

Atomic Sets and Extending Machinery from Low-Rank Optimization

Thomas Zhang

Contents

1 Atomic Sets: Definitions and Basic Properties	1
2 Inspiration from Low-Rank Matrix Optimization	3
2.1 Generalized Proof of Weak Submodularity over Rank-One Matrices	3
2.2 Proving Greedy Algorithm Gets Good Bounds	8

1 Atomic Sets: Definitions and Basic Properties

An atomic set \mathcal{A} is a (possibly infinite) collection of “atoms”, which can be anything from basis vectors, low-rank matrices to dirac measures. The reason we care about atomic sets is that they admit ideas about sparsity quite nicely, e.g. we can generalize the problem of regularization.

For the atomic norm to exist, we require that atomic sets are symmetric, i.e. if $a \in \mathcal{A}$, then $-a \in \mathcal{A}$. Given a member x of the same space that the atoms live in, we might wonder how to represent x in terms of the atoms. Say that

$$x = \sum_{\alpha \in \mathcal{A}_t} c_\alpha \alpha, \quad \mathcal{A}_t \subset \mathcal{A}$$

We immediately observe that since \mathcal{A} is symmetric, we can always represent x as a conic combination, i.e. $c_\alpha \geq 0$ for all $\alpha \in \mathcal{A}_t$. If for some choice $\tau \geq 0$, we have $\sum_{\alpha \in \mathcal{A}_t} c_\alpha \leq \tau$, we can also say

$$x \in \mathbf{conv}(\mathcal{A}_t, \tau).$$

Echoing some concepts from functional analysis, atomic sets induce a norm; namely, since atomic sets are symmetric, one can define a gauge functional with respect to the convex hull of a atomic set, and that turns out to be a norm:

$$\|x\|_{\mathcal{A}} := \inf \{t \mid x \in \mathbf{conv}(t\mathcal{A})\}.$$

Sometimes we might be interested in the dual atomic norm, which is defined in the following way

$$\|x\|_{\mathcal{A}}^* := \sup_{\|u\|_{\mathcal{A}} \leq 1} \langle u, x \rangle,$$

if we are concerned with recovering the minimizing direction/atom in an optimization problem.

What is the point of atomic sets?

Atomic sets in some sense are a general model designed to capture the sparsifying, yet convex, structure of things like ℓ^1 regularization and matrix completion. In the case of ℓ^1 regularization, the ultimate goal is to retrieve a solution that has a small support (which can be seen as ℓ^0), but that is not a convex problem (since ℓ^0 is not a norm). In other words, we want to solve some problem like this: $f(x)$ is convex

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & \text{supp}(x) \leq k. \end{aligned}$$

If we treat the elementary unit vectors as the atoms of \mathcal{A} , then we observe that the unit atomic norm ball, in other words $\text{conv}(\mathcal{A})$, is identical to the ℓ^1 norm ball. The properties of the ℓ^1 norm ball gives us an idea of why we would want to find a generalization that preserves them. One may recall coordinate descent for a strongly convex, differentiable objective, where the optimization subroutine at each iteration looks something like

$$\begin{aligned} \Delta x_{CD} &= \arg \min_{\|v\|_1=1} \langle \nabla f(x^{(t)}), v \rangle \\ &= \arg \min_{\|v\|_1=1} \sum_{i=1}^n \frac{\partial f}{\partial x_i^{(t)}} v_i \\ &= - \arg \max_i \left\{ \left| \frac{\partial f}{\partial x_i^{(t)}} e_i \right| \right\} \end{aligned}$$

where v_{\min} would be the corresponding e_i .

As one can see, imposing an ℓ^1 restriction on the subroutine induces the descent direction to be a coordinate direction, i.e. a sparse direction. More complex phenomena related to the ℓ^1 norm are described in the field of compressed sensing. However, from a geometric standpoint, many of these phenomena can be attributed to the “polytopal” structure of the ℓ^1 unit ball. In fact, many phenomena that could be described as “sparsifying” can be roughly ascribed to the underlying polytopal/ridge structures. In our above example, it is clear to see why if we had replaced the ℓ^1 unit ball with some other (I guess finite for now) atomic set, we would also get a similar result: since the atomic-norm ball is the convex hull of finitely many vectors, it is a polytope. Since the function $\langle \nabla f(x), \cdot \rangle$ is trivially a concave function, the minimum on a convex set is attained at a extreme point (a.k.a. vertex). Therefore, we could have easily replaced the constraint set $\|v\|_1 = 1$ with $\|v\|_{\mathcal{A}} = 1$ and instead of the minimizing direction being a coordinate direction, we would have gotten a minimizing direction that is an atom.

2 Inspiration from Low-Rank Matrix Optimization

In Khanna et al’s paper ([cite](#)) (among other papers [[cite](#)]), the connection is drawn between weak submodularity and the atomic set of rank-one matrices, and a greedy algorithm (Section 4.2) is proposed. However, the convergence properties of the algorithm relied heavily on the special properties of rank-one matrices and SVD. Our goal in this section is to first revise the proof of the weak submodularity of a restricted-strong-concave utility function ℓ (Theorem 2 in [[cite](#)]) so that it fits better under a more general atomic set narrative. Then, we will see if we can extend the OMPSel algorithm to atomic sets (under some conditions mentioned later).

2.1 Generalized Proof of Weak Submodularity over Rank-One Matrices

Background

The problem we address in this section is the following:

$$\begin{aligned} & (\arg) \max_X \ell(X) \\ & \text{s.t. } \text{rank}(X) \leq r. \end{aligned}$$

A core assumption is that the function we want to optimize (maximize) satisfies what is called “restricted strong concavity” (convexity for minimization) over the set we care about, in this case low-rank matrices.

Definition 2.1 (Low Rank Restricted Strong Concavity, Restricted Smoothness) *A function $\ell : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$ is restricted strong concave with parameter m_Ω and restricted smooth with parameter M_Ω if for all $X, Y \in \Omega \subset \mathbb{R}^{n \times d}$,*

$$\begin{aligned} -\frac{m_\Omega}{2} \|Y - X\|_F^2 & \geq \ell(Y) - \ell(X) - \langle \nabla \ell(X), Y - X \rangle \\ & \geq -\frac{M_\Omega}{2} \|Y - X\|_F^2. \end{aligned}$$

We remark that if $\Omega' \subseteq \Omega$, then by first principles

$$M_{\Omega'} \leq M_\Omega, \quad m_{\Omega'} \geq m_\Omega. \tag{1}$$

In related work ([cite](#)), it is shown that restricted strong concave/convex functions are closely related to the concept of weak submodularity, which generalizes the concept of submodularity of a set function.

Definition 2.2 (Weak Submodularity Ratio) *Let $S, L \subset [p]$ be two disjoint sets, and $f(\cdot) : [p] \rightarrow \mathbb{R}$. The submodularity ratio of L with respect to S is defined*

$$\gamma_{L,S} := \frac{\sum_{i \in S} (f(L \cup \{i\}) - f(L))}{f(L \cup S) - f(L)}.$$

Essentially, we are measuring the diminishing marginal returns property. Note that $\gamma_{L,S} \geq 1$ if f is submodular. We can also define the weak submodularity ratio of a set $U \subseteq [p]$ with respect to a number k :

$$\gamma_{U,k} := \min_{\substack{L,S:L \cap S = \emptyset; \\ L \subseteq U \\ |S|=k}} \gamma_{L,S}.$$

In order to translate the RSC-RS function utility function $\ell(\cdot)$ to a weakly submodular function f , We introduce a centered set function based on $\ell(\cdot)$, which we will want to show satisfies weak submodularity if $\ell(\cdot)$ satisfies RSC and RS.

Definition 2.3 Given atomic set L , we define the set function

$$f(L) = \max_{H \in \text{diag}(\mathbb{R}^{|L|})} \ell(U_L^\top H V_L) - \ell(0),$$

such that the low-rank maximization problem can be reformulated

$$\max_{|L| \leq r} f(L).$$

In its current state, given a restricted strong concave function $\ell(X)$, $X \in M_{m,n}$, the proof of its weak submodularity relies on the matrix structure of X and SVD. Namely, the proof defines and uses

$$\begin{aligned} \hat{H}^{(L)} &:= \arg \max_{H \in \mathbb{R}^{|L| \times |L|}} \ell(U_L H V_L^\top) \\ B^{(L)} &:= U_L \hat{H}^{(L)} V_L^\top, \end{aligned}$$

where L is a set of indices corresponding to vectors in \mathbb{R}^m and \mathbb{R}^n such that U_L is $|L|$ columns of \mathbb{R}^m vectors, and V_L^\top is $|L|$ rows of \mathbb{R}^n vectors. However, since H is an $\mathbb{R}^{|L| \times |L|}$ matrix, if we expand $B^{(L)}$, we get

$$B^{(L)} = \sum_{i=1}^{|L|} \alpha_{ii} u_i v_i^\top + \sum_{i \neq j} \alpha_{ij} u_i v_j^\top,$$

where α_{ij} are entries in H . This is where the reliance on the matrix structure of the problem comes in. If we are given another index set S corresponding to vectors orthogonal to L , establishing the weak submodularity bounds requires projections onto the span of S , and the expansion discussed does not necessarily exist for more abstract vector spaces. However, if we instead define $\hat{H}^{(L)}$ and $B^{(L)}$:

$$\begin{aligned} \hat{H}^{(L)} &:= \arg \max_{H \in \text{diag}(\mathbb{R}^{|L|})} \ell(U_L H V_L^\top) \\ B^{(L)} &:= U_L \hat{H}^{(L)} V_L^\top \\ &= \sum_{i=1}^{|L|} \alpha_i u_i v_i^\top, \end{aligned}$$

such that $\hat{H}^{(L)}$ is diagonal and hence no cross terms appear in $B^{(L)}$, we can treat $\{u_i v_i^\top\}_{i \in L}$ as a sparse atomic set. The only real barrier to generalizing the upcoming proof to general atomic sets is the fact that we are able to pick orthogonal $u_i v_i^\top$ for the atomic set. Note that this is also an underlying structure for ℓ^1 sparse optimization, where the atoms are the (orthogonal) elementary basis vectors. We will discuss a relaxation of this condition later (hopefully).

The Theorem and Proof

Theorem 2.4 (Weak Submodularity Ratio) *Let L be a set of k rank 1 atoms and S be a set of r rank 1 atoms that have been sequentially orthogonalized against L . If $\ell(\cdot)$ is m_i -strongly concave over matrices of rank i , and \tilde{M}_1 -smooth over the set*

$\tilde{\Omega} := \{(X, Y) : \text{rank}(X - Y) = 1\}$, *then we have the following submodularity ratio for $f(\cdot)$:*

$$\gamma_{L,S} := \frac{\sum_{i \in S} (f(L \cup \{i\}) - f(L))}{f(L \cup S) - f(L)} \geq \frac{m_{k+r}}{\tilde{M}_1}.$$

Proof: essentially this boils down to lower bounding the numerator and upper bounding the denominator. An important assumption is that the atoms in S are all mutually orthogonal and are orthogonal to the space spanned by L . We now lower bound the numerator.

We observe that

$$f(L \cup \{i\}) - f(L) = \ell(B^{(L \cup \{i\})}) - \ell(B^{(L)}),$$

and thus it suffices to deal with $\ell(\cdot)$. We denote B_{ii} , $i \in S$ as the projection of B onto the rank-one matrix $u_i v_i^\top$:

$$B_{ii} = \langle B, u_i v_i^\top \rangle u_i v_i^\top,$$

where $\|u_i\|_2 = \|v_i\|_2 = 1$ such that $u_i v_i^\top$ is normalized

$$\|u_i v_i^\top\|_F^2 = \text{tr}(v_i u_i^\top u_i v_i^\top) = 1.$$

We now have the following series of inequalities

$$\begin{aligned} \ell(B^{(L \cup \{i\})}) - \ell(B^{(L)}) &\geq \ell(B^{(L)} + \alpha_i B_{ii}^{(L \cup S)}) - \ell(B^{(L)}) \\ &\geq \langle \nabla \ell(B^{(L)}), (B^{(L)} + \alpha_i B_{ii}^{(L \cup S)}) - B^{(L)} \rangle \\ &\quad - \frac{\tilde{M}_1}{2} \left\| (B^{(L)} + \alpha_i B_{ii}^{(L \cup S)}) - B^{(L)} \right\|_F^2 \\ &= \langle \nabla \ell(B^{(L)}), \alpha_i B_{ii}^{(L \cup S)} \rangle - \alpha_i^2 \frac{\tilde{M}_1}{2} \left\| B_{ii}^{(L \cup S)} \right\|_F^2, \end{aligned}$$

where α_i is an arbitrary constant. Let us set

$$\alpha_i = \frac{\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \rangle}{\tilde{M}_1 \left\| B_{ii}^{(L \cup S)} \right\|_F^2}.$$

Plugging this back into the inequality, we get

$$\begin{aligned}
& \left\langle \nabla \ell(B^{(L)}), \alpha_i B_{ii}^{(L \cup S)} \right\rangle - \alpha_i^2 \frac{\tilde{M}_1}{2} \|B_{ii}^{(L \cup S)}\|_F^2 \\
&= \left\langle \nabla \ell(B^{(L)}), \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2} B_{ii}^{(L \cup S)} \right\rangle - \left(\frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2} \right)^2 \frac{\tilde{M}_1}{2} \|B_{ii}^{(L \cup S)}\|_F^2 \\
&= \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle^2}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2} - \frac{1}{2} \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2} \\
&= \frac{1}{2} \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle^2}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2}.
\end{aligned}$$

Therefore, we have established that

$$\ell(B^{(L \cup \{i\})}) - \ell(B^{(L)}) \geq \frac{1}{2} \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle^2}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2}.$$

The next step is to sum this over all $i \in S$:

$$\begin{aligned}
\sum_{i \in S} \ell(B^{(L \cup \{i\})}) - \ell(B^{(L)}) &\geq \sum_{i \in S} \frac{1}{2} \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle^2}{\tilde{M}_1 \|B_{ii}^{(L \cup S)}\|_F^2} \\
&= \frac{1}{2\tilde{M}_1} \sum_{i \in S} \frac{\left\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \right\rangle^2}{\|B_{ii}^{(L \cup S)}\|_F^2}.
\end{aligned}$$

The last expression looks awfully like something that results from an orthogonal projection. Let us denote $P_S(X)$ the projection of X onto the span of S , i.e. $\text{span}(\{u_i v_i^\top\})$. Let $\{c_i u_i v_i^\top\}$ be a basis for S (where we will draw a parallel between c_i and $B_{ii}^{(L \cup S)}$).

$$\begin{aligned}
P_S(\nabla \ell(B^{(L)})) &= \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), c_i u_i v_i^\top \rangle}{\langle c_i u_i v_i^\top, c_i u_i v_i^\top \rangle} c_i u_i v_i^\top \\
&= \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), c_i u_i v_i^\top \rangle}{c_i} u_i v_i^\top.
\end{aligned}$$

Note that $B_{ii}^{(L \cup S)} = \langle B^{(L \cup S)}, u_i v_i^\top \rangle u_i v_i^\top$, such that $c_i = \langle B^{(L \cup S)}, u_i v_i^\top \rangle = \|B_{ii}^{(L \cup S)}\|_F$. Plug-

ging this in, we get

$$\begin{aligned}
P_S(\nabla \ell(B^{(L)})) &= \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \rangle}{\|B_{ii}^{(L \cup S)}\|_F} u_i v_i^\top \\
\|P_S(\nabla \ell(B^{(L)}))\|_F^2 &= \left\langle \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \rangle}{\|B_{ii}^{(L \cup S)}\|_F} u_i v_i^\top, \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \rangle}{\|B_{ii}^{(L \cup S)}\|_F} u_i v_i^\top \right\rangle \\
&= \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \rangle^2}{\|B_{ii}^{(L \cup S)}\|_F^2}.
\end{aligned}$$

Therefore, we get

$$\sum_{i \in S} \ell(B^{(L \cup \{i\})}) - \ell(B^{(L)}) \geq \sum_{i \in S} \frac{\langle \nabla \ell(B^{(L)}), B_{ii}^{(L \cup S)} \rangle^2}{\|B_{ii}^{(L \cup S)}\|_F^2} \quad (2)$$

$$= \frac{1}{2\tilde{M}_1} \|P_S(\nabla \ell(B^{(L)}))\|_F^2. \quad (3)$$

We are done with the numerator. Now we want to bound the denominator. By the restricted strong concavity by $\ell(\cdot)$:

$$\ell(B^{(L)}) - \ell(B^{(L \cup S)}) + \langle \nabla \ell(B^{(L)}), B^{(L \cup S)} - B^{(L)} \rangle - \frac{m_{k+r}}{2} \|B^{(L \cup S)} - B^{(L)}\|_F^2 \geq 0.$$

Therefore, we can use this to bound the denominator

$$\begin{aligned}
\ell(B^{(L \cup S)}) - \ell(B^{(L)}) &\leq \langle \nabla \ell(B^{(L)}), B^{(L \cup S)} - B^{(L)} \rangle - \frac{m_{k+r}}{2} \|B^{(L \cup S)} - B^{(L)}\|_F^2 \\
&\leq \max_{\substack{X=U_{(L \cup S)} H V_{(L \cup S)}^\top \\ H \in \text{diag}(\mathbb{R}^{|L \cup S|})}} \langle \nabla \ell(B^{(L)}), X - B^{(L)} \rangle - \frac{m_{k+r}}{2} \|X - B^{(L)}\|_F^2.
\end{aligned}$$

Now we use the fact that $B^{(L)}$ is by construction the argmax of ℓ over $\text{span}(L)$. Therefore, by the restricted strong concavity of $\ell(\cdot)$, we have

$$\langle \nabla \ell(B^{(L)}), P_L(X) - B^{(L)} \rangle = 0 \quad \text{for all } X,$$

since $P_L(X) - B^{(L)} \in \text{span}(L)$. We want to know what the optimal X looks like. First, we observe that X must be of the form

$$X = B^{(L)} + X_S,$$

where X_S is in the span of S . This must be true because the L component of X doesn't change the inner product, but we want to minimize the second term $\frac{m_{k+r}}{2} \|X - B^{(L)}\|_F^2$. We

want X_L to cancel out $B^{(L)}$, since

$$\begin{aligned}
\|X - B^{(L)}\|_F^2 &= \langle X - B^{(L)}, X - B^{(L)} \rangle \\
&= \langle X_S + (X_L - B^{(L)}), X_S + (X_L - B^{(L)}) \rangle \\
&= \|X_S\|_F^2 + 2 \langle X_S, X_L - B^{(L)} \rangle + \|X_L - B^{(L)}\|_F^2 \\
&= \|X_S\|_F^2 + \|X_L - B^{(L)}\|_F^2.
\end{aligned}$$

This leaves us to determine X_S . Notice that our problem has been reduced to

$$\begin{aligned}
&\max_{X_S} \langle \nabla \ell(B^{(L)}), X_S \rangle - \frac{m_{k+r}}{2} \|X_S\|_F^2 \\
&= \max_{X_S} \langle P_S(\nabla \ell(B^{(L)})), X_S \rangle - \frac{m_{k+r}}{2} \|X_S\|_F^2,
\end{aligned}$$

since $X_S \in \text{span}(S)$. We observe that the optimal X_S will be collinear with $P_S(\nabla \ell(B^{(L)}))$, in which case this reduces to a calculus problem:

$$\max_c c \|P_S(\nabla \ell(B^{(L)}))\|_F^2 - c^2 \frac{m_{k+r}}{2} \|P_S(\nabla \ell(B^{(L)}))\|_F^2.$$

It is simple to verify that the maximizer is

$$X^* = c^* P_S(\nabla \ell(B^{(L)})) = \frac{1}{m_{k+r}} P_S(\nabla \ell(B^{(L)})).$$

Plugging this back into the inequality for the denominator, we get

$$\begin{aligned}
\ell(B^{(L \cup S)}) - \ell(B^{(L)}) &\leq \left\langle \nabla \ell(B^{(L)}), \frac{1}{m_{k+r}} P_S(\nabla \ell(B^{(L)})) \right\rangle - \frac{m_{k+r}}{2} \left\| \frac{1}{m_{k+r}} P_S(\nabla \ell(B^{(L)})) \right\|_F^2 \\
&\leq \frac{1}{2m_{k+r}} \|P_S(\nabla \ell(B^{(L)}))\|_F^2.
\end{aligned}$$

Putting our inequalities for the numerator and denominator together, we get

$$\begin{aligned}
\gamma_{L,S} &:= \frac{\sum_{i \in S} (f(L \cup \{i\}) - f(L))}{f(L \cup S) - f(L)} \\
&= \frac{\sum_{i \in S} (\ell(B^{(L \cup \{i\})}) - \ell(B^{(L)}))}{\ell(B^{(L \cup S)}) - \ell(B^{(L)})} \\
&\geq \frac{m_{k+r}}{\tilde{M}_1} \quad \blacksquare
\end{aligned}$$

2.2 Proving Greedy Algorithm Gets Good Bounds

We now want to show that we can use weak submodularity to obtain approximation guarantees on a simple greedy algorithm that depend only on the underlying rank of the problem. Since we rely very little on the matrix structure, the hope is that these approximation guarantees will transfer easily to simple greedy algorithms for more general sparse atomic

optimization. Given the definitions from above, we recall that the problem we are trying to solve is

$$\begin{aligned} \max f(L) \\ \text{s.t. } |L| \leq r, \end{aligned}$$

where L are sets of rank-one matrices $\{u_i v_i^\top\}$. We propose the following algorithm known as GECO (Greedy Efficient Component Optimization) (**cite**), where the parameters are $\text{GECO}(n, d, k, \beta)$: n is the dimension of the left singular vectors, d the dimension of the right vectors ($\mathbb{R}^{n \times d}$ is the ambient space), k is the sparsity parameter, and β is the precision parameter for the linear oracle in the algorithm.

Algorithm 1 $\text{GECO}(n, d, k, \beta)$

```

1:  $L_0 \leftarrow \emptyset$ 
2: for  $t = 1, \dots, k$  do
3:    $u_t v_t^\top \leftarrow \text{OMPSe1}(\beta, L_{t-1})$ 
4:    $L_t \leftarrow L_{t-1} \cup \{u_t v_t^\top\}$ 
5: end for
6: return  $L_k, B^{(L_k)}, f(L_k)$ 
```

$\text{OMPSe1}(\beta, L)$ is a linear oracle that given precision parameter β and atomic set L returns atom $\hat{u}\hat{v}^\top$ such that

$$\langle \nabla \ell(B^{(L)}), \hat{u}\hat{v}^\top \rangle \geq \beta \max_{uv^\top \perp L} \langle \nabla \ell(B^{(L)}), uv^\top \rangle. \quad (4)$$

We choose this particular oracle because it can be easily computed at least in the matrix case by computing the left and right singular vectors corresponding to the largest singular value of $\nabla \ell(B^{(L)}) \leftrightarrow \|\nabla \ell(B^{(L)})\|_2$. We now present the main result of this subsection

Theorem 2.5 (Approximation Guarantee for GECO) *Let L_k be the greedy solution obtained by GECO after k iterations, and let L^* be the optimal atomic set of size at most r satisfying the rank constraint of the problem. If $\ell(\cdot)$ is m_{k+r} restricted strong concave on the set of matrices with rank less than or equal to $k+r$, and \tilde{M}_1 restricted smooth on the set of matrices $\tilde{\Omega}$. Then*

$$f(L_k) \geq \left(1 - \frac{1}{e^w}\right) f(L^*),$$

where $w = \beta^2 \frac{m_{k+r}}{\tilde{M}_1} \frac{k}{r}$.

Note that the above theorem implies that we can attain a constant factor approximation of the true optimal value after iterating the algorithm r times, where r is the underlying rank of the optimal solution. This is an improvement on certain results (**cite**) that show a global linear convergence, except the constant factor approximation depends on n or d times, i.e. the dimension of the problem.

Toward proving the above theorem, we introduce some notation and prove a key lemma. Define $D(t) = f(L_t) - f(L_{t-1})$, $D(0) := 0$ as the improvement after the t -th iteration of the

algorithm. Let us also define $R(t) = f(L^*) - f(L_t)$ as the *remaining* distance from the optimal value after the t -th iteration. The following key lemma relates the weak submodularity constant from the previous section to the incremental improvement $D(t)$.

Lemma 2.6 *At iteration t , we have the following lower bound on the incremental gain*

$$D(t) \geq \beta^2 \frac{m_{t+r}}{\tilde{M}_1} \frac{1}{r} R(t) \geq \beta^2 \frac{m_{k+r}}{\tilde{M}_1} \frac{1}{r} R(t),$$

where the second inequality comes for free from the properties of weak submodularity (1).

Proof of Lemma 2.6: given L^* , which is an atomic set of size at most r , and $L_t =: L$, which has size t , we sequentially orthogonalize L^* with respect to L_t , and call the resulting set S . Observe that S can have at most r members. By Theorem 2.4, we have

$$\gamma_{L_t, S} := \frac{\sum_{i \in S} (f(L \cup \{i\}) - f(L))}{f(L \cup S) - f(L)} \geq \frac{m_{t+r}}{\tilde{M}_1}.$$

We therefore have the following series of inequalities for $D(t)$:

$$\begin{aligned} D(t) &= f(L \cup \{t\}) - f(L) \\ &= \ell(B^{(L \cup \{t\})}) - \ell(B^{(L)}) \\ &\geq \ell(B^{(L)} + \alpha u_t v_t^\top) - \ell(B^{(L)}) \quad \text{for any } \alpha \\ &\geq \langle \nabla \ell(B^{(L)}), \alpha u_t v_t^\top \rangle - \alpha^2 \frac{\tilde{M}_1}{2} \quad \text{by restricted smoothness.} \end{aligned}$$

We now choose $\alpha = \frac{\beta}{\tilde{M}_1} \|\nabla \ell(B^{(L)})\|_2$, where $\|\cdot\|_2$ is the operator induced-2-norm, i.e. the largest singular value σ_1 . Expressing $\nabla \ell(B^{(L)})$ by its singular value decomposition

$$\nabla \ell(B^{(L)}) = \sum_{i=1}^{\min(n,d)} \sigma_i u_i v_i^\top.$$

It is therefore simple to verify that

$$\langle \nabla \ell(B^{(L)}), u v^\top \rangle = u^\top \nabla \ell(B^{(L)}) v \leq \sigma_1 = \|\nabla \ell(B^{(L)})\|_2.$$

This implies that the projection of $\nabla \ell(B^{(L)})$ onto the span of S satisfies

$$\begin{aligned} P_S(\nabla \ell(B^{(L)})) &= \sum_{j=1}^r \langle \nabla \ell(B^{(L)}), u_j v_j^\top \rangle u_j v_j^\top \\ \implies \|P_S(\nabla \ell(B^{(L)}))\|_F^2 &\leq r \|\nabla \ell(B^{(L)})\|_2. \end{aligned}$$

We also recall that by OMP $Se \perp 1$, $u_t v_t^\top$ satisfies

$$\langle \nabla \ell(B^{(L)}), u_t v_t^\top \rangle \geq \beta \|\nabla \ell(B^{(L)})\|_2.$$

Going back to the sequence of lower bounds for $D(t)$, we now have

$$\begin{aligned}
& \langle \nabla \ell (B^{(L)}) , \alpha u_t v_t^\top \rangle - \alpha^2 \frac{\tilde{M}_1}{2} \\
& \geq \frac{\beta^2}{2\tilde{M}_1} \|\nabla \ell (B^{(L)})\|_2 \\
& \geq \frac{\beta^2}{\tilde{M}_1} \frac{1}{r} \|P_S (\nabla \ell (B^{(L)}))\|_F^2 \\
& \geq \frac{\beta^2 m_{t+r}}{\tilde{M}_1} \frac{1}{r} (\ell (B^{(L \cup S)}) - \ell (B^{(L)})) \quad \text{from proof of Theorem 2.4} \\
& \geq \beta^2 \frac{m_{t+r}}{\tilde{M}_1} \frac{1}{r} (\ell (B^{(L^*)}) - \ell (B^{(L)})) \quad \text{since } \text{span}(L^*) \subset \text{span}(L \cup S) \\
& = \beta^2 \frac{m_{t+r}}{\tilde{M}_1} \frac{1}{r} R(t) \geq \beta^2 \frac{m_{k+r}}{\tilde{M}_1} \frac{1}{r} R(t). \quad \blacksquare
\end{aligned}$$

Proof of Theorem 2.5: we observe that the lemma gives us

$$D(t+1) = R(t) - R(t+1) \geq \beta^2 \frac{m_{k+r}}{\tilde{M}_1} \frac{1}{r} R(t) =: \frac{w}{k} R(t),$$

which gives us

$$\begin{aligned}
R(t+1) & \leq (1 - w/k) R(t) \leq (1 - w/k)^{t+1} R(0) = (1 - w/k)^{t+1} f(L^*) \\
R(t) & = f(L^*) - f(L_t) \leq (1 - w/k)^t R(0) = (1 - w/k)^t f(L^*) \\
f(L_t) & \geq (1 - (1 - w/k)^t) f(L^*) \\
f(L_k) & \geq \left(1 - \frac{1}{e^w}\right) f(L^*) \\
& = \left(1 - \frac{1}{e^{\beta^2 \frac{m_{k+r}}{\tilde{M}_1} \frac{k}{r}}}\right) f(L^*). \quad \blacksquare
\end{aligned}$$