

# Course Notes for EE227C (Spring 2018): Convex Optimization and Approximation

Instructor: Moritz Hardt

Email: `hardt+ee227c@berkeley.edu`

Graduate Instructor: Max Simchowitz

Email: `msimchow+ee227c@berkeley.edu`

April 14, 2018

## 21 Lecture 21: Non-convex constraints part 1

Recall that convex minimization refers to minimizing convex functions over convex constraints. Today we will begin to explore minimizing convex functions with non-convex constraints. It is difficult to analyze the impact of “non-convexity” in general, since that can refer to anything that is not convex, which is a very broad class of problems. So instead, we will focus on solving least squares with sparsity constraints:

$$\min_{\|x\|_0 \leq s} \|Ax - y\|_2^2$$

for  $y \in \mathbb{R}^n$ ,  $A \in \mathbb{R}^{n \times d}$ , and  $x \in \mathbb{R}^d$  where  $d < n$ . We will show that in general even this problem is hard to solve but that for a restricted class of problems there is an efficient convex relaxation.

Least squares with sparsity constraints can be applied to solving compressed sensing and sparse linear regression, which are important in a variety of domains. In compressed sensing,  $A$  is a measurement model and  $y$  are the measurements of some sparse signal  $x$ . Compressed sensing is applied to reduce the number of measurements needed for, say, an MRI because by including a sparsity constraint on  $x$  we are able to recover the signal  $x$  in fewer measurements.

In sparse linear regression,  $A$  is the data matrix and  $y$  is some outcome variable. The goal of sparse linear regression is to recover a weights  $x$  on a sparse set of features that are responsible for the outcome variable. In genetics,  $A$  could be the genes of a patient,

and  $y$  is whether they have a particular disease. Then the goal is to recover a weights  $x$  on a sparse set of genes that are predictive of having the disease or not.

When there is no noise in the linear equations, we can simplify the problem to

$$\begin{aligned} \min \|x\|_0 \\ Ax = y \end{aligned}$$

## 21.1 Hardness

Even this simplification is NP-hard, as we will show with a reduction to exact 3-cover, which is NP-complete. Our proof is from [FR13].

**Definition 21.1.** The *exact cover by 3-sets* problem is given a collection  $\{T_i\}$  of 3-element subsets of  $[n]$ , does there exist an exact cover of  $[n]$ , a set  $z \subseteq [d]$  such that  $\cup_{j \in z} T_j = [n]$  and  $T_i \cap T_j = \emptyset$  for  $j \neq i \in z$ ?

**Definition 21.2.** The support of a vector  $x$  is defined as

$$\text{supp}(x) = \{i \mid x_i \neq 0\}.$$

**Theorem 21.3.**  $l_0$ -minimization for general  $A$  and  $y$  is NP-hard.

*Proof.* Define matrix  $A$  as

$$A_{ij} = \begin{cases} 1 & \text{if } i \in T_j \\ 0 & \text{o.w} \end{cases}$$

and  $y$  as the all ones vector. Note that from our construction we have  $\|Ax\|_0 \leq 3\|x\|_0$ , since each column of  $A$  has 3 non-zeros. If  $x$  satisfies  $Ax = y$ , we thus have  $\|x\|_0 \geq \frac{\|y\|_0}{3} = \frac{n}{3}$ . Let us now run  $l_0$ -minimization on  $A, y$  and let  $\hat{x}$  be the output. There are two cases

1. If  $\|\hat{x}\|_0 = \frac{n}{3}$ , then  $y = \text{supp}(\hat{x})$  is an exact 3-cover.
2. If  $\|\hat{x}\|_0 > \frac{n}{3}$ , then no exact 3-cover can exist because it would achieve  $\|\hat{x}\|_0 = \frac{n}{3}$  and hence violate optimality of the solution derived through  $l_0$  minimization.

Thus, since we can solve exact 3-cover through  $l_0$  minimization,  $l_0$  minimization must also be NP-hard. ■

## 21.2 Convex relaxation

Although  $l_0$ -minimization is NP-hard in general, we will prove that for a restricted class of  $A$ , we can relax  $l_0$ -minimization to  $l_1$ -minimization. First, define the set of approximately sparse vectors with support  $S$  as those whose  $l_1$  mass is dominated by  $S$ . Formally,

**Definition 21.4.** The set of approximately sparse vectors with support  $S$  is

$$C(S) = \{\Delta \in \mathbb{R}^d \mid \|\Delta_{\bar{S}}\|_1 \leq \|\Delta_S\|_1\}$$

where  $\bar{S} = [d]/S$  and  $\Delta_S$  is  $\Delta$  restricted to  $S$ ,

$$(\Delta_S)_i = \begin{cases} \Delta_i & \text{if } i \in S \\ 0 & \text{o.w} \end{cases}$$

Recall that the nullspace of matrix  $A$  is the set  $\text{null}(A) = \{\Delta \in \mathbb{R}^d \mid A\Delta = 0\}$ . The nullspace is the set of "bad" vectors in our estimation problem. Consider a solution  $Ax = y$ . If  $\Delta \in \text{null}(A)$ , then  $x + \Delta$  is also a solution since  $A(x + \Delta) = Ax + A\Delta = Ax = b$ . Thus, we focus on matrices whose nullspace only contains zero on the set of sparse vectors that we care about.

**Definition 21.5.** The matrix  $A$  satisfies the restricted nullspace property (RNP) with respect to the support  $S$  if  $C(S) \cup \text{null}(A) = \{0\}$ .

With these definitions in place, we can now state our main theorem.

**Theorem 21.6.** Given  $A \in \mathbb{R}^{n \times d}$  and  $y \in \mathbb{R}^n$  we consider the solution to the  $l_0$ -minimization problem  $x^* = \text{argmin}_{Ax=y} \|x\|_0$ . Assume  $x^*$  has support  $S$  and let the matrix  $A$  satisfy the restricted nullspace property with respect to  $S$ . Then given the solutions of the  $l_1$ -minimization problem  $\hat{x} = \text{argmin}_{Ax=y} \|x\|_1$  we have  $\hat{x} = x^*$ .

*Proof.* We first note that by definition both  $x^*$  and  $\hat{x}$  satisfy our feasibility constraint  $Ax = y$ . Letting  $\Delta = \hat{x} - x^*$  be the error vector we have  $A\Delta = A\hat{x} - Ax^* = 0$ , which implies that  $\Delta \in \text{null}(A)$ .

Our goal now is to show that  $\Delta \in C(S)$  then we would have  $\Delta = 0$  from the restricted nullspace property. First, since  $\hat{x}$  is optimal in  $l_1$  it follows that  $\|\hat{x}\|_1 \leq \|x^*\|_1$ . We then have

$$\begin{aligned} \|x_S^*\|_1 &= \|x^*\|_1 \geq \|\hat{x}\|_1 \\ &= \|x^* + \Delta\|_1 \\ &= \|x_S^* + \Delta_S\|_1 + \|x_{\bar{S}}^* + \Delta_{\bar{S}}\|_1 && \text{by splitting the } l_1 \text{ norm,} \\ &= \|x_S^* + \Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1 && \text{by the support assumption of } \|x^*\|_1, \\ &\geq \|x_S^*\|_1 - \|\Delta_S\|_1 + \|\Delta_{\bar{S}}\|_1. \end{aligned}$$

Hence  $\|\Delta_S\|_1 \geq \|\Delta_{\bar{S}}\|_1$ , which implies  $\Delta \in C(S)$ . ■

So far, so good. We have shown that the  $l_1$ -relaxation works for certain matrices. A natural question however is what kinds of matrices satisfy the restricted nullspace property. In order to get a handle on this, we will study yet another nice property of matrices, the so called restricted isometry property (RIP). Later, we will then see that specific matrix ensembles satisfy RIP with high probability.

**Definition 21.7.** A matrix  $A$  satisfies the  $(s, \delta)$ -RIP if for all  $s$ -sparse vectors  $x$  ( $\|x\|_0 \leq s$ ), we have

$$(1 - \delta)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta)\|x\|_2^2.$$

The intuition is that  $A$  acts like an isometry on sparse vectors (a true isometry would have  $\delta = 0$ ). The RIP is useful since it implies that the difference between two  $s$ -sparse vectors cannot be mapped to 0. By looking at the singular values of  $A$  we can derive the following lemma.

**Lemma 21.8.** *If  $A$  has  $(s, \delta)$ -RIP, then*

$$\|A_S^\top A_S - I_S\|_2 \leq \delta$$

for all subsets  $S$  of size  $s$ . Where

$$(A_S)_{ij} = \begin{cases} A_{ij} & \text{if } j \in S \\ 0 & \text{o.w.} \end{cases}$$

We now show that the RIP implies the restricted nullspace property.

**Theorem 21.9.** *If the matrix  $A$  has the  $(2s, \delta)$ -RIP, then it also has the restricted nullspace property for all subsets  $S$  of cardinality  $|S| \leq s$ .*

*Proof.* Let  $x \in \text{null}(A)$  be arbitrary but not equal to 0. Then we have to show that  $x \notin C(S)$  for any  $S$  with  $|S| \leq s$ . In particular, let  $S_0$  be the set indices of the  $s$  largest coefficients in  $x$ . It suffices to show that  $\|x_{S_0}\|_1 < \|x_{\bar{S}_0}\|_1$  since it would then hold for any other subset.

We write

$$\bar{S}_0 = \bigcup_{j=1}^{\lceil \frac{d}{s} \rceil - 1} S_j$$

where

- $S_1$  is the subset of indices corresponding to the  $s$  largest entries in  $\bar{S}_0$
- $S_2$  is the subset of indices corresponding to the  $s$  largest entries in  $\bar{S}_0 \setminus S_1$
- $S_3$  is the subset of indices corresponding to the  $s$  largest entries in  $\bar{S}_0 \setminus S_1 \setminus S_2$

• etc...

So we have  $x = x_{S_0} + \sum_j x_{S_j}$ . We have decomposed  $x$  into blocks of size  $s$ . This is sometimes called shelling. From RIP, we have

$$\|x_{S_0}\|_2^2 \leq \frac{1}{1-\delta} \|Ax_{S_0}\|_2^2.$$

Since  $x \in \text{null}(A)$  by assumption we have

$$\begin{aligned} A(x_{S_0} + \sum_{j \geq 1} x_{S_j}) &= 0 \\ \implies Ax_{S_0} &= - \sum_{j \geq 1} Ax_{S_j}. \end{aligned}$$

Hence

$$\begin{aligned} \|x_{S_0}\|_2^2 &\leq \frac{1}{1-\delta} \|Ax_{S_0}\|_2^2 \\ &= \frac{1}{1-\delta} \langle Ax_{S_0}, Ax_{S_0} \rangle \\ &= \frac{1}{1-\delta} \sum_{j \geq 1} \langle Ax_{S_0}, Ax_{S_j} \rangle \\ &= \frac{1}{1-\delta} \sum_{j \geq 1} \langle Ax_{S_0}, -Ax_{S_j} \rangle \\ &= \frac{1}{1-\delta} \sum_{j \geq 1} (\langle Ax_{S_0}, -Ax_{S_j} \rangle - \langle x_{S_0}, x_{S_j} \rangle) && \text{since } \langle x_{S_0}, x_{S_j} \rangle = 0 \\ &= \frac{1}{1-\delta} \sum_{j \geq 1} \langle x_{S_0}, (I - A^\top A)x_{S_j} \rangle \\ &\leq \frac{1}{1-\delta} \sum_{j \geq 1} \|x_{S_0}\|_2 \delta \|x_{S_j}\|_2 && \text{from Lemma 21.8.} \end{aligned}$$

So we have

$$\|x_{S_0}\|_2 \leq \frac{\delta}{1-\delta} \sum_{j \geq 1} \|x_{S_j}\|_2. \quad (1)$$

By construction, for each  $j \geq 1$ , we have

$$\|x_{S_j}\|_\infty \leq \frac{1}{S} \|X_{S_{j-1}}\|_1$$

and hence

$$\|x_{S_j}\|_2 \leq \frac{1}{\sqrt{S}} \|X_{S_{j-1}}\|_1.$$

Plugging into Equation 1, we get

$$\begin{aligned}
\|x_{s_0}\|_1 &\leq \sqrt{S}\|x_{s_0}\|_2 \\
&\leq \frac{\sqrt{S}\delta}{1-\delta} \sum_{j \geq 1} \|x_{s_j}\|_2 \\
&\leq \frac{\delta}{1-\delta} \sum_{j \geq 1} \|x_{s_{j-1}}\|_1 \\
&\leq \frac{\delta}{1-\delta} (\|x_{s_0}\|_1 + \sum_{j \geq 1} \|x_{s_{j-1}}\|_1)
\end{aligned}$$

which is equivalent to

$$\|x_{s_0}\|_1 \leq \frac{\delta}{1-\delta} (\|x_{s_0}\|_1 + \|x_{\bar{s}_0}\|_1).$$

Simple algebra gives us  $\|x_{s_0}\|_1 \leq \|x_{\bar{s}_0}\|_1$  as long as  $\delta < \frac{1}{3}$ . ■

Now that we've shown that if  $A$  has the RIP then  $l_1$ -relaxation will work, we look at a few examples of naturally occurring matrices with this property.

**Theorem 21.10.** *Let  $A \in \mathbb{R}^{n \times d}$  be defined as  $a_{ij} \sim \mathcal{N}(0,1)$  iid. Then the matrix  $\frac{1}{\sqrt{n}}A$  has  $(s, \delta)$ -RIP for  $n$  at least  $\mathcal{O}\left(\frac{1}{\delta^2} s \log \frac{d}{s}\right)$ .*

The same holds for sub-Gaussians. We have similar results for more structured matrices such as subsampled Fourier matrices.

**Theorem 21.11.** *Let  $A \in \mathbb{R}^{n \times d}$  be a subsampled Fourier matrix. Then  $A$  has  $(s, \delta)$ -RIP for  $n$  at least  $\mathcal{O}\left(\frac{1}{\delta^2} s \log^2 s \log d\right)$ .*

This result is from [HR15] using work from [RV07, Bou14, CT06].  $\mathcal{O}\left(\frac{1}{\delta^2} s \log d\right)$  is conjectured but open.

There is a lot more work on convex relaxations. For sparsity alone people have studied many variations e.g.

- **Basic pursuit denoising (BPDN)**  $\min \|x\|_1$  such that  $\|Ax - y\|_2 \leq \epsilon$
- **Constrained LASSO**  $\min \|Ax - y\|_2^2$  such that  $\|x\|_1 \leq \lambda$
- **Lagrangian LASSO**  $\min \|Ax - y\|_2^2 + \lambda \|x\|_1$

There are also convex relaxations for other constraints. For example  $\min \text{rank}(X)$  such that  $A(X) = Y$  is hard, a simpler problem is to solve the nuclear norm minimization instead:  $\min \|X\|_*$  such that  $A(X) = Y$ . This can be applied to low-rank estimation for images or matrix completion.

## References

- [Bou14] Jean Bourgain. *An Improved Estimate in the Restricted Isometry Problem*, pages 65–70. Springer International Publishing, Cham, 2014.
- [CT06] Emmanuel J. Candès and Terence Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE Trans. Information Theory*, 52(12):5406–5425, 2006.
- [FR13] Simon Foucart and Holger Rauhut. *A mathematical introduction to compressive sensing*, volume 1. 2013.
- [HR15] Ishay Haviv and Oded Regev. The restricted isometry property of subsampled fourier matrices. *CoRR*, abs/1507.01768, 2015.
- [RV07] Mark Rudelson and Roman Vershynin. On sparse reconstruction from fourier and gaussian measurements. 2007.