

# AM221 Final Project Proposal

Taylor Killian & Leonhard Spiegelberg

May 1, 2016

## Abstract

To be filled with an interesting summary of our work and results

## 1 Dictionary Learning as 2-stage supermodular minimization

The introduced problem of Dictionary learning, which in its general form can be seen as

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

where  $f$  describes the objective function used to measure goodness of approximation of  $X$  through  $D, R$ ,  $g$  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  $D \in \mathbb{R}^{d \times n}, \mathcal{D} = [d_1, \dots, d_n]$  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  $\|X - D\mathcal{R}\|_F^2$  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

$$\min_{D, R} \|X - DR\|_F^2 \tag{1}$$

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

$$\|r_i\|_0 \leq t, \forall i = 1, \dots, k \tag{3}$$

Description of constraints here

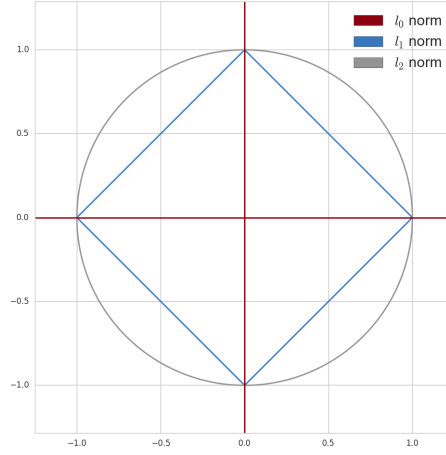


Figure 1: contour plot of the  $l_0, l_1, l_2$  norm at levelset  $L_1(f) := \{x \in \mathbb{R}^2 : f(x) = 1\}$ . Note that for any space  $\mathbb{R}^d$ , the image of the  $l_0$  consists of  $d + 1$  values ( $\{0, \dots, d\}$ ).

Using Lagrange multipliers this can be brought to the general form above

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

I.e. the functions are

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

How to bring this to the penalty function form?

The difficulty in solving this problem comes from the mathematical challenge the  $\|\cdot\|_0$  norm yields. To understand this better confer Figure 1. Let  $x \in \mathbb{R}^d$

$$l_0(x) := (\#i : x_i \neq 0)$$

here describe different norms... the goal is to find a norm that sort of converges to the infinite + of the l0 norm

Put here section about possible relaxations of l0 norm (related literature)

### 1.0.1 Two stage thoughts

We will now relax the problem by formulating it into two stages:

1. solve for fixed  $D, \theta$  the minimization to obtain optimal  $R, \lambda$ .
2. solve optimization for  $D, \theta$  with obtained values
3. repeat previous steps until done

This is actually the method of optimal directions! In the traditional method of optimal directions (ref here) the subproblem  $\min_{R, \lambda} f(D, R, X) + g(D, \theta) + h(R, \lambda)$  is solved either by various relaxation methods (i.e. LASSO, Matching Pursuit(maybe describe here)), however for our approach we want to embed the problem into a combinatorial optimization framework, which guarantees us bounds on optimality. Recent work done by [1] showed that the subproblem of finding an optimal representation with a given dictionary is equivalent to Sparse Multiple Linear Regression (SMLR) and can be reformulated as a  $\alpha$ -weakly supermodular function. We will now show how to obtain this formulation: Remember

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

then let  $S \subseteq [n]$  be a set of column indices for the dictionary. From now onwards we will assume that the given, fixed dictionary is normalized, i.e.  $\forall d_i : \|d_i\|_2 = 1$ . This allows us to remove the term  $g(D, \theta)$ . Define now  $D_S$  as the matrix obtained from  $D$  with all columns not indexed by  $S$  to be set to zero.

$$(D_S)_{ij} := \begin{cases} 0 & j \notin S \\ D_{ij} & j \in S \end{cases}$$

$D_S^+$  is the pseudoinverse which can be obtained i.e. via singular value decomposition when  $D_S = U\Sigma V^T$  through  $D_S^+ = V\Sigma^+U^T$ . Then the problem can

be equivalently stated as

$$\begin{aligned} \min_{R, \lambda} \|X - DR\|_F^2 + \sum_{i=1}^k \lambda_i (\|r_i\|_0 - t) \\ \iff \min_{S \subseteq [n], |S| \leq t} \|X - D_S D_S^+ X\|_F^2 \end{aligned}$$

[1] have shown that the objective function  $f(S) := \|X - D_S D_S^+ X\|_F^2$  of this optimization problem is  $\alpha$ -weakly supermodular with

$$\alpha = \max_{\tilde{S} \subseteq [n]} \|D_{\tilde{S}}^+\|_F^2$$

### 1.1 How to combine it

For the two stage formulation we can generally think of the problem as of adapting the idea of the method of optimal directions but instead of projecting the given representation onto the dictionary we strive for an additional minimization problem: Let  $f_1, \dots, f_m$  be instances of the SMLR problem with different dictionaries  $D_k$

$$f_k := \min_{D_k, R, \lambda} \|X - D_k R\|_F^2 + \sum_{i=1}^k \lambda_i (\|r_i\|_0 - t)$$

Then the dictionary problem can be seen as

$$\min_{D, \theta} q(f_1, \dots, f_m) + \sum_{j=1}^n \theta_j (\|d_j\|_2 - 1)$$

for some combination function  $q$ . We now ask, which properties such a combination function ideally should satisfy.

that is a great question! Can we formulate some ideas, thoughts on what it should look like? I am still unsure about the above formulation...

One choice for  $q$  could be

$$q(f_1(D, \theta), \dots, f_m(D, \theta)) := \sum_{i=1}^m f_i(D, \theta)$$

another one to minimize the product of the SMLR instances. Reformulating it with a logarithm might give even better results with increased computationally

tractability

$$q(f_1(D, \theta), \dots, f_m(D, \theta)) = \sum_{i=1}^m \log(f_i(D, \theta))$$

(redwe can omit the logical step of log'ing also the lambda term because it is a constraint...)

### Algorithm proposals

In the following we propose a general randomized greedy algorithm to solve the dictionary learning problem.

**Algorithm 1:** General randomized greedy algorithm for Sparse Dictionary Learning

```

input : A data matrix  $X \in \mathbb{R}^{d \times k}$ 
output: A dictionary  $D \in \mathbb{R}^{d \times n}$  and representation  $R \in \mathbb{R}^{n \times k}$ 
1  $D_0 \leftarrow \text{RandomDictionary}(-1, 1)$  // assert  $D_0 \neq 0$ 
2 normalize s.t.  $\|(D_0)_{*j}\|_2 = 1$ 
3  $S_0, R_0 \leftarrow \alpha\text{-GreedySolver}(X, D, t)$ 
4  $m \leftarrow 0$ 
5 while  $m < \text{MAXSTEPS}$  do
6    $D_{m+1} \leftarrow X_{S_m \setminus \{q\} \cup \{p\}}$ 
7   normalize s.t.  $\|(D_{m+1})_{*j}\|_2 = 1$ 
8    $S, R \leftarrow \alpha\text{-GreedySolver}(X, D_{m+1}, t)$ 
9   if  $\|X - D_{m+1}R\|_F^2 < \|X - D_m R_m\|_F^2$  then
10     $S_{m+1}, R_{m+1} \leftarrow S, R$ 
11  end
12   $m \leftarrow m + 1$ 
13 end

```

We start by initializing a random dictionary that is non-zero. To allow for all directions negative and positive values should be picked. To restrict the dictionary size, normalization along columns is done in lines 2 and 8. We use the greedy algorithm as proposed in [1] to solve for a given dictionary the SMLR problem by using  $t$  steps (yielding  $t$  columns for representation). The indices of the columns that best solve the SMLR are returned as set  $S$  along a representation  $R$ . We then iteratively refine the solution. Therefore one observation is crucial: Given a set of columns  $S$  that represents the indices of non-zero columns of  $R$ , only  $|S|$  corresponding columns in  $D$  obtained through matrix multiplication will contribute something to the norm. Thus, in order to

improve the solution we need to exchange one (or more) columns of the current set  $S$  with some other column that might give a better solution. To obtain a new dictionary candidate we construct it by using a similar argument via the data matrix  $X$ . I.e. choose randomly a column index  $q \in S$  that will be removed and add another random column index  $p \in \{1, \dots, k\} \setminus S$ . This step is performed in line 7, algorithm 1. We then obtain the dictionary candidate by  $l_2$ -column normalization and solve the resulting SMLR problem via the *alpha*-weakly supermodular greedy algorithm. If the new solution is better than the old one, we keep it, else we perform another round of local search.

## References

- [1] Christos Boutsidis, Edo Liberty, and Maxim Sviridenko. Greedy minimization of weakly supermodular set functions. *CoRR*, abs/1502.06528, 2015.