

Based on the notes from Ankur Moitra's class at MIT (see references for a link).

This lecture introduces submodular functions as a generalization of some functions we have previously seen for e.g. the cut function in graphs. We will see how we can use the ellipsoid method developed in the previous lecture to minimize an arbitrary submodular function.

## 1 Submodular Functions

**Definition 1** (Submodular Functions). *Let  $N$  be a set of  $n$  elements. A function  $f : 2^N \rightarrow \mathbb{R}$  is said to be submodular, if it satisfies following property of diminishing marginal returns: for every  $A \subseteq B \subseteq N$  and  $j \notin B$ ,  $f(A \cup \{j\}) - f(A) \geq f(B \cup \{j\}) - f(B)$ .*

One way to understand submodularity is to think of  $f$  as a *utility* functions. Then, the diminishing marginal returns property says that the marginal utility gained by adding a new item to a smaller set is no less than the marginal utility gained by adding a new item to a larger set.

The defining property of diminishing marginal returns is equivalent to the following (see the notes by Jeff Bilmes, in references, for a proof):

**Lemma 1.** *A function  $f : 2^N \rightarrow \mathbb{R}$  is submodular if and only if for every  $X, Y \subseteq N$ ,  $f(X) + f(Y) \geq f(X \cap Y) + f(X \cup Y)$ .*

## 2 Examples of Submodular Functions

Submodular functions generalize familiar quantities studied before in this class. For a graph  $G(V, E)$  on  $n$  vertices, let  $f(S) = E(S, \bar{S})$ , the number of edges that lie in the cut defined by the vertex set  $S$ . It is a short exercise to prove that  $f$  is a submodular function.

Another important example is a *coverage function*: let  $G(U, V, E)$  be a bipartite graph with the bipartition  $U$  and  $V$  of vertices with  $|U| = n$ . For every  $S \subseteq U$ , define  $f(S)$  to be the size of the neighborhood of vertices in  $S$ . Then  $f$  is a submodular function.  $f$  is also *monotone* — for any  $S$  and  $j \in U$ ,  $f(S \cup \{j\}) \geq f(S)$ . Observe that the cut function is not necessarily monotone.

Our final example is from information theory — let  $x_1, x_2, \dots, x_n$  be discrete random variables. For any  $A \subseteq [n]$ , let  $H(A)$  be the joint entropy of the  $\{x_i\}_{i \in A}$ , i.e.,  $H(A) = -\sum_a P(\{x_i\}_{i \in A} = a) \log P(\{x_i\}_{i \in A} = a)$ . Then,  $H$  is a submodular function.

### 3 Representation and Optimization

For the special cases above, we are familiar with the algorithmic tasks of optimization. For example, we know algorithms for finding the minimum size cut in a graph that run in polynomial time, along with a poly-time 0.878 approximation algorithm for the problem of finding a maximum size cut in a graph. Incidentally, while the above approximation factor requires the use of semidefinite programming, there's a simple algorithm that returns a random cut which yields a 0.5 approximation ratio. Incidentally, this trivial algorithm cannot be improved by even a sub-exponential size linear programs (see references) even though a simple semidefinite program does phenomenally better!

What about optimizing arbitrary submodular functions? This immediately leads to the question of how the submodular function we are interested in is represented. This is a tricky issue and in general, there exist submodular functions which do not have (even in an approximate sense) a small representation (see the paper by Balcan and Harvey in the references for details of one such construction). In most applications, we do have a concise representation. However, we would like our algorithms to not be tied to a particular representation, and we'd like to make minimal assumptions about it. We will see that we can give non-trivial algorithms by just having a oracle access to the submodular function — we will only need that given any  $x$ , some oracle returns to us the value of the underlying function at  $x$ .

As the case of min-cut might suggest, there's an efficient algorithm for minimizing an arbitrary submodular function given just oracle access to it! The maximization problem is of course NP-hard (Max-Cut is a special case) but it turns out that for a non-negative submodular function the oracle access is enough to yield a 0.5 approximation guarantee (this is by no means trivial though — see the paper by Buchbinder, Naor and Schwartz in references.) The naive algorithm that selects each element independently with probability half is known to give 0.25 approximation.

Our interest here is going to be in minimizing an arbitrary submodular function.

### 4 Minimizing Convex Functions over Convex Sets

Our main interest today is to give a polynomial time algorithm for minimizing an arbitrary submodular function. Even though submodular function is a discrete object, we will be able to minimize it by reduction to minimizing a *continuous convex function* over a convex set. Before describing this reduction, we recall how to use the ellipsoid method in order to minimize an arbitrary convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  over a convex set  $K$ .

We will need the following assumptions:

- 1) There's an efficient *evaluation oracle* and *gradient oracle* for  $f$  i.e., given a point  $x$ , we can compute  $f(x)$  and  $\nabla f(x)$ .
- 2) There's a known ellipsoid containing  $K$  and there's an efficient separation oracle for  $K$ .

First note that by doing a binary search over a parameter  $c$ , we can reduce minimizing  $f$  to the problem of finding a point (or proving emptiness) in  $S_c = K \cup \{x \mid f(x) \leq c\}$  given a (rough) interval where the minimum must lie. Notice that  $S_c$  is convex whenever  $f$  is a convex function.

We now want to use the ellipsoid method to find a point in  $S_c$  or to prove it is empty. Recall that we can use the ellipsoid method from previous classes to do this provided we can satisfy two conditions:

a) There's an ellipsoid that contains  $S_c$  (one usually knows a bound on the diameter of the set  $K$  in which case we can just set this ellipsoid to be the sphere of the same diameter.)

This is satisfied immediately from the assumption on  $K$  above.

b) There's an efficient separation oracle for  $S_c$ : We want to show that given a separation oracle for  $K$ , we can build an efficient separation oracle for  $S_c$ . Let  $x \in \mathbb{R}^n$ . Then, we first test if  $x \in K$ , and if  $x \notin K$  then obtain a separating hyperplane using the separation oracle for  $K$ . In case  $x \in K$ , we test if  $f(x) \leq c$  by using the evaluation oracle for  $f$ . If  $f(x) > c$ , then, observe that by convexity of  $f$ , for any minimizer of  $f$ , say  $x^*$ ,  $f(x) - f(x^*) \leq \nabla f(x)(x - x^*)$  which is a hyperplane (in  $x^*$  after we put  $c$  instead of  $f(x^*)$ ) that separates  $x$  and the convex set of all minimizers of  $f$ .

Thus, we have that we can minimize  $f$  over  $K$  efficiently.

## 5 Lovász Extension

We want to use the above idea for minimizing a convex function given oracle + gradient access over a convex set to minimize an arbitrary submodular function. To do this, we will define an *extension*  $\hat{f}$  of a given submodular function  $f$ ; here,  $\hat{f}$  will be a continuous function that *extends*  $f$  to all of  $[0, 1]^n$ . The definition itself makes sense for an arbitrary set function  $f : 2^N \rightarrow \mathbb{R}$ . However,  $\hat{f}$  will be convex if and only if  $f$  is submodular.

First, we think of  $f : 2^N \rightarrow \mathbb{R}$  as  $f : \{0, 1\}^n \rightarrow \mathbb{R}$  by associating a set  $S$  with its 0-1 indicator from  $\{0, 1\}^n$  (recall that  $n = |N|$ ).

**Definition 2.** Let  $f : \{0, 1\}^n \rightarrow \mathbb{R}$ . Define  $\hat{f}$ , the Lovász extension of  $f$ , so that for any  $z \in [0, 1]^n$ ,  $\hat{f}(z) = \mathbb{E}_{\lambda \sim [0, 1]}[f(\{i \mid z_i \geq \lambda\})]$  where  $\lambda \sim [0, 1]$  denotes a uniformly random sample from the interval  $[0, 1]$ .

Why is  $\hat{f}$  an extension of  $f$ ? Let  $z \in \{0, 1\}^n$ . Then notice that for any  $\lambda \in [0, 1]$ ,  $\{i \mid z_i \geq \lambda\} = S$  where  $S = \{i \mid z_i = 1\}$ . Thus,  $\hat{f}$  agrees with  $f$  over the hypercube and is some kind of an average of  $f$  at fractional points.

We can in fact explicitly find out this averaging representation of  $\hat{f}(z)$ . To do this, we define a “chain” of sets associated with any  $z \in \{0, 1\}^n$ . For simplicity of notation, assume that  $z_0 = 1 \geq z_1 \geq z_2 \dots \geq z_n \geq 0 = z_{n+1}$ . Let  $S_i$  for any  $i \in [n] \cup \{0\}$  equal  $\{1, 2, \dots, n\}$ . Then,  $S_0 = \emptyset \subseteq S_1 \subseteq S_2 \dots \subseteq S_n = [n]$ .

Further notice that if  $\lambda \in [z_i, z_{i+1})$  (the probability of which is equal to  $z_i - z_{i+1}$ ) for any  $i \in [n] \cup \{0\}$ , then  $\{j \mid z_j \geq \lambda\} = S_i$ . Thus,

$$\hat{f}(z) = \sum_{i=0}^n (z_i - z_{i+1}) f(S_i). \quad (1)$$

In particular, to evaluate  $\hat{f}$  at any  $z$ , we need only  $n + 1$  calls to the evaluation oracle for  $f$ .

## 6 Lovász Extension is Convex iff $f$ is submodular

The key result of today is the following:

**Theorem 2.** *Let  $\hat{f}$  be the Lovász Extension of  $f : \{0, 1\}^n \rightarrow \mathbb{R}$ . Then,  $\hat{f}$  is convex iff  $f$  is submodular.*

*Proof.* We will only show the “if” part — that is, if  $f$  is submodular then  $\hat{f}$  is convex. This is what is needed for our minimization algorithm. The proof is a neat application of LP Duality and is a reminder that duality is not only useful as in algorithm design but also a powerful method of arguing structural properties especially in situations where one can define the quantity of interest as the optimum value of a linear program.

We will assume that  $f(\emptyset) = 0$  and that  $1 \geq z_1 \geq z_2 \geq \dots \geq z_n \geq 0$ . This is WLOG as otherwise we can define a new function  $f' = f - f(\emptyset)$  which satisfies the first property and is submodular and rename the coordinates.

In Definition 2, we have a definition of the Lovász extension via an explicit formula. We now give a second definition of the same function as an optimum of a linear program. Specifically, for  $f$ , the submodular function above, define the associated *base polyhedron*  $\mathcal{P}_f$  (a polyhedron is just a possibly unbounded convex set defined as the intersection of a finite number of hyperplanes, as against a polytope which is a bounded polyhedron) by the following (exponentially many) inequalities:

$$x \in \mathcal{P}_f \text{ iff for every } S \subseteq [n] \sum_{i \in S} x_i \leq f(S)$$

$$\sum_{i \in [n]} x_i = f([n]).$$

We then define  $g(z) = \max z^\top x$  for  $x \in \mathcal{P}_f$ . This is a linear program — we will call it the primal (P) program. Note that we *don't have* constraints  $x_i \geq 0$  (since the submodular function could be negative), and hence this is a polyhedron.

Our first claim is that  $g(z) = \hat{f}(z)$  for every  $z \in [0, 1]^n$ . In other words, the LP optimum above is an alternate definition of  $\hat{f}(z)$ . To see why this helps, we observe that  $g$  is convex. The idea is simple and general, the maximum of linear functions is always convex. Concretely, for  $1 \geq \lambda_1 \geq 0$  and  $z_1, z_2 \in [0, 1]^n$ , observe that  $g(\lambda_1 z_1 + (1 - \lambda_1) z_2) = \max_{x \in \mathcal{P}_f} (\lambda_1 z_1^\top x + (1 - \lambda_1) z_2^\top x) \leq \lambda_1 \max_{x \in \mathcal{P}_f} z_1^\top x + (1 - \lambda_1) \max_{x \in \mathcal{P}_f} z_2^\top x = \lambda_1 g(z_1) + (1 - \lambda_1) g(z_2)$ .

Thus, proving  $\hat{f} = g$  completes the proof. To show this, we will use weak duality.

Verify that the dual linear program (D) to (P) is given by:

$$\min \sum_{S \subseteq [n]} y_S f(S)$$

$$\text{s.t. } \sum_{S \ni i} y_S = z_i \text{ for every } i$$

$$y_S \geq 0 \text{ for every } S \subsetneq [n].$$

How does this help us? We want to prove that the optimum value of (P) is  $\hat{f}(z)$ . We will do this by “guessing” the optimal solutions to (P) and (D) (you’ll see why it’s intuitive

to guess these solutions, especially when you know what we want to prove). Given such a guess, how do we verify that they are indeed optimal?

By weak duality, for any  $x \in \mathcal{P}_f$  and any  $y$  feasible for  $D$ ,  $z^\top x \leq \sum_{S \subseteq [n]} y_S f(S)$ . If for our guesses  $x^*$  for (P) and  $y^*$  for (D) it so happens that  $z^\top x^* = \sum_{S \subseteq [n]} y_S^* f(S)$ , then  $x^*$  and  $y^*$  have to be optimal for the respective linear programs as any better solution will violate weak duality!

We define  $x^*$  by:  $x_i^* = f(S_i) - f(S_{i-1})$  for every  $i \in [n]$  and  $y^*$  by:

$$y_S^* = \begin{cases} z_i - z_{i+1} & \text{for } S = S_i \text{ for } i < n \\ z_n & \text{if } S = N \\ 0 & \text{otherwise.} \end{cases}$$

We now make four claims about  $x^*$  and  $y^*$  which will complete the proof. First, we establish that  $z^\top x^* = \hat{f}(z)$ .

CLAIM 1 (Primal Objective):  $z^\top x^* = \hat{f}(z)$ .

Next, we claim that  $x^*$  and  $y^*$  are both optimal.

CLAIM 2 (Dual Objective):  $z^\top x^* = \sum_{S \subseteq [n]} y_S^* f(S)$ .

And finally,

CLAIM 3 (Dual Feasibility):  $y^*$  is feasible for (D).

CLAIM 4 (Primal Feasibility):  $x^*$  is feasible for (P).

Using our discussion from above, the above claims complete the proof. Let us now verify them one by one.

First,

$$z^\top x^* = \sum_i z_i \cdot (f(S_i) - f(S_{i-1})) \tag{2}$$

$$= \sum_{i=1}^{n-1} (z_i - z_{i+1}) f(S_i) + z_n f(S_n) \tag{3}$$

$$= \sum_{S \subseteq [n]} y_S^* f(S). \tag{4}$$

(3) along with (1) establishes Claim 1. (4) establishes Claim 2.

Let us now proceed to Claims 3 and 4. Observe that  $y^*$  has all non-negative entries because  $z_i \geq z_{i+1}$  for every  $i$ . Further, for any  $i$ ,  $\sum_{S: S \ni i} y_S^* = \sum_{j \geq i} y_{S_j}^* = (z_i - z_{i+1}) + (z_{i+1} - z_{i+2}) + \dots + (z_{n-1} - z_n) + z_n = z_i$  as required. This establishes that  $y^*$  is feasible for (D).

Let us now finally come to Claim 4. Observe that  $\sum_{i \in [n]} x_i^* = \sum_{i \leq n} (f(S_i) - f(S_{i-1})) = f([n]) - f(\emptyset) = f([n])$  using that  $f(\emptyset) = 0$ . Let  $S \subsetneq [n]$ . Then, we must show that  $\sum_{i \in S} x_i^* \leq f(S)$ . We will argue this by induction on the size of the set  $S$ . The base case trivially holds as  $f(\emptyset) = 0 \geq \sum_{i \in \emptyset} x_i^*$ . Let us now prove the inductive case (observe that so far we haven't used the submodularity of  $f$  — this is the part of the argument it comes in).

Let  $i$  be the largest index in  $S$ . By definition of submodularity,  $f(S) + f(S_{i-1}) \geq f(S \cup S_{i-1}) + f(S \cap S_{i-1}) = f(S_i) + f(S \setminus \{i\})$  which can be rearranged to

$$f(S) \geq f(S_i) - f(S_{i-1}) + f(S \setminus \{i\}) = x_i^* + f(S \setminus \{i\}) \geq \sum_{j \in S} x_j^*,$$

where we used the inductive hypothesis in the final inequality. This completes the proof of Claim 4. □

## 7 Wrapping Up

Now that we know that the Lovász extension of  $f$  is convex iff  $f$  is submodular, we can minimize the Lovász extension of  $f$  over  $[0, 1]^n$  whenever  $f$  is submodular. At first glance, it's not clear what this buys us: what happens if the minimizer isn't in  $\{0, 1\}^n$ ? How do we then turn this into a minimizer for  $f$ ?

The other magic property of the Lovász extension (which is especially surprising when it is convex), is that it always achieves its minimum at an extreme point (and moreover, a minimum on the extreme point can be found given any interior minimum).

To see this, consider any minimizer  $\vec{z}$ . For simplicity of notation, again assume that  $1 \geq z_1 \geq \dots \geq z_n \geq 0$ . Then  $\hat{f}(z) = \sum_{i=0}^n (z_i - z_{i+1})f(S_i)$  as stated previously. Therefore,  $\hat{f}(z)$  is a convex combination of  $f$  evaluated at  $n + 1$  different sets, and  $\hat{f}(z) \geq \min_i f(S_i)$ . In fact, as  $\vec{z}$  was the minimum, we must have  $\hat{f}(z) = f(S_i)$  for all  $i$  such that  $z_i - z_{i+1} > 0$ . Therefore, once we know the minimizer  $\vec{z}$  of  $\hat{f}(\cdot)$ , we can immediately output the minimizer  $S$  of  $f(\cdot)$ .

The “basic part” is the work immediately above: the minimizer of the Lovász extension is always an extreme point, and an extreme point minimizer can be deduced from a minimizer of the extension. Because this point is in fact a minimum even of the extension to  $[0, 1]^n$ , it is certainly the minimum over  $\{0, 1\}^n$ . Normally, functions which have minima on the extremes are not convex.

The “magic part” is that the Lovász extension happens to be convex when  $f$  is submodular. When  $f$  is submodular, we can actually minimize the extension in poly-time, enabling us to minimize  $f$ .

### BIBLIOGRAPHY

1. *Lecture Slides: “Submodular functions, their optimization and applications.”* Jeff Bilmes. Dept. of Elect. Eng., Univ. Washington, Seattle, Apr. 1, 2011 [Online]. Available: [http://melodi.ee.washington.edu/~bilmes/ee595a\\_spring\\_2011/lecture2.pdf](http://melodi.ee.washington.edu/~bilmes/ee595a_spring_2011/lecture2.pdf)
2. *A tight linear time (1/2)-approximation for unconstrained submodular maximization.* Niv Buchbinder, Moran Feldman, Joseph Naor, and Roy Schwartz. In FOCS, 2012.
3. *Approximating Rectangles by Juntas and a weakly exponential lower bound on linear programs for constraint satisfaction* Pravesh K. Kothari, Raghu Meka and Prasad Raghavendra. Available: <https://arxiv.org/abs/1610.02704>
4. *Learning Submodular Functions* Maria Florina Balcan and Nick Harvey. In STOC 2011.
5. *Lecture Notes: Submodular Functions and Lovász Extension* Ankur Moitra. Available at: <http://people.csail.mit.edu/moitra/docs/6854notes13.pdf>