

# AM221 Final Project

## Dictionary Selection Under a Supermodular Assumption

Taylor Killian & Leonhard Spiegelberg

May 1, 2016

### Abstract

We attempt to develop a method by which to solve dictionary selection under a supermodular assumption. This method is motivated primarily by the work of Singer et al. (2016) [10] where a two-stage submodular maximization generalization was developed. The development of a two-stage supermodular minimization algorithm to address Sparse Dictionary Selection utilizes approximation results derived by Boutsidis et al. (2015) [1] for supermodular optimization. We show that Sparse Dictionary Selection is functionally equivalent to Sparse Multi Linear Regression which allows us to leverage [1]. We analyze and experiment with an implementation of our two-stage optimization routine and compare it to current methods. With this project we have laid the ground for further exploration into the formalization of the two-stage method that we propose here.

## 1 Introduction

Dictionary selection and sparse regression (problems that we show to be functionally equivalent) are variants of representation learning. The goal of these kinds of problems is to determine a sparse representation of input data in the form of a linear combination of basis elements as well as the basis elements themselves. That is, one essentially factors the input, or design matrix, into a sparse catalog and a dictionary of basis elements, both of which need to be inferred from the data.

### 1.1 Problem Definition

Dictionary learning, in its general form, can be seen as

$$\min_{\mathcal{D}, \mathcal{R}, \theta, \lambda} f(\mathcal{D}, \mathcal{R}, X) + g(\mathcal{D}, \theta) + h(\mathcal{R}, \lambda)$$

where  $f$  describes the objective function used to measure goodness of approximation of  $X$  through  $\mathcal{D}, \mathcal{R}$ ,  $g$  describing suitable constraints on the dictionary,  $h$  on the representation respectively.

With an input dataset  $X = [x_1, \dots, x_k], x_i \in \mathbb{R}^d, X \in \mathbb{R}^{d \times k}$  we wish to find a dictionary  $\mathcal{D} \in \mathbb{R}^{d \times n}$ ,  $\mathcal{D} = [d_1, \dots, d_n]$  and a representation  $\mathcal{R} \in \mathbb{R}^{n \times k}$ ,  $\mathcal{R} = [r_1, \dots, r_k]$ ,  $r_i \in \mathbb{R}^n$  such that both  $\|X - \mathcal{D}\mathcal{R}\|_F^2$  is minimized and the representations  $r_i$  are “sparse enough”. To limit the dictionary becoming infinitely large (or small) we introduce a constraint on the dictionary’s columns. For this any sufficient norm can be used. An often used

norm is the  $l_2$ -norm. Thus for the dictionary learning problem we introduce the problem with constraints

$$\begin{aligned} \min_{\mathcal{D}, \mathcal{R}} \quad & \|X - \mathcal{D}\mathcal{R}\|_F^2 \\ \text{s.t.} \quad & \|d_j\|_2 \leq 1, \quad \forall j = 1, \dots, n \\ & \|r_i\|_0 \leq t, \quad \forall i = 1, \dots, k \end{aligned}$$

Here the constraints express the expectation that both  $\mathcal{D}$  is controllably determined and that  $\mathcal{R}$  is adequately sparse. The difficulty in solving this problem comes from the mathematical challenge the  $\|\cdot\|_0$  norm yields. This norm is defined as:

$$\begin{aligned} \text{Let } x &\in \mathbb{R}^d \\ \ell_0(x) &:= |\{x_i : x_i \neq 0\}| \end{aligned}$$

That is, any column of  $\mathcal{R}$  must have no more than  $t$  non-zero elements. By using Lagrange multipliers this can be brought to the general form above

$$\min_{\mathcal{D}, \mathcal{R}, \theta, \lambda} \|X - \mathcal{D}\mathcal{R}\|_F^2 + \sum_{j=1}^m \theta_j (\|d_j\|_2 - 1) + \sum_{i=1}^k \lambda_i (\|r_i\|_0 - t)$$

I.e. the functions are

$$\begin{aligned} f(\mathcal{D}, \mathcal{R}, X) &:= \|X - \mathcal{D}\mathcal{R}\|_F^2 \\ g(\mathcal{D}, \theta) &:= \sum_{j=1}^n \theta_j (\|d_j\|_2 - 1) \\ h(\mathcal{R}, \lambda) &:= \sum_{i=1}^k \lambda_i (\|r_i\|_0 - t) \end{aligned}$$

## 1.2 Related Work

In order to tackle the problem whose complexity mainly arises from the  $\|\cdot\|_0$  norm, one approach is to relax the norm to a more feasible one (i.e. LASSO, LARS for convex-relaxation or log penalty for non-convex relaxation). However, this might lead to over-penalization [9]. Another ansatz is to utilize a greedy approach which constantly adds or removes variables based on some measure [3]. Another method attempts to adjust the dictionary and sparse representation in an online fashion [8]. This can be stated as a maximization problem of a submodular function. A submodular function thereby can be viewed as a discrete analog to a convex function in the continuous case [6] which encodes the principle of “diminishing returns”. These submodular approximations to Sparse Dictionary Selection have become nearly ubiquitous in the literature with slight algorithmic variations that approach the theoretical bounds,  $(1 - 1/e)$ , set in [7].

Applications of Sparse Dictionary Learning—in particular the utilization of submodular approximations—are broad and varied; with uses found in classical Machine Learning (i.e. Sparse Linear Regression), Computer Vision (Image reconstruction, inpainting, denoising), Signal Processing [6] [9], Social Network Analysis, Statistics [4], and Economics [5] among many others.

There has been little work in Sparse Regression and Dictionary Learning that utilize the functional inverse of submodular approximations, known as supermodular functions. Supermodularity can facilitate the potential formulation of a dual to the native Dictionary Learning problem. Recently there has been some attempts at developing foundations for extending the suite of algorithms used in submodular maximization to supermodular formulations. In order to accomplish this, several caveats have been made on the functional representations themselves [1]. We leverage the gains made in submodular optimization and use recent results in approximating supermodular functions in constructing a two-stage supermodular minimization extension of the continuous greedy algorithm used in the submodular case [10].

### 1.3 Supermodular Assumption

Typical optimization problems that use submodular or supermodular functions are generally of the form:

Given a set of objects  $V = \{v_1, \dots, v_n\}$  and a function  $f : 2^V \rightarrow \mathbb{R}$  that returns a real value for any subset. Suppose we are interested in finding the subset that either maximizes or minimizes the function, e.g.,  $\arg \max_{S \subseteq V} f(S)$ , possible subject to some constraints on the size of  $S$ . In dictionary selection these constraints are generally on the number of items that the subset  $S$  can contain, where the function encodes the accuracy of representing  $V$  with  $S$ .

We begin by establishing key definitions and properties of sub-/supermodularity.

A function  $f : 2^V \rightarrow \mathbb{R}$  is *submodular* if for any  $A, B \subseteq V$  we have :

$$f(A) + f(B) \geq f(A \cup B) + f(A \cap B)$$

An alternate and equivalent definition that aids in developing intuition is, in the event that  $A \subseteq B \subseteq V$  and  $v \in V \setminus B$  we have:

$$f(A \cup \{v\}) - f(A) \geq f(B \cup \{v\}) - f(B)$$

This alternate definition of submodularity demonstrates a key feature of this class of functions; that is the marginal increase of any additional element diminishes with the overall size of the set. This is known as the “diminishing returns” property of submodular functions.

Now, a function  $f : 2^V \rightarrow \mathbb{R}$  is said to be *supermodular* if for any sets  $A, B \subseteq V$  we have :

$$f(A) + f(B) \leq f(A \cup B) + f(A \cap B)$$

This leads to a similar “increasing returns” property given sets  $A \subseteq B \subseteq V$  and  $v \in V \setminus B$  we have:

$$f(A \cup \{v\}) - f(A) \leq f(B \cup \{v\}) - f(B)$$

It is important to note that submodular and supermodular functions are inverses of each other. That relation motivates the excursion into the development of two-stage supermodular minimization, influenced by previous work in the development of similar methodology for submodular maximization.

## 2 Evaluation

### 2.1 Dictionary Selection as Two-Stage Supermodular Minimization

In [1], the authors step through a few examples to demonstrate the application of supermodular functions to relevant problems. One that we highlight is that of Sparse Multiple Linear Regression (SMLR). We focus on this problem due to its close relation to Dictionary Selection (in fact we aim to show that these problems are functionally equivalent) as both problems are under the umbrella of Representation Learning. The goal of demonstrating the equivalence of SMLR and Dictionary Selection is that we can extend the theoretical guarantees and algorithmic foundations in [1] in our future work.

SMLR is defined as follows:

Given two matrices  $X \in \mathbb{R}^{m \times n}$ ,  $Y \in \mathbb{R}^{m \times l}$  and an integer  $k$ , find a matrix  $W \in \mathbb{R}^{n \times l}$  that minimizes  $\|XW - Y\|_F^2$  subject to  $W$  having at more  $k$  non-zero rows. This problem is usually paired with a common assumption/simplification, where the columns of  $X$  are adjusted to have unit norm.

As outlined in Section 5, here is a general definition of the Dictionary Selection problem.

$$\underset{\mathcal{D}, \mathcal{R}}{\operatorname{argmin}} \|X - \mathcal{D}\mathcal{R}\|_F^2 + \lambda \sum_{i=1}^k \|r_i\|_0 \quad (1)$$

$$\text{s.t.} \quad \|d_j\|_2 \leq 1, \forall j = 1, \dots, n \quad (2)$$

Thus, given an input dataset  $X = [x_1, \dots, x_k]$ ,  $x_i \in \mathbb{R}^d$ ,  $X \in \mathbb{R}^{d \times k}$  we wish to find the dictionary  $\mathcal{D} \in \mathbb{R}^{d \times n}$ ,  $\mathcal{D} = [d_1, \dots, d_n]$  and a representation  $\mathcal{R} = [r_1, \dots, r_k]$ ,  $r_i \in \mathbb{R}^n$ ,  $\mathcal{R} \in \mathbb{R}^{n \times k}$ , such that both  $\|X - \mathcal{D}\mathcal{R}\|_F^2$  is minimized and the representations  $r_i$  are "sparse enough" (can be specified column by column as in [4], or be extended as defined in SMLR).

In SMLR, the data matrix  $X$  with the determined, sparse  $W$ , are used to approximate  $Y$  while in Dictionary Selection the matrices  $\mathcal{D}$ ,  $\mathcal{R}$  are determined to approximate the data  $X$ . In both problems, sparse representations are used to select a combination of data that closely replicates known data. For all intents and purposes these two problems are functionally equivalent with consideration being made to ensure that the sparsity requirements of Dictionary Selection (where there is a limit on the sparsity of each row/column) port into those of SMLR (where the sparsity requirements are placed on the matrix in full). This isn't of too much concern as this can be handled in the definition of each specific application of the Dictionary Selection problem.

Now, you'll note that the definition of the dictionary selection problem (Equation 6) requires the minimization over the matrices  $\mathcal{D}$  and  $\mathcal{R}$ . When taken together, this problem is combinatorially infeasible [11] and non-convex. However, if we separate the problem into a two-stage optimization (where we fix one matrix and solve for the other and then iterate) we can apply well understood convex solution strategies. [3], [2], [4], [10]. These two stages are:

1. Fix  $\mathcal{D}$ , find a sparse coding between it and  $X$ .
2. Solve the Dictionary Optimization problem:  $\mathcal{D} = X\mathcal{R}^+$  following which we renormalize  $\mathcal{D}$ .

Upon iteration we can expect convergence, using guarantees in the continuous greedy algorithm [2], [3], [6] as well as those made in [1]. The commonly known sparse regression can be seen as an instance of sparse multiple linear regression with  $Y$  existing only of one column. I.e.  $XW - Y$  is a vector.

## 2.2 Weakly- $\alpha$ Supermodularity

Weakly- $\alpha$ -supermodularity is a relaxation of supermodularity. With this relaxation, it can be shown that problems that are  $\alpha$ -weakly supermodular can be solved using a slightly adapted standard greedy algorithm. However, convergence typically requires more steps than if the problem was supermodular only.

A non-negative, non-increasing set function  $f(S) : 2^{[n]} \rightarrow \mathbb{R}^+$  is weakly- $\alpha$ -supermodular if there exists  $\alpha \geq 1$  such that for any two sets  $S, T \subseteq [n]$

$$f(S) - f(S \cup T) \leq \alpha |T \setminus S| \max_{i \in T \setminus S} f(S) - f(S \cup \{i\})$$

To understand the definition better, assume we are given two disjoint sets  $A, B \subseteq [n]$  with  $A \cap B = \emptyset$ . Then  $\alpha$ -weakly-supermodularity means that we can find an element  $i \in B$  s.t.

$$\frac{f(A) - f(A \cup B)}{\alpha |B|} \leq f(A) - f(A \cup \{i\})$$

In words, we can find an element  $i$  that is better (in terms of lowering the objective function) than  $\frac{1}{\alpha}$  times the average gain of an element of  $B$ . Note that for a supermodular function  $\alpha = 1$ .

In order to solve the SMLR problem

$$\min\{f(S) : |S| \leq k\}$$

for some  $k$ , the greedy extension algorithm needs at most

$$\left\lceil \alpha k \ln \left( f \left( \frac{S_0}{\epsilon} \right) \right) \right\rceil$$

steps given a start solution  $S_0$  (usually the empty set), a threshold error  $\epsilon$  (i.e.  $f(S_k) \leq (1 + \epsilon)f(S^*)$  for the optimal solution  $S^*$ ) and the weakness-parameter  $\alpha$  of the function  $f$  [1].

As in [1] the authors show that sparse regression can be seen as minimizing a  $\alpha$ -weakly supermodular function with  $\alpha = \|X^+\|_F$ , the extended greedy algorithm presented can be used to solve sparse regression.

## 3 Results

HERE WE WILL TALK ABOUT THE RESULTS WE GET FROM LEOs EXCURSION INTO THE ANALYSIS

## 4 Discussion

Our primary next step is in taking measures to develop the structure for an extension of the continuous greedy algorithm from submodular functions. After developing this algorithm, which we intend to pattern after that done in [10], we will need to develop guarantees for convergence, accuracy and efficiency. There are theoretical benchmarks that have been made in [2], [7] and [9] that we also need to be sure that we meet in the “dual”-like formulation of supermodular optimization for dictionary selection. Once we have this algorithm in place we will need to apply it some data sets like those in [10].

In another step, we would like to implement the extended greedy algorithm presented in [1] and compare its performance to other known approximative algorithms in order to better understand its limits or potential flaws. Eventually, this can also serve as a benchmark for the extension of the continuous greedy algorithm we seek to develop.

## 5 Model/Problem Refocusing

We don’t have any reason to refocus or redefine our project. There was initially some concern that our problem had already been solved in a satisfiable manner. We’re happy to report that there is significant room for contribution in the realm of supermodular optimization. We began by developing intuition in supermodular and submodular optimization, highlighted here in this section, and by reviewing recent research in the development of algorithms for both supermodular minimization [1] and for a two-stage formulation of submodular maximization [10].

## 6 Data

So far, we haven’t yet determined how we will use data to solidify our project as much of our work at this point is in understanding and developing the theory necessary to develop a new or modified algorithm. As we are looking to develop a novel algorithm to provide a new benchmark for approximate solutions to this class of problems, we will likely assimilate the data and experiments used in [3], [2], [4], [10], among others.

## 7 Completed Steps

In our proposal we outlined the following course of action:

Our first step in this project will be to understand the supermodular minimization framework in which sparse regression can be defined. Then, after having connected and identified sparse regression as supermodular minimization problem we want to show that dictionary selection is a 2-stage sparse regression.

In the following step we need to understand the greedy algorithm for submodular maximization and its improvement, the continuous greedy algorithm. After this, understanding 2-stage submodular optimization would ideally allow us to connect the pieces, gain insight and finally propose a new algorithm of dealing with the problem.

We have successfully completed the first of these stages (as will be shown in 2.1) and are on our way toward completing the second.

## 8 Next Steps

## References

- [1] Christos Boutsidis, Edo Liberty, and Maxim Sviridenko. Greedy minimization of weakly supermodular set functions. *CoRR*, abs/1502.06528, 2015.
- [2] Volkan Cevher and Andreas Krause. Greedy dictionary selection for sparse representation. *Selected Topics in Signal Processing, IEEE Journal of*, 5(5):979–988, 2011.
- [3] Abhimanyu Das and David Kempe. Submodular meets spectral: Greedy algorithms for subset selection, sparse approximation and dictionary selection. *arXiv preprint arXiv:1102.3975*, 2011.
- [4] F. Doshi-Velez and S. A. Williamson. Restricted Indian Buffet Processes. *ArXiv e-prints*, August 2015.
- [5] U. Feige. On maximizing welfare when utility functions are subadditive. In *STOC '06 Proceedings of the thirty-eighth annual ACM symposium on Theory of computing*, pages 41–50, 2006.
- [6] Zhuolin Jiang, Guangxiao Zhang, and L.S. Davis. Submodular dictionary learning for sparse coding. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3418–3425, June 2012.
- [7] Andreas Krause and Carlos Guestrin. Near-optimal nonmyopic value of information in graphical models. In *TWENTY-FIRST CONFERENCE ON UNCERTAINTY IN ARTIFICIAL INTELLIGENCE (UAI, 2005*.
- [8] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online dictionary learning for sparse coding. In *26th International Conference on Machine Learning*, 2009.
- [9] Jianping Shi, Xiang Ren, Guang Dai, Jingdong Wang, and Zhihua Zhang. A non-convex relaxation approach to sparse dictionary learning. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1809–1816, June 2011.
- [10] Yaron Singer, Eric Balkanski, and et. al. Learning sparse combinatorial representations via two-stage submodular maximization. (*preprint*) *Submitted to 2016 ICML*, 2016.
- [11] A.M. Tillmann. On the computational intractability of exact and approximate dictionary learning. *Signal Processing Letters, IEEE*, 22(1):45–49, Jan 2015.