c_i, S - b_i) - \sum_{i=1}^t d(b_i, S - b_i) \\ & = td(S) + \sum_{i=1}^t d(c_i, S) - \sum_{i=1}^t d(c_i, b_i) - \sum_{i=1}^t d(b_i, S - b_i) \\ & = td(S) + d(C, S) - \sum_{i=1}^t d(c_i, b_i) - d(A, B) - 2d(B). \end{aligned}$$

There are two cases. If  $t > 2$  then by Lemma 7, we have

$$\begin{aligned} & d(C, S) - \sum_{i=1}^t d(c_i, b_i) \\ & = d(A, C) + d(B, C) - \sum_{i=1}^t d(c_i, b_i) \\ & \geq d(A, C) + d(C). \end{aligned}$$

Furthermore, since  $d(S) = d(A) + d(B) + d(A, B)$ , we have

$2d(S) - d(A, B) - 2d(B) \geq d(A)$ . Therefore

$$\begin{aligned} & \sum_{i=1}^t d(S - b_i + c_i) \\ = & td(S) + d(C, S) - \sum_{i=1}^t d(c_i, b_i) - d(A, B) - 2d(B) \\ \geq & (t-2)d(S) + d(A, C) + d(C) + d(A) \\ \geq & (t-2)d(S) + d(O). \end{aligned}$$

If  $t = 2$ , then since the rank of the matroid is greater than two,  $A \neq \emptyset$ . Let  $z$  be an element in  $A$ , then we have

$$\begin{aligned} & 2d(S) + d(C, S) - \sum_{i=1}^t d(c_i, b_i) - d(A, B) - 2d(B) \\ = & d(A, C) + d(B, C) - \sum_{i=1}^t d(c_i, b_i) + 2d(A) + d(A, B) \\ \geq & d(A, C) + d(c_1, b_2) + d(c_2, b_1) + d(A) + d(z, b_1) + d(z, b_2) \\ \geq & d(A, C) + d(A) + d(c_1, c_2) \\ \geq & d(A, C) + d(A) + d(C) \\ = & d(O). \end{aligned}$$

Therefore

$$\begin{aligned} & \sum_{i=1}^t d(S - b_i + c_i) \\ = & td(S) + d(C, S) - \sum_{i=1}^t d(c_i, b_i) - d(A, B) - 2d(B) \\ \geq & (t-2)d(S) + d(O). \end{aligned}$$

This completes the proof.  $\square$

Now with the proofs of Lemma 5 and Lemma 7, we are ready to complete the proof of Theorem 2.

PROOF. Since  $S$  is a locally optimal solution, we have  $\phi(S) \geq \phi(S - b_i + c_i)$  for all  $i$ . Therefore, for all  $i$  we have

$$f(S) + \lambda d(S) \geq f(S - b_i + c_i) + \lambda d(S - b_i + c_i).$$

Summing up over all  $i$ , we have

$$tf(S) + \lambda td(S) \geq \sum_{i=1}^t f(S - b_i + c_i) + \lambda \sum_{i=1}^t d(S - b_i + c_i).$$

By Lemma 5, we have

$$tf(S) + \lambda td(S) \geq (t-2)f(S) + f(O) + \lambda \sum_{i=1}^t d(S - b_i + c_i).$$

By Lemma 7, we have

$$tf(S) + \lambda td(S) \geq (t-2)f(S) + f(O) + \lambda[(t-2)d(S) + d(O)].$$

Therefore,

$$2f(S) + 2\lambda d(S) \geq f(O) + \lambda d(O).$$

$$\phi(S) \geq \frac{1}{2}\phi(O),$$

this completes the proof.  $\square$

Theorem 2 shows that even in the more general case of a matroid constraint, we can still achieve the approximation ratio of 2. As is standard in such local search algorithms,

with a small sacrifice on the approximation ratio, the algorithm can be modified to run in polynomial time by requiring at least an  $\epsilon$ -improvement at each iteration rather than just any improvement. Note that the approximation ratio of two is tight for the local search algorithm by a similar example shown at the end of Section 4 with a small modification.

## 6. DYNAMIC UPDATE

In this section, we discuss dynamic updates for the maximum diversification problem with modular set functions. The setting is that we have initially computed a good solution with some approximation guarantee. The weights are changing over time, and upon seeing a change of weight, we want to maintain the quality (the same approximation ratio) of the solution by modifying the current solution without completely recomputing it. We use the number of updates to quantify the amount of modification needed to maintain the desired approximation. An *update* is a single swap of an element in  $S$  with an element outside  $S$ , where  $S$  is the current solution. We ask the following question:

Can we maintain a good approximation ratio with a limited number of updates?

Since the best known approximation algorithm achieves approximation ratio of 2, it is natural to ask whether it is possible to maintain that ratio through local updates. And if it is possible, how many such updates it requires. To simplify the analysis, we restrict to the following oblivious update rule. Let  $S$  be the current solution, and let  $u$  be an element in  $S$  and  $v$  be an element outside  $S$ . The marginal gain  $v$  has over  $u$  with respect to  $S$  is defined to be

$$\phi_{v \rightarrow u}(S) = \phi(S \setminus \{u\} \cup \{v\}) - \phi(S).$$

OBLIVIOUS (SINGLE ELEMENT SWAP) UPDATE RULE

Find a pair of elements  $(u, v)$  with  $u \in S$  and  $v \notin S$  maximizing  $\phi_{v \rightarrow u}(S)$ . If  $\phi_{v \rightarrow u}(S) \leq 0$ , do nothing; otherwise swap  $u$  with  $v$ .

Since the oblivious local search in Theorem 2 uses the same single element swap update rule, it is not hard to see that we can maintain the approximation ratio of 2. However, it is not clear how many updates are needed to maintain that ratio. We conjecture that the number of updates can be made relatively small (i.e., constant) by a non-oblivious update rule and carefully maintaining some desired configuration of the solution set. We leave this as an open question.

However, we are able to show that if we relax the requirement slightly, i.e., aiming for an approximation ratio of 3 instead of 2, and restrict slightly the magnitude of the weight-perturbation, we are able to maintain the desired ratio with a single update. Note that the weight restriction is only used for the case of a weight decrease (Theorem 4).

We divide weight-perturbations into four types: a weight increase (decrease) which occurs on an element, and a distance increase (decrease) which occurs between two elements. We denote these four types: (I), (II), (III), (IV); and we have a corresponding theorem for each case.

Before getting to the theorems, we first prove the following two lemmas. After a weight-perturbation, let  $S$  be the current solution set, and  $O$  be the optimal solution. Let  $S^*$  be the solution set after a single update using the oblivious update rule, and let  $\Delta = \phi(S^*) - \phi(S)$ . We again let  $Z = O \cap S$ ,  $X = O \setminus S$  and  $Y = S \setminus O$ .

LEMMA 8. *There exists  $z \in Y$  such that*

$$\phi_z(S \setminus \{z\}) \leq \frac{1}{|Y|}[f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

PROOF. If we sum up all marginal gain  $\phi_y(S \setminus \{y\})$  for all  $y \in Y$ , we have

$$\sum_{y \in Y} \phi_y(S \setminus \{y\}) = f(Y) + 2\lambda d(Y) + \lambda d(Z, Y).$$

By an averaging argument, there must exist  $z \in Y$  such that

$$\phi_z(S \setminus \{z\}) \leq \frac{1}{|Y|}[f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

□

Lemma 8 ensures the existence of an element in  $S$  such that after removing it from  $S$ , the objective function value does not decrease much. The following lemma ensures that there always exists an element outside  $S$  which can increase the objective function value substantially if we bring it in.

LEMMA 9. *If  $\phi(S^*) < \frac{1}{3}\phi(O)$ , then for all  $y \in Y$ , there exists  $x \in X$  such that*

$$\phi_x(S \setminus \{y\}) > \frac{1}{|X|}[2\phi(Z) + 3\phi(Y) + 3\lambda d(Z, Y) + 3\Delta].$$

PROOF. For any  $y \in Y$ , and by Lemma 1, we have

$$\begin{aligned} & f(X) + \lambda d(S \setminus \{y\}, X) \\ &= f(X) + \lambda d(Z, X) + \lambda d(Y \setminus \{y\}, X) \\ &\geq f(X) + \lambda d(Z, X) + \lambda d(X). \end{aligned}$$

Note that since  $\phi(S^*) = \phi(S) + \Delta < \frac{1}{3}\phi(O)$ , we have

$$\begin{aligned} \phi(O) &= \phi(Z) + f(X) + \lambda d(X) + \lambda d(Z, X) \\ &> 3\phi(Z) + 3\phi(Y) + 3\lambda d(Z, Y) + 3\Delta. \end{aligned}$$

Therefore,

$$\begin{aligned} & f(X) + \lambda d(S \setminus \{y\}, X) \\ &\geq f(X) + \lambda d(Z, X) + \lambda d(X) \\ &> 2\phi(Z) + 3\phi(Y) + 3\lambda d(Z, Y) + 3\Delta. \end{aligned}$$

This implies there must exist  $x \in X$  such that

$$\phi_x(S \setminus \{y\}) > \frac{1}{|X|}[2\phi(Z) + 3\phi(Y) + 3\lambda d(Z, Y) + 3\Delta].$$

□

Combining Lemma 8 and 9, we can give a lower bound for  $\Delta$ . We have the following corollary.

COROLLARY 1. *If  $\phi(S^*) < \frac{1}{3}\phi(O)$ , then we have  $|Y| > 3$  and furthermore*

$$\Delta > \frac{1}{|Y| - 3}[2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)].$$

PROOF. By Lemma 8, there exists  $y \in Y$  such that

$$\phi_y(S \setminus \{y\}) \leq \frac{1}{|Y|}[f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

Since  $\phi(S^*) < \frac{1}{3}\phi(O)$ , by Lemma 9, for this particular  $y$ , there exists  $x \in X$  such that

$$\phi_x(S \setminus \{y\}) > \frac{1}{|X|}[2\phi(Z) + 3\phi(Y) + 3\lambda d(Z, Y) + 3\Delta].$$

Since  $|X| = |Y|$ , we have

$$\Delta > \frac{1}{|Y|}[2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y) + 3\Delta].$$

If  $|Y| \leq 3$ , then it is a contradiction. Therefore  $|Y| > 3$ . Rearranging the inequality, we have

$$\Delta > \frac{1}{|Y| - 3}[2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)].$$

□

COROLLARY 2. *If  $p \leq 3$ , then for any weight or distance perturbation, we can maintain an approximation ratio of 3 with a single update.*

PROOF. This is an immediate consequence of Corollary 1 since  $p \geq |Y|$ . □

Given Corollary 2, we will assume  $p > 3$  for all the remaining results in this section. We first discuss weight-perturbations on elements.

THEOREM 3. [TYPE (1)] *For any weight increase, we can maintain an approximation ratio of 3 with a single update.*

PROOF. Suppose we increase the weight of  $s$  by  $\delta$ . Since the optimal solution can increase by at most  $\delta$ , if  $\Delta \geq \frac{1}{3}\delta$ , then we have maintained a ratio of 3. Hence we assume  $\Delta < \frac{1}{3}\delta$ . If  $s \in S$  or  $s \notin O$ , then it is clear the ratio of 3 is maintained. The only interesting case is when  $s \in O \setminus S$ . Suppose, for the sake of contradiction, that  $\phi(S^*) < \frac{1}{3}\phi(O)$ , then by Corollary 1, we have  $|Y| > 3$  and

$$\Delta > \frac{1}{|Y| - 3}[2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)].$$

Since  $\Delta < \frac{1}{3}\delta$ , we have

$$\delta > \frac{1}{|Y| - 3}[6\phi(Z) + 6f(Y) + 3\lambda d(Y) + 6\lambda d(Z, Y)].$$

On the other hand, by Lemma 8, there exists  $y \in Y$  such that

$$\phi_y(S \setminus \{y\}) \leq \frac{1}{|Y|}[f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

Now considering a swap of  $s$  with  $y$ , the loss by removing  $y$  from  $S$  is  $\phi_y(S \setminus \{y\})$ , while the increase that  $s$  brings to the set  $S \setminus \{y\}$  is at least  $\delta$  (as  $s$  is increased by  $\delta$ , and the original weight of  $s$  is non-negative). Therefore the marginal gain of the swap of  $s$  with  $y$  is  $\phi_{s \rightarrow y} \geq \delta - \phi_y(S \setminus \{y\})$  and hence

$$\phi_{s \rightarrow y}(S) \geq \delta - \frac{1}{|Y|}[f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

However,  $\phi_{s \rightarrow y}(S) \leq \Delta < \frac{1}{3}\delta$ . Therefore, we have

$$\frac{1}{3}\delta > \delta - \frac{1}{|Y|}[f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

This implies

$$\delta < \frac{1}{|Y|}\left[\frac{3}{2}f(Y) + 3\lambda d(Y) + \frac{3\lambda}{2}d(Z, Y)\right],$$

which is a contradiction. □

**THEOREM 4.** [TYPE (II)] *For a weight decrease of magnitude  $\delta$ , we can maintain an approximation ratio of 3 with*

$$\lceil \log_{\frac{p-2}{p-3}} \frac{w}{w-\delta} \rceil$$

*updates, where  $w$  is the weight of the solution before the weight decrease. In particular, if  $\delta \leq \frac{w}{p-2}$ , we only need a single update.*

**PROOF.** Suppose we decrease the weight of  $s$  by  $\delta$ . Without loss of generality, we can assume  $s \in S$ . Suppose, for the sake of contradiction, that  $\phi(S^*) < \frac{1}{3}\phi(O)$ , then by Corollary 1, we have  $|Y| > 3$  and

$$\begin{aligned} \Delta &> \frac{1}{|Y|-3} [2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)] \\ &\geq \frac{1}{p-3} \phi(S). \end{aligned}$$

Therefore

$$\phi(S^*) > \frac{p-2}{p-3} \phi(S).$$

This implies that we can maintain the approximation ratio with

$$\lceil \log_{\frac{p-2}{p-3}} \frac{w}{w-\delta} \rceil$$

number of updates. In particular, if  $\delta \leq \frac{w}{p-2}$ , we only need a single update.  $\square$

We now discuss the weight-perturbations between two elements. We assume that such perturbations preserve the metric condition. Furthermore, we assume  $p > 3$  for otherwise, by Corollary 1, the ratio of 3 is maintained.

**THEOREM 5.** [TYPE (III)] *For any distance increase, we can maintain an approximation ratio of 3 with a single update.*

**PROOF.** Suppose we increase the distance of  $(x, y)$  by  $\delta$ , and for the sake of contradiction, we assume that  $\phi(S^*) < \frac{1}{3}\phi(O)$ , then by Corollary 1, we have  $|Y| > 3$  and

$$\Delta > \frac{1}{|Y|-3} [2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)].$$

Since  $\Delta < \frac{1}{3}\delta$ , we have

$$\begin{aligned} \delta &> \frac{3}{|Y|-3} [2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)] \\ &\geq \frac{3}{p-3} \phi(S). \end{aligned}$$

If both  $x$  and  $y$  are in  $S$ , then it is not hard to see that the ratio of 3 is maintained. Otherwise, there are two cases:

1. Exactly one of  $x$  and  $y$  is in  $S$ , without loss of generality, we assume  $y \in S$ . Considering that we swap  $x$  with any vertex  $z \in S$  other than  $y$ . Since after the swap, both  $x$  and  $y$  are now in  $S$ , by the triangle inequality of the metric condition, we have

$$\Delta \geq (p-1)\delta - \phi(S) > \left(\frac{2}{3}p-2\right)\delta.$$

Since  $p > 3$ , we have

$$\Delta > \left(\frac{2}{3}p-2\right)\delta \geq \frac{2}{3}\delta > 2\Delta,$$

which is a contradiction.

2. Both  $x$  and  $y$  are outside in  $S$ . By Lemma 8, there exists  $z \in Y$  such that

$$\phi_z(S \setminus \{z\}) \leq \frac{1}{|Y|} [f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

Consider the set  $T = \{x, y\}$  with  $S \setminus \{z\}$ , by the triangle inequality of the metric condition, we have  $d(T, S \setminus \{z\}) \geq (p-1)\delta$ . Therefore, at least one of  $x$  and  $y$ , without loss of generality, assuming  $x$ , has the following property:

$$d(x, S \setminus \{z\}) \geq \frac{(p-1)\delta}{2}.$$

Considering that we swap  $x$  with  $z$ , we have:

$$\Delta \geq \frac{(p-1)}{2}\delta - \frac{1}{|Y|} [f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

Since  $\Delta < \frac{1}{3}\delta$ , we have

$$\frac{1}{3}\delta > \frac{(p-1)}{2}\delta - \frac{1}{|Y|} [f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

This implies that

$$\delta < \frac{6}{3p-5} \cdot \frac{1}{|Y|} [f(Y) + 2\lambda d(Y) + \lambda d(Z, Y)].$$

Since  $p > 3$ , we have

$$\delta < \frac{1}{|Y|} \left[ \frac{6}{7}f(Y) + \frac{12\lambda}{7}d(Y) + \frac{6\lambda}{7}d(Z, Y) \right],$$

which is a contradiction.

Therefore,  $\phi(S^*) \geq \frac{1}{3}\phi(O)$ ; this completes the proof.  $\square$

**THEOREM 6.** [TYPE (IV)] *For any distance decrease, we can maintain an approximation ratio of 3 with a single update.*

**PROOF.** Suppose we decrease the distance of  $(x, y)$  by  $\delta$ . Without loss of generality, we assume both  $x$  and  $y$  are in  $S$ , for otherwise, it is not hard to see the ratio of 3 is maintained. Suppose, for the sake of contradiction, that  $\phi(S^*) < \frac{1}{3}\phi(O)$ , then by Corollary 1, we have  $|Y| > 3$  and

$$\begin{aligned} \Delta &> \frac{1}{|Y|-3} [2\phi(Z) + 2f(Y) + \lambda d(Y) + 2\lambda d(Z, Y)] \\ &\geq \frac{1}{p-3} \phi(S). \end{aligned}$$

If  $\Delta \geq \delta$ , then the ratio of 3 is maintained. Otherwise,

$$\delta > \Delta \geq \frac{1}{p-3} \phi(S).$$

By the triangle inequality of the metric condition, we have

$$\phi(S) \geq (p-2)\delta > \frac{p-2}{p-3} \phi(S) > \phi(S),$$

which is a contradiction.  $\square$

Combining Theorem 3, 4, 5, 6, we have the following corollary.

**COROLLARY 3.** *If the initial solution achieves approximation ratio of 3, then for any weight-perturbation of TYPE (I), (III), (IV); and any weight-perturbation of TYPE (II) that is no more than  $\frac{1}{p-2}$  of the current solution for  $p > 3$  and arbitrary for  $p \leq 3$ , we can maintain the ratio of 3 with a single update.*

## 7. EXPERIMENTS

In this section, we present the results of some preliminary experiments. Note that all our results in this paper are of theoretical nature, and the purpose of the experiments in this section are to provide additional insights about our algorithms. They shall not be treated as experimental evidences for the performance of our algorithms.

We consider the max-sum diversification problem with modular set functions. In order to have gradual control of the parameters, we use the following different, but equivalent form of the objective function:

$$\alpha f(S) + (1 - \alpha) \sum_{u,v \in S} d(u,v),$$

where  $\alpha$  is a real number in  $[0, 1]$ . We conduct two sets of experiments:

1. We compare the performance (the approximation ratio) of the greedy algorithm proposed in [17] with the greedy algorithm that we propose in this paper.
2. We simulate three different dynamically changing environments, and record the worst approximation ratio occurring with a single application of the oblivious update rule.

Note that both experiments use synthetic data which is generated uniformly at random within a given range, and in order to compute the optimal solution (to evaluate the approximation ratio), we restrict the size of the input data. Therefore, our experiments are not representative for real data sets and large input cases, but nevertheless, they shed some light on the behavior of the proposed algorithm and the dynamic update rule.

### 7.1 Comparing Two Greedy Algorithms

The first experiment is designed to compare the greedy algorithm proposed in [17] with the greedy algorithm that we propose in this paper. We generate input instances of 50 vertices, with vertex weights chosen uniformly at random from  $[0, 1]$  and distances chosen uniformly at random from  $[1, 2]$  to ensure the triangle inequality. The target set size is chosen to be five. These instances are small enough that we can determine the optimal solution. We run both the greedy algorithm of [17], denoted as Greedy A, and our greedy algorithm, denoted as Greedy B, on the same instance for different values of  $\alpha$ . This is repeated 100 times for each value of  $\alpha$  and both approximation ratios are recorded for every instance. These approximation ratios are then averaged to represent the performance of the algorithm for each  $\alpha$  value. The results are shown in Fig. 1; the horizontal axis measures  $\alpha$  values, and the vertical axis measures the approximation ratio (the smaller the bar height, the better the ratio).

Although both algorithms have the same provable theoretical approximation ratio of 2, we observe the following phenomena:

1. Greedy B significantly outperforms Greedy A for every value of  $\alpha$ . The worst case for each algorithm occurs at  $\alpha = 0$  (i.e., no element weight), where Greedy A has an approximation ratio of 1.26 while Greedy B is still well below 1.05.
2. Both algorithms performs well below the theoretical approximation bound of 2.

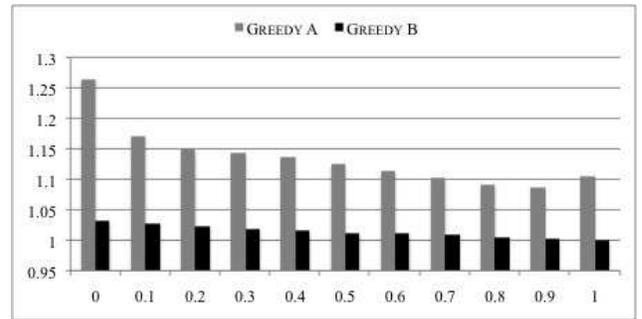


Figure 1: A Comparison of Two Greedy Algorithms

3. As  $\alpha$  increase, both algorithms tend to perform better. Note that Greedy B achieves the optimum when  $\alpha = 1$  as it then becomes the “standard” greedy algorithm [31, 14]; however, this is not the case for Greedy A.

Despite the fact that the experiment is only conducted for very limited cases, it gives some evidence that our greedy algorithm outperforms the one proposed in [17].

### 7.2 Approximation Ratio in Dynamic Updates

For dynamic update, we use the same type of random instance in the previous experiment. We have three different dynamically changing environments:

1. VPerturbation: each perturbation is a weight change on an element.
2. EPerturbation: each perturbation is a weight change between two elements.
3. MPerturbation: each perturbation is one of the above two with equal probability.

For each of the environments above and every value of  $\alpha$ , we start with a greedy solution (a 2-approximation) and run 20 steps of simulation, where each step consists of a weight change of the stated type, followed by a single application of the oblivious update rule. We repeat this 100 times and record the worst approximation ratio occurring during these 100 updates. The results are shown in Fig. 2; again the hori-

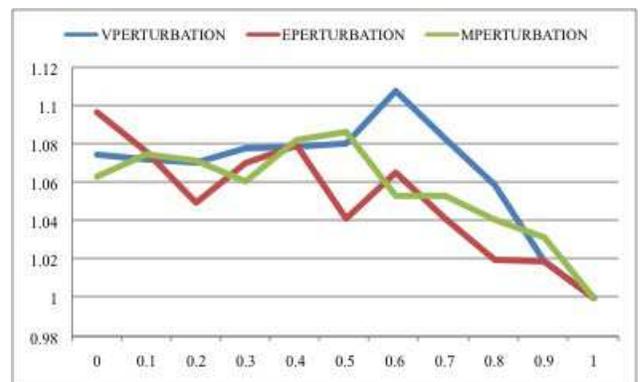


Figure 2: Approximation Ratio in Dynamic Updates

zontal axis measures  $\alpha$  values, and the vertical axis measures the approximation ratio.

We have the following observations:

1. In any dynamic changing environment, the maintained ratio is well below the provable ratio of 3. The worst observed ratio is about 1.11.
2. The maintained ratios are decreasing to 1 for increasing  $\alpha \geq 0.6$  approximately.

From the experiment, we see that oblivious update rule seems effective for maintaining a good approximation ratio in a dynamically changing environment.

## 8. CONCLUSION

We study the max-sum diversification with monotone submodular set functions and give a natural 2-approximation greedy algorithm for the problem. We further extend the problem to matroids and give a 2-approximation local search algorithm for the problem. We examine the dynamic update setting for modular set functions, where the weights and distances are constantly changing over time and the goal is to maintain a solution with good quality with a limited number of updates. We propose a simple update rule: the oblivious (single swap) update rule, and show that if the weight-perturbation is not too large, we can maintain an approximation ratio of 3 with a single update.

The diversification problem has many important applications and there are many interesting future directions. Although in this paper we restricted ourselves to the max-sum objective, there are many other well-defined notion of diversity that can be considered, see for example [7]. The max-sum case can be also viewed as the  $\ell_1$ -norm; what about other norms?

Another important open question is to find the tight approximation ratio for max-sum diversification with monotone submodular set functions (for both the uniform matroid case and the general matroid case). We know the ratio cannot be better than  $\frac{e}{e-1}$  assuming P is not equal to NP [29] and an approximation ratio of 2 is obtained in this paper. Is it possible to beat 2? In the general matroid case, the greedy algorithm given in Section 4 fails to achieve any constant approximation ratio, but how about other greedy algorithms? What if we start with the best pair?

In a dynamic update setting, we only considered the oblivious single swap update rule. It is interesting to see if it is possible to maintain a better ratio than 3 with a limited number of updates, by larger cardinality swaps, and/or by a non-oblivious update rule. We leave this as an interesting open question of the paper.

Finally, a crucial property used throughout our results is the triangle inequality. For a relaxed version of the triangle inequality can we relate the approximation ratio to the parameter of a relaxed triangle inequality?

## 9. ACKNOWLEDGMENTS

We thank MITACS and Thoor Inc. for its generous support and Justin Ward for many helpful discussions. We also thank anonymous referees for pointing out several small mistakes in the primary version of the paper and references to some early work in diversification.

## 10. REFERENCES

- [1] R. Agrawal, S. Gollapudi, A. Halverson, and S. Ieong. Diversifying search results. In *WSDM*, pages 5–14, 2009.
- [2] N. Bansal, K. Jain, A. Kazeykina, and J. Naor. Approximation algorithms for diversified search ranking. In *ICALP (2)*, pages 273–284, 2010.
- [3] C. Brandt, T. Joachims, Y. Yue, and J. Bank. Dynamic ranked retrieval. In *WSDM*, pages 247–256, 2011.
- [4] R. A. Brualdi. Comments on bases in dependence structures. *Bulletin of the Australian Mathematical Society*, 1(02):161–167, 1969.
- [5] G. Călinescu, C. Chekuri, M. Pál, and J. Vondrák. Maximizing a monotone submodular function subject to a matroid constraint. *SIAM J. Comput.*, 40(6):1740–1766, 2011.
- [6] J. Carbonell and J. Goldstein. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '98*, pages 335–336, New York, NY, USA, 1998. ACM.
- [7] B. Chandra and M. M. Halldórsson. Facility dispersion and remote subgraphs. In *Proceedings of the 5th Scandinavian Workshop on Algorithm Theory*, pages 53–65, London, UK, 1996. Springer-Verlag.
- [8] B. Chandra and M. M. Halldórsson. Approximation algorithms for dispersion problems. *J. Algorithms*, 38(2):438–465, 2001.
- [9] H. Chen and D. R. Karger. Less is more: probabilistic models for retrieving fewer relevant documents. In *SIGIR*, pages 429–436, 2006.
- [10] E. Demidova, P. Fankhauser, X. Zhou, and W. Nejdl. Divq: diversification for keyword search over structured databases. In *Proceeding of the 33rd international ACM SIGIR conference on Research and development in information retrieval, SIGIR '10*, pages 331–338. ACM, 2010.
- [11] Z. Dou, S. Hu, K. Chen, R. Song, and J.-R. Wen. Multi-dimensional search result diversification. In *WSDM*, pages 475–484, 2011.
- [12] M. Drosou and E. Pitoura. Diversity over continuous data. *IEEE Data Eng. Bull.*, 32(4):49–56, 2009.
- [13] M. Drosou and E. Pitoura. Search result diversification. *SIGMOD Record*, 39(1):41–47, 2010.
- [14] J. Edmonds. Matroids and the greedy algorithm. *Mathematical Programming*, 1:127–136, 1971.
- [15] E. Erkut. The discrete p-dispersion problem. *European Journal of Operational Research*, 46(1):48–60, May 1990.
- [16] E. Erkut and S. Neuman. Analytical models for locating undesirable facilities. *European Journal of Operational Research*, 40(3):275–291, June 1989.
- [17] S. Gollapudi and A. Sharma. An axiomatic approach for result diversification. In *World Wide Web Conference Series*, pages 381–390, 2009.
- [18] M. M. Halldórsson, K. Iwano, N. Katoh, and T. Tokuyama. Finding subsets maximizing minimum structures. In *Symposium on Discrete Algorithms*, pages 150–159, 1995.

- [19] P. Hansen and I. D. Moon. Dispersion facilities on a network. *Presentation at the TIMS/ORSA Joint National Meeting, Washington, D.C.*, 1988.
- [20] R. Hassin, S. Rubinstein, and A. Tamir. Approximation algorithms for maximum dispersion. *Oper. Res. Lett.*, 21(3):133–137, 1997.
- [21] D. Kempe, J. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining, KDD '03*, pages 137–146, 2003.
- [22] S. Khanna, R. Motwani, M. Sudan, and U. V. Vazirani. On syntactic versus computational views of approximability. *Electronic Colloquium on Computational Complexity (ECCC)*, 2(23), 1995.
- [23] M. J. Kuby. Programming models for facility dispersion: The p-dispersion and maximum dispersion problems. *Geographical Analysis*, 19(4):315–329, 1987.
- [24] H. Lin and J. Bilmes. Multi-document summarization via budgeted maximization of submodular functions. In *HLT-NAACL*, pages 912–920, 2010.
- [25] H. Lin and J. Bilmes. A class of submodular functions for document summarization. In *North American chapter of the Association for Computational Linguistics/Human Language Technology Conference (NAACL/HLT-2011)*, Portland, OR, June 2011. (long paper).
- [26] H. Lin, J. Bilmes, and S. Xie. Graph-based submodular selection for extractive summarization. In *Proc. IEEE Automatic Speech Recognition and Understanding (ASRU)*, Merano, Italy, December 2009.
- [27] Z. Liu, P. Sun, and Y. Chen. Structured search result differentiation. *PVLDB*, 2(1):313–324, 2009.
- [28] E. Minack, W. Siberski, and W. Nejdl. Incremental diversification for very large sets: a streaming-based approach. In *SIGIR*, pages 585–594, 2011.
- [29] G. Nemhauser, L. Wolsey, and M. Fisher. An analysis of the approximations for maximizing submodular set functions. *Mathematical Programming*, 1978.
- [30] F. Radlinski, R. Kleinberg, and T. Joachims. Learning diverse rankings with multi-armed bandits. In *ICML*, pages 784–791, 2008.
- [31] R. Rado. A note on independence functions. *Proceedings of the London Mathematical Society*, 7:300–320, 1957.
- [32] D. Rafiei, K. Bharat, and A. Shukla. Diversifying web search results. In *WWW*, pages 781–790, 2010.
- [33] S. S. Ravi, D. J. Rosenkrantz, and G. K. Tayi. Heuristic and special case algorithms for dispersion problems. *Operations Research*, 42(2):299–310, March-April 1994.
- [34] R. L. T. Santos, C. Macdonald, and I. Ounis. Intent-aware search result diversification. In *SIGIR*, pages 595–604, 2011.
- [35] A. Schrijver. *Combinatorial Optimization: Polyhedra and Efficiency*. Springer, 2003.
- [36] A. Slivkins, F. Radlinski, and S. Gollapudi. Learning optimally diverse rankings over large document collections. In *ICML*, pages 983–990, 2010.
- [37] M. R. Vieira, H. L. Razente, M. C. N. Barioni, M. Hadjieleftheriou, D. Srivastava, C. T. Jr., and V. J. Tsotras. Divdb: A system for diversifying query results. *PVLDB*, 4(12):1395–1398, 2011.
- [38] M. R. Vieira, H. L. Razente, M. C. N. Barioni, M. Hadjieleftheriou, D. Srivastava, C. T. Jr., and V. J. Tsotras. On query result diversification. In *ICDE*, pages 1163–1174, 2011.
- [39] D. W. Wang and Y.-S. Kuo. A study on two geometric location problems. *Inf. Process. Lett.*, 28:281–286, August 1988.
- [40] C. Yu, L. Lakshmanan, and S. Amer-Yahia. It takes variety to make a world: diversification in recommender systems. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology, EDBT '09*, pages 368–378, 2009.
- [41] Y. Yue and T. Joachims. Predicting diverse subsets using structural svms. In *ICML*, pages 1224–1231, 2008.
- [42] C. Zhai, W. W. Cohen, and J. D. Lafferty. Beyond independent relevance: methods and evaluation metrics for subtopic retrieval. In *SIGIR*, pages 10–17, 2003.
- [43] F. Zhao, X. Zhang, A. K. H. Tung, and G. Chen. Broad: Diversified keyword search in databases. *PVLDB*, 4(12):1355–1358, 2011.
- [44] X. Zhu, A. B. Goldberg, J. V. Gael, and D. Andrzejewski. Improving diversity in ranking using absorbing random walks. In *HLT-NAACL*, pages 97–104, 2007.