

# Streaming Algorithms for Maximising a Submodular Function

November 18, 2019

## Introduction

A set function  $f : 2^E \rightarrow \mathbb{R}_+$  on a ground set  $E$  is *submodular* if it satisfies the *diminishing marginal return property*, i.e., for any subsets  $S \subseteq T \subsetneq E$  and  $e \in E \setminus T$ ,

$$f(S \cup \{e\}) - f(S) \geq f(T \cup \{e\}) - f(T).$$

A function is *monotone* if  $f(S) \leq f(T)$  for any  $S \subseteq T$ . Submodular functions play a fundamental role in combinatorial optimization, as they capture rank functions of matroids, edge cuts of graphs, and set coverage, just to name a few examples. Besides their theoretical interests, submodular functions have attracted much attention from the machine learning community because they can model various practical problems such as online advertising [1, 2, 7], sensor location [3], text summarization [5, 6], and maximum entropy sampling [4].

Many of the aforementioned applications can be formulated as the maximization of a monotone submodular function under a knapsack constraint. In this problem, we are given a monotone submodular function  $f : 2^E \rightarrow \mathbb{R}_+$ , a size function  $c : E \rightarrow \mathbb{N}$ , and an integer  $K \in \mathbb{N}$ , where  $\mathbb{N}$  denotes the set of positive integers. The problem is defined as

$$\text{maximize } f(S) \quad \text{subject to } c(S) \leq K, \quad S \subseteq E, \quad (1)$$

where we denote  $c(S) = \sum_{e \in S} c(e)$  for a subset  $S \subseteq E$ . Note that, when  $c(e) = 1$  for every item  $e \in E$ , the constraint coincides with a cardinality constraint:

$$\text{maximize } f(S) \quad \text{subject to } |S| \leq K, \quad S \subseteq E. \quad (2)$$

## Internship Project

This project involves designing *streaming* algorithms for the above two problems (which are known to be approximable within the factor of  $1 - e^{-1}$  in the *offline* setting). The *streaming* setting means: each item in the ground set  $E$  arrives sequentially, and we can keep only a small number of the items in memory at any point.

Designing streaming algorithms for submodular function optimisation is a relatively unexplored area. There is much gap between the upper bound and the lower bound (in terms of approximation ratios). The student working on this project will have much chance of making new progress.

## References

- [1] N. Alon, I. Gamzu, and M. Tennenholtz. Optimizing budget allocation among channels and influencers. In *Proceedings of the 21st International Conference on World Wide Web (WWW)*, pages 381–388, 2012.

- [2] D. Kempe, J. Kleinberg, and É. Tardos. Maximizing the spread of influence through a social network. In *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 137–146, 2003.
- [3] A. Krause, A. P. Singh, and C. Guestrin. Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning Research*, 9:235–284, 2008.
- [4] J. Lee. *Maximum Entropy Sampling*, volume 3 of *Encyclopedia of Environmetrics*, pages 1229–1234. John Wiley & Sons, Ltd., 2006.
- [5] H. Lin and J. Bilmes. Multi-document summarization via budgeted maximization of submodular functions. In *Proceedings of the 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 912–920, 2010.
- [6] H. Lin and J. 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 (ACL-HLT)*, pages 510–520, 2011.
- [7] T. Soma, N. Kakimura, K. Inaba, and K. Kawarabayashi. Optimal budget allocation: Theoretical guarantee and efficient algorithm. In *Proceedings of the 31st International Conference on Machine Learning (ICML)*, pages 351–359, 2014.