

Curvature and Optimal Algorithms for Learning and Minimizing Submodular Functions

Authored by:

Jeff A. Bilmes
Stefanie Jegelka
Rishabh K. Iyer

Abstract

We investigate three related and important problems connected to machine learning, namely approximating a submodular function everywhere, learning a submodular function (in a PAC like setting [26]), and constrained minimization of submodular functions. In all three problems, we provide improved bounds which depend on the ‘curvature’ of a submodular function and improve on the previously known best results for these problems [9, 3, 7, 25] when the function is not too curved – a property which is true of many real-world submodular functions. In the former two problems, we obtain these bounds through a generic black-box transformation (which can potentially work for any algorithm), while in the case of submodular minimization, we propose a framework of algorithms which depend on choosing an appropriate surrogate for the submodular function. In all these cases, we provide almost matching lower bounds. While improved curvature-dependent bounds were shown for monotone submodular maximization [4, 27], the existence of similar improved bounds for the aforementioned problems has been open. We resolve this question in this paper by showing that the same notion of curvature provides these improved results. Empirical experiments add further support to our claims.

1 Paper Body

We investigate three related and important problems connected to machine learning: approximating a submodular function everywhere, learning a submodular function (in a PAC-like setting [28]), and constrained minimization of submodular functions. We show that the complexity of all three problems depends on the ‘curvature’ of the submodular function, and provide lower and upper bounds that refine and improve previous results [2, 6, 8, 27]. Our proof techniques are fairly generic. We either use a black-box transformation of the function (for approximation and learning), or a transformation of algorithms to use an appropriate surrogate function (for minimization). Curiously, curvature

has been known to influence approximations for submodular maximization [3, 29], but its effect on minimization, approximation and learning has hitherto been open. We complete this picture, and also support our theoretical claims by empirical results.

1

Introduction

Submodularity is a pervasive and important property in the areas of combinatorial optimization, economics, operations research, and game theory. In recent years, submodularity's use in machine learning has begun to proliferate as well. A set function $f : 2V \rightarrow \mathbb{R}$ over a finite set $V = \{1, 2, \dots, n\}$ is submodular if for all subsets $S, T \subseteq V$, it holds that $f(S) + f(T) \geq f(S \cup T) + f(S \cap T)$. Given a set $S \subseteq V$, we define the gain of an element $j \notin S$ as $f(S \cup \{j\}) - f(S)$. A function f is submodular if it satisfies diminishing marginal returns, namely $f(j \mid S) \geq f(j \mid T)$ for all $S \subseteq T, j \notin T$, and is monotone if $f(j \mid S) \geq 0$ for all $j \notin S, S \subseteq V$. While submodularity, like convexity, occurs naturally in a wide variety of problems, recent studies have shown that in the general case, many submodular problems of interest are very hard: the problems of learning a submodular function or of submodular minimization under constraints do not even admit constant or logarithmic approximation factors in polynomial time [2, 7, 8, 10, 27]. These rather pessimistic results however stand in sharp contrast to empirical observations, which suggest that these lower bounds are specific to rather contrived classes of functions, whereas much better results can be achieved in many practically relevant cases. Given the increasing importance of submodular functions in machine learning, these observations beg the question of qualifying and quantifying properties that make sub-classes of submodular functions more amenable to learning and optimization. Indeed, limited prior work has shown improved results for constrained minimization and learning of sub-classes of submodular functions, including symmetric functions [2, 25], concave functions [7, 18, 24], label cost or covering functions [9, 31]. In this paper, we take additional steps towards addressing the above problems and show how the generic notion of the curvature – the deviation from modularity – of a submodular function determines both upper and lower bounds on approximation factors for many learning and constrained optimization problems. In particular, our quantification tightens the generic, function-independent bounds in [8, 2, 27, 7, 10] for many practically relevant functions. Previously, the concept of curvature has been used to

tighten bounds for submodular maximization problems [3, 29]. Hence, our results complete a unifying picture of the effect of curvature on submodular problems. Moreover, curvature is still a fairly generic concept, as it only depends on the marginal gains of the submodular function. It allows a smooth transition between the ‘easy’ functions and the ‘really hard’ subclasses of submodular functions.

2

Problem statements, definitions and background

Before stating our main results, we provide some necessary definitions and

introduce a new concept, the curve normalized version of a submodular function. Throughout this paper, we assume that the submodular function f is defined on a ground set V of n elements, that it is nonnegative and $f(\emptyset) = 0$. We also use normalized modular (or P -additive) functions $w : 2^V \rightarrow \mathbb{R}$ which are those that can be written as a sum of weights, $w(S) = \sum_{i \in S} w(i)$. We are concerned with the following three problems: Problem 1. (Approximation [8]) Given a submodular function f in form of a value oracle, find an approximation \hat{f} (within polynomial time and representable within polynomial space), such that for all $X \subseteq V$, it holds that $\hat{f}(X) \leq f(X) \leq (1 + \epsilon) \hat{f}(X)$ for a polynomial $\epsilon = \epsilon(n)$. Problem 2. (PMAC-Learning [2]) Given i.i.d training samples $\{(X_i, f(X_i))\}_{i=1}^m$ from a distribution D , learn an approximation $\hat{f}(X)$ that is, with probability $1 - \delta$, within a multiplicative factor of $(1 + \epsilon)$ from f . Problem 3. (Constrained optimization [27, 7, 10, 16]) Minimize a submodular function f over a family C of feasible sets, i.e., $\min_{X \in C} f(X)$. In its general form, the approximation problem was first studied by Goemans et al. [8], who approximate any monotone submodular function to within a factor of $O(\sqrt{n \log n})$, with a lower bound of $\Omega(\sqrt{n / \log n})$. Building on this result, Balcan and Harvey [2] show how to PMAC-learn a monotone submodular function within a factor of $\epsilon(n) = O(\sqrt{n})$, and prove a lower bound of $\Omega(\sqrt{n}/3)$ for the learning problem. Subsequent work extends these results to sub-additive and fractionally sub-additive functions [1]. Better learning results are possible for the subclass of submodular shells [23] and Fourier sparse set functions [26]. Both Problems 1 and 2 have numerous applications in algorithmic game theory and economics [2, 8] as well as machine learning [2, 22, 23, 26, 15]. Constrained submodular minimization arises in applications such as power assignment or transportation problems [19, 30, 13]. In machine learning, it occurs, for instance, in the form of MAP inference in high-order graphical models [17] or in size-constrained corpus extraction [21]. Recent results show that almost all constraints make it hard to solve the minimization even within a constant factor [27, 6, 16]. Here, we will focus on the constraint of imposing a lower bound on the cardinality, and on combinatorial constraints where C is the set of all s - t paths, s - t cuts, spanning trees, or perfect matchings in a graph. A central concept in this work is the total curvature κ_f of a submodular function f and the curvature $\kappa_f(S)$ with respect to a set $S \subseteq V$, defined as [3, 29]

$$\begin{aligned} \kappa_f &= 1 - \min_{j \in V} \frac{f(j) - f(\emptyset)}{f(j)} \\ \kappa_f(S) &= 1 - \min_{j \in S} \frac{f(j) - f(S - j)}{f(j)} \end{aligned} \quad (1)$$

Without loss of generality, assume that $f(j) \geq 0$ for all $j \in V$. It is easy to see that $\kappa_f(S) \leq \kappa_f(V) = \kappa_f$, and hence $\kappa_f(S)$ is a tighter notion of curvature. A modular function has curvature $\kappa_f = 0$, and a matroid rank function has maximal curvature $\kappa_f = 1$. Intuitively, κ_f measures how far away f is from being modular. Conceptually, curvature is distinct from the recently proposed submodularity ratio [5] that measures how far a function is from being submodular. Curvature has served to tighten bounds for submodular maximization problems, e.g., from $(1 - 1/e)$ to $(1 - \epsilon) \kappa_f$ for monotone submodular maximization subject to a

cardinality constraint [3] or matroid constraints [29], and these results are tight. For submodular minimization, learning, and approximation, however, the role of curvature has not yet been addressed (an exception are the upper bounds in [13] for minimization). In the following sections, we complete the picture of how curvature affects the complexity of submodular maximization and minimization, approximation, and learning. The above-cited lower bounds for Problems 1-3 were established with functions of maximal curvature ($\kappa = 1$) which, as we will see, is the worst case. By contrast, many practically interesting functions have smaller curvature, and our analysis will provide an explanation for the good empirical results

observed with such functions [13, 22, 14]. An example for functions with $\kappa \leq 1$ is the class of concave over modular functions that have been used in speech processing [22] and computer vision [17]. This class comprises, for instance, functions of the form $f(X) = \sum_{i=1}^n (w_i(X))^a$, for some $a \in [0, 1]$ and nonnegative weight vectors w_i . Such functions may be defined over clusters $C_i \subseteq V$, in which case the weights $w_i(j)$ are nonzero only if $j \in C_i$ [22, 17, 11].

Curvature-dependent analysis. To analyze Problems 1-3, we introduce the concept of a *curve-normalized polymatroid*. Specifically, we define the κ -curve-normalized version of f as $P_\kappa f(X) = (1 - \kappa) \sum_{j \in X} f(j) + \kappa f(X)$. If $\kappa = 0$, then we set $P_0 f = f$. We call $P_\kappa f$ the *curve-normalized version* of f because its curvature is κ . The function $P_\kappa f$ allows us to decompose a submodular function f into a *difficult* polymatroid function and an *easy* modular part as $f(X) = P_\kappa f(X) + (1 - \kappa) \sum_{j \in X} f(j)$. Moreover, we may modulate the curvature of given any function g with $\kappa_g = 1$, by constructing a function $f(X) = c g(X) + (1 - c) \sum_{j \in X} g(j)$ with curvature $\kappa_f = c$ but otherwise the same polymatroidal structure as g . Our curvature-based decomposition is different from decompositions such as that into a totally normalized function and a modular function [4]. Indeed, the curve-normalized function has some specific properties that will be useful later on (proved in [12]):

Lemma 2.1. If f is monotone submodular with $\kappa_f \leq 1$, then $P_\kappa f(X) \geq \sum_{j \in X} P_\kappa f(j)$ and $P_\kappa f(X) \geq (1 - \kappa) \sum_{j \in X} f(j)$.

Lemma 2.2. If f is monotone submodular, then $P_\kappa f(X)$ in Eqn. (2) is a monotone non-negative submodular function. Furthermore, $P_\kappa f(X) \geq \sum_{j \in X} P_\kappa f(j)$.

The function $P_\kappa f$ will be our tool for analyzing the hardness of submodular problems. Previous information-theoretic lower bounds for Problems 1-3 [6, 8, 10, 27] are independent of curvature and use functions with $\kappa = 1$. These curvature-independent bounds are proven by constructing two essentially indistinguishable matroid rank functions h and f_R , one of which depends on a random set $R \subseteq V$. One then argues that any algorithm would need to make a super-polynomial number of queries to the functions for being able to distinguish h and f_R with high enough probability. The lower bound will be the ratio $\max_{X \subseteq V} h(X)/f_R(X)$. We extend this proof technique to functions with a fixed given curvature. To this end, we define the functions $f_R(X) = P_\kappa f_R(X) + (1 - \kappa) \sum_{j \in X} f_R(j)$

and

$$h(X) = P_\kappa h(X) + (1 - \kappa) \sum_{j \in X} h(j).$$

(3)

Both of these functions have curvature $\frac{1}{2}$. This construction enables us to explicitly introduce the effect of curvature into information-theoretic bounds for all monotone submodular functions. Main results. The curve normalization (2) leads to refined upper bounds for Problems 1-3, while the curvature modulation (3) provides matching lower bounds. The following are some of our main results: for approximating submodular functions (Problem 1), we replace the known bound $\frac{1}{2} n \log n$. We of $\frac{1}{2} n = O(n \log n)$ [8] by an improved curvature-dependent $O(\frac{1}{2}(n \log n)(1+\frac{1}{2}f))$ complement this with a lower bound of $\frac{1}{2}(1+\frac{1}{2}n)(1+\frac{1}{2}f)$. For learning submodular functions $\frac{1}{2}$ (Problem 2), we refine the known bound of $\frac{1}{2} n = O(n)$ [2] in the PMAC setting to a curvature $\frac{1}{2} n^{1/3}$ dependent bound of $O(\frac{1}{2}(1+\frac{1}{2}n)(1+\frac{1}{2}f))$, with a lower bound of $\frac{1}{2}$. Finally, $1+\frac{1}{2}(1+\frac{1}{2}n)(1+\frac{1}{2}f)$ Table 1 summarizes our curvature-dependent approximation bounds for constrained minimization (Problem 3). These bounds refine many of the results in [6, 27, 10, 16]. In general, our new curvature-dependent upper and lower bounds refine known theoretical results whenever $f \geq 1$, in many cases replacing known polynomial bounds by a curvature-dependent constant factor $1/(1+\frac{1}{2}f)$. Besides making these bounds precise, the decomposition and the curve-normalized version (2) are the basis for constructing tight algorithms that (up to logarithmic factors) achieve the lower bounds. 1

A polymatroid function is a monotone increasing, nonnegative, submodular function satisfying $f(\emptyset) = 0$.

3

Constraint

Modular approx. (MUB)

Card. LB

$$\frac{k+1}{1+(m^?)(1^{??})} \cdot \frac{1+(n^?)(1^{??})}{n} \cdot \frac{2+(n^?)(1^{??})}{n} \cdot \frac{1+(n^?)(1^{??})}{m}$$

Spanning Tree Matchings s-t path s-t cut

Ellipsoid approx. (EA)

$$\begin{aligned} & \log n \cdot O(1 + (\log n)^2) = O(n^2 \log n) \\ & m \log m \cdot O(1 + (m \log m)^2) = O(m^3 \log^3 m) \\ & m \log m \cdot O(1 + (m \log m)^2) = O(m^3 \log^3 m) \\ & \log m \cdot O(1 + (\log m)^2) = O(\log^3 m) \end{aligned}$$

Lower bound $1/2$) ?) ?($1+(n1/2n?1)(1?? f)$) $n ? ?() 1+(n?1)(1??f)$)

$$n^{-1/2} (1 + (n^{-1/2} f))^{2/3} n^{-1/2} (2/3 + (n^{-1/2} f))^{2/3} (1 + (n^{-1/2} f))^{2/3}$$

1)(1f) ?

$$n \cdot \frac{1}{n} = 1$$

Table 1: Summary of our results for constrained minimization (Problem 3).

3

Approximating submodular functions everywhere

We first address improved bounds for the problem of approximating a monotone submodular function everywhere. Previous work established $\frac{1}{2}$ -approximations to a submodular function f satisfying $g(S) \geq \frac{1}{2} f(S)$ for all $S \subseteq V$ [8]. We begin with a theorem showing how any algorithm computing such an approximation may be used to obtain a curvature-specific, improved approximation.

Note that the curvature of a monotone submodular function can be obtained within $2n + 1$ queries to f . The key idea of Theorem 3.1 is to only approximate the curved part of f , and to retain the modular part exactly. The full proof is in [12].

Theorem 3.1. Given a polymatroid function f with $\kappa_f \leq 1$, let f° be its curve-normalized version defined in Equation (2), and let f^{mod} be a submodular function satisfying $f^{\text{mod}}(X) \leq f^\circ(X) \leq P(n)f^{\text{mod}}(X)$, for some $X \subseteq V$. Then the function $f^{\text{mod}}(X) + \frac{1}{P(n)} \sum_{j \in X} f(j)$ satisfies $f^{\text{mod}}(X) \leq f(X) \leq$

$$\frac{1}{P(n)} f^{\text{mod}}(X) + \frac{1}{P(n)} \sum_{j \in X} f(j) \leq f(X) \leq \frac{1}{P(n)} f^{\text{mod}}(X) + \sum_{j \in X} f(j) \quad (4)$$

Theorem 3.1 may be directly applied to tighten recent results on approximating submodular functions everywhere. An algorithm by Goemans et al. [8] computes an approximation to a polymatroid function f in polynomial time by approximating the submodular polyhedron via an ellipsoid. This approximation p (which we call the ellipsoidal approximation) satisfies $\kappa_p = O(n \log n)$, and has the form $wf(X)$ for a certain weight vector w . Corollary 3.2 states that a tighter approximation is possible for functions with $\kappa_f \leq 1$.

Corollary 3.2. Let f be a polymatroid function with $\kappa_f \leq 1$, and let $wf^\circ(X)$ be the ellipsoidal approximation to the κ -curve-normalized version $f^\circ(X)$ of f . Then the function $f^{\text{mod}}(X) = p$ satisfies

$\kappa_{f^{\text{mod}}} \leq n \log n$. $f^{\text{mod}}(X) \leq f(X) \leq O(f^{\text{mod}}(X))$. (5) $1 + (n \log n - 1)(1 - \kappa_f)$ If $\kappa_f = 0$, then the approximation is exact. This is not surprising since a modular function can be inferred exactly within $O(n)$ oracle calls. The following lower bound (proved in [12]) shows that Corollary 3.2 is tight up to logarithmic factors. It refines the lower bound in [8] to include κ_f .

Theorem 3.3. Given a submodular function f with curvature κ_f , there does not exist a (possibly randomized) polynomial-time algorithm that computes an approximation to f within a factor of $n^{1/2}$, for any $\kappa_f \in (0, 1 + (n^{1/2} - 1)(1 - \kappa_f))$.

The simplest alternative approximation to f one might conceive is the modular function $f(X) = \sum_{j \in X} f(j)$ which can easily be computed by querying the n values $f(j)$.

Lemma 3.1. Given a monotone submodular function f , it holds that $f(X) \leq \sum_{j \in X} f(j) \leq f(X) + \sum_{j \in X} (f(j) - f(X))$.

In [12], we show this result with a stronger notion of curvature: $\kappa_f(X) = 1 - \frac{f(X)}{\sum_{j \in X} f(j)}$.

$$\frac{1}{P(n)} \sum_{j \in X} f(j) \leq f(X) \leq \frac{1}{P(n)} \sum_{j \in X} f(j) + \sum_{j \in X} (f(j) - f(X)) \quad (6)$$

The form of Lemma 3.1 is slightly different from Corollary 3.2. However, there is a straightforward correspondence: given f^{mod} such that $f^{\text{mod}}(X) \leq f(X) \leq P(n)f^{\text{mod}}(X)$, by defining $f^{\text{mod}}(X) = \frac{1}{P(n)} f^{\text{mod}}(X)$, we get that $f(X) \leq f^{\text{mod}}(X) \leq P(n)f^{\text{mod}}(X)$. Lemma 3.1 for the modular approximation is complementary to Corollary 3.2: First, the modular approximation is better whenever $\kappa_f(X) \leq n$. Second, the bound in Lemma 3.1 depends on the curvature $\kappa_f(X)$ with respect to the set

X , which is stronger than \sqrt{f} . Third, f^* is extremely simple to compute. For sets of larger cardinality, however, the ellipsoidal approximation of Corollary 3.2 provides a better approximation, in fact, the best possible one (Theorem 3.3). In a similar manner, Lemma 3.1 is tight for any modular approximation to a submodular function: Lemma 3.2. For any $\epsilon > 0$, there exists a monotone submodular function f with curvature ϵ such that no modular upper bound on f can approximate $f(X)$ to a factor better than $1 + (\epsilon - 1)(1/\epsilon)$. The improved curvature dependent bounds immediately imply better bounds for the class of concave over modular functions used in [22, 17, 11]. Corollary 3.4. Given weight vectors $w_1, \dots, w_k \geq 0$, and a submodular function $f(X) = \sum_{i=1}^k w_i f_i(X)$, $w_i \geq 0$, for a $\alpha \in (0, 1)$, it holds that $f(X) \leq \sum_{j \in X} f_j(j) \leq \sqrt{X}$. In particular, when $\alpha = 1/2$, the modular upper bound approximates the sum of square-root over modular functions by a factor of \sqrt{X} .

4

Learning Submodular functions

We next address the problem of learning submodular functions in a PMAC setting [2]. The PMAC (Probably Mostly Approximately Correct) framework is an extension of the PAC framework [28] to allow multiplicative errors in the function values from a fixed but unknown distribution D over 2^V . We are given training samples $\{(X_i, f(X_i))\}_{i=1}^m$ drawn i.i.d. from D . The algorithm may take time polynomial in $n, 1/\epsilon, 1/\delta$ to compute a (polynomially-representable) function \hat{f} that is a good approximation to f with respect to D . Formally, \hat{f} must satisfy that $\mathbb{E}_i \Pr_{X_1, X_2, \dots, X_m \sim D} [\hat{f}(X) \leq \epsilon f(X) + (1/\delta) \sum_{j \in X} \hat{f}(j)] \leq \delta$ for some approximation factor $\epsilon(n)$. Balcan and Harvey [2] propose an algorithm that PMAC-learns any monotone, nonnegative submodular function within a factor $\epsilon(n) = n + 1$ by reducing the problem to that of learning a binary classifier. If we assume that we have an upper bound on the curvature ϵ , or that we can estimate it ϵ , and have access to the value of the singletons $f(j), j \in V$, then we can obtain better learning results with non-maximal curvature: Lemma 4.1. Let f be a monotone submodular function for which we know an upper bound on its curvature and the singleton weights $f(j)$ for all $j \in V$. For every $\epsilon > 0$ there is an algorithm that uses a polynomial number of training examples, runs in time polynomial in $(n, 1/\epsilon, 1/\delta)$ and ϵ PMAC-learns f within a factor of $1 + (\epsilon + 1)(1/\epsilon)$. If D is a product distribution, then there exists

\log

1

an algorithm that PMAC-learns f within a factor of $O(1 + (\log 1/\delta)(1/\epsilon))$.

)

f

The algorithm of Lemma 4.1 uses the reduction of Balcan and Harvey [2] to learn the ϵ -curvenormalized version f^* of f . From the learned function $\hat{f}^*(X)$, we construct the final estimate $\hat{f}(X) = \hat{f}^*(X) + (1 - \epsilon) \sum_{j \in X} f(j)$. Theorem 3.1 implies Lemma 4.1 for this $\hat{f}(X)$. Moreover, no polynomial-time algorithm can be guaranteed to PMAC-learn f within a factor of $O(n^{1/3})$, for

any $0 \leq \epsilon \leq 1$ [12]. We end this section by showing how we can learn with a $1/(1-\epsilon)$ construction analogous to that in Lemma 3.1. Lemma 4.2. If f is a monotone submodular function with known curvature (or a known upper bound) $\kappa f(X)$, $X \subseteq V$, then for every $\epsilon \in (0, 1]$ there is an algorithm that uses a polynomial number of training examples, runs in time polynomial in $(n, 1/\epsilon)$ and PMAC learns $f(X)$ within a factor $1 + \epsilon \kappa$ of $f(X)$. Note that κf can be estimated from a set of $2n + 1$ samples $\{(j, f(j))\}_{j \in V}$, $\{(V, f(V))\}$, and $\{(V \setminus j, f(V \setminus j))\}_{j \in V}$ included in the training samples.

Compare this result to Lemma 4.1. Lemma 4.2 leads to better bounds for small sets, whereas Lemma 4.1 provides a better general bound. Moreover, in contrast to Lemma 4.1, here we only need an upper bound on the curvature and do not need to know the singleton weights $\{f(j), j \in V\}$. Note also that, while κf itself is an upper bound of $\kappa f(X)$, often one does have an upper bound on $\kappa f(X)$ if one knows the function class of f (for example, say concave over modular). In particular, an immediate corollary is that the class of concave over modular f with $\kappa f(X) \leq \alpha$, for $\alpha \in (0, 1]$ can be learnt within a factor of $\min\{n + 1, 1/\alpha\}$.

5

Constrained submodular minimization

Next, we apply our results to the minimization of submodular functions under constraints. Most algorithms for constrained minimization use one of two strategies: they apply a convex relaxation [10, 16], or they optimize a surrogate function f^* that should approximate f well [6, 8, 16]. We follow the second strategy and propose a new, widely applicable curvature-dependent choice for surrogate functions. A suitable selection of f^* will ensure theoretically optimal results. Throughout this section, we refer to the optimal solution as $X^* = \arg\min_{X \subseteq C} f(X)$. Lemma 5.1. Given a submodular function f , let f^1 be an approximation of f such that $f^1(X) \leq b_1 \arg\min_{X \subseteq C} f^1(X)$ of f^* satisfies $f(X) \leq (n) f^1(X)$, for all $X \subseteq V$. Then any minimizer X^1 of f^1 satisfies $f(X^1) \leq (n) f(X^*)$. Likewise, if an approximation of f is such that $f(X) \leq f^2(X) \leq \gamma(X) f(X)$ for all $X \subseteq V$, then its minimizer X^2 satisfies $f(X^2) \leq \gamma(X^*) f(X^*)$ for a set-specific factor $\gamma(X)$, then its minimizer X^2 satisfies $f(X^2) \leq \gamma(X^*) f(X^*)$. If only γ -approximations are possible for minimizing f^1 or f^2 over C , then the final bounds are (n) and $\gamma(X^*)$ respectively. For Lemma 5.1 to be practically useful, it is essential that f^1 and f^2 be efficiently optimizable over C . We discuss two general curvature-dependent approximations that work for a large class of combinatorial constraints. In particular, we use Theorem 3.1: we decompose f into f^* and a modular part f_m , and then approximate f^* while retaining f_m , i.e., $f^* = f^{**} + f_m$. The first approach uses a simple modular upper bound (MUB) and the second relies on the Ellipsoidal approximation (EA) we used in Section 3. MUB: The simplest approximation to a submodular function is the modular approximation $P f_m(X), j \in X f(j) \leq f(X)$. Since here, f^{**} happens to be equivalent to f_m , we obtain the overall approximation $f^* = f_m$. Lemmas 5.1 and 3.1 directly imply a set-dependent approximation factor for f_m : $b \leq C$ be a γ -approximate solution for minimizing P Corollary 5.1. Let $X^j \subseteq C$ be a γ -approximate solution for minimizing P Corollary 5.1. Let $X^j \subseteq C$ be a γ -approximate solution for minimizing P Corollary 5.1. Let $X^j \subseteq C$ be a γ -approximate solution for minimizing P Corollary 5.1.

over C , i.e. $P \leq f(j) \leq \min f(j)$. Then $X \in C$ b $j \in X$ $j \in X$ $f(X)$

$$\begin{aligned} & 1 + \\ & (-X \leq - \\ & ? - X \leq - f(X \leq ?) \cdot (1 \leq ?) f(X \leq ?)) \\ & (8) \end{aligned}$$

Corollary 5.1 has also been shown in [13]. Similar to the algorithms in [13], MUB can be extended to an iterative algorithm yielding performance gains in practice. In particular, Corollary 5.1 implies improved approximation bounds for practically relevant concave over modular functions, such as those used in [17]. For instance, for $f(X) = \sum_{i=1}^n w_i \cdot x_i$, we obtain a worst-case p -approximation bound of $-X \leq - \leq n$. This is significantly better than the worst case factor of $-X \leq -$ for general submodular functions. EA: Instead of employing a modular upper bound, we can approximate p using the construction f by Goemans et al. [8], as in Corollary 3.2. In that case, $f(X) = \sum w_i x_i + (1 - \sum w_i) f_m(X)$ has a special form: a weighted sum of a concave function and a modular function. Minimizing such a function over constraints C is harder than minimizing a merely modular function, but with the algorithm in [24] we obtain an FPTAS for minimizing f over C whenever we can minimize a nonnegative linear function over C . 4 5

A p -approximation algorithm for minimizing a function g finds set $X : g(X) \leq \min_{X \in C} g(X)$. The FPTAS will yield a $(1 + \epsilon)$ -approximation through an algorithm polynomial in $1/\epsilon$.

6

Corollary 5.2. For a submodular function with curvature $\rho \leq 1$, algorithm EA will return a b that satisfies solution X

$\leq n \log n \cdot b \cdot f(X) \leq O(\rho \cdot f(X \leq ?)) \cdot (n \log n \leq 1)(1 \leq ?) + 1$ Next, we apply the results of this section to specific optimization problems, for which we show (mostly tight) curvature-dependent upper and lower bounds. We just state our main results; a more extensive discussion along with the proofs can be found in [12]. Cardinality lower bounds (SLB). A simple constraint is a lower bound on the cardinality of the solution, i.e., $C = \{X \leq V : -X \leq k\}$. Svitkina and Fleischer [27] prove that for monotone submodular functions of arbitrary curvature, it is impossible to find a polynomial-time algorithm with an approximation factor better than $n/\log n$. They show an algorithm which matches this approximation factor. Corollaries 5.1 and 5.2 immediately imply curvature-dependent approximation $\log n$ bounds of $1 + (k \leq 1)(1 \leq ?)$ and $O(1 + (n \log n \leq 1)(1 \leq ?))$. These bounds are improvements over the results of [27] whenever $\rho \leq 1$. Here, MUB is preferable to EA whenever k is small. Moreover, the bound of EA is tight up to poly-log factors, in that no polynomial time algorithm can achieve a general approximation factor better than $1 + (n \log n \leq 2)$ for any $\epsilon > 0$. In the following problems, our ground set V consists of the set of edges in a graph $G = (V, E)$ with two distinct nodes $s, t \in V$ and $n = |V|$, $m = |E|$. The submodular function is $f : 2^E \rightarrow \mathbb{R}$. Shortest submodular s - t path (SSP). Here, we aim to find an s - t path X of minimum (submodular) length $f(X)$. Goel et al. [6] show a $O(n^{2/3})$ -approximation with matching curvature-independent lower bound $\Omega(n^{2/3})$. By Corollary 5.1, the

curvature-dependent worst-case bound for MUB is $n^{1+(n^{-1})(1-\alpha)^{\alpha}}$ since any minimal s-t path has at most n edges. Similarly, the factor for EA is $n^{1+(n^{-1})(1-\alpha)^{\alpha}}$.

$m \log m$ $O(1 + (m \log m)^{\alpha})$. The bound of EA will be tighter for sparse graphs while MUB provides $m^{1+(n^{-1})(1-\alpha)^{\alpha}}$.

better results for dense ones. Our curvature-dependent lower bound for SSP is for any $\alpha > 0$, which reduces to the result in [6] for $\alpha = 1$.

$$n^{2/3\alpha}, 1 + (n^{2/3\alpha - 1})(1-\alpha)^{\alpha}$$

Minimum submodular s-t cut (SSC): This problem, also known as the cooperative cut problem [16, 17], asks to minimize a monotone submodular function f such that the solution $X \subseteq E$ is a set of edges whose removal disconnects s from t in G . Using curvature refines the We can also show a $n^{1/2\alpha}$ lower bound of [16] to $1 + (n^{1/2\alpha - 1})(1-\alpha)^{\alpha}$, for any $\alpha > 0$. Corollary 5.1 implies an approximation $(1-\alpha)^{\alpha}$.

$m \log m$ m factor of $O((m \log m)^{\alpha})$ for EA and a factor of $1 + (m^{1/2\alpha - 1})(1-\alpha)^{\alpha}$ for MUB, where $m = |E - f| + 1$ is the number of edges in the graph. Hence the factor for EA is tight for sparse graphs. Specifically for cut problems, there is yet another useful surrogate function that is exact on local neighborhoods. Jegelka and Bilmes [16] demonstrate how this approximation may be optimized via a generalized maximum flow algorithm that maximizes a polymatroidal network flow [20]. This algorithm still applies to the combination $f = \alpha f + (1 - \alpha)f$, where we only approximate f . We refer to this approximation as Polymatroidal Network Approximation (PNA). Corollary 5.3. Algorithm PNA achieves a worst-case approximation factor of $2 + (n^{1/2\alpha - 1})(1-\alpha)^{\alpha}$ for the f cooperative cut problem.

For dense graphs, this factor is theoretically tighter than that of the EA approximation. Minimum submodular spanning tree (SST). Here, C is the family of all spanning trees in a given graph G . Such constraints occur for example in power assignment problems [30]. Goel et al. [6] show a curvature-independent optimal approximation factor of $O(n)$ for this problem. Corollary 5.1 refines this bound to $1 + (n^{1/2\alpha - 1})(1-\alpha)^{\alpha}$ when using MUB; Corollary 5.2 implies a slightly worse bound f for EA. We also show that the bound of MUB is tight: no polynomial-time algorithm can guarantee a $n^{1/2\alpha}$, for any $\alpha > 0$. factor better than $1 + (n^{1/2\alpha - 1})(1-\alpha)^{\alpha}$. Minimum submodular perfect matching (SPM): Here, we aim to find a perfect matching in a graph that minimizes a monotone submodular function. Corollary 5.1 implies that an MUB n approximation will achieve an approximation factor of at most $2 + (n^{1/2\alpha - 1})(1-\alpha)^{\alpha}$. Similar to the f spanning tree case, the bound of MUB is also tight [12]. 7

4 2 50
100 150 200 250 n
 $\alpha = 0.7$ $\alpha = 0.5$ $\alpha = 0.3$ $\alpha = 0.1$
4 3 2 50
100 150 200 250 n
(c) with varying α and $\alpha = 1$ 100 80 60
 $\alpha = n/2$ $\alpha = n^{3/4}$ $1/2$
 $\alpha = n$
40 20 50

100 n
150
200
(d) varying γ with $\beta = n/2$ and $\alpha = 1$ emp. approx. factor
6
(b) varying γ , $\beta = 0.1$ 5
emp. approx. factor
 $\beta = 0.1$ $\beta = 0.2$ $\beta = 0.3$ $\beta = 0.4$
emp. approx. factor
emp. approx. factor
(a) varying γ , $\beta = 0$ 8
 $\beta = 0.9$ $\beta = 0.6$ $\beta = 0.3$ $\beta = 0.1$
8 6 4 2 50
100 n
150
200

Figure 1: Minimization of g_γ for cardinality lower bound constraints. (a) fixed $\beta = 0$, $\beta = n/2 + 2$, $\beta = n/2$ for varying γ ; (b) fixed $\gamma = 0.1$, but varying β ; (c) different choices of β for $\gamma = 1$; (d) varying γ with $\beta = n/2$, $\alpha = 1$. Dashed lines: MUB, dotted lines: EA, solid lines: theoretical bound. The results of EA are not visible in some instances since it obtains a factor of 1. 5.1

Experiments

We end this section by empirically demonstrating the performance of MUB and EA and their precise dependence on curvature. We focus on cardinality lower bound constraints, $C = \{X \subseteq V : |X| \geq \beta\}$ and the ‘worst-case’ class of functions that has been used throughout this paper to prove lower $\beta = V$ R and $R \subseteq V$ is random set such that $\beta + \gamma$, $|X| \geq \beta$ where R bounds, $f_R(X) = \min\{|X| - \beta, 1/2 + 2 - |R| + \gamma\}$. We adjust $\beta = n$ and $\gamma = n$ by a parameter γ . The smaller γ is, the harder the problem. This function has curvature $\gamma f = 1$. To obtain a function with specific curvature γ , we define $f_\gamma R(X) = \gamma f(X) + (1 - \gamma)|X| - \beta$ as in Equation (3). In all our experiments, we take the average over 20 random draws of R. We first set $\gamma = 1$ and vary β . Figure 1(a) shows the empirical approximation factors obtained using EA and MUB, and the theoretical bound. The empirical factors follow the theoretical results very closely. Empirically, we also see that the problem becomes harder as β decreases. Next we fix $\gamma = 0.1$ and vary the curvature γ in $f_\gamma R$. Figure 1(b) illustrates that the theoretical and empirical approximation factors improve significantly as γ decreases. Hence, much better approximations than the previous theoretical lower bounds are possible if γ is not too large. This observation can be very important in practice. Here, too, the empirical upper bounds follow the theoretical bounds very closely. Figures 1(c) and (d) show results for larger β and $\gamma = 1$. In Figure 1(c), as β increases, the empirical factors improve. In particular, as predicted by the theoretical bounds, EA outperforms MUB for large β and, for $\beta \geq n/3$, EA finds the optimal solution. In addition, Figures 1(b) and (d) illustrate the theoretical and empirical effect of curvature: as n grows, the bounds saturate and approximate a constant $1/(1 - \gamma)$ — they do not grow

polynomially in n . Overall, we see that the empirical results quite closely follow our theoretical results, and that, as the theory suggests, curvature significantly affects the approximation factors.

6

Conclusion and Discussion

In this paper, we study the effect of curvature on the problems of approximating, learning and minimizing submodular functions under constraints. We prove tightened, curvature-dependent upper bounds with almost matching lower bounds. These results complement known results for submodular maximization [3, 29]. Given that the functional form and effect of the submodularity ratio proposed in [5] is similar to that of curvature, an interesting extension is the question of whether there is a single unifying quantity for both of these terms. Another open question is whether a quantity similar to curvature can be defined for subadditive functions, thus refining the results in [1] for learning subadditive functions. Finally it also seems that the techniques in this paper could be used to provide improved curvature-dependent regret bounds for constrained online submodular minimization [15]. Acknowledgments: Special thanks to Kai Wei for pointing out that Corollary 3.4 holds and for other discussions, to Bethany Herwaltdt for reviewing an early draft of this manuscript, and to the anonymous reviewers. This material is based upon work supported by the National Science Foundation under Grant No. (IIS-1162606), a Google and a Microsoft award, and by the Intel Science and Technology Center for Pervasive Computing. Stefanie Jegelka's work is supported by the Office of Naval Research under contract/grant number N00014-11-1-0688, and gifts from Amazon Web Services, Google, SAP, Blue Goji, Cisco, Clearstory Data, Cloudera, Ericsson, Facebook, General Electric, Hortonworks, Intel, Microsoft, NetApp, Oracle, Samsung, Splunk, VMware and Yahoo!. 8

2 References

- [1] M. F. Balcan, F. Constantin, S. Iwata, and L. Wang. Learning valuation functions. COLT, 2011.
- [2] N. Balcan and N. Harvey. Submodular functions: Learnability, structure, and optimization. In Arxiv preprint, 2012.
- [3] M. Conforti and G. Cornuejols. Submodular set functions, matroids and the greedy algorithm: tight worstcase bounds and some generalizations of the Rado-Edmonds theorem. Discrete Applied Mathematics, 7(3): 251?274, 1984.
- [4] W. H. Cunningham. Decomposition of submodular functions. Combinatorica, 3(1):53?68, 1983.
- [5] A. Das and D. Kempe. Submodular meets spectral: Greedy algorithms for subset selection, sparse approximation and dictionary selection. In ICML, 2011.
- [6] G. Goel, C. Karande, P. Tripathi, and L. Wang. Approximability of combinatorial problems with multi-agent submodular cost functions. In FOCS, 2009.
- [7] G. Goel, P. Tripathi, and L. Wang. Combinatorial problems with discounted price functions in multi-agent systems. In FSTTCS, 2010.
- [8] M. Goemans, N. Harvey, S. Iwata, and V. Mirrokni. Approximating submodular functions everywhere. In SODA, pages 535?544, 2009.
- [9] R. Hassin, J. Monnot,

and D. Segev. Approximation algorithms and hardness results for labeled connectivity problems. *J Combinatorial Optimization*, 14(4):437–453, 2007. [10] S. Iwata and K. Nagano. Submodular function minimization under covering constraints. In *FOCS*, pages 671–680. IEEE, 2009. [11] R. Iyer and J. Bilmes. Algorithms for approximate minimization of the difference between submodular functions, with applications. In *UAI*, 2012. [12] R. Iyer, S. Jegelka, and J. Bilmes. Curvature and Optimal Algorithms for Learning and Optimization of Submodular Functions: Extended arxiv version, 2013. [13] R. Iyer, S. Jegelka, and J. Bilmes. Fast semidifferential based submodular function optimization. In *ICML*, 2013. [14] S. Jegelka. Combinatorial Problems with submodular coupling in machine learning and computer vision. PhD thesis, ETH Zurich, 2012. [15] S. Jegelka and J. Bilmes. Online submodular minimization for combinatorial structures. *ICML*, 2011. [16] S. Jegelka and J. A. Bilmes. Approximation bounds for inference using cooperative cuts. In *ICML*, 2011. [17] S. Jegelka and J. A. Bilmes. Submodularity beyond submodular energies: coupling edges in graph cuts. In *CVPR*, 2011. [18] P. Kohli, A. Osokin, and S. Jegelka. A principled deep random field for image segmentation. In *CVPR*, 2013. [19] A. Krause and C. Guestrin. Near-optimal nonmyopic value of information in graphical models. In *Proceedings of Uncertainty in Artificial Intelligence*. *UAI*, 2005. [20] E. Lawler and C. Martel. Computing maximal λ -polymatroidal network flows. *Mathematics of Operations Research*, 7(3):334–347, 1982. [21] H. Lin and J. Bilmes. Optimal selection of limited vocabulary speech corpora. In *Interspeech*, 2011. [22] H. Lin and J. Bilmes. A class of submodular functions for document summarization. In *The 49th Meeting of the Assoc. for Comp. Ling. Human Lang. Technologies (ACL/HLT-2011)*, Portland, OR, June 2011. [23] H. Lin and J. Bilmes. Learning mixtures of submodular shells with application to document summarization. In *UAI*, 2012. [24] E. Nikolova. Approximation algorithms for offline risk-averse combinatorial optimization, 2010. [25] J. Soto and M. Goemans. Symmetric submodular function minimization under hereditary family constraints. *arXiv:1007.2140*, 2010. [26] P. Stobbe and A. Krause. Learning fourier sparse set functions. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2012. [27] Z. Svitkina and L. Fleischer. Submodular approximation: Sampling-based algorithms and lower bounds. In *FOCS*, pages 697–706, 2008. [28] L. G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984. [29] J. Vondrák. Submodularity and curvature: the optimal algorithm. *RIMS Kokyuroku Bessatsu*, 23, 2010. [30] P.-J. Wan, G. Calinescu, X.-Y. Li, and O. Frieder. Minimum-energy broadcasting in static ad hoc wireless networks. *Wireless Networks*, 8:607–617, 2002. [31] P. Zhang, J.-Y. Cai, L.-Q. Tang, and W.-B. Zhao. Approximation and hardness results for label cut and related problems. *Journal of Combinatorial Optimization*, 21(2):192–208, 2011.