

Scribed by Wenlong Mou

Lecture 18

In which we show how to use semidefinite programs to certify non-existence of sparse vectors in a random subspace.

1 Planted Sparse Vector, Tensor Decomposition and Polynomial Optimization

In the past few lectures, we discussed algorithms for finding planted sparse vectors in random subspaces, and tensor decomposition. Though the two problems seem to be completely unrelated, powerful algorithms exist for both problems based on similar ideas. Let's first review the two problems:

Planted Sparse Vector in Random Subspace. Given a description of the subspace $\text{span}(\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_d) \subseteq \mathbb{R}^n$, \mathbf{v}_0 is k -sparse and $\mathbf{v}_1, \dots, \mathbf{v}_d \sim \text{i.i.d.} \mathcal{N}(0, I_n)$, find \mathbf{v}_0 .

Tensor Decomposition. Given tensor $T = \sum_{i=1}^n \mathbf{a}_i \otimes \mathbf{a}_i \otimes \mathbf{a}_i \in \mathbb{R}^{d \times d \times d}$, $\{\mathbf{a}_i\}_{i=1}^n$ linearly independent (or randomly drawn) find $\mathbf{a}_1, \dots, \mathbf{a}_n$.

In previous lectures, we use LP relaxation for the first problem, and guarantees recovery of $\frac{cn}{\sqrt{d \log n}}$ -sparse vector for $d < \frac{n}{2}$. The tensor decomposition problem can be solved with only linear independence assumptions for $n \leq d$, by reduction to matrix decomposition.

On the other hand, better bounds can be achieved under stronger assumptions for both problems. Specifically, if we assume $d = O\left(\frac{\sqrt{n}}{\sqrt{\log n}}\right)$, cn -sparse planted vectors can be recovered by a polynomial-time algorithm for some small constant c . For tensor decomposition with $\mathbf{a}_i \sim \text{i.i.d.} \mathcal{N}(0, I_d)$, polynomial-time algorithm exists in the regime of $n \leq \frac{d^{3/2}}{\text{polylog}(n)}$. An interesting fact is that methods that achieve these bounds are based on similar ideas, for which the following polynomial optimization problem provides a unified view:

$$\max\{p(\mathbf{x})\} \quad \text{s.t.} \|\mathbf{x}\|_2 = 1, \tag{1}$$

where $p(\cdot)$ is a polynomial.

Sparsity of a vector can be relaxed using ratio between its $\|\cdot\|_p$ norms, which leads to polynomial optimization. For tensor decomposition, the problem of finding top singular vector of a 3-tensor is naturally a degree-3 polynomial problem:

$$\max \left\{ \sum_{j,h,k} T_{j,h,k} x_j x_h x_k \right\} \quad s.t. \|\mathbf{x}\|_2 = 1, \quad (2)$$

In general, the problem is hard for polynomials with degree greater than 2. But efficient algorithms or certifiable upper bounds exists for the random cases. In this lecture, we will illustrate the technique by certifying non-existence of cn -sparse vector in random subspace - which is usually the first step for finding planted solution.

Theorem 1 *There exists absolute constant $c, c' > 0$, such that there is a polynomial-time algorithm that certifies $\text{span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d)$ contains no cn -sparse vector with $1 - o(1)$ probability, if $d \leq \frac{c'\sqrt{n}}{\sqrt{\log n}}$ and $\mathbf{v}_1, \dots, \mathbf{v}_d \sim \text{i.i.d.} \mathcal{N}(0, I_n)$.*

In the following we will prove this theorem.

2 SDP Relaxation

We consider continuous relaxation for the discrete problem $\min_{\mathbf{x} \in S} \|\mathbf{x}\|_0$.

Let $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ for $p \geq 1$. Note that for k -sparse vector \mathbf{x} , by Cauchy-Schwartz inequality, we have:

$$\|\mathbf{x}\|_2^4 = \left(\sum_{i:x_i \neq 0} |x_i|^2 \right)^2 \leq \left(\sum_{i:x_i \neq 0} 1^2 \right) \left(\sum_{i:x_i \neq 0} x_i^4 \right) = k \|\mathbf{x}\|_4^4. \quad (3)$$

So we have:

$$\min_{\mathbf{x} \in S} \|\mathbf{x}\|_0 \geq \min_{\mathbf{x} \in S} \frac{\|\mathbf{x}\|_2^4}{\|\mathbf{x}\|_4^4} = \left(\max_{\mathbf{x} \in S, \|\mathbf{x}\|_2=1} \|\mathbf{x}\|_4^4 \right)^{-1}. \quad (4)$$

Let matrix $A = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d] \in \mathbb{R}^{n \times d}$, we have $A_{ij} \sim \text{i.i.d.} \mathcal{N}(0, 1)$. The above optimization problem is equivalent to

$$\max_{\|\mathbf{Ax}\|_2=1} \|\mathbf{Ax}\|_4^4 \quad s.t. \|\mathbf{Ax}\|_2 = 1. \quad (5)$$

Therefore, we have:

$$\max_{\|\mathbf{Ax}\|_2=1} \|\mathbf{Ax}\|_4^4 \leq O(1/n) \iff \left(\max_{\mathbf{x} \in S, \|\mathbf{x}\|_2=1} \|\mathbf{x}\|_4^4 \right)^{-1} \geq \Omega(n) \implies \min_{\mathbf{x} \in S} \|\mathbf{x}\|_0 \geq \Omega(n) \quad (6)$$

The above arguments lead to the first lemma:

Lemma 2 *For subspace $S \subseteq \mathbb{R}^n$, to certify non-existence of cn -sparse vector in S , it suffices to certify that $\max_{\|\mathbf{Ax}\|_2=1} \|\mathbf{Ax}\|_4^4 \leq O(1/n)$.*

We first deal with the constraint $\|A\mathbf{x}\|_2 = 1 \iff \mathbf{x}^T A^T A \mathbf{x} = 1$. Note that for $d < n/2$, $A^T A$ is sample covariance of Gaussian random vectors, which will be concentrated around identity in a multiplicative factor. So, asking for $\|A\mathbf{x}\|_2 = 1$ is roughly asking for $\|\mathbf{x}\|_2 = 1$. We use the following fact from random matrix theory:

Lemma 3 *For $A \sim \mathcal{N}(0, 1)^{n \times d}$, we have the following with high probability:*

$$\frac{n}{2}I_d \preceq A^T A \preceq 2nI_d. \quad (7)$$

Conditioned on $\frac{n}{2}I_d \preceq A^T A$, the constraint $\|A\mathbf{x}\|_2 = 1$ implies $\|\mathbf{x}\|_2 \leq \frac{2}{n}$. So we have:

$$\max_{\|A\mathbf{x}\|_2=1} \|A\mathbf{x}\|_4^4 \leq \max_{\|\mathbf{x}\|_2 \leq \frac{2}{n}} \|A\mathbf{x}\|_4^4 \quad (8)$$

Thus we can make $A^T A \succeq \frac{n}{2}I$ the first part of our certificate (which can be certified in polynomial time), and the revised goal is to certify with high probability that:

$$\max_{\|\mathbf{x}\|_2 \leq \frac{2}{n}} \|A\mathbf{x}\|_4^4 = O\left(\frac{1}{n}\right). \quad (9)$$

3 Certifying the Upper Bound

By now, we have reduce the problem of certifying non-existence of sparse vector to certification of upper bound for 9, which is a random instance of 4-th order polynomial optimization problem 1. The connection to tensor decomposition also lies here.

Rewrite the random matrix as:

$$A = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d] = [\mathbf{a}(1), \mathbf{a}(2), \dots, \mathbf{a}(n)]^T \quad (10)$$

The optimization objective in 9 can be written as:

$$\|A\mathbf{x}\|_4^4 = \sum_{i=1}^n \sum_{j,k,h,l} a_j(i)a_k(i)a_h(i)a_l(i)x_j x_k x_h x_l. \quad (11)$$

Instead of dealing with the 4-th order polynomial, it is more convenient to use quadratic forms. For $\mathbf{x} \in \mathbb{R}^n$, define $\mathbf{x} \otimes \mathbf{x} \in \mathbb{R}^{d^2}$ as a vector indexed by $(j, k) \in \{1, 2, \dots, d\}^2$, with $(\mathbf{x} \otimes \mathbf{x})_{jk} = x_j x_k$. We can write the problem as:

$$\|A\mathbf{x}\|_4^4 = (\mathbf{x} \otimes \mathbf{x})^T M (\mathbf{x} \otimes \mathbf{x}), \quad (12)$$

where $M \in \mathbb{R}^{d^2 \times d^2}$ is a matrix indexed by $(j, k), (h, l)$, defined as:

$$(M)_{(j,k)(h,l)} = \sum_{i=1}^n a_j(i)a_k(i)a_h(i)a_l(i) \quad (13)$$

By definition, the vector $\mathbf{x} \otimes \mathbf{x}$ satisfies:

$$\|\mathbf{x} \otimes \mathbf{x}\|_2 = \left(\sum_{j=1}^d \sum_{k=1}^d x_j^2 x_k^2 \right)^{1/2} = \|\mathbf{x}\|_2^2 \leq \frac{2}{n}. \quad (14)$$

A first attempt is to show that $\|M\|_2 \leq O(n)$ with high probability, and use $\|M\|_2$ as a certificate. However, as we will see in following arguments, this is not true.

Fortunately, this does not mean the above relaxation is hopeless, since we restrict the vector to be the form of $\mathbf{x} \otimes \mathbf{x}$, which is the flattening of a rank-1 matrix. The good news is that the eigenvectors corresponding to large eigenvalues of M are flattened high-rank matrices, which are far from $\mathbf{x} \otimes \mathbf{x}$.

To understand eigen-pairs of random matrix M , we first consider its expectation:

$$\mathbb{E}(M)_{(j,k),(h,l)} = n \mathbb{E} g_j g_h g_k g_l, \quad (15)$$

where $g_1, g_2, \dots, g_d \sim \text{i.i.d.} \mathcal{N}(0, 1)$.

There are only two cases that make $\mathbb{E}(M)_{(j,k),(h,l)}$ non-zero: $j = k = h = l$, or they are two pairs of values but are different.

$$B_{(j,k),(h,l)} = \mathbb{E} g_j g_h g_k g_l = \begin{cases} \mathbb{E} g^4 = 3 & j = k = h = l \\ (\mathbb{E} g^2)^2 = 1 & (j, k, h, l) \text{ are two pairs} \\ 0 & \text{otherwise} \end{cases} \quad (16)$$

For example, let $\mathbf{y} = (\mathbf{1}_{\{j=k\}})_{1 \leq j, k \leq d}$ be flattened identity matrix, we have:

$$\mathbb{E}(\mathbf{y}^T M \mathbf{y}) = n(3d + d(d-1)) = (d+2)n\|\mathbf{y}\|^2. \quad (17)$$

That is why bounding the spectral norm $\|M\|_2$ does not work. And we have to exploit the structure of $\mathbf{x} \otimes \mathbf{x}$:

$$(\mathbf{x} \otimes \mathbf{x})^T B(\mathbf{x} \otimes \mathbf{x}) = 3 \sum_{j=1}^d x_j^4 + 6 \sum_{j,h} x_j^2 x_h^2 = 3\|\mathbf{x}\|_4^4 + 6\|\mathbf{x}\|_2^4 \leq 9\|\mathbf{x}\|_2^4 = O\left(\frac{1}{n^2}\right). \quad (18)$$

So we have: $\mathbb{E}(\mathbf{x} \otimes \mathbf{x})^T M(\mathbf{x} \otimes \mathbf{x}) = O\left(\frac{1}{n}\right)$.

If $\|M - \mathbb{E} M\|_2 = O(n)$, we can easily construct high-probability certificate. However, this is not true either. On the few directions with large eigenvalues, it is possible for M to deviate a lot from the expectation, in terms of spectral norm.

This problem can be circumvented by using multiplicative approximation ratio, instead of the norms of additive error, since the overall guarantee only requires multiplicative factor. It is possible that on some direction \mathbf{y} , we have $\mathbf{y}^T (M - \mathbb{E}(M)) \mathbf{y}$ is large, but on these directions $\mathbf{y}^T \mathbb{E}(M) \mathbf{y}$ is also large. And we only care about directions in the form of $\mathbf{x} \otimes \mathbf{x}$, which guarantees $\mathbf{y}^T \mathbb{E}(M) \mathbf{y}$ to be small. Actually, we can prove the following lemma:

Lemma 4 *With high probability, $M \preceq 2 \mathbb{E}(M)$.*

If this lemma is true, it will guarantee that

$$(\mathbf{x} \otimes \mathbf{x})^T M (\mathbf{x} \otimes \mathbf{x}) \leq 2(\mathbf{x} \otimes \mathbf{x})^T \mathbb{E}(M) (\mathbf{x} \otimes \mathbf{x}) = O\left(\frac{1}{n}\right), \quad \text{w.h.p} \quad (19)$$

According to previous arguments, this is what we need to certify upper bounds for polynomial optimization problem 9. We can verify this within polynomial time, and combine it with the first part, making a valid certificate for non-existence of sparse vector in S that holds with high probability. The entire certificate is:

$$\begin{cases} M \preceq 2 \mathbb{E}(M) \\ A^T A \succeq \frac{n}{2} I \end{cases} \quad (20)$$

The certificate 20 can be computed in polynomial time, and certifies non-existence of cn -sparse vector for some $c > 0$. To prove Theorem 1, it remains to prove Lemma 4.

Note that:

$$M = \sum_{i=1}^n M_i = \sum_{i=1}^n (\mathbf{a}_i \otimes \mathbf{a}_i)(\mathbf{a}_i \otimes \mathbf{a}_i)^T. \quad (21)$$

M is the sum of i.i.d. random matrices, which is convenient to apply matrix concentration inequalities. However, standard matrix concentration techniques such as matrix Chernoff bounds, matrix Bernstein bounds, etc., deal only with additive error in terms of spectral norm, which is not possible in this setting. The solution to that problem is by whitening, i.e., to normalize each matrix M_i so that the expectation becomes identity. If each \hat{M}_i has expectation $I_{d \times d}$, the additive error makes no differences with multiplicative error. Furthermore, the multiplicative error bounds for \hat{M}_i can be translated back to guarantees for M_i , and the lemma can be proved with standard matrix concentration inequalities with normalization trick.

Concretely, we normalize M_i as:

$$\hat{M}_i = B^{-\frac{1}{2}} M_i B^{-\frac{1}{2}} = B^{-\frac{1}{2}} (\mathbf{a}_i \otimes \mathbf{a}_i)(\mathbf{a}_i \otimes \mathbf{a}_i)^T B^{-\frac{1}{2}} \in \mathbb{R}^{d^2 \times d^2}. \quad (22)$$

Here we define the square root of a positive semidefinite matrix as follows: for PSD matrix B with eigen-decomposition $B = P D P^T$, with orthogonal matrix P and diagonal matrix $D = \text{diag}(\lambda_1, \lambda_2, \dots)$ with $\lambda_i \geq 0$, we define $B^{\frac{1}{2}} = P \cdot \text{diag}\left(\lambda_1^{\frac{1}{2}}, \lambda_2^{\frac{1}{2}}, \dots\right) \cdot P^T$. Consequently, if B is positive definite, each $\lambda_i > 0$, we have $B^{-\frac{1}{2}} = \left(B^{\frac{1}{2}}\right)^{-1} = P \cdot \text{diag}\left(\lambda_1^{-\frac{1}{2}}, \lambda_2^{-\frac{1}{2}}, \dots\right) \cdot P^T$.

Apparently, since $M_i = (\mathbf{a}_i \otimes \mathbf{a}_i)(\mathbf{a}_i \otimes \mathbf{a}_i)^T$ is always positive semidefinite, $B = \mathbb{E} M_i$ is also a PSD matrix. To illustrate the proof idea, let us temporarily ignore the invertibility issue, and proceed the analysis. By definition, we have:

$$\mathbb{E} \hat{M}_i = B^{-\frac{1}{2}} \mathbb{E}(M_i) B^{-\frac{1}{2}} = B^{-\frac{1}{2}} \cdot B \cdot B^{-\frac{1}{2}} = I_{d^2} \quad (23)$$

Let $\hat{M} = \sum_{i=1}^n \hat{M}_i$. If we can prove that $\|\hat{M} - \mathbb{E} \hat{M}\|_2 \leq n$ with high probability, this will imply $\hat{M} \preceq \mathbb{E} \hat{M} + n I_{d^2}$. Then we will have the following with high probability:

$$2 \mathbb{E}(M) - M = B^{\frac{1}{2}} \left(2 \mathbb{E}(\hat{M}) - \hat{M}\right) B^{\frac{1}{2}} \succeq B^{\frac{1}{2}} \left(2 \mathbb{E}(\hat{M}) - \mathbb{E}(\hat{M}) - n I_{d^2}\right) B^{\frac{1}{2}} = 0 \quad (24)$$

This will prove Lemma 4 and therefore Theorem 1.

However, the first step of above arguments is not technically correct: actually, the matrix B is never invertible, due to the symmetry among paired quadruples (j, k, h, l) . Concretely, the matrix $B \in \mathbb{R}^{d^2 \times d^2}$ can be written as:

$$B = \begin{bmatrix} 2I_d + J_d & 0_{d \times \frac{d(d-1)}{2}} & 0_{d \times \frac{d(d-1)}{2}} \\ 0_{\frac{d(d-1)}{2} \times d} & I_{\frac{d(d-1)}{2}} & I_{\frac{d(d-1)}{2}} \\ 0_{\frac{d(d-1)}{2} \times d} & I_{\frac{d(d-1)}{2}} & I_{\frac{d(d-1)}{2}} \end{bmatrix}, \quad (25)$$

where the first d rows and columns correspond to (j, k) with $j = k$; the rows and columns at location $d + 1 \cdots d + \frac{d(d-1)}{2}$ correspond to the case of $j < k$; and the last $\frac{d(d-1)}{2}$ locations are for $j > k$. Apparently, the $d^2 \times d^2$ matrix B has rank $\frac{d(d+1)}{2}$.

There are two approaches towards solving this problem: we could either reduce the matrix by exploiting symmetry in our construction of X , or to use pseudo-inverse.

For the first perspective, we can modify the formulation of this SDP, so that $M \in \mathbb{R}^{\frac{d(d+1)}{2} \times \frac{d(d+1)}{2}}$, indexed by $(j, k), (h, l)$ with $j \leq k$ and $h \leq l$. The vector will also be indexed with (j, k) with $j \leq k$, by taking corresponding entries of $\mathbf{x} \otimes \mathbf{x}$. In this setup we will have $B \succeq I_{\frac{d(d+1)}{2}} \succ 0$. And the normalization trick goes through.

Another perspective is by taking pseudo-inverse of B . Recall that the pseudo-inverse B^\dagger of a matrix $B = P \cdot \text{diag}(\lambda_1, \lambda_2, \dots) \cdot P^T$ is defined as $B^\dagger = P \cdot \text{diag}(\nu_1, \nu_2, \dots) \cdot P^T$, where each eigenvalue is $\nu_i = \begin{cases} \lambda_i^{-1} & \lambda_i \neq 0 \\ 0 & \lambda_i = 0 \end{cases}$. This is equivalent to inverting the matrix B restricted to subspace which is orthogonal to B 's nullspace, and the normalization argument can be also carried out in this subspace.

Thus, we can resolve the issue in either ways, and to prove Theorem 1, it remains to show that $\|\hat{M} - \mathbb{E} \hat{M}\|_2 \leq n$ with high probability. In the next lecture, we will use standard tools of matrix concentration inequalities to prove this fact in the regime of $d = O\left(\frac{\sqrt{n}}{\sqrt{\log n}}\right)$.