
“BRING YOUR OWN GREEDY”+MAX: Near-Optimal $1/2$ -Approximations for Submodular Knapsack

Grigory Yaroslavtsev
Indiana University
The Alan Turing Institute

Samson Zhou
Indiana University
Carnegie Mellon University

Dmitrii Avdiukhin
Indiana University

Abstract

The problem of selecting a small-size representative summary of a large dataset is a cornerstone of machine learning, optimization and data science. Motivated by applications to recommendation systems and other scenarios with query-limited access to vast amounts of data, we propose a new rigorous algorithmic framework for a standard formulation of this problem as a submodular maximization subject to a linear (knapsack) constraint. Our framework is based on augmenting all partial GREEDY solutions with the best additional item. It can be instantiated with negligible overhead in any model of computation, which allows the classic GREEDY algorithm and its variants to be implemented. We give such instantiations in the offline (GREEDY+MAX), multi-pass streaming (SIEVE+MAX) and distributed (DISTRIBUTED SIEVE+MAX) settings. Our algorithms give $(1/2 - \epsilon)$ -approximation with most other key parameters of interest being near-optimal. Our analysis is based on a new set of first-order linear differential inequalities and their robust approximate versions. Experiments on typical datasets (movie recommendations, influence maximization) confirm scalability and high quality of solutions obtained via our framework. Instance-specific approximations are typically in the 0.6-0.7 range and frequently beat even the $(1 - 1/e) \approx 0.63$ worst-case barrier for polynomial-time algorithms.

1 Introduction

A fundamental problem in many large-scale machine learning, data science and optimization tasks is finding a small representative subset of a big dataset. This problem arises from applications in recommendation systems [Leskovec et al. \(2007\)](#); [El-Arini and Guestrin \(2011\)](#); [Bogunovic et al. \(2017\)](#); [Mitrović et al. \(2017\)](#); [Yu et al. \(2018\)](#); [Avdiukhin et al. \(2019\)](#), exemplar-based clustering [Gomes and Krause \(2010\)](#), facility location [Lindgren et al. \(2016\)](#), image processing [Iyer and Bilmes \(2019\)](#), viral marketing [Hartline et al. \(2008\)](#), principal component analysis [Khanna et al. \(2015\)](#), and document summarization [Lin and Bilmes \(2011\)](#); [Wei et al. \(2013\)](#); [Sipos et al. \(2012\)](#) and can often be formulated as constrained monotone submodular optimization under various constraints such as cardinality [Badanidiyuru et al. \(2014\)](#); [Bateni et al. \(2018\)](#); [Kazemi et al. \(2019\)](#), knapsack [Huang et al. \(2017\)](#), matchings [Chakrabarti and Kale \(2014\)](#), and matroids [Călinescu et al. \(2011\)](#); [Anari et al. \(2019\)](#) due to restrictions demanded by space, budget, diversity, fairness or privacy. As a result, constrained submodular optimization has been recently and extensively studied in various computational models, including centralized [Nemhauser et al. \(1978\)](#), distributed [Mirzasoleiman et al. \(2013\)](#); [Kumar et al. \(2015\)](#); [da Ponte Barbosa et al. \(2015\)](#); [Mirrokni and Zadimoghaddam \(2015\)](#); [Mirzasoleiman et al. \(2016\)](#); [da Ponte Barbosa et al. \(2016\)](#); [Liu and Vondrák \(2019\)](#), streaming [Badanidiyuru et al. \(2014\)](#); [Buchbinder et al. \(2015\)](#); [Norouzi-Fard et al. \(2018\)](#); [Agrawal et al. \(2019\)](#); [Kazemi et al. \(2019\)](#), and adaptive [Golovin and Krause \(2011\)](#); [Balkanski and Singer \(2018\)](#); [Balkanski et al. \(2019\)](#); [Fahrbach et al. \(2019\)](#); [Ene and Nguyen \(2019b\)](#); [Chekuri and Quanrud \(2019\)](#) among others.

In this paper we focus on monotone submodular maximization *under a knapsack constraint*, which captures the scenario when the representative subset should have a small cost or size. While a number of algorithmic techniques exist for this problem, there are few that robustly

scale to large data and can be easily implemented in various computing frameworks. This is in contrast with a simpler *cardinality-constrained* version in which only the number of elements is restricted. In this setting the celebrated GREEDY algorithm of Nemhauser et al. (1978) enjoys both an optimal approximation ratio and a simplicity that allows easy adaptation in various environments. For knapsack constraints, such a simple and universal algorithm is unlikely. In particular, GREEDY does not give any approximation guarantee.

We develop a framework that augments solutions constructed by GREEDY and its variations and gives almost $1/2$ -approximations¹ in various computational models. For example, in the multi-pass streaming setting we achieve optimal space and almost optimal number of queries and running time. We believe that our framework is robust to the choice of the computational model as it can be implemented with essentially the same complexity as that of running GREEDY and its variants.

Preliminaries and our contributions. A set function $f: 2^U \rightarrow \mathbb{R}$ is *submodular* if for every $S \subseteq T \subseteq U$ and $e \in U$ it holds that $f(e \cup T) - f(T) \leq f(e \cup S) - f(S)$. Moreover, f is *monotone* if for every $S \subseteq T \subseteq U$ it holds that $f(T) \geq f(S)$. In the monotone submodular maximization problem subject to a knapsack constraint, each item e has cost $c(e)$. Given a parameter $K > 0$, the task is to maximize a non-negative monotone submodular function $f(S)$ under the constraint $c(S) := \sum_{e \in S} c(e) \leq K$. Without loss of generality, we assume that $\min_{e \in S} c(e) \geq 1$, which can be achieved by rescaling the costs and taking all items with cost 0. Then $\tilde{K} = \min(n, K)$ is an upper bound on the number of elements in any feasible solution.

Any algorithm for submodular maximization requires *query access* to f . As query access can be expensive, the number of queries is typically considered one of the performance metrics. Furthermore, in some critical applications of submodular optimization such as recommendation systems, another constraint often arises from the fact that only queries to feasible sets are allowed (e.g. when click-through rates can only be collected for sets of ads which can be displayed to the users). Practical algorithms for submodular optimization hence typically only make such queries, an assumption commonly used in the literature (see e.g. Norouzi-Fard et al. (2018)). For *any* algorithm that only makes queries on feasible sets, it is easy to show that $\Omega(n^2)$ queries are required to go beyond $1/2$ -approximation under various assumptions on f (Theorem 2.9). Hence it is natural to ask whether we can get a $1/2$ -approximation, while keeping other performance metrics of interest nearly

optimal and hence not compromising on practicality. We answer this question positively.

We first state the following simplified result in the most basic *offline model* (i.e. when an algorithm can access any element at any time) to illustrate the main ideas and then improve parameters in our other results. In this model, we are given an integer knapsack capacity $K \in \mathbb{Z}^+$ and a set E of elements e_1, \dots, e_n from a finite universe U .²

Theorem 1.1 (Offline GREEDY+MAX)

Let $\tilde{K} = \min(n, K)$. There exists an offline algorithm GREEDY+MAX (Algorithm 1) that gives a $1/2$ -approximation for the submodular maximization problem under a knapsack constraint with query complexity and running time $\mathcal{O}(\tilde{K}n)$ (Theorem 2.4).

In the *single-pass streaming model*, the algorithm is given K and a stream E consisting of elements $e_1, \dots, e_n \in U$, which arrive sequentially. The objective is to minimize the auxiliary space used by algorithm throughout the execution. In the *multi-pass streaming model*, the algorithm is further allowed to make multiple passes over E . This model is typically used for modeling storage devices with sequential access (e.g. hard drives) while using a small amount of RAM. In this setting minimizing the number of passes becomes another key priority. Note that since $\Omega(\tilde{K})$ is a trivial lower bound on space and $\Omega(n)$ is a trivial lower bound on time and query complexity of any approximation algorithm that queries feasible sets, our next result is almost optimal in most parameters of interest.

Theorem 1.2 (Multi-pass streaming algorithm SIEVE+MAX)

Let $\tilde{K} = \min(n, K)$. There exists a multi-pass streaming algorithm SIEVE+MAX (Algorithm 2) that uses $\mathcal{O}(\tilde{K})$ space and $\mathcal{O}(1/\epsilon)$ passes over the stream and outputs a $(1/2 - \epsilon)$ -approximation to the submodular maximization problem under a knapsack constraint, with query complexity and running time³ $\mathcal{O}(n(1/\epsilon + \log \tilde{K}))$.

We also give an algorithm in the massively-parallel computation (MPC) model Karloff et al. (2010) used to model MapReduce/Spark-like systems. We use the most restrictive version, which only allows linear total memory, running time and communication per

²W.l.o.g. for all e we have $1 \leq c(e) \leq K$ as one can rescale the capacity and costs and filter out all items with cost more than K (in all our results this means replacing K with the aspect ratio $K / \min_{e \in E} c(e)$).

³Note that when $\frac{1}{\epsilon} \ll K$, in terms of running time our streaming algorithm is more efficient than our offline algorithm. Hence, in the offline setting one can use the best of the two algorithms depending on the parameters.

¹Algorithm gives an α -approximation if it outputs S such that $f(S) \geq \alpha f(\text{OPT})$.

round [Andoni et al. \(2014\)](#). In this model, the input set E of size n is arbitrarily distributed across m machines, each with $s = \mathcal{O}(n/m)$ memory so that the overall memory is $\mathcal{O}(n)$. A standard setting of parameters for submodular optimization is $m = \sqrt{n/\tilde{K}}$ and $s = \mathcal{O}(\sqrt{n\tilde{K}})$ (see e.g. [Liu and Vondrák \(2019\)](#); [Avdiukhin et al. \(2019\)](#)). One of the machines is designated as the *central machine* and outputs the solution in the end. The machines communicate to each other in a number of synchronous rounds. In each round, each machine receives an input of size $\mathcal{O}(\sqrt{n\tilde{K}})$, performs a local linear-time computation, and sends an output of size $\mathcal{O}(\sqrt{n\tilde{K}})$ to other machines before the next round begins. The primary objective in this model is minimizing the number of rounds. Our main result in this model is given below.

Theorem 1.3 (MPC algorithm DISTRIBUTED SIEVE+MAX) *Let $\tilde{K} = \min(n, K)$. There exists an MPC algorithm that runs in $\mathcal{O}(1/\epsilon)$ rounds on $\sqrt{n/\tilde{K}}$ machines, each with $\mathcal{O}(\sqrt{n\tilde{K}})$ memory. Each machine uses query complexity and runtime $\mathcal{O}(\sqrt{n\tilde{K}})$ per round. The algorithm outputs a $(1/2 - \epsilon)$ -approximation to the submodular maximization problem under a knapsack constraint.*

In particular, our algorithm uses execution time $\mathcal{O}(\sqrt{n\tilde{K}}/\epsilon)$ and total communication, CPU time and number of queries $\mathcal{O}(n/\epsilon)$. Details are given in the supplementary material.

Relationship to previous work. The classic version of the problem considered in this work sets $c(e) = 1$ for all $e \in U$ and is known as monotone submodular maximization under a cardinality constraint and has been extensively studied. The celebrated result of [Nemhauser et al. \(1978\)](#) gives a $1 - 1/e \approx 0.63$ -approximation using GREEDY, which is optimal unless $P \neq NP$, [Feige \(1998\)](#). The problem of maximizing a monotone submodular function under a knapsack constraint was introduced by [Wolsey \(1982\)](#), who gave an algorithm with ≈ 0.35 -approximation. [Khuller et al. \(1999\)](#) gave a simple GREEDY+MAX algorithm with $1 - 1/\sqrt{e} \approx 0.39$ -approximation as well as a more complicated algorithm PARTIALENUM+GREEDY which requires a partial enumeration over an initial seed of three items and hence runs in $\mathcal{O}(\tilde{K}n^4)$ time. PARTIALENUM+GREEDY was later analyzed by [Sviridenko \(2004\)](#) who showed a $(1 - 1/e) \approx 0.63$ -approximation, matching the hardness of [Feige \(1998\)](#). While faster algorithms with $(1 - 1/e - \epsilon)$ -approximation exist [Badanidiyuru and Vondrák \(2014\)](#); [Ene and Nguyen \(2019a\)](#), they are self-admittedly impractical due to their expo-

nential dependence on large polynomials in $1/\epsilon$.

Compared to the well-studied cardinality-constrained case, streaming literature on monotone submodular optimization under a knapsack constraint is relatively sparse. A summary of results in the streaming setting is given in Figure 1. Prior to our work, the best results in streaming are by [Huang et al. \(2017\)](#); [Huang and Kakimura \(2019\)](#). While the most recent work of [Huang and Kakimura \(2019\)](#) achieves the $(1/2 - \epsilon)$ -approximation, its space, runtime and query complexities are far from optimal and depend on large polynomials of $1/\epsilon$, making it impractical for large data. Compared to this result, our Theorem 1.2 gives an improvement on all main parameters of interest, leading to near-optimal results. On the other hand, for the cardinality-constrained case, an optimal single-pass $(1/2 - \epsilon)$ -approximation has very recently been achieved by [Kazemi et al. \(2019\)](#). While using different ideas, our multi-pass streaming result matches theirs in terms of approximation, space and improves slightly on the number of queries and runtime (from $\mathcal{O}(n \log \tilde{K}/\epsilon)$ to $\mathcal{O}(n(1/\epsilon + \log \tilde{K}))$) only at the cost of using a constant number of passes for constant ϵ .

In the distributed setting, [Mirzasoleiman et al. \(2013\)](#) give an elegant two round protocol for monotone submodular maximization subject to a knapsack constraint that achieves a subconstant guarantee. [Kumar et al. \(2015\)](#) later give algorithms for both cardinality and matroid constraints that achieve a constant factor approximation, but the number of rounds is $\Theta(\log \Delta)$, where Δ is the maximum increase in the objective due to a single element, which is infeasible for large datasets since Δ even be significantly larger than the size of the entire dataset. [da Ponte Barbosa et al. \(2015, 2016\)](#) subsequently give a framework for both monotone and non-monotone submodular maximization under cardinality, matroid, and p -system constraints. Specifically, the results of [da Ponte Barbosa et al. \(2016\)](#) achieves almost $1/2$ -approximation for these settings using two rounds, a result subsequently matched by Liu and Vondrák [Liu and Vondrák \(2019\)](#) without requiring the duplication of items, as well as a $(1 - 1/e - \epsilon)$ approximation using $\mathcal{O}(1/\epsilon)$ rounds. [da Ponte Barbosa et al. \(2015\)](#) also gives a two-round algorithm for a knapsack constraint that achieves roughly 0.17-approximation in expectation.

For extensions to other constraints, non-monotone objectives and other generalizations see e.g. [Chakrabarti and Kale \(2014\)](#); [Chekuri et al. \(2015\)](#); [Chan et al. \(2017\)](#); [Elenberg et al. \(2017\)](#); [Epasto et al. \(2017\)](#); [Mirzasoleiman et al. \(2018\)](#); [Feldman et al. \(2018\)](#); [Chekuri and Quanrud \(2019\)](#).

Reference	Approx.	Passes	Space	Runtime and Queries
Huang et al. (2017)	$1/3 - \epsilon$	1	$\mathcal{O}(\frac{1}{\epsilon} K \log K)$	$\mathcal{O}(\frac{1}{\epsilon} n \log K)$
Huang et al. (2017)	$4/11 - \epsilon$	1	$\mathcal{O}(\frac{1}{\epsilon} K \log K)$	$\mathcal{O}(\frac{1}{\epsilon} n \log K)$
Huang et al. (2017)	$2/5 - \epsilon$	3	$\mathcal{O}(\frac{1}{\epsilon^2} K \log^2 K)$	$\mathcal{O}(\frac{1}{\epsilon} n \log K)$
Huang and Kakimura (2019)	$1/2 - \epsilon$	$\mathcal{O}(1/\epsilon)$	$\mathcal{O}(\frac{1}{\epsilon^2} K \log^2 K)$	$\mathcal{O}(\frac{1}{\epsilon^8} n \log^2 K)$
SIEVE+MAX (Alg. 2)	$1/2 - \epsilon$	$\mathcal{O}(1/\epsilon)$	$\mathcal{O}(K)$	$\mathcal{O}(n(\frac{1}{\epsilon} + \log K))$

Fig. 1: Monotone submodular maximization under a knapsack constraint in the streaming model.

Our techniques. Let $f(e|S) = f(e \cup S) - f(S)$ be the *marginal gain* and $\rho(e|S) = f(e|S)/c(e)$ be the *marginal density* of e with respect to S . GREEDY starts with an empty set G and repeatedly adds an item that maximizes $\rho(e|G)$ among the remaining items that fit. While by itself this does not guarantee any approximation, the classic result of Khuller et al. (1999) shows that GREEDY+MAX algorithm, which takes the best of the greedy solution and the single item with maximum value, gives a 0.39-approximation but cannot go beyond 0.44-approximation. Our algorithm GREEDY+MAX (Algorithm 1) instead attempts to augment every partial greedy solution with the item giving the largest marginal gain. For each i , let \mathcal{G}_i be the set of the first i items taken by greedy. We augment this solution with the item s_i which maximizes $f(s_i|\mathcal{G}_i)$ among the remaining items that fit. GREEDY+MAX then outputs the best solution among such augmentations.

Our main technical contribution lies in the analysis of this algorithm and its variants, which shows a $1/2$ -approximation (this analysis is tight, see the supplementary material). Let o_1 be the item from OPT with the largest cost. The main idea is to consider the last partial greedy solution such that o_1 still fits. Since o_1 has the largest cost in OPT, we can augment the partial solution with any element from OPT, and all of them have a non-greater marginal density than the next selected item. While GREEDY+MAX augments partial solutions with the best item, for the sake of analysis it suffices to consider only augmentations with o_1 (note that the item itself is unknown to the algorithm).

To simplify the presentation, in the analysis we rescale f and the costs so that $f(\text{OPT}) = 1$ and $K = 1$. Suppose that at some point, the partial greedy solution has collected elements with total cost $x \in [0, 1]$. We use a continuous function $g(x)$ to track the performance of GREEDY. We also introduce a function $g_1(x)$ to track the performance of augmentation with o_1 and then show that g and g_1 satisfy a differential inequality $g_1(x) + (1 - c(o_1))g'(x) \geq 1$ (Lemma 2.3), where g' denotes the right derivative. To give some intuition about the proof, consider the case when there exists

a partial greedy solution of cost exactly $1 - c(o_1)$. If $g_1(1 - c(o_1)) \geq 1/2$, then the augmentation with o_1 gives a $1/2$ -approximation. Otherwise, by the differential inequality, $g'(1 - c(o_1)) \geq 1/2(1 - c(o_1))$. Since $g(0) = 0$ and g' is non-increasing, $g(1 - c(o_1)) \geq (1 - c(o_1))g'(1 - c(o_1)) \geq 1/2$. See full analysis for how to handle the cases when there is no partial solution of cost exactly $1 - c(o_1)$.

Our streaming algorithm SIEVE+MAX and distributed algorithm DISTRIBUTED SIEVE+MAX approximately implement GREEDY+MAX in their respective settings. SIEVE+MAX makes $\mathcal{O}(1/\epsilon)$ passes over the data, and for each pass it selects items with marginal density at least a threshold $\frac{cf(\text{OPT})}{K(1+\epsilon)^i}$ in the i -th pass for some constant $c > 0$. This requires having a constant-factor approximation of $f(\text{OPT})$ which can be computed using a single pass. DISTRIBUTED SIEVE+MAX combines the thresholding approach with the sampling technique developed by Liu and Vondrák (2019) for the cardinality constraint. The differential inequality which we develop for GREEDY+MAX turns out to be robust to various sources of error introduced through thresholding and sampling. As we show, it continues to hold with functions and derivatives replaced with their $(1 + \epsilon)$ -approximations, which results in $(1/2 - \epsilon)$ -approximation guarantees for both algorithms.

2 Algorithms and analysis

Offline algorithm GREEDY+MAX. We introduce the main ideas by first describing our offline algorithm GREEDY+MAX which is then adapted to the streaming and distributed settings. Recall that GREEDY starts with an empty set G and in each iteration selects an item e with the highest marginal density $\rho(e|G)$ that still fits into the knapsack. GREEDY+MAX is based on augmenting each partial solution constructed by GREEDY with the item of the largest marginal value (as opposed to density) and taking the best among such augmentations. Recall that \mathcal{G}_i is the set of the first i items in the greedy solution. GREEDY+MAX finds for each i an augmenting item s_i which maximizes $f(s_i \cup \mathcal{G}_i)$ among all items that still fit. The final

output is the best among all such augmented solutions. Implementation is given as Algorithm 1. In the rest

Algorithm 1: Offline algorithm GREEDY+MAX

Input: Set of elements $E = e_1, \dots, e_n$, knapsack capacity K , cost function $c(\cdot)$, non-negative monotone submodular function f ;

Output: $\frac{1}{2}$ -approximation for submodular maximization under knapsack constraint;

$G \leftarrow \emptyset, S \leftarrow \emptyset$;

while $E \neq \emptyset$ **do**

$s \leftarrow \operatorname{argmax}_{e \in E} f(e | G)$;

if $f(S) < f(G \cup s)$ **then**

$S \leftarrow G \cup s$;

$a \leftarrow \operatorname{argmax}_{e \in E} \rho(e | G)$;

$G \leftarrow G \cup a$;

$K \leftarrow K - c(a)$;

 Remove all elements $e \in E$ with $c(e) > K$;

return S

of this section we give an outline of the key technical lemmas behind the proof that GREEDY+MAX gives $\frac{1}{2}$ -approximation.

W.l.o.g. and only for analysis of approximation we rescale the function values and costs so that $f(\text{OPT}) = 1$ and $c(\text{OPT}) = K = 1^4$. We first define a greedy performance function $g(x)$ which allows us to track the performance of the greedy solution in a continuous fashion. Let g_1, g_2, \dots, g_m be the elements in G in the order they were added and recall that $\mathcal{G}_i = \{g_1, \dots, g_i\}$. For a fixed x , let i be the smallest index such that $c(\mathcal{G}_i) > x$.

Definition 2.1 (Greedy performance function)

For $x \in [0, 1]$ we define $g(x)$ as:

$$g(x) = f(\mathcal{G}_{i-1}) + (x - c(\mathcal{G}_{i-1}))\rho(g_i | \mathcal{G}_{i-1}).$$

For the function g_1 we only consider adding o_1 , the largest item from OPT, to the current partial greedy solution. Consider the last item added by the greedy solution before the cost of this solution exceeds $1 - c(o_1)$. We define c^* so that $1 - c(o_1) - c^*$ is the cost of the greedy solution before this item is taken.

Definition 2.2 (GREEDY+MAX performance lower bound) For $x \in [0, 1 - c(o_1) - c^*]$ we define $g_1(x) = g(x) + f(o_1 | \mathcal{G}_{i-1})$.

Lemma 2.3 (GREEDY+MAX inequality) Let g' denote the right derivative of g . Then for all $x \in [0, 1 - c(o_1) - c^*]$, the following differential inequality

⁴Note that if $c(\text{OPT}) < K$ then we can set $K = c(\text{OPT})$ first as this does not affect $f(\text{OPT})$.

holds:

$$g_1(x) + (1 - c(o_1))g'(x) \geq 1$$

We defer the proof of Lemma 2.3 to the supplementary material.

Theorem 2.4 Recall that $\tilde{K} = \min(n, K)$ is an upper bound on the number of elements in feasible solutions. Then GREEDY+MAX gives a $\frac{1}{2}$ -approximation to the submodular maximization problem under a knapsack constraint and runs in $\mathcal{O}(\tilde{K}n)$ time.

Proof : By applying Lemma 2.3 at the point $x = 1 - c(o_1) - c^*$, we have:

$$g_1(1 - c(o_1) - c^*) + (1 - c(o_1))g'(1 - c(o_1) - c^*) \geq 1$$

If $g_1(1 - c(o_1) - c^*) \geq \frac{1}{2}$, then we have $\frac{1}{2}$ -approximation, because $g_1(1 - c(o_1) - c^*)$ is a lower bound on the value of the augmented solution when the cost of the greedy part is $1 - c(o_1) - c^*$. Otherwise:

$$\begin{aligned} g'(1 - c(o_1) - c^*) &\geq \frac{1 - g_1(1 - c(o_1) - c^*)}{1 - c(o_1)} \\ &> \frac{1}{2(1 - c(o_1))}. \end{aligned}$$

Note that since $g(0) = 0$ and g' is non-increasing by the definition of GREEDY, for any $x \in [0, 1]$ we have $g(x) \geq g'(x) \cdot x$:

$$g(x) \geq \int_{\chi=0}^x g'(\chi) d\chi \geq \int_{\chi=0}^x g'(x) d\chi = g'(x) \cdot x,$$

Therefore, applying this inequality at $x = 1 - c(o_1) - c^*$:

$$\begin{aligned} g(1 - c(o_1) - c^*) &\geq (1 - c(o_1) - c^*)g'(1 - c(o_1) - c^*) \\ &\geq \frac{1 - c(o_1) - c^*}{2(1 - c(o_1))}. \end{aligned}$$

Recall that $1 - c(o_1) - c^*$ was the last cost of the greedy solution when we could still augment it with o_1 ; therefore, the next element e that the greedy solution selects has the cost at least $(1 - c(o_1)) - (1 - c(o_1) - c^*) = c^*$. Thus, the function value after taking e is at least

$$\begin{aligned} &g(1 - c(o_1) - c^*) + c^*g'(1 - c(o_1) - c^*) \\ &\geq \frac{1 - c(o_1) - c^*}{2(1 - c(o_1))} + \frac{c^*}{2(1 - c(o_1))} = \frac{1}{2} \end{aligned}$$

Hence, Algorithm 1 gives a $\frac{1}{2}$ -approximation to the submodular maximization problem under a knapsack constraint. It remains to analyze the running time and

query complexity of Algorithm 1. Since \tilde{K} is the maximum size of a feasible set, Algorithm 1 makes at most \tilde{K} iterations. In each iteration, it makes $\mathcal{O}(n)$ oracle queries, so the total number of queries and runtime is $\mathcal{O}(\tilde{K}n)$. \square

Streaming algorithm SIEVE+MAX. Our multi-pass streaming algorithm is given as Algorithm 2. To simplify the presentation, we first give the algorithm under the assumption that it is given a parameter λ , which is a constant-factor approximation of $f(\text{OPT})$. We then show how to remove this assumption using standard techniques in the supplementary material. As discussed in the description of our techniques SIEVE+MAX uses $\mathcal{O}(1/\epsilon)$ passes over the data to simulate the execution of GREEDY+MAX approximately.

Algorithm 2: Multi-pass streaming algorithm SIEVE+MAX

Input: Stream e_1, \dots, e_n , knapsack capacity K , cost function $c(\cdot)$, non-negative monotone submodular function f , λ which is an α -approximation of $f(\text{OPT})$ for some fixed constant $\alpha > 0$, $\epsilon > 0$;

Output: $(1/2 - \epsilon)$ -approx. for submodular maximization under a knapsack constraint;

$T \leftarrow \emptyset$, $\tau \leftarrow \frac{\lambda}{\alpha K}$;
while $\tau > \frac{\lambda}{2K}$ **do** // Thresholding stage
 Take a new pass over the stream;
 for each read item e **do**
 if $\rho(e|T) \geq \tau$ **and** $c(e \cup T) \leq K$ **then**
 $T \leftarrow T \cup \{e\}$;
 $\tau \leftarrow \tau/(1 + \epsilon)$;

For each i , let G_i be the first i items selected in the construction of T above and initialize $s_i = \emptyset$ for the best augmenting item for G_i ;

Take a pass over the stream;

for each read item e **do** // Augmentation stage
 if $e \notin T$ **then**
 $j = \max\{i | c(G_i) + c(e) \leq K\}$;
 if $f(G_j \cup s_j) < f(G_j \cup e)$ **then**
 $s_j \leftarrow \{e\}$;
return $\arg\max f(G_i \cup s_i)$

Let T be the set of items constructed SIEVE+MAX (as in Algorithm 2) and let t_1, t_2, \dots be the order that they are collected. We refer to the part of the algorithm which constructs T as “thresholding” and the rest as “augmentation” below. We use \mathcal{T}_i to denote the set containing the i items $\{t_1, t_2, \dots, t_i\}$. We again use o_1 to denote the item with highest cost in OPT. Similar to the above, we define two functions representing the values of our thresholding algorithm, and augmented solutions given the utilized proportion of the knapsack.

Definition 2.5 (Thresholding performance function) For any $x \in [0, 1]$, let i be the smallest index such that $c(\mathcal{T}_i) > x$. We define $t(x) = f(\mathcal{T}_{i-1}) + (x - c(\mathcal{T}_{i-1}))\rho(t_i | \mathcal{T}_{i-1})$.

We define a function $t_1(x)$ that lower bounds the performance of SIEVE+MAX when the thresholding solution collects a set of cost x :

Definition 2.6 (SIEVE+MAX performance function and lower bound) For any fixed x , let i be the smallest index such that $c(\mathcal{T}_i) > x$. Then we define $t_1(x) = t(x) + f(o_1 | \mathcal{T}_{i-1})$, where $o_1 = \arg\max_{e \in \text{OPT}} c(e)$.

In order to analyze the output of the algorithm, we prove a differential inequality for t_1 . If $c(T) \geq 1 - c(o_1)$ then let $c^* \geq 0$ be defined so that $1 - c(o_1) - c^*$ is the cost of the thresholding solution before the algorithm takes the item which makes the cost exceed $1 - c(o_1)$.

Lemma 2.7 (SIEVE+MAX Inequality) If $c(T) \geq 1 - c(o_1)$ then for all $x \in [0, 1 - c(o_1) - c^*]$, then t and t_1 satisfy the following differential inequality:

$$t_1(x) + (1 + \epsilon)(1 - c(o_1))t'(x) \geq 1.$$

Theorem 2.8 There exists an algorithm that uses $\mathcal{O}(\tilde{K})$ space and $\mathcal{O}(1/\epsilon)$ passes over the stream, makes $\mathcal{O}(n/\epsilon + n \log \tilde{K})$ queries, and outputs a $(1/2 - \epsilon)$ -approximation to the submodular maximization problem under a knapsack constraint.

The proofs are similar to the proofs of Lemma 2.3 and Theorem 2.4 and are provided in the supplementary material.

Query lower bound. We show a simple query lower bound under the standard assumption Norouzi-Fard et al. (2018); Kazemi et al. (2019) that the algorithm only queries f on feasible sets.

Theorem 2.9 For $\alpha > 1/2$, any α -approximation algorithm for maximizing a function f under a knapsack constraint that succeeds with constant probability and only queries values of the function f on feasible sets (i.e. sets of cost at most K) must make at least $\Omega(n^2)$ queries if f is either: 1) non-monotone submodular, 2) monotone and submodular on the feasible sets, 3) monotone subadditive.

We defer the proof of Theorem 2.9 to the supplementary material.

3 Experimental results

We compare our offline algorithm GREEDY+MAX and our streaming algorithm SIEVE+MAX with baselines, answering the following questions: (1) What are the approximation factors we are getting on real data? (2) How do the objective values compare? (3) How do the runtimes compare? (4) How do the numbers of queries compare? We compare GREEDY+MAX to the offline baselines GREEDY, GREEDYORMAX (Khuller et al. (1999)), as well as PARTIALENUM+GREEDY (Khuller et al., 1999; Sviridenko, 2004), which enumerates knapsacks containing all subsets of d items, and then runs the GREEDY algorithm upon each of those knapsacks. PARTIALENUM+GREEDY uses $\Omega(Kn^{d+1})$ runtime and queries, so only $d = 1$ and the smallest ego-Facebook dataset are feasible in our experiments. In streaming we compare SIEVE+MAX to SIEVE (Badanidiyuru et al. (2014)) and SIEVEORMAX (Huang et al. (2017)), which are similar thresholding-based algorithms. We also implemented a single-pass BRANCHINGMRT by Huang et al. (2017) that uses thresholding along with multiple branches and gives a $4/11 \approx 0.36$ -approximation. We did not implement Huang and Kakimura (2019) as their algorithms are orders of magnitude slower than BRANCHINGMRT which is already several orders of magnitude slower than other algorithms.

Our code is available at <https://github.com/aistats20submodular/aistats20submodular>. We used two types of datasets:

Graph coverage. For a graph $G(V, E)$ and $Z \subset V$, the objective is to maximize the neighborhood vertex coverage function $f(Z) := |Z \cup N(Z)|/|V|$, where $N(Z)$ is the set of neighbors of Z . We ran experiments on two graphs from SNAP Leskovec and Krevl (2014): ego-Facebook and com-DBLP.

Movie ratings. We also analyze a dataset of movies to model the scenario of movie recommendation. The objective function, defined as in Avdiukhin et al. (2019), is maximized for a set of movies that is similar to a user’s interests. For more details about the settings, see the supplementary material. We analyze the ml-20 MovieLens dataset GroupLens (2015).

Approximation and runtimes. We first give instance-specific approximation factors for different values of K for offline (Fig. 2) and streaming (Fig. 3) algorithms. These approximations are computed using upper bounds on $f(\text{OPT})$ which can be obtained using the analysis of GREEDY. GREEDY+MAX and SIEVE+MAX typically perform at least 20% better than their $1/2$ worst-case guarantees. In the supplementary material we show that the value can be improved by up to 50%, both by GREEDY+MAX upon GREEDY and

by SIEVE+MAX upon SIEVE.

In the supplementary material we show that runtimes of GREEDY+MAX and GREEDYORMAX are similar and at most 20% greater than the runtime of GREEDY. On the other hand, even though PARTIALENUM+GREEDY does not outperform GREEDY+MAX, it is only feasible for $d = 1$ and the ego-Facebook dataset and uses on average almost 500 times as much runtime for $K = 10$ across ten iterations of each algorithm.

For the streaming algorithms, Figure 4 shows that runtimes of SIEVE+MAX, SIEVEORMAX, and SIEVE performs generally similar; however in the case of the com-dblp dataset, the runtime of SIEVE+MAX grows faster with K . This can be explained by the fact that oracle calls on larger sets typically require more time and augmenting sets typically contain more elements than sets encountered during execution of SIEVE. On the other hand, the runtime of BRANCHINGMRT was substantially slower, and we did not include its runtime for scaling purposes. E.g. for $K = 5$, the runtime of BRANCHINGMRT was already a factor 80K more than SIEVE.

In the supplementary material, we also compare the number of oracle queries performed by the algorithms. GREEDY+MAX, GREEDYORMAX and GREEDY require the same amount of oracle calls, since computing marginal gains and finding the best element for augmentation can be done using the same queries. On the other hand, PARTIALENUM+GREEDY requires 544x more calls than GREEDY for $K = 8$. For the streaming algorithms, the number of oracle calls made by SIEVE, SIEVE+MAX, and SIEVE, never differed by more than a factor of two, while BRANCHINGMRT requires a factor 125K more oracle calls than SIEVE for $K = 8$.

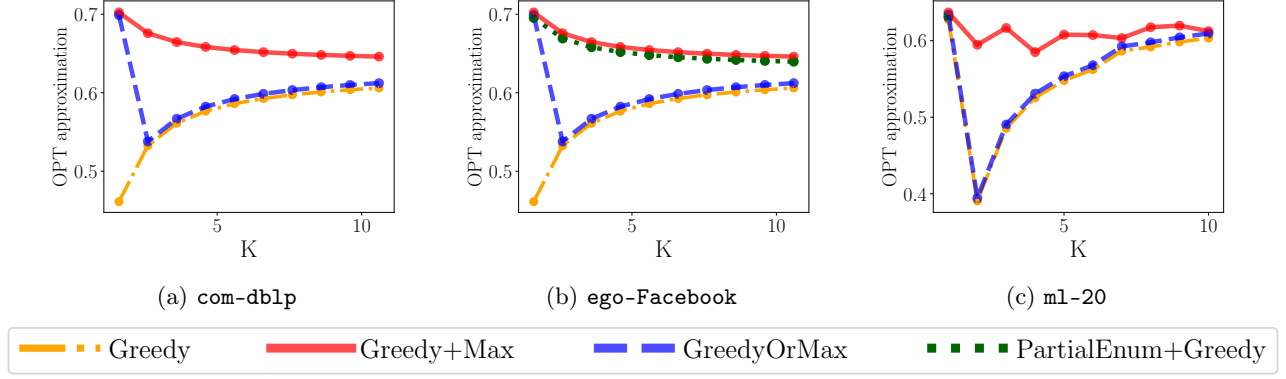


Fig. 2: Instance-specific approximations for different K . GREEDY+MAX performs substantially better than its worst-case $1/2$ -approximation guarantee and typically beats even the $(1 - 1/e) \approx 0.63$ bound. Despite much higher runtime, PARTIALENUM+GREEDY does not beat GREEDY+MAX even on the only dataset where its runtime is feasible (ego-Facebook).

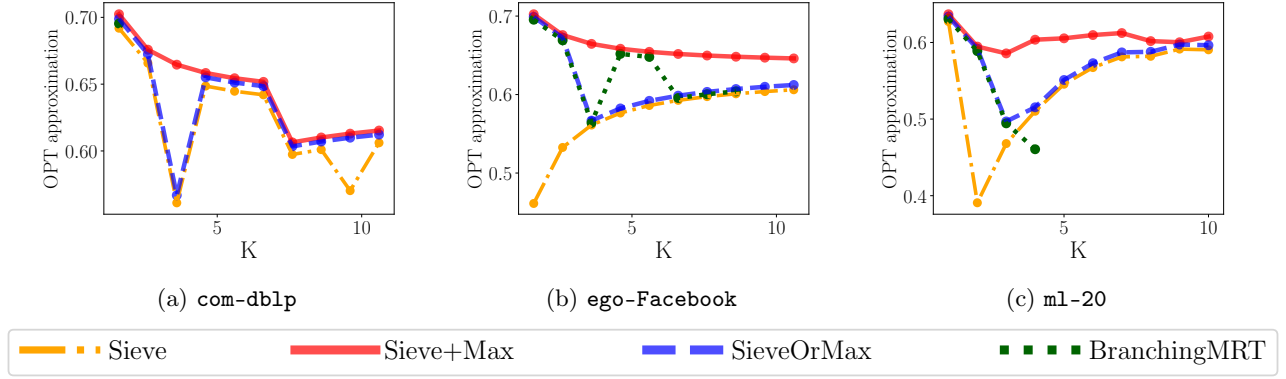


Fig. 3: Instance-specific approximations for different K . SIEVE+MAX performs substantially better than its worst-case $(1/2 - \epsilon)$ -approximation guarantee and robustly dominates all other approaches. It can improve by up to 40% upon SIEVE. Despite much higher runtime, BRANCHINGMRT does not beat SIEVE+MAX

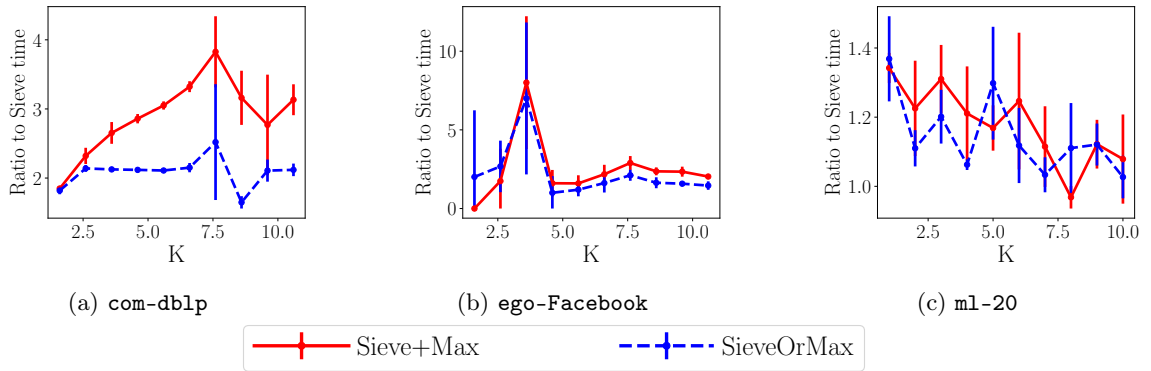


Fig. 4: Ratios between average runtimes of streaming algorithms and average runtimes of SIEVE over ten executions. Error bars show standard deviation across runs.

References

- Shipra Agrawal, Mohammad Shadravan, and Cliff Stein. Submodular secretary problem with shortlists. In *10th Innovations in Theoretical Computer Science Conference, ITCS*, 2019.
- Nima Anari, Nika Haghtalab, Seffi Naor, Sebastian Pokutta, Mohit Singh, and Alfredo Torrico. Structured robust submodular maximization: Offline and online algorithms. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pages 3128–3137. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/anari19a.html>.
- Alexandr Andoni, Aleksandar Nikolov, Krzysztof Onak, and Grigory Yaroslavtsev. Parallel algorithms for geometric graph problems. In *Symposium on Theory of Computing, STOC*, pages 574–583, 2014.
- Dmitrii Avdiukhin, Slobodan Mitrovic, Grigory Yaroslavtsev, and Samson Zhou. Adversarially robust submodular maximization under knapsack constraints. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD.*, pages 148–156, 2019.
- Ashwinkumar Badanidiyuru and Jan Vondrák. Fast algorithms for maximizing submodular functions. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 1497–1514, 2014.
- Ashwinkumar Badanidiyuru, Baharan Mirzasoleiman, Amin Karbasi, and Andreas Krause. Streaming submodular maximization: Massive data summarization on the fly. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 671–680. ACM, 2014.
- Eric Balkanski and Yaron Singer. The adaptive complexity of maximizing a submodular function. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC*, pages 1138–1151, 2018.
- Eric Balkanski, Aviad Rubinfeld, and Yaron Singer. An exponential speedup in parallel running time for submodular maximization without loss in approximation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 283–302, 2019.
- MohammadHossein Bateni, Hossein Esfandiari, and Vahab S. Mirrokni. Optimal distributed submodular optimization via sketching. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD*, pages 1138–1147, 2018.
- Ilija Bogunovic, Slobodan Mitrović, Jonathan Scarlett, and Volkan Cevher. Robust submodular maximization: A non-uniform partitioning approach. In *Proceedings of the 34th International Conference on Machine Learning, ICML*, pages 508–516, 2017.
- Niv Buchbinder, Moran Feldman, and Roy Schwartz. Online submodular maximization with preemption. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 1202–1216, 2015.
- Gruia Călinescu, Chandra Chekuri, Martin Pál, and Jan Vondrák. Maximizing a monotone submodular function subject to a matroid constraint. *SIAM J. Comput.*, 40(6):1740–1766, 2011.
- Amit Chakrabarti and Sagar Kale. Submodular maximization meets streaming: Matchings, matroids, and more. In *Integer Programming and Combinatorial Optimization - 17th International Conference, IPCO. Proceedings*, pages 210–221, 2014.
- T.-H. Hubert Chan, Zhiyi Huang, Shaofeng H.-C. Jiang, Ning Kang, and Zhihao Gavin Tang. Online submodular maximization with free disposal: Randomization beats 0.25 for partition matroids. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 1204–1223, 2017.
- Chandra Chekuri and Kent Quanrud. Submodular function maximization in parallel via the multilinear relaxation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 303–322, 2019.
- Chandra Chekuri, Shalmoli Gupta, and Kent Quanrud. Streaming algorithms for submodular function maximization. In *Automata, Languages, and Programming - 42nd International Colloquium, ICALP, Proceedings, Part I*, pages 318–330, 2015.
- Rafael da Ponte Barbosa, Alina Ene, Huy L. Nguyen, and Justin Ward. The power of randomization: Distributed submodular maximization on massive datasets. In *Proceedings of the 32nd International Conference on Machine Learning, ICML*, pages 1236–1244, 2015.
- Rafael da Ponte Barbosa, Alina Ene, Huy L. Nguyen, and Justin Ward. A new framework for distributed submodular maximization. In *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS*, pages 645–654, 2016.
- Khalid El-Arini and Carlos Guestrin. Beyond keyword search: discovering relevant scientific literature. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 439–447. ACM, 2011.

- Ethan R. Elenberg, Alexandros G. Dimakis, Moran Feldman, and Amin Karbasi. Streaming weak submodularity: Interpreting neural networks on the fly. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*, pages 4047–4057, 2017.
- Alina Ene and Huy L. Nguyen. A nearly-linear time algorithm for submodular maximization with a knapsack constraint. In *46th International Colloquium on Automata, Languages, and Programming, ICALP*, pages 53:1–53:12, 2019a.
- Alina Ene and Huy L. Nguyen. Submodular maximization with nearly-optimal approximation and adaptivity in nearly-linear time. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 274–282, 2019b.
- Alessandro Epasto, Silvio Lattanzi, Sergei Vassilvitskii, and Morteza Zadimoghaddam. Submodular optimization over sliding windows. In *Proceedings of the 26th International Conference on World Wide Web, WWW*, pages 421–430, 2017.
- Matthew Fahrbach, Vahab S. Mirrokni, and Morteza Zadimoghaddam. Submodular maximization with nearly optimal approximation, adaptivity and query complexity. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 255–273, 2019.
- Uriel Feige. A threshold of $\ln n$ for approximating set cover. *J. ACM*, 45(4):634–652, 1998.
- Moran Feldman, Amin Karbasi, and Ehsan Kazemi. Do less, get more: Streaming submodular maximization with subsampling. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS*, pages 730–740, 2018.
- Daniel Golovin and Andreas Krause. Adaptive submodularity: Theory and applications in active learning and stochastic optimization. *Journal of Artificial Intelligence Research*, 42:427–486, 2011.
- Ryan Gomes and Andreas Krause. Budgeted nonparametric learning from data streams. In *Proceedings of the 27th International Conference on Machine Learning, ICML*, pages 391–398, 2010.
- GroupLens. <https://grouplens.org/datasets/movielens>, 2015. MovieLens Datasets.
- Jason Hartline, Vahab Mirrokni, and Mukund Sundararajan. Optimal marketing strategies over social networks. In *Proceedings of the 17th International Conference on World Wide Web, WWW '08*, pages 189–198, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-085-2. doi: 10.1145/1367497.1367524. URL <http://doi.acm.org/10.1145/1367497.1367524>.
- Chien-Chung Huang and Naonori Kakimura. Multi-pass streaming algorithms for monotone submodular function maximization. In *Algorithms and Data Structures - 16th International Symposium, WADS*, 2019. (to appear). <https://arxiv.org/abs/1802.06212>.
- Chien-Chung Huang, Naonori Kakimura, and Yuichi Yoshida. Streaming algorithms for maximizing monotone submodular functions under a knapsack constraint. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM*, pages 11:1–11:14, 2017.
- Rishabh Iyer and Jeffrey Bilmes. Near optimal algorithms for hard submodular programs with discounted cooperative costs. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pages 276–285. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/iyer19a.html>.
- Howard J. Karloff, Siddharth Suri, and Sergei Vassilvitskii. A model of computation for mapreduce. In *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 938–948, 2010.
- Ehsan Kazemi, Marko Mitrovic, Morteza Zadimoghaddam, Silvio Lattanzi, and Amin Karbasi. Submodular streaming in all its glory: Tight approximation, minimum memory and low adaptive complexity. In *Proceedings of the 36th International Conference on Machine Learning, ICML*, pages 3311–3320, 2019.
- Rajiv Khanna, Joydeep Ghosh, Russell Poldrack, and Oluwasanmi Koyejo. Sparse Submodular Probabilistic PCA. In Guy Lebanon and S. V. N. Vishwanathan, editors, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pages 453–461, San Diego, California, USA, 09–12 May 2015. PMLR. URL <http://proceedings.mlr.press/v38/khanna15.html>.
- Samir Khuller, Anna Moss, and Joseph Naor. The budgeted maximum coverage problem. *Inf. Process. Lett.*, 70(1):39–45, 1999.
- Ravi Kumar, Benjamin Moseley, Sergei Vassilvitskii, and Andrea Vattani. Fast greedy algorithms in mapreduce and streaming. *ACM Transactions on Parallel Computing (TOPC)*, 2(3):14, 2015.
- Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie

- Glance. Cost-effective outbreak detection in networks. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 420–429. ACM, 2007.
- Hui Lin and Jeff Bilmes. A class of submodular functions for document summarization. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 510–520. Association for Computational Linguistics, 2011.
- Erik M. Lindgren, Shanshan Wu, and Alexandros G. Dimakis. Leveraging sparsity for efficient submodular data summarization. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems*, pages 3414–3422, 2016.
- Paul Liu and Jan Vondrák. Submodular optimization in the mapreduce model. In *2nd Symposium on Simplicity in Algorithms, SOSA@SODA 2019*, pages 18:1–18:10, 2019.
- Vahab S. Mirrokni and Morteza Zadimoghaddam. Randomized composable core-sets for distributed submodular maximization. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC*, pages 153–162, 2015.
- Baharan Mirzasoleiman, Amin Karbasi, Rik Sarkar, and Andreas Krause. Distributed submodular maximization: Identifying representative elements in massive data. In *Advances in Neural Information Processing Systems*, pages 2049–2057, 2013.
- Baharan Mirzasoleiman, Morteza Zadimoghaddam, and Amin Karbasi. Fast distributed submodular cover: Public-private data summarization. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems*, pages 3594–3602, 2016.
- Baharan Mirzasoleiman, Stefanie Jegelka, and Andreas Krause. Streaming non-monotone submodular maximization: Personalized video summarization on the fly. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18)*, pages 1379–1386, 2018.
- Slobodan Mitrović, Ilija Bogunovic, Ashkan Norouzi-Fard, Jakub Tarnawski, and Volkan Cevher. Streaming robust submodular maximization: A partitioned thresholding approach. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*, pages 4560–4569, 2017.
- George L. Nemhauser, Laurence A. Wolsey, and Marshall L. Fisher. An analysis of approximations for maximizing submodular set functions - I. *Math. Program.*, 14(1):265–294, 1978.
- Ashkan Norouzi-Fard, Jakub Tarnawski, Slobodan Mitrovic, Amir Zandieh, Aidasadat Mousavifar, and Ola Svensson. Beyond 1/2-approximation for submodular maximization on massive data streams. In *Proceedings of the 35th International Conference on Machine Learning, ICML*, pages 3826–3835, 2018.
- Ruben Sipos, Adith Swaminathan, Pannaga Shivashwamy, and Thorsten Joachims. Temporal corpus summarization using submodular word coverage. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pages 754–763. ACM, 2012.
- Maxim Sviridenko. A note on maximizing a submodular set function subject to a knapsack constraint. *Oper. Res. Lett.*, 32(1):41–43, 2004.
- Kai Wei, Yuzong Liu, Katrin Kirchhoff, and Jeff Bilmes. Using document summarization techniques for speech data subset selection. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 721–726, 2013.
- Laurence A. Wolsey. Maximising real-valued submodular functions: Primal and dual heuristics for location problems. *Math. Oper. Res.*, 7(3):410–425, 1982.
- Qilian Yu, Easton Li Xu, and Shuguang Cui. Streaming algorithms for news and scientific literature recommendation: Monotone submodular maximization with a $\$d\$$ -knapsack constraint. *IEEE Access*, 6: 53736–53747, 2018.