
Sparse Recovery-I

1 Introduction

In the previous three lectures we covered some standard techniques used in convex optimization like Projected Gradient Descent and Frank-Wolfe and their convergence analysis when the objective function and the constraint set satisfied certain 'nice' properties. The previous lecture was devoted to analyzing convergence of Projected Gradient Descent in convex settings and it ended with an introduction of Non-Convex Optimization and Sparse Recovery.

Non-Convex optimization in itself is NP- Hard but if the function and constraint set satisfy certain properties it can be performed and even shown to have a good convergence rate. The problem of Sparse recovery is one such example. It is a popular technique used in Statistics and Compressive Sensing for obtaining a sparse representation of the data.

In this lecture, we will study the structure of the Sparse Recovery problem and analyse its solution using the PGD algorithm on non-convex sets.

2 Recap

$$\min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$$

The optimization problem in Non-Convex Optimization has the above form where optimization objective $f(\mathbf{x})$ and the constraint set \mathcal{C} may be non-convex.

The problem of Sparse Recovery can be expressed in terms of finding a sparse linear model $\mathbf{x} \in \mathcal{B}_0(k)$ for a positive semi-definite design or sketch matrix $A \in \mathbb{R}^{n \times d}$ (which may be hand-crafted or the set of observations) and a vector $\mathbf{b} \in \mathbb{R}^d$. Here $\mathcal{C} = \mathcal{B}_0(k) = \{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|_0 \leq k\}$, which is the set of k -sparse vectors. The objective function takes the following form for sparse recovery -

$$f(\mathbf{x}) = \|A\mathbf{x} - \mathbf{b}\|_2^2$$

2.1 α - RSC and β - RSS

We repeat the definitions of α -Restricted Set Convexity and β -Restricted Set Smoothness which are essential for analysis of functions over possibly non-convex sets. Note that these restricted properties are defined with respect to an objective function f and a set \mathcal{C} .

A function f is said to be α -RSC and β -RSS over a set \mathcal{C} if $\forall \mathbf{x}, \mathbf{y} \in \mathcal{C}$,

$$f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \in \left(\frac{\alpha}{2}, \frac{\beta}{2} \right) \|\mathbf{x} - \mathbf{y}\|_2^2$$

As shown in previous lecture , for the sparse recovery settings, the objective function can be represented as a sum of its linear approximation and a quadratic term.

$$f(\mathbf{x}) = f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{1}{2} (\mathbf{x} - \mathbf{y})^\top A^\top A (\mathbf{x} - \mathbf{y})$$

From the above expression it is evident that for f , $\alpha = \lambda_{\min}(A^\top A)$ and $\beta = \lambda_{\max}(A^\top A)$. The convexity condition can be easily violated if the smallest eigenvalue is 0. For an ill-posed problem ($n < d$), A can never be full rank and thus $A^\top A$ has minimum eigenvalue 0.

Remark 17.1. The values of α and β can be different from the minimum and maximum eigenvalues of $A^\top A$ if the eigenvectors corresponding to these vectors don't lie in the constraint set \mathcal{C} .

3 Ensuring α – RSC and β – RSS

We first state the exact optimization problem for sparse recovery.

$$\min_{\mathbf{x} \in \mathcal{B}_0(k)} \frac{1}{n} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

Here, $k \ll d$ and the factor of $\frac{1}{n}$ ensures that eigenvalues don't scale with n .

Sparse Sets offer very rich structure with 'nice' properties in a lot of different scenarios. As seen in earlier lectures , k -sparse sets had small Rademacher complexity and regression over these sets required only $\mathcal{O}(k \log d)$ samples . Similarly, α -RSC and β -RSS can be guaranteed with high probability for the sparse recovery objective over the set of k -sparse vectors using sampling techniques with low sample complexity.

Theorem 17.1. If $A \in \mathbb{R}^{n \times d}$, such that $n = \Omega\left(\frac{1}{(1-\alpha)^2} k \log d\right)$ and each element $A_{i,j}$ is sampled i.i.d from a $\mathcal{N}(0, 1)$ (or $\{-1, 1\}$ Rademacher distribution), then $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$ is α -RSC over the set $\mathcal{B}_0(k)$ with high probability $(1 - \exp(-\Omega(d)))$

A similar result exists for β -RSS.

Theorem 17.2. If $A \in \mathbb{R}^{n \times d}$, such that $n = \Omega\left(\frac{1}{(1-\beta)^2} k \log d\right)$ and each element $A_{i,j}$ is sampled i.i.d from a $\mathcal{N}(0, 1)$ (or $\{-1, 1\}$ Rademacher distribution), then $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$ is β -RSS over the set $\mathcal{B}_0(k)$ with high probability $(1 - \exp(-\Omega(d)))$

Here, $k \log d \ll n$. Thus, with sufficient number of samples we can make the restricted minimum eigenvalue and in turn α as large as we want.

3.1 Implications of α -RSC

α -RSC in the sparse settings leads to uniqueness in the algorithm output for sparse recovery, i.e, when f is α -RSC over $\mathcal{B}_0(k)$ then $\forall \mathbf{x}, \mathbf{y} \in \mathcal{B}_0(k)$

$$\begin{aligned} (\mathbf{x} - \mathbf{y})^\top A^\top A (\mathbf{x} - \mathbf{y}) &\geq \alpha \|\mathbf{x} - \mathbf{y}\|_2^2 \\ \implies \|A(\mathbf{x} - \mathbf{y})\|_2^2 &\geq \alpha \|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned}$$

. The contrapositive of the above statement says that for some distinct $\hat{\mathbf{x}}, \hat{\mathbf{y}} \in \mathcal{B}_0(k)$

$$\begin{aligned} \|A(\hat{\mathbf{x}} - \hat{\mathbf{y}})\|_2^2 &= 0 \\ \implies A\hat{\mathbf{x}} &= A\hat{\mathbf{y}} \end{aligned}$$

This leads to an ill-posed problem, where the algorithm cannot distinguish between two distinct vectors and thus cannot give a unique solution . α -RSC ensures that this never happens.

Exercise 17.1. Prove that the optima $\mathbf{x}^* \in \mathcal{B}_0(k)$ for an α -RSC f is unique.

3.2 Deterministic Construction of A

The design/sketch matrix can be obtained by 2 very different methods. In the settings of Sparse Regression, each row of A corresponds to the input feature vector for a datapoint. However, for compressive sensing applications, each row of A corresponds to a 'carefully' handcrafted functional. The degree of control over A in the later case is very high which motivates us for constructing these design matrices deterministically. It turns out that the best algorithms for deterministic construction of A are quadratic or nearly quadratic in the sparsity level k as opposed to sampling based techniques which are linear in k . The interested reader can refer to Jain and Kar (2017) for a discussion on this.

4 PGD for Sparse Recovery

4.1 Algorithm

Algorithm 1: Projected Gradient Descent

Input: Data A, \mathbf{b} , step lengths η_t , constraint set \mathcal{C}
Output: $\mathbf{x} \in \mathcal{C}$

```

1:  $\mathbf{x}^1 \leftarrow \mathbf{0}$ 
2: for  $t = 1, 2, \dots, T - 1$  do
3:    $\mathbf{z}^{t+1} \leftarrow \mathbf{x}^t - \eta_t \cdot A^\top (A\mathbf{x}^t - \mathbf{b})$ 
4:    $\mathbf{x}^{t+1} \leftarrow \Pi_{\mathcal{C}}(\mathbf{z}^{t+1})$  //Projection onto Non-convex set
5: end for
6: return  $\mathbf{x}^\top$ 
```

The algorithm we will be using to solve this optimization problem is exactly similar to PGD. The only difference here is that projection step projects onto a non-convex set \mathcal{C} . Projection onto a non-convex set is NP-Hard but simplifies for the cases of structured non-convex sets like sparse vectors ($\mathcal{B}_0(k)$) or low rank matrices($\mathcal{B}_{rank}(k)$).

Exercise 17.2. Show that projection onto the set of sparse vectors $\mathcal{B}_0(k)$ is Hard Thresholding.

Exercise 17.3. Show that projection onto the set of low rank matrices $\mathcal{B}_{rank}(k)$ is SVD.

4.2 Analysis

The analysis scheme involves the following steps

1. Apply α -RSC between \mathbf{x}^t and \mathbf{x}^*
2. Apply β -RSS between \mathbf{x}^{t+1} and \mathbf{x}^t
3. Apply Projection Property 0

Remark 17.2. Projection Property 0 ($\|\hat{\mathbf{x}} - \mathbf{z}\|_2 \leq \|\mathbf{x} - \mathbf{z}\|_2 \forall \mathbf{x} \in \mathcal{C}$) is always satisfied by the definition of projection.

Remark 17.3. Projection Properties I and II may be violated for a non-convex set \mathcal{C} .

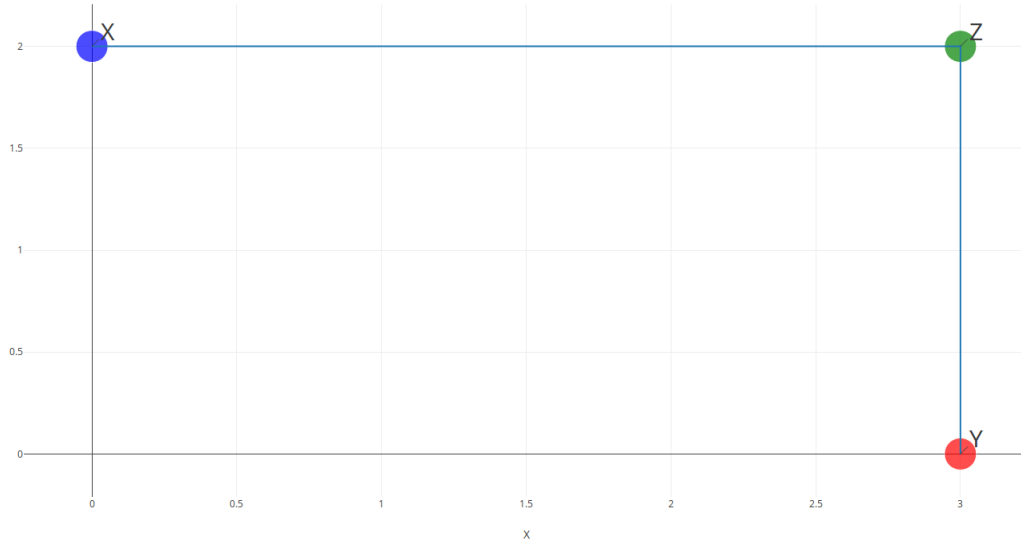


Figure 1: Violation of Projection Property II for 1-sparse vectors

Example 17.1. In \mathbb{R}^2 , $\mathcal{B}_0(1)$ is the union of the 2 axes. Let $\mathbf{z} \in \mathbb{R}^2$ and \mathbf{y} be its projection on $\mathcal{B}_0(1)$. From the figure, $\|\mathbf{x} - \mathbf{y}\| > \|\mathbf{z} - \mathbf{x}\|$ but $\mathbf{x} \in \mathcal{B}_0(1)$

Assumption 17.1. $\nabla f(\mathbf{x}^*) = 0$

This assumption is essential for easy analysis for non-convex sets and the result can be proved without this assumption. Moreover this can be ensured for the noiseless setting ($\mathbf{b} = A\mathbf{x}^*$) when global minima lies in the set \mathcal{C} . This assumption is required to obtain $\langle \nabla f(\mathbf{x}^*), \hat{\mathbf{x}} - \mathbf{x}^* \rangle = 0 \forall \hat{\mathbf{x}} \in \mathcal{C}$

Claim 17.3. If \mathcal{C} is a convex set, then $\forall \hat{\mathbf{x}} \in \mathcal{C} \langle \nabla f(\mathbf{x}^*), \hat{\mathbf{x}} - \mathbf{x}^* \rangle \geq 0$

While analysing PGD, we did not combine the α -RSS and β -RSC conditions together as done in previous lectures (shown below) since the other Projection Properties are invalid. Inability to apply Projection Property I gets us stuck when following this approach.

$$\Phi_{t+1} \leq \langle \nabla f(\mathbf{x}^t), \mathbf{x}^{t+1} - \mathbf{x}^* \rangle + \frac{\beta}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2 - \frac{\alpha}{2} D_t^2$$

However, such a convergence analysis might also be successful but it might not be optimal due to very lossy intermediate steps. We separately deal with the α -RSC and the β -RSS inequalities and combine them only when they have suitable forms.

By β -RSS between \mathbf{x}^{t+1} and \mathbf{x}^t we get

$$f(\mathbf{x}^{t+1}) - f(\mathbf{x}^t) \leq \frac{1}{\eta} \langle \mathbf{x}^t - \mathbf{z}^{t+1}, \mathbf{x}^{t+1} - \mathbf{x}^t \rangle + \frac{\beta}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2$$

Using the identity $2ab = (a + b)^2 - a^2 - b^2$

$$\leq \frac{1}{2\eta} \left(\|\mathbf{x}^{t+1} - \mathbf{z}^{t+1}\|_2^2 - \|\mathbf{x}^t - \mathbf{z}^{t+1}\|_2^2 - \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2 \right) + \frac{\beta}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2$$

Step length is constant and equal to $\frac{1}{\beta}$ for sake of simplicity. Analysis also holds for $\eta \leq \frac{1}{\beta}$

$$= \frac{\beta}{2} \left(\|\mathbf{x}^{t+1} - \mathbf{z}^{t+1}\|_2^2 - \|\mathbf{x}^t - \mathbf{z}^{t+1}\|_2^2 \right)$$

Applying Projection Property 0 with $\hat{x} = x^{t+1}$ and $\mathbf{x} = \mathbf{x}^*$

$$\begin{aligned} &\leq \frac{\beta}{2} \left(\|\mathbf{x}^* - \mathbf{z}^{t+1}\|_2^2 - \|\mathbf{x}^t - \mathbf{z}^{t+1}\|_2^2 \right) \\ &= \frac{\beta}{2} \left(\|\mathbf{x}^* - \mathbf{x}^t + \mathbf{x}^t - \mathbf{z}^{t+1}\|_2^2 - \|\mathbf{x}^t - \mathbf{z}^{t+1}\|_2^2 \right) \end{aligned}$$

Using $(a + b)^2 = a^2 + b^2 + 2ab$

$$\begin{aligned} &\leq \frac{\beta}{2} (D_t^2 + 2 \langle \mathbf{x}^* - \mathbf{x}^t, \mathbf{x}^t - \mathbf{z}^{t+1} \rangle) \\ &= \frac{\beta}{2} D_t^2 + \langle \mathbf{x}^* - \mathbf{x}^t, \nabla f(\mathbf{x}^t) \rangle \end{aligned}$$

Applying α -RSC between \mathbf{x}^* and \mathbf{x}^t

$$\begin{aligned} f(\mathbf{x}^*) - f(\mathbf{x}^t) &\geq \langle \nabla f(\mathbf{x}^t), \mathbf{x}^* - \mathbf{x}^t \rangle + \frac{\alpha}{2} D_t^2 \\ -\frac{\alpha}{2} D_t^2 - \Phi_t &\geq \langle \nabla f(\mathbf{x}^t), \mathbf{x}^* - \mathbf{x}^t \rangle \\ \implies f(x^{t+1}) - f(x^t) &\leq \frac{(\beta - \alpha)}{2} D_t^2 - \Phi_t \\ \implies \Phi_{t+1} - \Phi_t &\leq \frac{(\beta - \alpha)}{2} D_t^2 - \Phi_t \end{aligned}$$

To express D_t in terms of Φ_t , we again use α -RSC between \mathbf{x}^t and \mathbf{x}^*

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \geq \langle \nabla f(\mathbf{x}^*), \mathbf{x}^* - \mathbf{x}^t \rangle + \frac{\alpha}{2} D_t^2$$

Here, we use the assumption that $\nabla f(\mathbf{x}^*) = 0$

$$\begin{aligned} \implies \frac{2\Phi_t}{\alpha} &\leq D_t^2 \\ \implies \Phi_{t+1} - \Phi_t &\leq \left(\frac{\beta}{\alpha} - 1 \right) \Phi_t - \Phi_t \\ \implies \Phi_{t+1} &\leq \Phi_t \left(\frac{\beta}{\alpha} - 1 \right) \\ \implies \Phi_{t+1} &\leq \Phi_t (\kappa - 1) \end{aligned}$$

Here $\kappa = \frac{\beta}{\alpha}$ is the condition number. By definition, $\kappa \geq 1$. We can guarantee exponential decay if $\kappa - 1 \in [0, 1)$. If $\kappa < 2$,

$$\Phi_T \leq \Phi_0 \exp(-T(2 - \kappa))$$

Thus, the value of condition number determines our rate of convergence. Note that here the condition number is defined using the parameters of Restricted Convexity and Restricted Smoothness. Moreover, based on the value of κ , we can guarantee decrease in the potential function by a constant fraction. For exponential decay we need a small condition number. Thus we need a large value of α which translates to more number of samples. For ill-posed problems with high condition numbers under more assumptions convergence can be guaranteed although the analysis becomes much more complicated.

The subsequent lectures will deal with other methods of solving the sparse recovery problem.

References

Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *Found. Trends Mach. Learn.*, 10(3-4):142–336, December 2017. ISSN 1935-8237. doi: 10.1561/22000000058. URL <https://doi.org/10.1561/22000000058>.